



# Adatbányászati algoritmusok

Dr. Bodon Ferenc

2010. február 28.

Copyright © 2002-2010 Dr. Bodon Ferenc

Ezen dokumentum a Free Software Foundation által kiadott *GNU Free Documentation license* 1.2-es, vagy bármely azt követő verziójának feltételei alapján másolható, terjeszthető és/vagy módosítható. Nincs *Nem Változtatható Szakasz*, nincs *Címlap-szöveg*, nincs *Hátlap-szöveg*. A licenc magyar nyelvű fordítása a [http://hu.wikipedia.org/wiki/A\\_GNU\\_Szabad\\_Dokumentációs\\_Licenc\\_szövege](http://hu.wikipedia.org/wiki/A_GNU_Szabad_Dokumentációs_Licenc_szövege) oldalon található.

Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 (<http://www.gnu.org/copyleft/fdl.html>) or any later version published by the Free Software Foundation; with no *Invariant Sections*, no *Front-Cover Texts*, and no *Back-Cover Texts*. A copy of the license is included in the section entitled "GNU Free Documentation License".

# Köszönetnyilvánítás

Ezúton szeretnék köszönetet mondani **Rónyai Lajosnak**, a Budapesti Műszaki és Gazdaságtudományi Egyetem tanárának az egész munka során nyújtott segítségével, hasznos ötleteiért, útmutatásaiért, de legfőképpen azért, mert megismertetett az adatbányászattal. Köszönöm **Molnár-Sáska Gábornak**, **Pintér Mártának**, **Szabó Jácintnak**, **Hum Katalinnak**, **Biro Istvánnak** és **Fekete Zsoltnak** az MTA-SZTAKI dolgozóinak valószínűségszámítással kapcsolatos tanácsaikat.

Köszönöm **Buza Krisztiánnak** hasznos megjegyzéseit, ötleteit, szemléletes példáit és a kidolgozott ábráit, amelyekkel hozzájárult a tanulmány sikeréhez.

Külön köszönet illeti **Czibula Veronikát** a tanulmány többszöri, alapos átnézéséért és a felfedezett hibák kijavításáért. **Marx Dániel** rengeteg információval látott el a  $\text{\LaTeX}$ , emacs, Xfig hatékony használatát illetően. Köszönöm neki a fáradozásait.

**Friedl Katának**, **ifjabb Benczúr Andrásnak**, **Lukács Andrásnak**, **Maricza Istvánnak**, **Sarlós Tamásnak** és **Bereczki Tamásnak** köszönöm az értékes észrevételeit, megjegyzéseit.

Értékes észrevételeik és konstruktív javaslataiért köszönet illeti a BME diákjait, többek között (névsorrendben) Erős Pétert, Fekete Gábort, Hajnacs Zoltánt, Lajkó Pétert, Petróczi Attilát, Schlotter Ildikót, Szántó Ádámot, Szőke Mónikát és Varga Dánielt.

Végezetül külön köszönöm Kedvesemnek, **Móninak**, hogy esténként megteremtette az íráshoz és gondolkodáshoz szükséges nyugodt körülményt, továbbá elfogadta és támogatta ezt a sok időt felemésztő hobbimat.

# Tartalomjegyzék

Előszó . . . . .	3
<b>1. Bevezetés</b>	<b>6</b>
1.1. Legjelentősebb adatbányászati feladatok . . . . .	8
1.2. A tudásfeltárás folyamata . . . . .	10
1.3. Adatbányászat kontra statisztika . . . . .	13
1.4. Sikeres alkalmazások . . . . .	15
1.5. Szabványok . . . . .	18
1.6. Adatbányászati rendszer architektúrája . . . . .	19
1.7. Adatbányászat és az etika . . . . .	20
1.8. Az adatbányászat feltételei . . . . .	22
<b>2. Alapfogalmak, jelölések</b>	<b>24</b>
2.1. Halmazok, relációk, függvények, sorozatok . . . . .	24
2.2. Lineáris algebra . . . . .	26
2.3. Gráfelmélet . . . . .	26
2.4. Matematika logika . . . . .	27
2.4.1. Ítéletlogika . . . . .	27
2.4.2. Elsőrendű logika . . . . .	27
2.5. Valószínűségszámítás . . . . .	27
2.5.1. Nevezetes eloszlások . . . . .	27
2.5.2. Egyenlőtlenségek . . . . .	28
2.5.3. Entrópia . . . . .	28
2.6. Statisztika . . . . .	29
2.6.1. Hipotézisvizsgálat . . . . .	29
2.6.2. A binomiális próba . . . . .	30
2.6.3. Az $F$ -próba . . . . .	30
2.6.4. A $\chi^2$ -próba . . . . .	30
2.6.5. Függetlenségvizsgálat . . . . .	31
2.6.6. Student t-próba . . . . .	32
2.7. Algoritmus-elmélet . . . . .	32
2.8. Adatstruktúrák . . . . .	32
2.8.1. Szófák . . . . .	32
2.8.2. Piros-fekete fák . . . . .	34
2.8.3. Hash-tábla . . . . .	34
2.9. Számítógép-architektúrák . . . . .	34
2.9.1. Többszintű memória, adatlokalitás . . . . .	35

2.9.2.	Csóvezetékés feldolgozás, elágazás-előrejelzés . . . . .	35
<b>3.</b>	<b>Előfeldolgozás, hasonlósági függvények</b>	<b>37</b>
3.1.	Attribútum típusok . . . . .	37
3.2.	Hasonlósági mértékek . . . . .	39
3.2.1.	Bináris attribútum . . . . .	39
3.2.2.	Kategória típusú attribútum . . . . .	40
3.2.3.	Sorrend típusú attribútum . . . . .	40
3.2.4.	Intervallum típusú attribútum . . . . .	41
3.2.5.	Vegyes attribútumok . . . . .	42
3.2.6.	Speciális esetek . . . . .	42
3.3.	Előfeldolgozás . . . . .	43
3.3.1.	Hiányzó értékek kezelése . . . . .	43
3.3.2.	Attribútum transzformációk . . . . .	44
3.3.3.	Hibás bejegyzések, a zaj eltávolítása . . . . .	45
3.3.4.	Adatok elrontása, összezagyválása . . . . .	46
3.3.5.	Diszkrétizálás . . . . .	46
3.3.6.	Normalizálás . . . . .	48
3.3.7.	Mintavételezés . . . . .	48
3.3.8.	Dimenziócsökkentés . . . . .	55
<b>4.</b>	<b>Gyakori elemhalmazok</b>	<b>62</b>
4.1.	A gyakori elemhalmaz fogalma . . . . .	62
4.2.	Az APRIORI algoritmus . . . . .	66
4.2.1.	Jelöltek előállítása . . . . .	67
4.2.2.	Jelöltek támogatottságának meghatározása . . . . .	67
4.2.3.	Ritka jelöltek törlése . . . . .	71
4.2.4.	Zsákutca nyesés . . . . .	71
4.2.5.	A bemenet tárolása . . . . .	71
4.2.6.	Tranzakciók szűrése . . . . .	72
4.2.7.	Equisupport nyesés . . . . .	73
4.2.8.	Borgelt-féle támogatottság-meghatározás . . . . .	75
4.2.9.	Futási idő és memóriaigény . . . . .	76
4.3.	Az ECLAT algoritmus . . . . .	80
4.3.1.	kdc1 . . . . .	82
4.3.2.	lcm . . . . .	82
4.4.	Az FP-GROWTH algoritmus . . . . .	82
4.4.1.	Az FP-GROWTH* algoritmus . . . . .	86
4.4.2.	Patricia . . . . .	87
4.5.	Elavult technikák . . . . .	87
4.6.	Mintavételező algoritmus elemzése . . . . .	87
4.6.1.	Mintavétel nagysága . . . . .	87
4.7.	Elemhalmazok Galois lezárja . . . . .	88
4.7.1.	A zárt elemhalmazok fogalma . . . . .	89
4.8.	Kényszerek kezelése . . . . .	90
4.8.1.	ExAnte . . . . .	90

4.9.	Többszörös támogatottsági küszöb . . . . .	91
4.9.1.	MSApriori algoritmus . . . . .	92
<b>5.</b>	<b>Asszociációs szabályok</b>	<b>93</b>
5.1.	Az asszociációs szabály fogalma . . . . .	94
5.1.1.	Maximális következményű asszociációs szabály . . . . .	95
5.1.2.	Egzakt asszociációs szabályok bázisa . . . . .	95
5.2.	Érdekességi mutatók . . . . .	96
5.3.	Szabályok függetlensége . . . . .	97
5.3.1.	lift érték . . . . .	97
5.3.2.	Empirikus kovariancia, empirikus korreláció . . . . .	98
5.3.3.	A $\chi^2$ -statisztika . . . . .	99
5.3.4.	A binomiális próba . . . . .	101
5.3.5.	Fisher-féle egzakt próba . . . . .	101
5.3.6.	További mutatószámok . . . . .	104
5.3.7.	Asszociációs szabályok rangsora . . . . .	105
5.4.	Általánosság, specialitás . . . . .	106
5.5.	Asszociációs szabályok általánosítása . . . . .	107
5.5.1.	Hierarchikus asszociációs szabályok . . . . .	107
5.5.2.	Kategória asszociációs szabályok . . . . .	109
5.6.	A korreláció nem jelent implikációt . . . . .	110
5.7.	Asszociációs szabályok és az osztályozás . . . . .	111
<b>6.</b>	<b>Osztályozás és regresszió</b>	<b>113</b>
6.1.	Bevezetés . . . . .	113
6.2.	Az osztályozás és a regresszió feladata . . . . .	114
6.2.1.	Az elméleti regressziós görbe . . . . .	115
6.2.2.	Maximum likelihood osztályozó . . . . .	116
6.3.	$k$ -legközelebbi szomszéd módszere . . . . .	116
6.3.1.	Dimenzióátlak - Curse of dimensionality . . . . .	118
6.3.2.	A legközelebbi szomszéd érzékenysége . . . . .	118
6.3.3.	Az osztályozás felgyorsítása . . . . .	120
6.4.	Lineárisan szeparálható osztályok . . . . .	122
6.4.1.	Perceptron tanulási szabály . . . . .	124
6.4.2.	Winnow módszer . . . . .	124
6.4.3.	Rocchio-eljárás . . . . .	125
6.4.4.	Lineáris regresszió . . . . .	126
6.4.5.	Logisztikus regresszió . . . . .	127
6.5.	Mesterséges neurális hálózatok . . . . .	131
6.6.	Döntési szabályok . . . . .	134
6.6.1.	Szabályhalmazok és szabálysorozatok . . . . .	135
6.6.2.	Döntési táblázatok . . . . .	136
6.6.3.	Az 1R algoritmus . . . . .	137
6.6.4.	A Prism módszer . . . . .	138
6.7.	Döntési fák . . . . .	139
6.7.1.	Döntési fák és döntési szabályok . . . . .	140

6.7.2.	A döntési fa előállítása . . . . .	141
6.7.3.	Az ID3 algoritmus . . . . .	143
6.7.4.	Feltételek a csomópontokban . . . . .	144
6.7.5.	Vágási függvények . . . . .	144
6.7.6.	Továbbfejlesztések . . . . .	146
6.7.7.	Súlyozott divergenciafüggvények alapján definiált vágási függvények . . .	147
6.7.8.	Döntési fák nyesése . . . . .	149
6.7.9.	Döntési fák ábrázolása . . . . .	149
6.7.10.	Hanyag döntési fák . . . . .	150
6.8.	Bayesi hálózatok . . . . .	150
6.8.1.	Naív Bayes-hálók . . . . .	150
6.8.2.	Naív Bayes-hálók és a logisztikus regresszió kapcsolata . . . . .	152
6.8.3.	Bayes hihetőségi hálók . . . . .	154
6.9.	Osztályozók kombinálása . . . . .	154
6.9.1.	Bagging . . . . .	154
6.9.2.	Randomizálás . . . . .	154
6.9.3.	Boosting . . . . .	154
6.10.	Osztályozók kiértékelése . . . . .	154
6.10.1.	Értekezés . . . . .	157
6.10.2.	Hiba mérése regresszió esetében . . . . .	157
6.10.3.	Hiba mérése valószínűségi döntési rendszerek esetén . . . . .	158
6.10.4.	Osztályozók hatékonyságának mutatószámai . . . . .	159
6.11.	Osztályozók összehasonlítása . . . . .	161
6.11.1.	Binomiális próbán alapuló összehasonlítás . . . . .	161
6.11.2.	Student próbán alapuló összehasonlítás . . . . .	162
6.11.3.	Osztályozó módszerek összehasonlítása . . . . .	162

## 7. Klaszterezés 164

7.1.	Egy lehetetlenség-elmélet . . . . .	165
7.2.	Hasonlóság mértéke, adatábrázolás . . . . .	167
7.3.	A klaszterek jellemzői . . . . .	168
7.4.	A klaszterezés „jósága” . . . . .	169
7.4.1.	Klasszikus mértékek . . . . .	170
7.4.2.	Konduktancia alapú mérték . . . . .	171
7.5.	Klaszterező algoritmusok típusai . . . . .	173
7.6.	Particionáló eljárások . . . . .	174
7.6.1.	Forgy $k$ -közép algoritmus . . . . .	175
7.6.2.	A $k$ -medoid algoritmusok . . . . .	175
7.7.	Hierarchikus eljárások . . . . .	177
7.7.1.	Single-, Complete-, Average Linkage Eljárások . . . . .	177
7.7.2.	Ward módszere . . . . .	178
7.7.3.	A BIRCH algoritmus . . . . .	178
7.7.4.	A CURE algoritmus . . . . .	179
7.7.5.	A Chameleon algoritmus . . . . .	181
7.8.	Sűrűség-alapú módszerek . . . . .	181
7.8.1.	A DBSCAN algoritmus . . . . .	181

<b>8. Idősorok elemzése</b>	<b>183</b>
<b>9. Webes adatbányászat</b>	<b>184</b>
9.1. Oldalak rangsorolása . . . . .	184
9.1.1. Az egyszerű Page Rank . . . . .	185
9.1.2. Az igazi Page Rank . . . . .	188
9.2. Webes keresés . . . . .	188
9.2.1. Gyűjtőlapok és Tekintélyek – a HITS algoritmus . . . . .	188
9.2.2. A SALSA módszer (Jakabfy Tamás) . . . . .	192
9.2.3. Gyűjtőlapok, Tekintélyek és véletlen séták (Jakabfy Tamás) . . . . .	194
9.2.4. Automatikus forrás előállító - Gyűjtőlapok és Tekintélyek módosításai . . . . .	195
9.2.5. Gyűjtőlapok és Tekintélyek módszerének hátrányai . . . . .	196
<b>10. Gyakori minták kinyerése</b>	<b>198</b>
10.1. A gyakori minta definíciója . . . . .	199
10.1.1. Hatékonysági kérdések . . . . .	200
10.2. További feladatok . . . . .	201
10.2.1. Nem bővíthető és zárt minták . . . . .	201
10.2.2. Kényszerek kezelése . . . . .	202
10.2.3. Többszörös támogatottsági küszöb . . . . .	203
10.2.4. Dinamikus gyakori mintakinyerés . . . . .	204
10.3. Az algoritmusok jellemzői . . . . .	204
10.4. Az APRIORI módszer . . . . .	204
10.4.1. Jelöltek előállítása . . . . .	205
10.4.2. Zárt minták kinyerése, az APRIORI-CLOSE algoritmus . . . . .	207
10.5. Sorozat típusú bemenet . . . . .	208
10.5.1. APRIORI . . . . .	208
10.5.2. Zaki módszere . . . . .	210
10.5.3. Mintanövelő algoritmusok . . . . .	212
10.5.4. Kétlépcsős technikák . . . . .	214
10.5.5. A zárt minták „törékenysége” . . . . .	216
10.5.6. Dinamikus gyakori mintabányászat . . . . .	217
<b>11. Gyakori sorozatok, bool formulák és epizódok</b>	<b>219</b>
11.1. Gyakori sorozatok kinyerése . . . . .	219
11.1.1. A Gyakori Sorozat Fogalma . . . . .	220
11.1.2. APRIORI . . . . .	220
11.1.3. Elemhalmazokat tartalmazó gyakori sorozatok . . . . .	221
11.1.4. Sorozat típusú minta általánosítása . . . . .	225
11.2. Gyakori bool formulák . . . . .	226
11.3. Gyakori epizódok . . . . .	226
11.3.1. A támogatottság definíciója . . . . .	227
11.3.2. APRIORI . . . . .	228

<b>12. Gyakori fák és feszített részgráfok</b>	<b>231</b>
12.1. Az izomorfia problémája . . . . .	231
12.2. A gyakori gráf fogalma . . . . .	233
12.3. gyakori gyökeres fák . . . . .	234
12.3.1. TreeMinerH . . . . .	235
12.3.2. TreeMinerV . . . . .	236
12.4. Gyakori részfák . . . . .	238
12.5. A gyakori feszített részgráfok . . . . .	238
12.5.1. Az AcGM algoritmus . . . . .	238
12.6. A gyakori részgráfok keresése . . . . .	241
12.6.1. Az FSG algoritmus . . . . .	241
12.6.2. gSpan . . . . .	242
<b>13. Adatbányászat a gyakorlatban</b>	<b>245</b>
13.1. Felhasználási területek . . . . .	245
13.1.1. Az ügyfél életciklusa . . . . .	245
13.1.2. Kereskedelem . . . . .	246
13.1.3. Pénzügy . . . . .	247
13.1.4. Biológia és Orvostudomány . . . . .	248
13.2. Az adatbányászat bölcsője: az elektronikus kereskedelem (e-commerce) . . . . .	249
13.3. Adatbányász szoftverek . . . . .	250
13.3.1. Adatbányászati rendszerek tulajdonságai . . . . .	251
13.3.2. Esettanulmányok röviden . . . . .	252
<b>Függelék</b>	<b>258</b>
Függelék A . . . . .	258



# Jelölésjegyzék

Az alábbi táblázat tartalmazza a jegyzetben használt legfontosabb jelöléseket.

jelölés	jelentés
$c$	kategóriák száma
$i$	elem (gyakori elemhalmaz keresésénél)
$\ell$	döntési fákban a gyermekek száma
$n$	attribútumok, elemek száma
$m,  \mathcal{T} $	tanítópontok száma

# Előszó

A 90-es években a tárolókapacitások méretének igen erőteljes növekedése, valamint az árak nagymértékű csökkenése<sup>1</sup> miatt az elektronikus eszközök és adatbázisok a hétköznapi életben is mind inkább elterjedtek. Az egyszerű és olcsó tárolási lehetőségek a nyers, feldolgozatlan adatok tömeges méretű felhalmozását eredményezték, ezek azonban a közvetlen visszakeresésen és ellenőrzésen kívül nem sok egyéb haszonnal jártak. A ritkán látogatott adatokból „adat temetők” (data tombs) alakultak ki [55], amelyek tárolása haszon helyett költséget jelentett. Ekkor még nem álltak rendelkezésre olyan eszközök, amivel az adatokba ágyazott értékes információt ki tudtak nyerni. Következésképpen a fontos döntések a döntéshozók megérzésein alapultak, nem pedig az információ-gazdag adatokon. Jól jellemzi ezt a helyzetet John Naisbitt híres mondása, miszerint „We are drowning in information, but starving for knowledge” (Megtulladunk az információtól, miközben tudásra éhezünk).

Egyre több területen merült fel az igény, hogy az adathalmazokból a hagyományosnál árnyaltabb szerkezetű információkat nyerjenek ki. A hagyományos adatbázis-kezelő rendszerek – a közvetlen keresőkérdéseken kívül, illetve az alapvető statisztikai funkciókon túl (átlag, szórás, maximális és minimális értékek meghatározása) – komplexebb feladatokat egyáltalán nem tudtak megoldani, vagy az eredmény kiszámítása elfogadhatatlanul hosszú időbe telt. A szükség egy új tudományterületet keltett életre, az adatbányászatot, amelynek célja: „hasznos, látens információ kinyerése az adatokból”. Az adatbányászati algoritmusokat immár arra tervezték, hogy képesek legyenek az árnyaltabb információ kinyerésére akár óriási méretű adatbázisok esetén is.

Az adatbányászat, mint önálló tudományterület létezéséről az 1980-as évek végétől beszélhetünk. Kezdetben a különböző heurisztikák, a matematikailag nem elemzett algoritmusok domináltak. A 90-es években megjelent cikkek többségét legfeljebb elhinni lehetett, de semmiképpen sem kétely nélkül meggyőződni az egyes írások helytállóságáról. Az algoritmusok futási idejéről és memóriaigényéről általában felszínes elemzéseket és tesztelési eredményeket olvashattunk. Az igényes olvasóban mindig maradt egy-két kérdés, amire említés szintjén sem talált választ. Bizonyos káosz uralkodott, amiben látszólag mindenre volt megoldás, ám ezek a megoldások többnyire részlegesek voltak, tele a legkülönbözőbb hibákkal.

A XXI. századba való belépéssel a kutatók körében egyre nagyobb népszerűségnek kezdett örvendeni az adatbányászat. Ennek két oka van. Egyrészt a növekvő versenyhelyzet miatt a piaci élet szereplőinek óriási az igénye az adatbázisokban megbújó hasznos információkra. A növekvő igény növekvő kutatói beruházásokat indukált. Másrészt, az adatbányászat a maga nehézségével, multi-diszciplináris voltával a kutatni, gondolkodni és újszerű problémákat megoldani vágyó igényét tökéletesen kielégíti.

---

<sup>1</sup>A tárolókapacitás növekedése még Moore jóslatát is jócskán felülmúlja. Az utóbbi 15 év alapján ugyanis a tárolókapacitás 9 hónaponként duplázódik meg [110]

Sorra születtek meg a színvonalas munkák, elemzések, összehasonlítások, mint tiszta irányvonalak rajzolódtak ki a káoszban. A megoldatlan, nyitott problémákra még mindig keressük a választ, így valószínűleg az adatbányászat diadalmenete még sokáig töretlen marad.

Ez a jegyzet a jelenlegi adatbányászati problémákról és az azokat megoldó algoritmusokról szól. A területek áttekintése mellett az algoritmusok mélyebb szintű megismerése is a cél. Az írás informatikus beállítottságú olvasóknak készült. Feltételezzük, hogy az olvasó tisztában van algoritmus- [79] és adatbázis-elméleti alapokkal, továbbá nem ismeretlen terület számára a valószínűségszámítás [7, 40] és a lineáris algebra [115] sem.

A jegyzet célja az, hogy az adatbányászati apparátus olyan megismerését nyújtsa, melynek segítségével az olvasó sikerrel oldja meg az egyre több területen felbukkanó újabb és újabb adatbányászati problémákat. Algoritmikus adatbányászatról írunk, ezért azon mesterséges intelligencia területéhez tartozó eszközök (mesterséges neurális hálózatok, genetikus algoritmusok és fuzzy rendszerek), amelyekről azt tartják, hogy az adatbányászatban is használhatók, kevés hangsúlyt kapnak.

A jegyzet legfrissebb változata letölthető a

<http://www.cs.bme.hu/~bodon/magyar/adatbanyaszat>

címen található oldalról.

A jegyzet nem végleges! Folyamatosan bővül, változik. Egyes részek kisebb súlyt kapnak, mások viszont jobban részletezettek. Örömmel fogadok bármilyen észrevételt, javaslatot akár helyesírási, stilisztikai vagy tipográfiai hibára vonatkozóan. Ezeket kérném, hogy a

`bodon@cs.bme.hu`

címre küldjék.

A tanulmány a Budapesti Műszaki és Gazdaságtudományi Egyetem műszaki informatikusok számára kiírt *Adatbányászati algoritmusok* című tárgy hivatalos jegyzete. Adatbányászatból laborgyakorlatok is vannak, amelynek során a hallgatók a *weka* szabadon hozzáférhető szoftvert ismerik meg. Ezért találkozunk a jegyzetben lépten-nyomon *weka* használati utasításokkal.

Az írás L<sup>A</sup>T<sub>E</sub>X-ben készült, eleinte a *kile*, későbbiekben az *emacs* szövegszerkesztő segítségével. Egyes ábrák *Xfig*-el, mások a *pst-node* csomaggal lettek rajzolva. Az egész munkához az UHU-linux operációs rendszer (<http://www.uhu.linux.hu>) nyújtotta a stabil és biztonságos hátteret.

## Ajánlott irodalom

Először azt kell tisztáznunk, hogy mitől jó egy adatbányászatról szóló könyv. A rengeteg kutatás, projekt, konferencia és folyóirat hatására sok adatbányászati módszert fejlesztettek ki. Mintha elmozdultunk volna a „Megfulladunk az információtól, miközben tudásra éhezünk” kórból a „Megfulladunk az elemző eszközöktől, miközben tudásra éhezünk”. Egy rossz adatbányászati könyv pusztán a módszerek ismertetéséről szól. Olyan érzésünk támad, mintha a kutatók már mindent megoldottak volna és bővelkedünk a jobbnál-jobb eszközökben. Ugyanakkor a megoldások lógnak a levegőben.

Egy jó könyv ezzel szemben keretbe foglalja az eljárásokat, megmutatja hogyan jutunk el az egyik módszerből a másikba, mi a közös és mitől különböznek egymástól a módszerek. Mivel

nincsen tökéletes adatbányászati eljárás, ezért ki kell térni a feladatok nehézségére a módszerek korlátaira és hátrányaira is.

Ezen szempontok alapján osztályozzuk (egyőtől ötig) a következő két részben felsorolt könyveket. A pontok szubjektívek és e tanulmány szerzőjének véleményét tükrözik.

## Magyar nyelvű irodalom

Adatbányász témában az első magyar nyelvű könyv Pieter Adriaans and Dolf Zantinge *Adatbányászat* (1 pont) című könyve [1] volt. Mára a könyv elavult ezért nem ajánljuk senkinek. 2004-ben jelent meg a magyar nyelvű fordítása [54], *ADATBÁNYÁSZAT – Konceptiók és technikák* (3 pont) címmel *Jiawei Han és Micheline Kamber* nagy sikerű könyvének [55]. Azóta megjelent az angol nyelvű könyv második kiadása, ezért ha tehetjük inkább ezt olvassuk.

A legjobb magyar nyelvű adatbányászatról szóló könyvnek a Dr. Abonyi János által szerkesztett *Adatbányászat a hatékonyság eszköze* (4 pont) című könyvet [67] tekintjük. Remek kiegészítése a jelen tanulmánynak. A könyvben helyet kapnak olyan témák, amelyekről ebben a tanulmányban nem esik szó (pl. adattárházak, idősorok, regressziós technikák) habár fontos lenne. Nagyon hasznos, hogy a módszerek bemutatása után a szerzők kitérnek arra, hogy a *weka* szoftvert hogyan kell beállítani a módszer használatához. Mi is az ő példájukat követjük.

Az adatbányászat rokonterületéről írt könyvet Tikk Domonkos *Szövegbányászat* (5 pont) címmel. Kitűnő írás, ajánljuk mind informatikus hallgatóknak és kutatóknak, mind a téma iránt érdeklődőknek.

## Angol nyelvű irodalom

Eibe Frank és Ian H Witten írta az egyik legnépszerűbb adatbányászati könyvet *Data Mining: Practical Machine Learning Tools and Techniques* (5 pont) címmel [142]. Fontos megemlítenünk, hogy Eibe Frank a *weka* egyik főfejlesztője, ennek megfelelően a könyv egy része a *weka* használatát tárgyalja. A könyv egyszerűsége törekszik, kerüli a képleteket, a leírások érthetőek és világosak. Az adatbányászati cikkekben gyakran az ellenkezője figyelhető meg; egyszerű elméleteket és megoldásokat elbonyolítanak, új terminológiát vezetnek be, túlzott formalizmust használnak és elvesznek a figyelemelterelő részletekben, mindez azért, hogy ne lássuk a fától az erdőt, a sok sortól a lényegét. Ebben a könyvben az ellenkező törekvés figyelhető meg, legfontosabb a lényeg megértetése. Ha erre egy példa a legjobb eszköz, akkor el is hagyják a formalizmust, a precíz képleteket. Ajánljuk a könyvet ezért azoknak is, akik nem annyi járatosak a matematikában, viszont alkalmazni szeretnék az adatbányászati eszközöket.

Másik kiemelkedő munka Trevor Hastie, Robert Tibshirani és Jerome Friedman által szerkesztett *The Elements of Statistical Learning: Data Mining, Inference and Prediction* (5 pont) című rendkívül igényes könyv [57]. Az előző könyvvel szemben ez a könyv már komoly matematikai felkészültséget feltételez. Aki viszont rendelkezik statisztikai alapokkal, annak kétségkívül hasznos lesz e olvasmány.

# 1. fejezet

## Bevezetés

A számítógép, korunk legdicsebb találmánya, rohamléptekkel hódít teret magának az élet minden területén. Egy generáció alatt nélkülözhetetlenné vált, amit szüleink még el sem tudtak képzelni, számunkra már elválaszthatatlanná vált munkánktól és szórakozásunktól egyaránt.

Az Internet elterjedésével még intenzívebben érzékelhető a számítógép térhódítása. A világon az egyik legnagyobb problémát, a távolságot hidalta át. Üzleti és magáncélú érintkezések váltak lehetővé rövidebb idő alatt és hatékonyabban, mint valaha. Adatok millióit kezelik és szállítják a számítógépes rendszerek. Az információkon alapuló döntéshozatal ideje lerövidült, hiszen a hozzáférés könnyebbé és gyorsabbá vált. Az üzleti élet szereplőinek élete is felgyorsult.

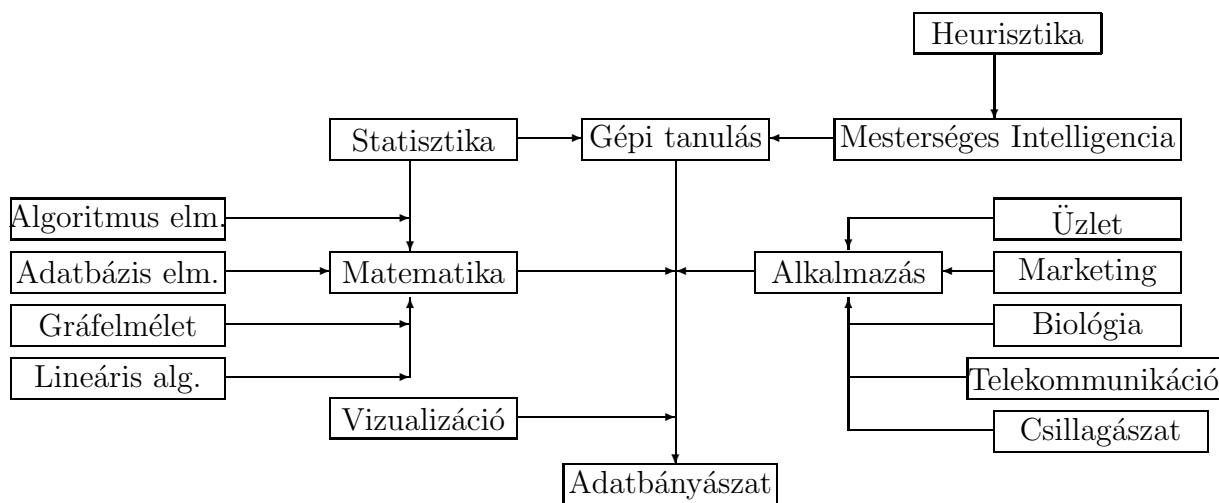
Ma a vállalatok léte múlhat az információk gyors és pontos begyűjtésén, elemzésén, a rugalmas fejlődésen, valamint az innováción. Egyre több felső vezető ismeri fel, hogy az Internet, az adatok elektronikus tárolása a vállalat szolgálatába állítható. Az adatok azonban önmagukban nem hasznosak, hanem a belőlük kinyerhető, a vállalat igényeihez igazodó, azt kielégítő információkra lenne szükség. Ez egy újabb szükségletet teremt: egy olyan eszköz iránti igényt, ami képes arra, hogy információszerzés céljából elemezze a nyers adatokat. Ez az új eszköz az **adatbányászat**.

„ Az angol tudósok azt állapították meg, hogy aki sokat jár disco-ba, annak nagyobb valószínűséggel alakul ki asztmája.” Forrás: Sláger rádió, 2007. október 2., 8 óra 26 perc

Adatbányászati (data mining) algoritmusokat az adatbázisból történő tudásfeltárás (knowledge discovery in databases) során alkalmaznak. A tudáskinyerés adatbázisokból egy olyan folyamat, melynek során érvényes, újszerű, lehetőleg hasznos és végső soron érthető mintákat fedezünk fel az adatokban. Ezt gyakran megtehetjük különböző lekérdezések eredményeinek vizsgálatával, azonban ez a megoldás lassú, drága és nem elég átfogó. Nem is beszélve arról, hogy az emberi szubjektivitás sokszor hibás, továbbá az adatbázisok olyan nagyok lehetnek, hogy egyes lekérdezések elfogadhatatlanul lassan futnak le. Jogos tehát az igény, hogy a legismertebb, leggyakoribb elemzéstípusokhoz speciális módszereket, algoritmusokat fejlesszenek ki, amelyek gyorsan és pontosan szolgáltatnak egy objektív képet az adatbázisokban található „kincsről”.

Sokféleképpen definiálták az adatbányászatot. Felsorolunk néhányat a legismertebbek közül kiemelve a kulcsszavakat:

- „The nontrivial **extraction** of implicit, previously unknown, and potentially **useful information** from **data**” (Piatetsky Shapiro)



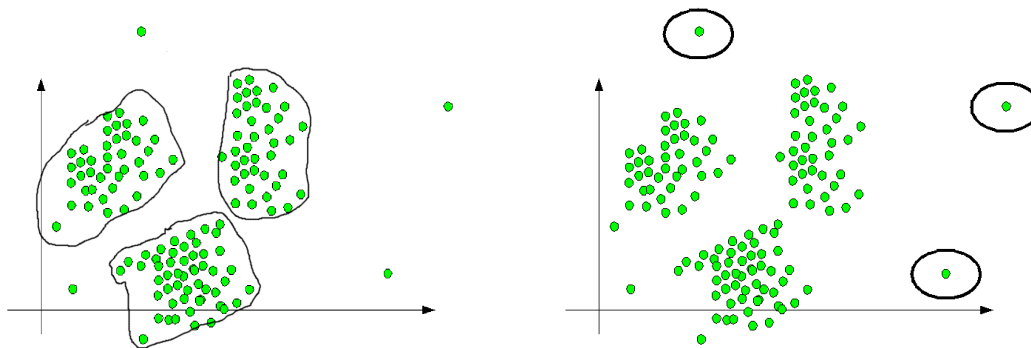
1.1. ábra. Az adatbányászat kialakulása

- „... the automated or convenient **extraction** of **patterns** representing **knowledge** implicitly stored or captured in **large databases, data warehouses, the Web, ... or data streams.**” (Han [55], xxi oldal)
- „... the process of **discovering patterns** in **data**. The process must be automatic or (more usually) semiautomatic. The **patterns** discovered must be **meaningful...**” (Witten [142], 5. oldal)
- „... **finding** hidden **information** in a **database.**” (Dunham [35], 3. oldal)
- „... the process of employing one or more computer learning techniques to automatically **analyze and extract knowledge** from **data contained within a database.**” (Roiger, 4. oldal)

Egyesek szerint az adatbányászat, mint megnevezés némiképp szerencsétlen [54]. Ha szénbányászatról beszélünk, a szén bányászására gondolunk. Ezzel ellentétben adatbányászat esetén *nem* adatot bányászunk, hanem — amint a példákban is láttuk — a rejtett és számunkra hasznos *tudást* (információt), összefüggéseket keressük egy nagy adathalmazban („adathegyben”).

Az adatbányászatot az üzleti élet és a marketing keltette életre. Még ma is ezek az adatbányászat fő mozgató rugói. Szerencsére az adatbányászat lehetőségeit egyre több területen ismerik fel, melynek eredményeként az alapkutatásoknak is egy fontos eszköze lett. Alkalmazzák az orvosbiológiában, genetikában, távközlésben, csillagászatban, ...

Az adatbányászat egy multi-diszciplináris terület. Az 1.1 ábrán látható, hogy mely tudományterületek eszközeit használja az adatbányászat. Az adatbányászat több hangsúlyt fektet az algoritmusokra, mint a statisztika, és többet a modellekre, mint a gépi tanulás eszközei (pl. neurális hálózatok). Mára az adatbányászat akkora területté nőtte ki magát, hogy szinte lehetetlen átlátni magas színvonalon az egészet.



1.2. ábra. Klaszterezés (bal oldali ábra) és különc pontok keresése (jobb oldali ábra)

## 1.1. Legjelentősebb adatbányászati feladatok

Feltehetjük, hogy az adatbázis valamilyen objektumok (ügyfelek, betegségek, vásárlók, telekommunikációs események, ...) különböző tulajdonságait írja le. A tulajdonság helyett gyakran használjuk majd az attribútum szót<sup>1</sup>. Az adatbányászat feladata a rejtett összefüggések, kapcsolatok felderítése. Az összefüggések típusa szerint a következő adatbányászati alapproblémákról beszélhetünk:

**Gyakori minták kinyerése:** Adott objektumok egy sorozata. Célunk megtalálni a gyakran előforduló (rész-) objektumokat. Az objektumok lehetnek elemhalmazok vagy sorozatok, esetleg epizódok (részben rendezések), gráfok stb.

**Attribútumok közötti kapcsolatok:** Gyakran hasznos, ha az objektumokra úgy tekintünk, mint az attribútumok megvalósulásaira és keressük az összefüggéseket az attribútumok között. Többféle összefüggés létezik. Ilyenek például az asszociációs-, korrelációs szabályok, a funkcionális függőségek és hasonlóságok. Az *osztályozás* is attribútumok közötti összefüggések felfedezésére szolgál. Az osztályozásnál egy kitüntetett attribútum értékét kell megjósolnunk a többi attribútum értéke alapján. Ezt egy modell felépítésével teszi. Leggyakrabban a modell egy döntési fa, de lehet if-then szabályok sorozata, valamilyen matematikai formula, vagy akár egy neurális hálózat stb. is.

**Klaszterezés:** Objektumokat előre nem definiált csoportokba (klaszterekbe) kell sorolnunk úgy, hogy az egy csoportba tartozó objektumok hasonlóak legyenek, míg a különböző csoportba kerültek különbözzenek egymástól. Két pont hasonlóságát egy előre megadott (távolságszerű) függvény segítségével szokás értelmezni. Klaszterezésre mutat példát az 1.2 ábra első fele.

**Sorozatelemzés:** A sorozatelemzésbe többféle adatbányászati feladat tartozik. Kereshetünk egymáshoz hasonlító (akár rész-) sorozatokat. Ezen kívül elemezhetjük a sorozat alakulását, és különböző regressziós módszerekkel próbálhatjuk megjósolni a jövőbeli valószínűleg előforduló eseményeket.

<sup>1</sup>A közgazdászok a tulajdonság helyett *ismérvet*, valamilyen tulajdonság konkrét értéke helyett *ismérv változatot* mondanak.

**Eltéréselemzés:** Azokat az elemeket, amelyek nem felelnek meg az adatbázis általános jellemzőinek, tulajdonságaik nagy mértékben eltérnek az általánostól *különc* pontoknak nevezük. A legtöbb adatbányászati algoritmus az ilyen külön pontoknak nem tulajdonít nagy jelentőséget, legtöbbször zajnak vagy kivételnek kezeli őket. Azonban az élet egyre több területén merül fel az igény, hogy éppen az ilyen külön pontokat találjuk meg. Eltéréselemzés főbb alkalmazási területe a másolás-, koppintáskeresés, továbbá a csalások, visszaélések, vírusok, hackertámadások kiszűrése. Különc pontok kezelésére mutat példát az 1.2 ábra második fele.

**Webes adatbányászat:** Az Interneten óriási adattömeg található, így az Interneten alapuló információ-kinyerő algoritmusok is az adatbányászat területéhez tartoznak. A jegyzetben szó lesz intelligensebb keresésről, oldalak rangsorolásáról, illetve hasonló tartalmú oldalak megtalálásáról.

Előfordulhat, hogy az adatbányászati rendszer, még megfelelően megválasztott paraméterek mellett is, túl sok szabályt, összefüggést tár fel. Az egyik legnehezebb kérdés az, hogy ezek közül melyek az érdekesek. Érdekességi mutatókról általánosságban nem sok mondható el, mert a különböző felhasználási területeken más-más minta lehet hasznos. Megkülönböztetünk szubjektív és objektív érdekességi mutatókat. Egy minta mindenképpen érdekes, ha meglepő, azaz eddigi tudásunknak ellentmond, vagy újszerű, azaz tudásunkat kiegészíti. Ugyanakkor egy információ csak akkor érdekes, ha felhasználható, azaz tudunk valamit kezdeni vele [127]. Azt, hogy egy szabály mennyire meglepő – több-kevesebb sikerrel – tudjuk formalizálni. Az újszerűségről és a felhasználhatóságról azonban csak a terület szakértője tud nyilatkozni.

Annak ellenére, hogy az adatbányászat egy új terület, a fentiekből látható, hogy régi, már ismert problémákat is magába foglal. Gondoljunk itt arra, hogy klaszterező algoritmusokat már a 60-as években is javasoltak, vagy arra, hogy az osztályozás feladatát függvény approximációként is felfoghatjuk, aminek irodalmával több könyvespolcot is meg lehetne tölteni<sup>2</sup>. Tehát az adatbányászatban gyakran nem maga a probléma új, hanem az adatok mérete, továbbá az a követelmény, hogy az egyes algoritmusok futási ideje olyan rövid legyen, hogy az eredmények a gyakorlatban elfogadható időn belül érkezzenek. Az alkalmazásokban nem ritkák a giga- sőt terrabájt nagyságú adathalmazok. A [36] írásban például egy beszámolót olvashatunk egy bank adatbázisának elemzéséről adatbányászati eszközökkel, ahol az ügyfelek száma elérte a 190 milliót az adatok mérete pedig a 4 TB-ot. Ilyen méretek mellett már kvadrátikus lépésigényű algoritmusokat sem engedhetünk meg. Látni fogjuk, hogy a legtöbb adatbányászati algoritmus a teljes adatbázist kevés alkalommal olvassa végig.

Skálázható (scalable) és hatékony (efficient) algoritmusokat keresünk, amelyek megbirkóznak nagy méretű adatbázisokkal. Elvárjuk, hogy az adatbázis fontosabb paramétereinek ismeretében az algoritmusok futási ideje megjósolható legyen. Az óriási memóriaméretek miatt a legtöbb elemzendő adatbázis – megfelelő átalakításokkal – valószínűleg elfér a memóriában, de mégis sokszor azt feltételezzük, hogy az adat a háttértáron található.

„Magyar kutatók szerint a mobil pusztítja a spermiumokat.”  
 Forrás: <http://www.origo.hu/tudomany/élet/20040628amobiltelefon.html>

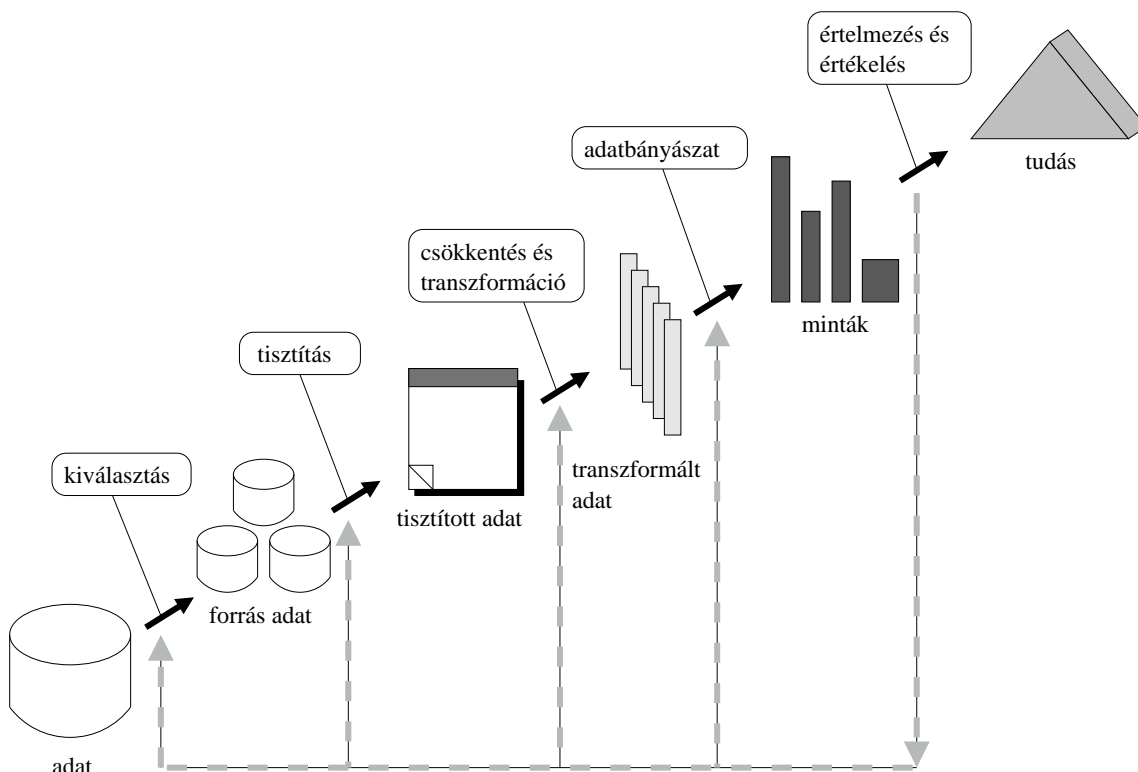
<sup>2</sup>Vannak olyan eredmények is, amelyeket egymástól függetlenül megkaptak az adatbányászat és a statisztika kutatói is. Például döntési fák előállításáról írt négy statisztikus egy közismert könyvet [21]. Eközben egy jeles adatbányász kutató J. Ross Quinlan döntési fa előállító szoftvert készített. A két kutatásban sok közös módszer lelhető fel.



Az adatbázisok méretének növekedése miatt egyre fontosabbak a párhuzamosítható algoritmusok (lásd például partíciós algoritmus rész). Ezek az adatbázist részekre osztják, majd az egyes részeket külön memóriával és háttértárral rendelkező egységek dolgozzák fel, és végül egy kitüntetett egység egyesíti a részeredményeket. Szintén a méretnövekedés az oka azon algoritmusok népszerűségének, amelyek futási ideje nagy mértékben csökkenthető valamilyen előzetes információk (például korábbi futási eredmények) ismeretében (lásd asszociációs szabályok karbantartása rész).

## 1.2. A tudásfeltárás folyamata

A tudáskinyerés folyamata során 6-10 fázist szokás elkülöníteni [39, 55] attól függően, hogy mely lépéseket vonjuk össze (tekinthetjük például az 1.3 ábrát):



1.3. ábra. A tudásfeltárás folyamata

- I. Az alkalmazási terület feltárása és megértése, fontosabb előzetes ismeretek begyűjtése, és a felhasználási célok meghatározása.
- II. Céladatbázis létrehozása: kiválasztani a használni kívánt adatbázist, (vagy annak csak egy részét), amiből a tudást ki akarjuk nyerni.
- III. Adattisztítás: itt olyan alapvető operációkat értünk, mint a téves bejegyzések eltávolítása, hiányos mezők pótlása, zajok szűrése stb. Zajon az adatba épült véletlen hibát értünk. Vannak zajok, amelyeket egyszerű felfedezni és javítani. Például sztring érték ott, ahol

számot várunk, vagy felsorolás típusú attribútumnál érvénytelen érték található. Sajnos sok esetben a hiba észrevétlen marad (például 0.53 helyett 0.35 érték gépelése).

- IV. Adatintegráció: a feldolgozás számára fontos, esetleg elosztott adatbázisok egyesítése. A harmadik és negyedik lépést együtt gyakran nevezik az adatok előfeldolgozásának. A különböző forrásból vett adatok integrációja során sok problémába ütközhetünk. A különböző osztályok különböző módon tárolják adataikat, különböző konvenciókat követnek, különböző mértékegységeket, elsődleges kulcsokat és elnevezést használhatnak és különféle hibák lehetnek jelen. Az egész céget átfogó adatintegrációt adattárházban tárolják, mely egy speciális, az elemzést támogató adatbázis.<sup>3</sup>
- V. Adattér csökkentés: az adatbázisból a cél szempontjából fontos attribútumok kiemelése.
- VI. Adatbányászati algoritmus típusának kiválasztása: eldönteni, hogy a megoldandó feladat klaszterezés, vagy szabály-, illetve mintakeresés, esetleg osztályozás.
- VII. A megfelelő adatbányászati algoritmus meghatározása. Előnyeinek, hátrányainak, paramétereinek vizsgálata, futási idő- és memóriaigény elemzése.
- VIII. Az algoritmus alkalmazása.
- IX. A kinyert információ értelmezése, esetleg visszatérés az előző lépésekhez további finomítások céljából.
- X. A megszerzett tudás megerősítése: összevetés elvárásokkal, előzetes ismeretekkel. Eredmények dokumentálása és átadása a felhasználónak. Egy adatbányászati elemzés eredménye akkor „nem felel meg az elvárásainknak”, ha nem sikerül semmilyen új, hasznos és természetesen valós összefüggést feltárni. Ennek nyilván több oka is lehet, a következőkben két példát mutatunk [25].
1. Előfordulhat, hogy rosszul választottuk meg az elemzéshez (adatbányászathoz) használt algoritmust vagy ennek paramétereit, és egy másik eljárással (vagy más paraméterekkel) találni fogunk valamilyen érdekes összefüggést. Szemléletesen szólva: más oldalról ránézve az adathegyre, lehet, hogy látunk rajta valami érdekeset.
  2. Természetesen az is lehetséges, hogy az adatok egyáltalán nem rejtenek semmiféle új, a gyakorlatban hasznosítható összefüggést. Ekkor — sajnos — teljesen előlről kell kezdeni a folyamatot, új adatok gyűjtésével.

---

<sup>3</sup>A „hétköznapi” működést támogató operatív adatbázis, és az adattárházak közötti különbségre egy szemléletes példa az alábbi [25]: Ha tudni szeretnénk Kis János aktuális számlaegyenlegét, akkor ezt egy operatív adatbázis alapján pontosan és gyorsan meg tudjuk határozni. Egy „átfogóbb” kérdés — például: „Hogyan alakultak az ügyfelek bankban elhelyezett megtakarításai az elmúlt 12 hónapban?” — megválaszolása egy operatív adatbázis esetén bonyolult lehet, és sok ideig tarthat. Egy adattárház az utóbbi kérdésre gyors választ tud adni, támogatva ezáltal a döntéshozókat. A válasz azonban nem teljesen pontos: ha délután 4-kor kérdezzük le az utóbbi 12 hónapbeli megtakarításokat, abban még nem biztos, hogy benne lesz Kis János aznap délelőtt lekötött betétje. Az adattárház adatai tehát nem feltétlenül abszolút frissek, nyilván szükséges azonban a periodikus frissítésük. Adattárházak alkalmazásakor a trendek, folyamatok elemzése a cél. Az, hogy nem az aktuálisan legfrissebb adatokkal dolgozunk, általában nem okoz gondot, feltéve, hogy a legutóbbi frissítés óta nem következett be radikális változás. Ugyanakkor Kis János nyilván nem örülne, ha a betét elhelyezése után este lekérdezve számláját „nem látná” a pénzét, mert a periodikus frissítés csak hetente egyszer esedékes: az ő igényeinek nyilván az operatív adatbázis felel meg.

A sikeres adatbányászati projektekben az első 5 lépés teszi ki az idő- és pénzráfordítások legalább 80%-át. Ha a célok nem kellőképpen átgondoltak és a bányászandó adatok nem elég minőségiek, akkor könnyen előfordulhat, hogy az adatbányász csak vaktában dolgozik és a kinyert információnak tulajdonképpen semmi haszna sincs. A tudásfeltárás során elengedhetetlen, hogy az adatbányász és az alkalmazási terület szakértője szorosan együttműködjön, a projekt minden fázisában ellenőrizték a betartandó irányvonalakat. Nézzünk erre egy példát: ha adatbányászati eszközökkel sikerül kimutatni, hogy X betegséggel gyakran együtt jár Y betegség is, a kutatóorvos képes eldönteni azt, hogy ez valóban így van-e: megvizsgálhatja, hogy ugyanezen összefüggés más adathalmaz esetén is fennáll-e (esetleg direkt ebből a célból gyűjt adatot). Ha igen, akkor kiderítheti azt, hogy az egyik betegség során keletkezik-e olyan kémiai anyag, vagy elszaporodott-e olyan kórokozó, mely hozzájárul a másik betegség kialakulásához. Ezek alapján azt mondhatjuk, hogy az adatbányász „tippeket” ad a kutatóorvosoknak. Ezen „tippek” jelentőségét nem szabad alábecsülnünk: ezek óvhatják meg a kutatóorvost attól, hogy — szemléletesen fogalmazva — „rossz helyen tapogatózzon”. Az adatbányászat tehát első sorban új, ígéretes hipotézisekkel járulhat hozzá a közegészségügyi kutatásokhoz.

A következő valós példa is az adatbányász és a kutatóorvos szerepét szemlélteti. Egy adatbányász az életmódra és a megbetegedésekre vonatkozó adatokat elemezve juthat arra a következtetésre, hogy a prosztatatarák összefügg a szenesedésig sült hús fogyasztásával. Ezzel „irányt mutat” a kutatóorvosnak, aki a háttérben rejlő kémiai reakciókat és azok biológiai következményeit tárja fel. Ez a konkrét esetben lényegében így is történt: előbb tárták fel a jól átsült hús fogyasztása és a prosztatatarák gyakorisága közötti összefüggést, majd megtalálták a hús sütésakor keletkező PhIP vegyületet és kimutatták, hogy hatására prosztatatarák alakulhat ki [62].

Ez a jegyzet a 6. és 7. lépéseket veszi szemügyre: rendelkezésünkre áll egy adatbázis, tudjuk, milyen jellegű információra van szükségünk, és az adatbányász feladata, hogy ennek megoldására minél gyorsabb és pontosabb algoritmust adjon.

Általánosabban kétféle adatbányászati tevékenységet különítünk el:

**Feltárás:** A feltárás során az adatbázisban található mintákat keressük meg. A minták legtöbbször az általános trendeket/szokásokat/jellemzőket írják le, de vannak olyan alkalmazások is (például csalásfelderítés), ahol éppen az általánostól eltérő/nem várt mintákat keressük.

**Előrejelzés:** Az előrejelzésnél a feltárt minták alapján próbálunk következtetni a jövőre. Például egy elem ismeretlen értékeit próbáljuk előrejelezni az ismert értékek és a feltárt tudás alapján.

Négy fontos elvárásunk van a megszerzett tudással kapcsolatban: (1) legyen könnyen érthető, (2) érvényes, (3) hasznos és (4) újszerű. Az érvényesség eldöntése a terület szakértője mellett az adatbányász (esetleg statisztikus) feladata is. Előfordulhat, hogy helyes modellt adunk, az algoritmus is jól működött, mégis a kinyert szabály nem fedti a valóságot. Bonferroni tétele arra figyelmeztet bennünket, hogy amennyiben a lehetséges következtetések száma túl nagy, akkor egyes következtetések tényleges valóságtartalom nélkül igaznak mutatkoznak, tisztán statisztikai megfontolások alapján. Az egyik legjobb példa a valóságtartalom nélküli szabály kinyerésére az alábbi megtörtént eset. Amerikában a Dow Jones átlag becsléséhez keresni kezdték azt a terméket, amely árának alakulása leginkább hasonlított a Dow Jones átlag alakulásához. A kapott termék a bangladesi gyapot volt.

Az adatok illetve a kinyert információk megjelenítésének módja legalább annyira fontos, mint az összefüggések meghatározása. A végfelhasználókat (akik általában vezetők) jobban megragadja egy jól elkészített ábra, mint különböző matematikai struktúrák nyers tálalása. A megjelenítés tehát fontos része az adatbányászatnak. Ezt jól igazolja, hogy nagy sikert könyvelnek el az olyan adatbányászati szoftverek, amelyek adatbányászati algoritmusokat nem is futtatnak, pusztán az adatokat jelenítik meg intelligens módon (háromdimenziós, színes, forgatható ábrák). Ezeknél a rendszereknél az összefüggéseket, mintázatokat, közös tulajdonsággal rendelkező csoportokat maguk a felhasználók veszik észre. Az adatbányászati szoftverekről részletesebben a 13. fejezetben olvashatunk.

### 1.3. Adatbányászat kontra statisztika

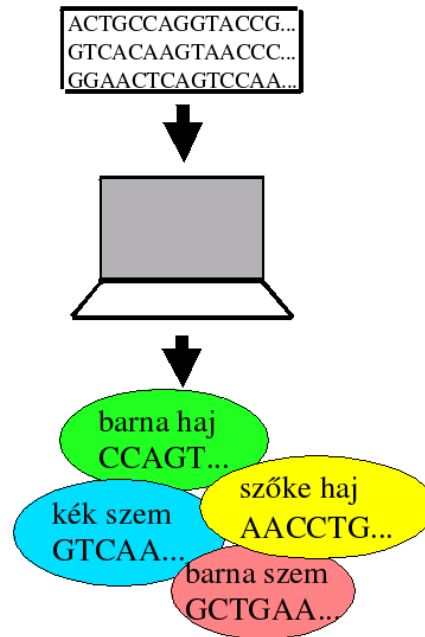
Nehéz definiálni, hogy egy feladat és annak megoldása mikor tartozik a statisztika és mikor az adatbányászat felségterülete alá. A statisztika több hangsúlyt fektet hipotézisek vizsgálatára, míg az adatbányászatban a hipotézisek megtalálásának módja áll a középpontban. Az adatbányászat egy gyakorlatorientált terület, kevesebb súlyt kapnak (sajnos) az elméleti elemzések. Viszont központi kérdés egy algoritmus futási ideje és memóriaigénye. Az adatbányászati algoritmusok bemutatása során kitérünk az adatstruktúrális és akár implementációs kérdésekre is.

Sok kutató az adatbányászatot nem különbözteti meg a gépi tanulástól. Elvégre a gépi tanuláshoz is adatok alapján tanul meg egy koncepciót a gép. Cinikusok szerint az adatbányászat nem más, mint statisztika plusz egy kis marketing. Valóban, nincs éles határ köztük. Úgy általában beszélhetünk adat elemző technikákról. Egyes adat elemző technikákat inkább adatbányászati módszernek mondunk, másokat pedig a statisztikához vagy a gépi tanuláshoz sorolunk.

A 20. század második felétől egyre jellemzőbb a tudományra, hogy bizonyos klasszikus elméletet kiragadnak és új kutatási területnek kiáltják ki. Ugyanígy van ezzel a marketing; ugyanazt a terméket egyszer csak új, hangzatosabb névvel kezdik el értékesíteni. A tudományban is a kutatási feladatokat el kell adni a pályázatokat bíráló zsűrieknek és az új névvel ellátott tudományterület új irányokat sugall; az új irányzatok és élbéli kutatások pedig nagy támogatást kapnak. Ez a tény jelentősen hozzájárult az adatbányászat elterjedéséhez és az egyes adatelemző feladatok "adatbányászati" címkével való ellátásához.

Adatbányászathoz soroljuk a klaszterezés, osztályozás, asszociációs szabálykinyerés és az idősorelemzés nem klasszikus (pl. regressziószámítás, simítás) feladatait. A következőkben néhány példán keresztül szemléltjük az adatbányászat és a statisztika közötti különbséget és egyben a két terület rokonságát is [25].

- I. Tegyük fel, hogy egy adatbázisban sokmillió ember DNS-szekvenciáit és tulajdonságait tároljuk (1.4 ábra). Egy jellegzetes statisztikai kérdés lehet az, hogy például a kék szemű emberek mekkora részére jellemző egy adott DNS-szekvencia. Természetesen olyan kérdést is feltehetünk, melynek megválaszolása ennél kifinomultabb eszköztárat igényel: ha azt szeretnénk tudni, van-e szignifikáns függés egy adott DNS-szekvencia megléte és a „kék szem” tulajdonság között, statisztikai próbát alkalmazhatunk ennek eldöntésére. Egy adatbányász nem kérdezné rá egy konkrét szekvencia és egy konkrét tulajdonság közötti összefüggésre, hanem egy általánosabb kérdést tenne fel, például azt, hogy mi-



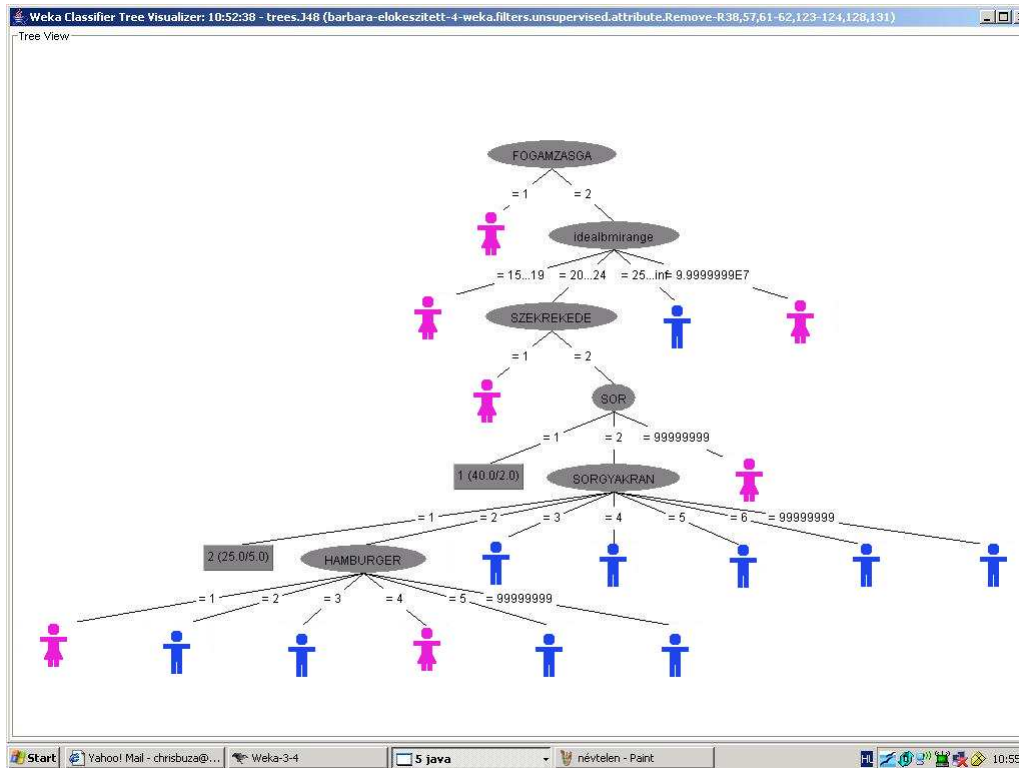
1.4. ábra. Egy jellegzetes adatbányászati feladat: DNS-szekvenciák elemzése

Ilyen összefüggés van a tulajdonságok és szekvenciák között, melyik tulajdonságért melyik szekvencia felelős?

II. Egy másik példa az adatbányászat és statisztika közötti különbségre az alábbi: egy statisztikai elemzés során megvizsgálhatjuk, hogy a nők illetve férfiak hány százaléka dohányzik, fogyaszt rendszeresen nagy mennyiségben alkoholt, van-e szignifikáns eltérés a két csoport között. Egy adatbányászati elemzés során itt is általánosabb kérdést tennénk fel, például azt, hogy milyen jellegzetes csoportok vannak az alkoholfogyasztásra és dohányzásra nézve? Tehát azt nem mondjuk meg előre, hogy az egyik csoportba a nők, a másikba pedig a férfiak tartoznak. Az adatbányász feladata, hogy úgy csoportosítsa az embereket (rekordokat), hogy a hasonlóak egy csoportba, a különbözők pedig különböző csoportba kerüljenek. (Ez egy klaszterezési feladat.) Az adatbányászatban az ilyen feladatokat nem hosszas emberi munka és intuíció árán oldjuk meg, hanem törekszünk a minél nagyobb fokú automatizálásra kifinomult szoftverek alkalmazásával. Eredményként könnyen lehet, hogy nem a nemek szerinti csoportosítást kapjuk, hanem egy olyat, melyben ugyanazon csoportokba férfiak és nők is kerültek, akik — egyéb tulajdonságaik alapján — hasonlóak.<sup>4</sup>

III. Az előbbi példában természetesen más irányba is „általánosíthatjuk” a statisztikai elemzés során feltett kérdésünket: lehet, hogy arra vagyunk kíváncsiak, hogy mi a különbség a férfiak és a nők között. Ismerjük tehát a két csoportot, de nem tudjuk, hogy mely tulajdonságok vagy tulajdonságkombinációk jellemzőek egy-egy csoportra. Ekkor egy

<sup>4</sup>Ahhoz, hogy egy ilyen elemzés sikeres legyen, nagyon fontos a hasonlósági mérték megfelelő megválasztása, valamint az elemzésbe bevont attribútumok (adattábla-oszlopok) „ügyes” kiválasztása. Ha például az alkoholfogyasztásra és dohányzásra vonatkozó adatok mellett „túl sok” további attribútumot vonunk be a vizsgálatba, akkor lehet, hogy a csoportosítás nem az alkoholfogyasztásra és dohányzásra vonatkozó jellegzetes csoportokat tartalmazza, hanem „általános” csoportokat kapunk.



1.5. ábra. Döntési fa: nők és férfiak közötti különbségek a Semmelweis Egyetem hallgatóinak körében végzett felmérés alapján.

osztályozási feladattal állunk szemben, a csoportokat osztályoknak nevezzük. Ezt a kérdést egyébként fel is tettük a Semmelweis Egyetem hallgatóinak körében végzett egyik felmérés adatbázisán. Az eredmény az 1.5. ábrán látható. Ez egy döntési fa. A levelek az osztályoknak (nők illetve férfiak) felelnek meg. A fa közbülső csomópontjaiban egy-egy attribútum (adattáblabeli oszlop) neve látható. A fa egy csomópontjából kiinduló ágak az adott csomóponthoz tartozó attribútum egy-egy lehetséges értékének felelnek meg. Egy döntési fa azt mutatja meg, hogy ha nem ismernénk, hogy egy rekord melyik osztályba tartozik, akkor hogyan dönthetnénk ezt el. Például a fogamzásgátlót szedő hallgatók nők (pontosabban: azon rekordok, amelyek FOGAMZASGA attribútuma „1” értékű, a női hallgatók osztályába tartoznak).<sup>5</sup>

## 1.4. Sikeres alkalmazások

Az „adat bányászata” eredetileg statisztikusok által használt kifejezés, az adatok nem kellőképpen megalapozott felhasználására, amely során valaki helytelen következtetést von le. Igaz ugyanis, hogy tetszőleges adathalmazban felfedezhetünk valamilyen struktúrát, ha elég sokáig nézzük az adatot. Ismét utalunk a lehetséges következtetések nagy számából eredő veszélyre. A helytelen következtetésre az egyik leghíresebb példa az alábbi: Az 50-es években

<sup>5</sup>A döntési fa építéskor általában nem követelmény, hogy egy levélbeli összes rekord ugyanazon osztályba tartozzon, elég, ha „nagy részük” azonos osztályba tartozik. Ebben a konkrét példában az összes fogamzásgátlót szedő hallgató nő volt.

David Rhine parapszichológus diákokat vizsgált meg azzal a céllal, hogy parapszichológiai képességgel rendelkezőket találjon. Minden egyes diáknak 10 lefedett kártya színét kellett megtippelnie (piros vagy fekete). A kísérlet eredményeként bejelentette, hogy a diákok 0,1%-a parapszichológiai képességgel rendelkezik (a teljesen véletlenszerűen tippelők között a helyesen tippelők várható száma statisztikailag nagyjából ennyi, hiszen annak valószínűsége, hogy valaki mind a tíz kártyát eltalálja  $\frac{1}{2^{10}} = \frac{1}{1024}$ ). Ezekkel a diákokkal újra elvégezte a kísérletet, ám ezúttal a diákok eredménye teljesen átlagos volt. Rhine következtetése szerint az, aki parapszichológiai képességgel rendelkezik és erről nem tud, elveszti eme képességét, miután tudomást szerez róla.

A fenti példa ellenére mára az adatbányászat szó elvesztette jelentésének negatív tartalmát, a számos sikeres alkalmazásnak köszönhetően. A teljesség igénye nélkül felsorolunk belőlük néhányat.

- A bankok egyre gyakrabban alkalmaznak olyan automatikusan előállított döntési fákat, amelyek alapján egy program javaslatot tesz egy hitel megítéléséről. Ezt a kérelmezők személyes, továbbá előzetes hitelfelvételi és törlesztési adatai alapján teszi (osztályozás) [132]. Tesztek például igazolták, hogy a hitelbírálat minősége javult az USA-ban, amikor a bankok áttértek a kötelezően alkalmazott, írásban rögzített szabályok alkalmazására [132]. Ezeket a szabályokat pedig az adatbányászat segítségével állították össze.
- A vásárlói szokások felderítése szupermarketekben, illetve nagy vevőkörrel rendelkező áruházakban hasznos lehet az áruház terméktérképének kialakításánál, akciók, eladáshelyi reklámok (Point of Sales, Point of Purchase), leárazások szervezésénél... (asszociációs szabályok).
- Az ember genotípusának elemzéséhez a gének nagy száma miatt szintén adatbányászati algoritmusok szükségesek. Az eddigi sikeres kísérletek célja olyan géncsoportok feltárása volt, amelyek a cukorbetegség bizonyos változataiért felelősek. A teljes emberi génrendszer feltárásával ez a terület egyre fontosabb lesz.
- Az on-line áruházak a jövőben egyre elfogadottabbak és elterjedtebbek lesznek. Mivel az on-line kereskedelemben nem használhatóak a megszokott személyes marketing eszközök a forgalom (és a profit) személyre szabott vásárlási ajánlatokkal növelhető. Az ajánlatokat az eddigi vásárlási adatok és a rendelkezésre álló demográfiai adatok elemzése alapján tehetjük meg (epizód kutatás, asszociációs szabályok).
- A csillagászatban az égitestek óriási száma miatt a hagyományos klaszterező algoritmusok még a mai számítási kapacitások mellett sem képesek racionális időn belül különbséget tenni galaxisok, közeli csillagok és más égi objektumok között. Az újabb, kifinomultabb algoritmusok futási ideje jóval kevesebb, ami lehetővé teszi a klaszterezést (klaszterezés).
- Utazás szervezéssel kapcsolatos minták kinyerésével hatékonyabban (és ennek következtében nagyobb nyereséggel) megszervezhetőek a nagy költségfaktorú tényezők, pl. szállodai szobák, repülőjegyek leárazása, vagy áremelése (epizód kutatás, gyakori minta).
- Kifinomult gyártási folyamatok során gyakran a beállítási paraméterek finomhangolására van szükség. A kőolaj és a földgáz szétválasztása az olajfinomítás szükséges előfeltétele, de az elválasztási folyamat kontrollálása nem könnyű feladat. A British Petroleum

olajvállalat a gépi tanulás technikáját használta a paraméter-beállítás szabályainak megalkotására. Most ez tíz percet vesz igénybe, míg korábban szakértők több, mint egy napi munkáját vette igénybe.

- A Westinghouse cég nukleáris tüzelőanyag-cellák gyártása során ütközött problémákba, és szintén a gépi tanulás segítségével hoztak létre folyamatkontrollálási szabályokat. Ezzel 10 millió USD-t sikerült megspórolniuk az 1984-es évben. A Tennessee állambeli R.R. Donelly nyomdaipari cég ugyanezt az ötletet alkalmazta a retogravúr nyomdagépek irányítására, így csökkentve a hibás paraméter-beállítások következtében keletkező selejtes nyomatok számát évi 500-ról 30-ra.
- A vírusölő programok az ismert vírusokat lenyomataik alapján detektálják, az ismeretleneket pedig többnyire valamilyen heurisztikus módon próbálják kiszűrni. Osztályozó algoritmusok felhasználásával az ismert vírusok tulajdonságai alapján olyan modellt lehet felállítani, ami jól leírja a vírusok tulajdonságait [120, 121]. A modellt sikeresen alkalmazták új ismeretlen vírusok kiszűrésére (osztályozás).

További esettanulmányokról a 13.3.2 részben olvashatunk.

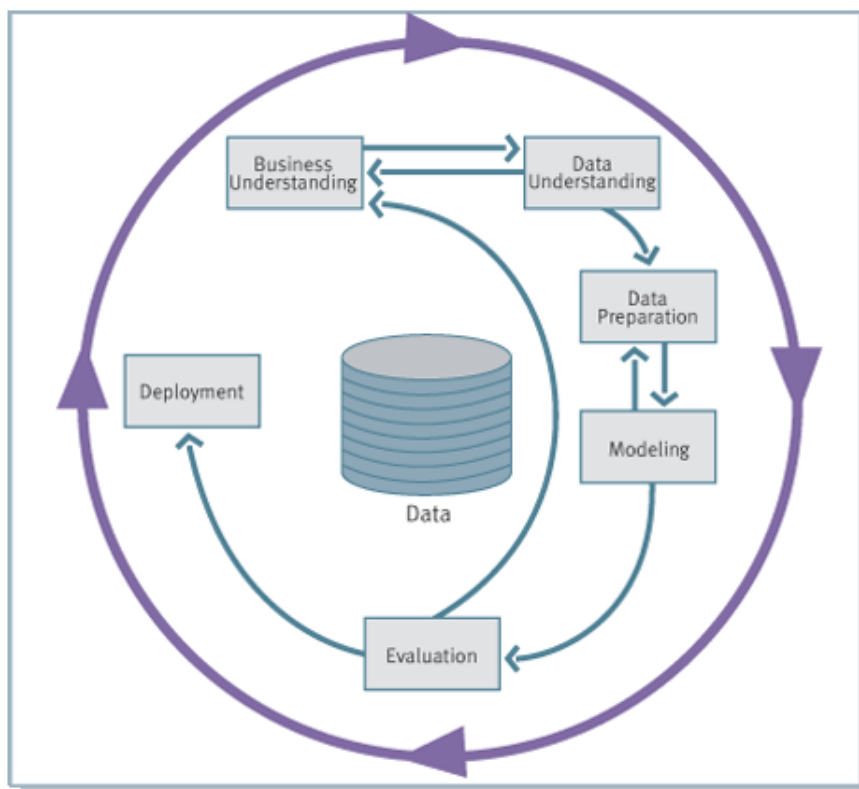
A fentiekben a sikeres alkalmazásokat ismertettük. A következőben további alkalmazásokat mutatunk be. Célunk, hogy szemléltessük a diszciplína kiterjedtségét és aktuális állását.

- Az emberi mesterséges megtermékenyítés során petesejtek sokaságát gyűjtik össze a női petefészekből. Ezeket a partner, vagy donor spermáival megtermékenyítve számos embrió fejlődik ki. Közülük néhányat kiválasztanak, és az anyaméhbe ültetnek. A problémát a leginkább életképes, legjobb túlélési esélyekkel rendelkező embriók kiválasztása jelenti. A kiválasztás az embriók körülbelül hatvan rögzített jellegzetességén – a magzat morfológiáján, oocita-, tüszősejt- és spermamintákon – alapszik. A jellemzők számossága elegendően nagy ahhoz, hogy túl bonyolult legyen az embriológusoknak valamennyit párhuzamosan megbecsülni és összefüggést találni a múltbéli esetek kezdeti jellemzői és azok kimenetele között, azaz, hogy az embrióból végül életképes csecsemő született-e vagy sem. Egy angol kutatási projekt arra irányuló kutatást folytat, hogy hogyan lehet a kiválasztást gépi tanulással – az embriók rögzített adatait tanítóhalmazként használva – megvalósítani.
- Az új-zélandi tejgazdaságoknak minden évben kemény üzleti döntést kell meghozniuk: ki kell választani, hogy a szarvasmarha állomány mely egyedeit tartják meg, és melyeket értékesítik vágóhidaknak. Tipikusan minden gazdaság ötödik egyede kerül mészárszékre a fejési idény végén, ahogy az élelmezési tartalékok kiapadnak. A döntést az egyes példányok tenyészádatai és múltbéli tejtermelékenységi mutatója befolyásolja. További kritikus faktorok az egyed kora (egy példány kb. 8 évesen éri el produktív korszakának végét), kórtörténete, szülési komplikációk, nemkívánatos jellemvonások (agresszivitás, kerítés átugrása), illetve az, hogy a következő szezonban vemhes-e. Több millió szarvasmarha egyedenként több mint 700 tulajdonságát rögzítették az évek során. A kutatók azt vizsgálják, hogyan használható fel a gépi tanulás annak megállapítására, hogy a sikeres farmerek mely faktorokat veszik számításba a szelektálásnál. Ezzel nem a döntési folyamat gépesítése a céljuk, hanem a sikerstratégia kitanulása, és annak közkinccsé tétele.



## 1.5. Szabványok

Kezdetben sok adatbányászati projektre jellemző volt, hogy az adatbányászok megkapták az adatokat és némi információt az alkalmazási területről és cserébe várták tőlük a kincset érő információkat. A szoros együttműködés hiánya azonban csak olyan információkhoz vezetett, amelyekkel az alkalmazási terület embererei nem sok mindent tudtak kezdeni. Az adatbányászat elterjedésével (és a minőségbiztosítási elvárásokkal) fellépett az igény, hogy legyen egy szabvány, egy útmutató az adatbányászati projektek lebonyolításáról. Így született meg a CRISP-DM (CROSS Industry Standard Process for Data Mining) [28], amely adatbányászati eszköztől és felhasználási területtől függetlenül leírja, hogy miként kellene kinéznie egy adatbányászati projektnek, illetve ismerteti a kulcsfontosságú lépéseket, és a potenciális veszélyeket. A CRISP-DM szerint a tudáskinyerés az 1.6 ábra szerinti módon jön létre.



1.6. ábra. A tudásfeltárás folyamata a CRISP-DM szerint

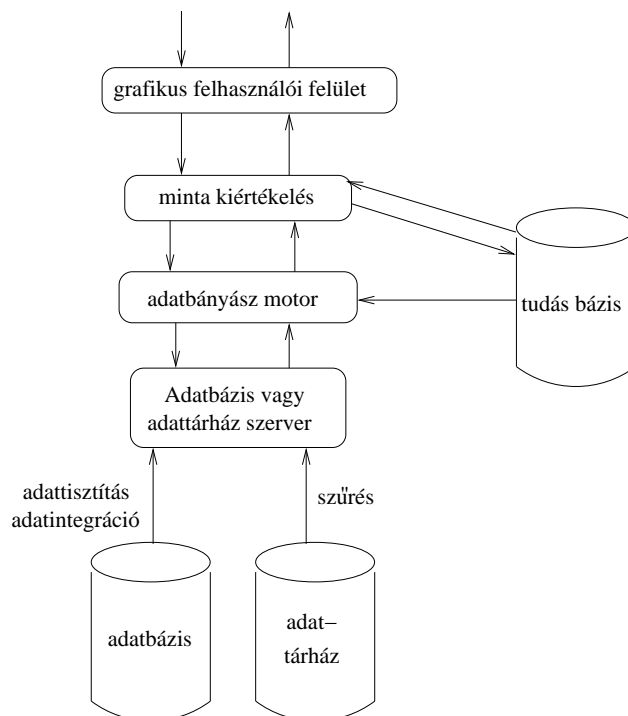
Az adatbányászati folyamat szabványosítása mellett egyre nagyobb az igény a folyamat egyes lépéseiben felmerülő megoldások, problémák, eszközök szabványosítására. Ezek közül a legismertebbek:

- az XML alapú PMML (Predictive Modeling Markup Language), amely az adatbányászati eredmények szabványos leírását szolgálja,
- a Microsoft analysis server adatbányászati funkciókkal kibővített szabványa (OLE DB for data mining),

- az ISO törekvései multimédia és alkalmazás specifikus SQL típusok és a hozzá tartozó eljárások definiálására (SQL/MM)
- java adatbányászati API (JDMAPI)

## 1.6. Adatbányászati rendszer architektúrája

Egy adatbányászati rendszernek kapcsolatban kell lennie az adatbázissal, a felhasználóval és esetleg valami tudásalapú rendszerrel. Ezek alapján egy tipikus adatbányászati architektúra az 1.7. ábrán látható.



1.7. ábra. Tipikus adatbányászati rendszer architektúrája

**Adatbázis, adattárház vagy más információ raktár:** Itt található a tényleges adatok, ami lehet egy adatbázis, vagy adattárház, akár egy munkalap vagy bármilyen tárolt információ. Az adattisztítás és integráció közvetlenül az adatokon is elvégezhető.

**Adatbázis vagy adattárház szerver:** A szerver felelős a felhasználó által kért adat kézbesítéséért.

**Tudás bázis:** A területre jellemző, valamilyen szinten formalizálható tudás található itt. Fontos szerepe lehet ennek a keresési tér szűkítésénél, a kinyert minták érdekességének meghatározásánál, különböző paraméterek és küszöbszámok meghatározásánál.

**Adatbányász motor:** Az adatbányász motorban futnak a különböző adatbányászati algoritmusok.

**Minta kiértékelő modul:** Ez a modul felelős a kinyert minta vagy összefüggések kiértékeléséért a területre jellemző érdekességi mutatók alapján. Sokszor látni fogjuk, hogy minél jobban egybe tudjuk építeni az adatbányászatot a minta kiértékelésével, annál hatékonyabb és gyorsabb lehet a tudásfeltárás.

**Grafikus felhasználói felület:** Itt zajlik a kommunikáció a felhasználó és az adatbányászati rendszer között. A felhasználó itt adhatja meg, hogy melyik adatbázisban milyen jellegű összefüggéseket keres és ezen a rétegen keresztül láthatja a végeredményt. Az összefüggések átlátható, értelmes tálalása rendkívül fontos, hiszen ennek hiánya elriaszthatja a felhasználót az adatbányásztól.

## 1.7. Adatbányászat és az etika

Az internet széleskörű terjedésével, és a modern technikák megjelenésével a jogi szabályozás sokszor képtelen lépést tartani. Ilyen terület az adatbányászat is, mint felfutó tudományág. Bár a személyes adatok védelméről már születtek jogszabályok, ezek nem minden esetben teljesszerűek, illetve még így is sok etikai problémát hagynak nyitva. Elmondható, hogy a törvény nagyrészt azt szabályozza, ami egyben etikátlan is, így ebben a részben<sup>6</sup> éljünk azal a feltételezéssel, hogy minden törvényesen történik az adatbányászati projekt során.

Először is sokak által kifogásolt terület az adatok begyűjtése. A mai ember lépten-nyomon információkat szór szét magáról: minden üzletben kamerák vannak elhelyezve, szörföl az interneten, bárhol használhat hitelkártyát, törzsvásárlói kártyát. Ezen tevékenységek során rengeteg adatbázisban hagy magáról információt. A személyes információk begyűjtésénél - a törvény szerint - az érintetteket tájékoztatni kell arról, hogy ki, milyen célból fogja azt feldolgozni. Ezt a célt az adatbányászatnál nem lehet előre azonosítani, mert az elemzés előtt nem tudjuk megmondani, hogy milyen információkat fogunk kinyerni az adatbázisból, és néhány esetben még azt sem tudjuk garantálni, hogy az eredmény egyáltalán felhasználásra kerül. Ez is mutatja, hogy az adatok begyűjtésére új szabályozásokat kellene alkotni. Erre egy lehetséges megoldás, ha az érintett azt döntheti el, hogy – általánosan – adatbányászati célra felhasználhatóak-e az adatai. Érezhető azonban, hogy amíg az ügyfél nincs tisztában azzal, hogy az adatokból mire tudnak egyáltalán következtetni, és mire tudják ezt a következtetést felhasználni, addig kevésbé fog élni a lehetőségével, hogy letiltsa adatainak felhasználását vagy tárolását. Az pedig még nyilvánvalóbb, hogy a cégeknek nem áll érdekükben önként megosztani ezeket az információkat az ügyféllel, vagy akár a versenytársakkal [100].

Az osztályozási feladatok kapcsán másfajta morális problémák is felléphetnek: egyrészt magát a módszertant támadják az ellenzői, másrészt könnyen felhasználható diszkriminációra.

Sokan azt a hozzáállást tartják aggasztónak, hogy előzetes személyes ismeretség nélkül soroljanak be egyéneket bizonyos csoportokba. Például az amerikai állampolgárok egy felmérés során kifogásolták a hitelbírálatban bevezetett automatikus döntéshozást, mert úgy érzik, így kevésbé foglalkoznak velük, személytelenebbé vált a rendszer, csak egy adathalmazként tekintenek rájuk. A személyes ismeretség nélküli besorolás során fennáll a veszélye, hogy rossz eredményt ad a rendszer, főként, ha kevés attribútummal dolgozunk. Próbáljuk meg például az embereket gyerek-felnőtt csoportokba besorolni a kor attribútum ismerete nélkül. A legkézenfekvőbb megoldás a testsúlyuk és testmagasságuk szerint osztályozni őket. Ez a módszer viszont sok esetben

---

<sup>6</sup>Ezt a részt Huczman Zsuzsanna írta.

téves eredményhez vezet.

A diszkriminációnak az asszociációs szabályok kinyerése adhat teret. Érzékeny adatokat – mint például a vallás, faji hovatartozás, szexuális irányultság – tilos feldolgozni, egy egyszerű adatbányászati algoritmussal viszont nagyon könnyen megoldhatónak tűnik a vásárlói szokásokból a nemi hovatartozás, vagy akár egy megfelelő kérdőív esetében a faji hovatartozás megállapítása. Az, hogy ezek az adatok ne kerüljenek felhasználásra, a projekt vezetőjének felelőssége. Az előbb említettek „kényesebb” területek, de beláthatjuk, hogy már az is diszkriminációnak számít, ha egy bizonyos vásárlói csoportot előnyökhöz juttatunk egy olyannal szemben, amely sokkal kevesebb profitot ígér elemzéseink alapján. A pozitív diszkriminációra jó példa, amikor a telefonszolgáltatók ajándékot ajánlanak fel új előfizetőik számára [142].

Ugyanakkor Európában az automatikus döntési fák alkalmazása a hitelbírálatban pont a diszkrimináció kiküszöbölése a rendszerből, hiszen a jogszabály szerint a matematikai háttérrel az ügyfél kérésére fel kell fedni. (Ez az „USA-ban” például nem kötelessége a banknak, így lehetősége van diszkriminálni.)

További etikailag megkérdőjelezhető, érdekes felvetések:

- A külön pontok kezelése egy másik adatbányászati terület, ahol felmerül, hogy jogában áll-e bárkinek meghatározni, hogy mi tekinthető normális viselkedésnek, illetve előkerülhet a kvantumelméletből ismert tétel: miszerint a megfigyelt rendszer viselkedése a megfigyelés tényétől is megváltozik.
- Hazánkban az adatbányászatot – annak ára miatt – főként a marketing területén használják fel: célzott reklámok küldésére (direkt marketing); az akciós termékek meghatározására gyakran vett termékcsoporthoz alapján; új tarifacsomagok bevezetésére. Felmerül a kérdés, hogy a személyes adatokat, amiket elvileg azért tárolnak pl. telefonszolgáltató esetén, hogy számlázni tudjanak [131], etikus-e profitszerzésre használni?
- Ha egy adatot kivesznek egy adott adathalmazból - az egyén kérésére, - aminek mérete pl. a vásárlói kosarakból való adatbányászat esetén akár  $10^9$  rekord is lehet, ettől még az elemzés elvégezhető, és az egyén érdekeit nem sérti, de ő is részese lesz a következményeknek, pl. direkt-marketing ajánlatokat kaphat.

Végül két példa az adatbányászat etikai vonatkozásaira:

Az emberi DNS által hordozott információ kinyerése tipikusan adatbányászati feladat. Rendkívül érzékeny adatról lévén szó, az adatbázis, amivel eddig az adatbányászok dolgoztak, nem egy konkrét személy DNS szekvenciáját tartalmazták, hanem több rövidebb szakaszból álló DNS-láncokat. Nemrégiben viszont egy kutató előállt a saját DNS-szekvenciájának kulcsával, vagyis felfedte, hogy mi van belekódolva. Eddig ez technikailag és törvényesen nehezen volt megvalósítható. Beláthatjuk, hogy nem sok értelme van vizsgálni, hogy több töredék DNS milyen információt hordoz, hiszen pont azon lenne a hangsúly, hogy létező előfordulásokat és összefüggéseket lehessen vizsgálni [85]. A kutató ezzel a lépésével nagyban hozzájárult a DNS-szerkezet megismeréséhez, példájából láthatjuk, hogy néha a privacy feladása kell egy-egy eredmény eléréséhez.

Adatbányászati projekteknél egyébként is gyakori, hogy az adathalmaz, amin az elemzést el kellene végezni, annyira titkos, vagy érzékeny adatokat tartalmaz, hogy a cégek nem is adják

„A krónikus fertőzések depressziót és skizofréniát okozhatnak – állítják német kutatók.” Forrás: [http://hvg.hu/egeszseg/20070512\\_depresszio\\_skizofrenia.aspx](http://hvg.hu/egeszseg/20070512_depresszio_skizofrenia.aspx)

ki a kezükből. Ilyenek pl. a távközlési vállalatok vagy bankok adatai, amik a konkurencia számára nagy értékkel bírnak. Ezekben az esetekben megoldás lehet, hogy a cég maga generál a rendelkezésére álló adatokból egy új adatbázist, - amiknek nyilván bizonyos követelményeknek eleget kell tenniük - és ezt kapja meg a kutatócsoport elemzésre. Itt a cég felelőssége, hogy ellenőrizze, hogy az eredeti adathalmazban ugyanazok az összefüggések fennállnak-e, mint a vizsgált adathalmazban.

A másik példa az adatbányászat és az etika ütközésére az a szinte már utópisztikus projekt, amit a NASA és a Northwest Airlines indított el: a Washington Times 2002 nyarán számolt be egy új információ-technológiai megvalósításról: a két „szervezet” olyan alkalmazást fejlesztett, és üzemeltet, amely képes előrejelezni, megjósolni az utasok várható viselkedését a repülőtéren, illetve a repülőn. A technológia alapja egy különleges műszer és egy pszichológusok segítségével készített program ötvözött alkalmazása nagy adatbázisokon. A kapu képes érzékelni a rajta áthaladó személy elektromágneses agyhullámain, a szívritmusát, a pislogását és a testhőmérsékletét: gyakorlatilag egy „szuper hazugságvizsgáló” berendezés. Ezeket az adatokat analizálva és összevetve különböző adatbázisokkal, mint például bűnügyi nyilvántartással, a gép jelez, ha kockázatot érzékel. Ilyenkor a biztonsági őrök még mérlegelhetnek, mint ahogy egyébként is tennék, így a rendszert védők arra hivatkozhatnak, hogy csak döntéstámogatást végeznek. Az is mellettük szól, hogy az Egyesült Államokban a repülőtereken a titkos megfigyelés végzése nem törvénybe ütköző cselekedet. Az USA-ban a személyes adatok védelmére kisebb hangsúlyt fektetnek a terrortámadások óta. A lakosság elfogadta azt a nézetet, hogy a biztonságuk érdekében áldozzák fel a személyes adataikat.

## 1.8. Az adatbányászat feltételei

Tagadhatatlan, hogy a sikertelen adatbányászati projektek száma nagy, és az adatbányászat nagyon sok esetben nem váltotta be a hozzá fűzött reményeket. Ennek oka egyrészt az adatbányászati szakemberhiány (a jó adatbányászati szakember ritka, mint a fehér holló), másrészt az, hogy alapvető feltételek nem teljesültek a projektek során. A sikeres adatbányászati projekt egyik legfontosabb feltétele az adatbányász és a terület szakértőjének szoros együttműködése. A további feltételek az alábbiak:

**Nagy mennyiségű adat:** A nagy mennyiségű adat a kinyert szabályok statisztikai jelentőségét növeli. Minél nagyobb az adatmennyiség, annál biztosabban tudjuk kizárni bizonyos összefüggések esetiségét, azaz annál kisebb az esélye, hogy a talált összefüggés csak a véletlen eredménye. Sajnos sok adatot sokáig tart feldolgozni, sőt az algoritmusok egy jelentős része érzékeny arra, hogy az adatbázis elfér-e a memóriában.

**Sok attribútum:** Ha az objektumokat leíró attribútumok száma kicsi, akkor hagyományos eszközökkel (grafikonok, egyszerű táblázatok, kis dimenziós, forgatható, színes ábrák, stb.) is fel tudjuk tární a tudást. Kevés attribútum esetén a kinyerhető tudás sem lehet túl sokféle. Az adatbányászat ereje akkor mutatkozik meg, amikor az attribútumszám olyan nagy, hogy a hagyományos módszereknek nincs esélyük.

**Tiszta adat:** Az adatok jó minősége az adatbányászat egyik alapfeltétele. A zajok, a hibás bejegyzések jó esetben csak nehezítik az adatbányászatot (például amikor ismerjük az adatokban található zaj, ill. bizonytalanság fokát), rosszabb esetben azonban hamis

eredményekhez vezetnek. Az ilyen rossz minőségű adatokra remek példa hazánk orvosi adatbázisa (rengeteg hibás bejegyzés, kitöltetlen mező, eltérő mértékegység alapú bejegyzések, szöveges bejegyzések), pedig az ezekből kinyert információk értékesek lennének. A "szeméthalmazban" való kutakodást tréfásan GIGO-nak (garbage in, garbage out<sup>7</sup>) nevezik.

**Torzítatlan adat:** Az adatbányászat sikeressége múlhat az adatok nem megfelelő kiválasztásán. Ide tartozó fogalom az ún. BIBO (bias in, bias out<sup>8</sup>), amely arra hívja fel a figyelmünket, hogy ha egy részsokaság alapján akarunk következtetni az alapsokaságra, akkor figyelembe kell vennünk a részsokaság kiválasztásának szempontjait, illetve az abból adódó (esetleges) torzításokat. Például, ha a lakosságot az anyagi helyzet szerint akarjuk csoportokba sorolni, de csak nyugat-magyarországi adatok állnak rendelkezésünkre, akkor tudnunk kell, hogy a kapott eredmény (a csoportok leírása) torz lesz, hiszen a részsokaság átlag életszínvonala jobb az alapsokaságénál.

**Alkalmazási terület akcióképessége:** Gyakran előfordul, hogy a tudást csak kinyerik, de a felhasználása elmarad. Gyakran a felhasználási területek túl merevek, vagy a változtatás túlságosan magas költségekkel járna. A legtöbb adatbányászati esettanulmányban a tudás kinyerésének módjáról esik szó, a tudás felhasználásáról pedig ritkán hallunk.

**A befektetés megtérülésének (Return On Investment) mérhetősége:** Egy adatbányászati projektről akkor állíthatjuk biztosan, hogy sikeres, ha a befektetés hatását mérni, vagy viszonylag pontosan becsülni tudjuk.

A jegyzet fejezeteiben a legkevésbé ismert, de napjainkban egyre nagyobb teret nyerő területeket járjuk körül: a gyakori minták kinyerését, az attribútumok közötti összefüggések meghatározását, a sorozatelemzést, a klaszterezést és a webes adatbányászatot. Minden esetben az algoritmusok gyakorlati felhasználását példákon keresztül szemléltetjük; emellett megadjuk a problémák formális definícióit, és bemutatjuk a legismertebb, leghatékonyabb algoritmusokat is. A jegyzet további célja, hogy összefoglalja az eddig nem, vagy csak kis hatékonysággal megoldott problémákat, továbbá a jelenlegi kutatási területeket.

---

<sup>7</sup>szemét be, szemét ki

<sup>8</sup>torzítás be, torzítás ki

## 2. fejezet

# Alapfogalmak, jelölések

Ebben a részben tisztázzuk a jegyzet során használt fogalmak jelentését. Célszerű akkor átnéznünk e fejezet egyes részeit, amikor az olvasás során olyan részbe ütközünk, ami nem teljesen tiszta.

### 2.1. Halmazok, relációk, függvények, sorozatok

A *halmaz* különböző objektumok együttese, amelyeket a halmaz *elemeinek* hívunk. Ha  $x$  eleme a  $H$  halmaznak, akkor azt így jelöljük:  $x \in H$ , a halmaz elemeinek számát (rövidebben *elemszámát*) pedig  $|H|$ -val. A jegyzetben a természetes számok halmazát  $(\{0,1,\dots\})$   $\mathbb{N}$ -el jelöljük, a valós számok halmazát  $\mathbb{R}$ -el, az egész számok halmazát  $\mathbb{Z}$ -vel, az üres halmazt (egyetlen elemet sem tartalmazó halmaz)  $\emptyset$ -val. Két halmaz akkor egyezik meg, ha ugyanazok az elemeik.  $X$  részhalmaza  $Y$ -nak ( $X \subseteq Y$ ), ha  $X$  minden eleme  $Y$ -nak is eleme. Ha  $X \subseteq Y$ , de  $X \neq Y$ , akkor  $X$  *valódi részhalmaza*  $Y$ -nak. A valódi jelzőt gyakran fogjuk használni, és a valódi részhalmaz analógiájára azt értjük rajta, hogy az egyenlőséget kizárjuk. Sajnos a *super*-set angol szónak nincsen általánosan elfogadott fordítása, pedig sokszor szeretnénk használni. Azt fogjuk mondani, hogy  $Y$  *bővebb*  $X$ -nél, ha  $(X \subseteq Y)$ . A halmazműveletek jelölése és pontos jelentésük: metszet:  $X \cap Y = \{z : z \in X \text{ és } z \in Y\}$ , unió:  $X \cup Y = \{z : z \in X \text{ vagy } z \in Y\}$ , különbség:  $X \setminus Y = \{z : z \in X \text{ és } z \notin Y\}$ .

Két halmaz  $(X, Y)$  *Descartes-szorzata*  $(X \times Y)$  az összes olyan rendezett párból álló halmaz, amelynek az első komponense (tagja)  $X$ -ben, a második  $Y$ -ban van. Az  $X, Y$  halmazokon értelmezett *bináris reláció* az  $X \times Y$  részhalmaza. Ha  $(x, y)$  eleme a  $\phi$  relációnak, akkor azt így is jelölhetjük:  $x\phi y$ . A  $\preceq$  reláció *részben rendezés* (vagy *parciális rendezés*), ha *reflexív* ( $x \preceq x$ ), *antiszimmetrikus* ( $x \preceq y$  és  $y \preceq x$  feltételekből következik, hogy  $x = y$ ), *tranzitív* ( $x \preceq y$  és  $y \preceq z$  feltételekből következik, hogy  $x \preceq z$ ). Ha az előző 3 feltételben az antiszimmetrikus helyett szimmetrikusat ( $x \preceq y$ -ből következik, hogy  $y \preceq x$ ) mondunk, akkor *ekvivalencia-relációról* beszélünk. A továbbiakban, tetszőleges  $\preceq$  rendezés esetén, ha  $x \neq y$  és  $x \preceq y$ , akkor azt így jelöljük  $x \prec y$ . Legyen  $X$  részhalmaza  $X'$ . A  $X'$  halmaznak  $y \in X$  egy *alsó korlátja*, ha  $y \preceq x$  minden  $x \in X'$ -re. Az  $y$  *legnagyobb alsó korlát*, ha minden  $y'$  alsó korlátra  $y' \preceq y$ . Az  $y$  *maximális alsó korlátja*  $X'$ -nak, ha nem létezik olyan  $y$ -tól különböző  $y'$  alsó korlát, amire  $y \preceq y'$ . Hasonlóan értelmezhető a felső, legkisebb felső, minimális felső korlát fogalmak is. A  $\prec$  rendezés *teljes rendezés*, ha minden  $x \neq y$  elemre  $x \prec y$ ,  $y \prec x$  közül az egyik fennáll. Az  $(X, \preceq)$  párost *hálónak* nevezzük, ha  $\preceq$  az  $X$ -en értelmezett parciális rendezés, és tetszőleges  $x, y \in X$  elemeknek létezik

legnagyobb alsó (jelölésben:  $x \wedge y$ ) és legkisebb felső korlátjuk ( $x \vee y$ ).

Központi fogalom lesz a *lexikografikus rendezés*. Nézzük először ennek a matematikai definícióját. Legyen  $X$  és  $Y$  két halmaz, amelyeken értelmezve van egy-egy parciális rendezés ( $\prec_X, \prec_Y$ ). Azt mondjuk, hogy a  $(x_1, y_1) \in X \times Y$  lexikografikusan megelőzi  $(x_2, y_2) \in X \times Y$  párt, ha  $x_1 \prec_X x_2$ , vagy  $x_1 = x_2$  és  $y_1 \prec_Y y_2$ . A lexikografikus rendezést tetszőleges számú halmaz Descartes-szorzatára is kiterjeszthetjük rekurzív módon az alábbiak alapján:  $X \times Y \times Z = X \times (Y \times Z)$ . Látható, hogy a lexikografikus rendezést Descartes szorzatokon értelmezzük, vagy más szóval olyan összetett struktúrákon, amelyeknek ugyanannyi tagjuk van ( $n$ -eseknek is hívják ezeket). Mi ezt szeretnénk általánosítani, hiszen például szavak sorba rendezésénél is előfordulnak eltérő hosszúságú szavak. Ha a rövidebb szó megegyezik a hosszabb szó első felével (például komp és kompenzál szavak), akkor megegyezés alapján a rövidebb szó előzi meg lexikografikusan a hosszabbikat. Ezek alapján mindenki tudja definiálni a lexikografikus rendezést eltérő számú halmazok Descartes szorzatára. A legtöbb esetben a Descartes szorzat tagjainak halmaza és a rajtuk definiált rendezések megegyeznek (pl.:  $X = Y$  és  $\prec_X = \prec_Y$ ). Ilyenre, adott rendezés szerinti lexikografikus rendezésként hivatkozunk.

Az  $X, Y$  halmazokon értelmezett  $f$  bináris reláció *függvény*, ha bármely  $x \in X$  esetén pontosan egy olyan  $y \in Y$  létezik, hogy  $(x, y) \in f$ . Ez jelölésben  $f : X \rightarrow Y$ , és, ha  $(x, y) \in f$ , akkor  $y = f(x)$ . Az  $X$  halmazt a  $f$  értelmezési tartományának hívjuk (vagy máshogy:  $f$  az  $X$ -en értelmezett),  $Y$ -t az  $f$  képhalmazának, az  $f(X)$  halmazt pedig az  $f$  értékkészletének. Azt a függvényt, amelyet úgy kapunk, hogy először a  $f$ , majd az  $g$  függvényt alkalmazzuk  $g \circ f$ -el jelöljük. *Predikátum* egy függvény, ha az értékkészlete az  $\{igaz, hamis\}$  halmaz. *Szürjektív* egy függvény, ha a képhalmaza megegyezik az értékkészletével, *injektív* (vagy más néven egy-egy értelmű leképezés), ha az értelmezési tartomány bármely két különböző eleméhez különböző értéket rendel és *bijektív* (másképpen a függvény egy *bijekció*), ha szürjektív és injektív is egyben.

Legyen  $H$  tetszőleges halmaz. Az  $f : \overbrace{H \times \dots \times H}^n \rightarrow H$  függvényt  $n$  változós *műveletnek* nevezzük. A  $H$  halmazon értelmezett kétváltozós  $\star$  műveletet *asszociatívnak* nevezzük, ha tetszőleges  $a, b, c \in H$  esetén  $(a \star b) \star c = a \star (b \star c)$ . A  $(H, \star)$  párt *félcsoportnak* nevezzük, ha  $\star$  a  $H$ -n értelmezett asszociatív művelet. A  $(H, \star)$  félcsoport elemein a  $H$  elemeit értjük. Ha a  $(H, \star)$  félcsoport elemei között létezik olyan  $e$  elem, amelyre  $e \star a = a \star e = a$  minden  $a \in H$  elemre, akkor  $e$ -t *egységelemnek* hívjuk és egységelemes félcsoportról beszélünk. Ha egy egységelemes félcsoportban minden elemnek létezik inverze, akkor *csopotról* beszélünk. Az  $a$  inverzére ( $a^{-1}$ ) teljesüljön, hogy  $a \star a^{-1} = a^{-1} \star a = e$ . A csoport *Ábel-csoport*, ha a  $\star$  művelet *kommutatív* ( $a \star b = b \star a$ ) is. A  $(H, \star, +)$  hármas egy *gyűrű*, amennyiben  $(H, \star)$  Ábel csoport,  $(H, +)$  félcsoport és a  $\star, +$  műveletek *disztributívek* egymásra nézve, azaz  $(a + b) \star c = a \star c + b \star c$ . A  $\star$  és a  $+$  műveletek egységelemeit az 1 és a 0 szimbólumok jelölik. *Testnek* hívjuk az olyan kommutatív gyűrűt, ahol az  $1 \neq 0$  és a 0-an kívül a  $H$  minden elemének van inverze.

A  $H$  halmaz felett értelmezett *multihalmaznak* vagy *zsáknak* nevezzük azt a halmazt, amelynek elemei olyan párok, amelyek első tagja  $H$  egy eleme, második tagja pedig egy pozitív egész szám. Egy multihalmazt szokás úgy ábrázolni mintha olyan halmaz lenne, amely egy elemet többször is tartalmazhat. Ilyenkor a pár első tagját annyiszor írjuk le, amennyi a pár második tagja. Például a  $\{(A, 1), (C, 3)\}$ -at  $\{A, C, C, C\}$ -vel ábrázoljuk. A multihalmaz *méretén* a párok második tagjainak összegét, elemszámán pedig a párok számát értjük.

Sokat fogjuk használni a *sorozat* fogalmát. Legyen  $S$  egy halmaz. Az  $f : \mathbb{N} \rightarrow S$  függvényt az  $S$  felett értelmezett sorozatnak hívjuk. Leírására az  $f(0), f(1), \dots$  helyett a  $\langle s_0, s_1, \dots \rangle$



jelölést fogjuk használni. *Véges sorozatok* esetében az  $f$  értelmezési tartománya (általában az  $\{1, 2, \dots, n\}$ ) véges halmaz. Véges sorozat *hossza* az értelmezési tartományának elemszáma. Az  $S = \langle s_1, s_2, \dots, s_n \rangle, S' = \langle s'_1, s'_2, \dots, s'_{n'} \rangle$  sorozat konkatenációján az  $\langle s_1, s_2, \dots, s_n, s'_1, s'_2, \dots, s'_{n'} \rangle$  sorozatot értjük, és  $\langle S, S' \rangle$ -el jelöljük.

## 2.2. Lineáris algebra

Legyen  $H$  egy test, amelynek elemeit *skalároknak* hívjuk. A  $H$  felett értelmezett *vektortér* egy  $V$  halmaz (amelynek elemei a *vektorok*) és két bináris operátor (vektor összeadás:  $+$  és skalárral való szorzás:  $\cdot$ ), amelyekre teljesül néhány axióma (1.  $u, v, w \in V$ -re  $u + (v + w) = (u + v) + w$ , 2.  $u + v = v + u$ , stb.). A  $W \subseteq V$  halmazt *altérnek* nevezzük, ha zárt a vektorösszeadás és skalárszorzás műveletekre. Adott vektorhalmazt tartalmazó alterek metszetét a *vektorhalmaz által feszített altérnek* nevezzük. Ha a halmazból nem távolíthatunk el elemet a feszített altér megváltoztatása nélkül, akkor a vektorhalmazt *lineárisan függetlennek* hívjuk. A  $V$  altér egy *bázisa* egy olyan lineárisan független vektorhalmaz, amelynek feszített altere  $V$ .

A hagyományoknak megfelelően az  $A$  mátrix  $i$ -edik sorából képzett vektort  $A^i$ -vel jelöljük,  $\|v\|$ -vel a  $v$  vektor euklideszi normáját ( $\sqrt{\sum_i v_i^2}$ ) és  $v^T w$ -vel a  $v^T, w$  vektorok skaláris szorzatát ( $\sum_i v_i^T w_i$ ).

## 2.3. Gráfelmélet

*Irányított gráf* egy  $G = (V, E)$  pár, ahol  $V$  *csúcsok* (vagy *pontok*) véges halmaza,  $E$  pedig egy bináris reláció  $V$ -n.  $E$  elemeit *éleknek* nevezzük. Ha  $(u, v) \in E$ , akkor az  $u, v$  csúcsok egymás *szomszédai*. *Irányítatlan gráfról* beszélünk, ha az  $E$  reláció szimmetrikus. A *címkezett* (vagy *súlyozott*) gráfnál a csúcsokhoz, *címkezett élű* (vagy *élsúlyozott*) gráfnál pedig az élekhez rendelünk címkéket. A címkezett élű gráfot *súlyozott gráfnak* hívjuk, ha a címkék számokkal kifejezhető súlyokat jelentenek. A gráf méretén ( $|G|$ ) a csúcsok számát értjük. Egy csúcs *fokán* a csúcsot tartalmazó éleket értjük. Irányított gráfoknál megkülönböztetünk *kifokot* és *befokot*. A  $G$  irányítatlan gráf *k-reguláris*, ha minden csúcs foka pontosan  $k$ .

A  $G' = (V', E')$  gráf a  $G = (V, E)$  *részgráfja*, ha  $V' \subseteq V$  és  $E' \subseteq E$ . A  $G = (V, E)$  gráf  $V' \subseteq V$  *által feszített részgráfja* (induced subgraph) az a  $G' = (V', E')$  gráf, ahol  $E' = \{(u, v) \in E : u, v \in V'\}$ . A  $G_1(V_1, E_1)$  *izomorf* a  $G_2(V_2, E_2)$  gráffal, jelölésben  $G_1 \cong G_2$ , ha létezik  $\phi: V_1 \rightarrow V_2$  bijekció, amelyre  $(u, v) \in E_1$  esetén  $(\phi(u), \phi(v)) \in E_2$  is fennáll. Címkezett gráfoknál emellett megköveteljük, hogy az  $u$  csúcs címkéje megegyezzen a  $\phi(u)$  címkéjével minden  $u \in V_1$ -re, címkezett élű gráfnál pedig az  $(u, v)$  címkéje egyezzen meg a  $(\phi(u), \phi(v))$  él címkéjével. Ha  $G \cong G$ , akkor *automorfizmusról* beszélünk.

A gráfok ábrázolásának elterjedt módja a *szomszédossági mátrix* (adjacency matrix) és a *szomszédosság lista*. Az  $|G| \times |G|$  méretű  $A$  szomszédossági mátrix  $a_{ij}$  eleme 1 (élcímkezett esetben az él címkéje), ha a  $G$  gráf  $i$ -edik csúcsából indul él a  $j$ -edik csúcsba, különben 0. Természetesen a szomszédossági mátrixot a gráfon kívül az határozza meg, hogy melyik csúcsot hívjuk az elsőnek, másodiknak, ... A szomszédossági mátrixot tehát a gráf és az  $f: V \rightarrow \{1, \dots, |V|\}$  bijekció adja meg. Hurokél nélküli, címkezett gráfban a szomszédossági mátrix  $a_{ii}$  eleme az  $i$  csúcs címkéjét tárolja. A szomszédossági lista  $|G|$  darab lista, ahol az  $i$ -edik lista tárolja az  $i$ -edik csúcs szomszédait.

Az  $u$  csúcsot az  $u'$  csúccsal összekötő  $k$ -hosszú úton csúcsoknak egy olyan (véges)  $\langle v_0, v_1, \dots, v_k \rangle$  sorozatát értjük, amelyre  $u = v_0$ ,  $u' = v_k$ , és  $(v_{i-1}, v_i) \in E$  ( $i = 1, 2, \dots, k$ ). Egy út *egyszerű*, ha a benne szereplő csúcsok páronként különbözők. A  $\langle v_0, v_1, \dots, v_k \rangle$  út *kör*, ha  $v_0 = v_k$ , és az út legalább egy élt tartalmaz. Egy gráfot *összefüggőnek* hívunk, ha bármely két csúcsa összeköthető úttal. A körmenetes, irányítás nélküli gráfot *erdőnek* hívjuk. Ha az erdő összefüggő, akkor pedig *fának*. Az olyan fát, amely tartalmazza egy  $G$  gráf minden csúcsát, a  $G$  *feszítőfájának* hívjuk.

A *gyökeres fában* az egyik csúcsnak kitüntetett szerepe van. Ezt a csúcsot *gyökérnek* nevezzük. A gyökérből egy tetszőleges  $x$  csúcsba vezető (egyértelműen meghatározott) út által tartalmazott bármely  $y$  csúcsot az  $x$  *ősének* nevezzük. Azt is mondjuk ekkor, hogy  $x$  az  $y$  *leszármazottja*. Ha  $x \neq y$ , akkor *valódi ősről* és *valódi leszármazotról* beszélünk. Ha az úton  $x$  1 élen keresztül érhető el  $y$ -ből, akkor  $x$  az  $y$  *gyereke* és  $y$  az  $x$  *szülője*. Ha két csúcsnak ugyanaz a szülője, akkor *testvéreknek* mondjuk őket.

A  $G = (V, E)$  gráf  $S, V \setminus S$  *vágásán* a  $V$  halmaz kétrészes partícióját értjük. Az  $(u, v) \in E$  él *keresztezi* az  $S, V \setminus S$  vágást, ha annak egyik végpontja  $S$ -ben a másik  $V \setminus S$ -ben van. Egy vágás *súlya* – súlyozott gráfok esetében – megegyezik a vágást keresztező élek összsúlyával.

## 2.4. Matematika logika

### 2.4.1. Ítéletlogika

### 2.4.2. Elsőrendű logika

## 2.5. Valószínűségszámítás

Feltételezzük, hogy az olvasó tisztában van a *diszkrét/folytonos valószínűségi változó*, valószínűségi változó *eloszlásának*, *sűrűségfüggvényének*, *eloszlásfüggvényének*, a valószínűségi változó *várható értékének* ( $\mathbb{E}[X] = \mu = \sum x \cdot p(x)$ ), annak fontos tulajdonságait ( $\mathbb{E}[aX + bY] = a\mathbb{E}[X] + b\mathbb{E}[Y]$ ,  $\mathbb{E}[X] = \mathbb{E}[\mathbb{E}(X|Y)]$ ) és *szórásának* ( $D^2[X] = \sigma_X^2 = \mathbb{E}[(X - \mu)^2]$ ) vagy általánosan az  $n$ -edik *centrális momentumok* fogalmával ( $D^n[X] = \mathbb{E}[(X - \mu)^n]$ ), továbbá ismeri két valószínűségi változó közötti kovarianciát ( $Cov(X, Y) = \mathbb{E}[(X - \mu)(Y - \nu)]$ ) és korrelációt ( $Corr(X, Y) = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}$ ). A sűrűségfüggvény (vagy diszkrét eloszlás) maximumhelyét az eloszlás móduszának hívjuk. Az  $F$  eloszlásfüggvény  $p$ -*kvartilisé*t az a  $K$  szám adja, amelyre  $F(K) < p$  és  $F(K + 0) \geq p$ . Az 1/2-hez tartozó kvartilist *mediánnak* nevezzük.

### 2.5.1. Nevezetes eloszlások

A következő nevezetes eloszlásokkal fogunk találkozni tanulmányaink során.

#### Binomiális és Poisson eloszlás

Legyen  $(\Omega, F, P)$  Kolmogorov-féle valószínűségi mező,  $A \in F$  pozitív valószínűségű esemény,  $p = P(A) > 0$ . Hajtsunk végre  $n$ -szeres független kísérletsorozatot és legyen  $X$  értéke annyi, ahányszor  $A$  bekövetkezett a kísérletsorozatban.  $X$ -et ekkor  $n, p$  paraméterű binomiális eloszlású valószínűségi változónak nevezzük, jele  $X \in B(n, p)$ .  $X$  eloszlása  $p_k = P(X = k) = \binom{n}{k} p^k (1-p)^{n-k}$ , várható értéke  $\mathbb{E}(X) = np$ , szórása  $\sigma^2(X) = np(1-p)$ .

A Poisson-eloszlás a binomiális eloszlás határesetete.  $\lim_{n \rightarrow \infty, p \rightarrow 0, np = \lambda} \binom{n}{k} p^k q^{n-k} = \frac{\lambda^k}{k!} e^{-\lambda}$ . A Moivre-Laplace tétel szerint, az  $n$ -ed rendű  $p$  paraméterű binomiális eloszlás standardizáltja  $n$  minden határon túl való növelése esetén normális eloszlású:  $\forall x \in \mathbb{R}: \lim_{n \rightarrow \infty} \sum_{\frac{k-np}{\sqrt{npq}} < x} \binom{n}{k} p^k q^{n-k} = \Phi(x)$

### Hipergeometrikus eloszlás

Tegyük fel, hogy van  $N$  különböző elemünk, amelyből  $R$  darab rossz. A hipergeometrikus eloszlás adja meg annak az esélyét, hogy  $r$  darab rossz elem lesz, ha az  $N$  elemből  $n$  darabot kivesszünk véletlenszerűen. Elemi kombinatorikus úton a valószínűség kiszámítható ( $0 \leq r \leq \min\{n, R\}$ ):

$$\mathbb{P}(r, N, R, n) = \frac{\binom{R}{r} \binom{N-R}{n-r}}{\binom{N}{n}}$$

A fenti sűrűségfüggvénnyel rendelkező diszkrét valószínűségi eloszlást hívjuk *hipergeometrikus eloszlásnak*.

Amennyiben  $n \ll N$ , akkor a hipergeometrikus eloszlást közelíthetjük az  $n, R/N$  paraméterű binomiális eloszlással.

### Normális eloszlás

#### $\chi^2$ eloszlás

Legyenek  $\xi_1, \xi_2, \dots, \xi_n$  egymástól független, standard normális eloszlású valószínűségi változók. Ekkor az  $\sum_{i=1}^n \xi_i^2$  valószínűségi változó eloszlását  $n$  paraméterű  $\chi^2$  eloszlásnak ( $\chi_n^2$ ) nevezzük.

### 2.5.2. Egyenlőtlenségek

Legyen  $X$  egy  $E[X]$  várható értékű valószínűségi változó. A Markov egyenlőtlenség szerint  $\mathbb{P}(|X| \geq a) \leq \frac{E[|X|]}{a}$ , ahol  $a > 0$ . A Hoeffding-korlát a mintavételzéssel kapcsolatos állítások alapja.

**2.1. lemma.** *Legyen  $X_i, 1 \leq i \leq n$   $\mu$  várható értékű, független, azonos eloszlású valószínűségi változók és  $a \leq X_i \leq b$  minden  $i$ -re. Ekkor tetszőleges  $\lambda > 0$ -ra fennáll a következő egyenlőtlenség:*

$$\mathbb{P}\left[\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| \geq \lambda\right] \leq 2e^{-2\lambda^2 n / (b-a)^2}.$$

### 2.5.3. Entrópia

Legyen  $X$  egy diszkrét valószínűségi változó, amely értékeit egy  $\mathcal{X}$  halmazból veheti fel. Az  $l_X = -\log_2 p(X)$  valószínűségi változót az  $X$  *entrópiásűrűségének* nevezzük.  $X$  entrópiáját –  $H(X)$ -et – ezen változó várható értékével definiáljuk:

$$H(X) = - \sum_{x \in \mathcal{X}} p(x) \log_2 p(x).$$

Az entrópia valamiképpen a változó *bizonytalanságát* fejezi ki. Ha  $\mathcal{X}$  elemszáma rögzített és az  $X$  változó csak egy értéket vehet fel (mert az egyik érték valószínűsége 1), akkor  $H(X)$  értéke 0 (nincs bizonytalanság), ha pedig  $X$  eloszlása egyenletes eloszlást követ, akkor az entrópia a maximumát veszi fel,  $\log_2(|\mathcal{X}|)$ -t.

Legyen  $X$  és  $Y$  két diszkrét értékű valószínűségi változó. Az  $X$ -nek az  $Y$  feltétellel vett feltételes entrópiája:

$$H(X|Y) = - \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} p(x, y) \log_2 p(x|y),$$

vagy egy kicsit átalakítva kapjuk, hogy

$$H(X|Y) = - \sum_{y \in \mathcal{Y}} p(y) \sum_{x \in \mathcal{X}} p(x|y) \log_2 p(x|y).$$

Be lehet bizonyítani, hogy  $H(X|Y) = H(XY) - H(Y)$ , ami informálisan úgy lehet megfogalmazni, hogy a feltételes entrópia megadja, hogy mennyi bizonytalanság marad  $X$ -ben, ha elveszünk az  $Y$  bizonytalanságát.

A feltételes entrópia számos tulajdonsága közül mi csak az alábbi fogjuk felhasználni:

$$0 \leq H(X|Y) \leq H(X).$$

## 2.6. Statisztika

A statisztikában általában  $X_1, X_2, \dots, X_n$  független, azonos eloszlású valószínűségi változók vannak megadva, amiket *mintáknak* nevezünk. Az eloszlást nem ismerjük pontosan, de rendelkezésünkre állnak megfigyelések.

Legyenek  $X_1, X_2, \dots, X_n$  független, azonos eloszlású valószínűségi változók. Ekkor a  $\bar{X} = \frac{X_1 + X_2 + \dots + X_n}{n}$  valószínűségi változót *empirikus középnek*, vagy *mintátlagnak*, a  $s_n^{*2} = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$  valószínűségi változót pedig *korrigált empirikus szórásnégyzetnek* nevezzük.

A  $\chi^2$  eloszlás definíciójából következik, hogy az  $\frac{(n-1)s^{*2}}{\sigma^2}$  valószínűségi változó eloszlása  $\chi_n^2$ , amennyiben a  $s^{*2}$   $\sigma$  szórású, normális eloszlású valószínűségi változók korrigált empirikus szórásnégyzetét jelöli

**2.2. definíció.** *Legyenek  $X$  és  $Y$  két olyan valószínűségi változó, amelyek eloszlása rendre  $\chi_n^2$  és  $\chi_m^2$ . Ekkor a  $Z = \frac{X/n}{Y/m}$  valószínűségi változó eloszlását  $F_{n,m}$  eloszlásnak hívjuk.*

### 2.6.1. Hipotézisvizsgálat

A hipotézisvizsgálat feladata mindig valamilyen állítás helyességének vizsgálata. Ezt az állítást *nullhipotézisnek* nevezzük, jele  $H_0$ . A nullhipotézis általában egy valószínűségi változó valamely paraméterére vagy a változó viselkedésére vonatkozó állítás. Az állítás igazolásához vagy elvetéséhez kísérletezgetések, minták állnak rendelkezésünkre. Ha a minták alapján a nullhipotézist elvetjük, holott az igaz, akkor *elsőfajú hibát* követünk el. Ellenkező esetben – amikor a nullhipotézis hamis, de mi elfogadjuk – *másodfajú hibáról* beszélünk. Pusztán minták segítségével nem tudunk teljesen biztos választ adni. A gyakorlatban egy paraméterrel ( $\alpha$ )

rögzítik az elsőfajú hiba elkövetésének megengedett valószínűségét. Az  $1 - \alpha$  értéket a *próba szintjének* hívjuk.

Összefoglalva tehát, adott egy állítás, egy paraméter ( $\alpha$ ) és minták sorozata. Feladatunk, hogy a minták alapján cáfoljuk vagy igazoljuk az állítást úgy, hogy bizonyíthatóan  $\alpha$ -nál kisebb legyen annak valószínűsége, hogy az állítás igaz, holott mi cáfoljuk. A hipotézisvizsgálatnál a minták eredményeit felhasználva kiszámítunk egy ún. *próbat statisztika* értéket, és ezt vetjük össze egy ismert eloszlással. Az  $\alpha$ -nak célszerű kis (0.1 és 0.01 közötti) értéket választani<sup>1</sup>.

## 2.6.2. A binomiális próba

### 2.6.3. Az $F$ -próba

Az  $F$ -próba arra szolgál, hogy két független, normális eloszlású valószínűségi változó ( $X, Y$ ) szórásának egyenlőségét eldöntsük.

$$H_0 : \sigma_X = \sigma_Y.$$

Tudjuk, hogy  $\frac{(n_X-1)s_X^{*2}}{\sigma_X^2}$  és  $\frac{(n_Y-1)s_Y^{*2}}{\sigma_Y^2}$   $\chi^2$  eloszlásúak  $(n_X - 1)$  illetve  $(n_Y - 1)$  paraméterrel. Ha a nullhipotézis fennáll, akkor az

$$F = \frac{s_X^{*2}}{s_Y^{*2}}$$

próbat statisztika  $F$ -eloszlású  $(n_X - 1, n_Y - 1)$  paraméterrel. Azonban  $\frac{1}{F}$  is  $F$ -eloszlású  $(n_Y - 1, n_X - 1)$  paraméterrel, ezért a gyakorlatban  $F^* = \max\{F, 1/F\} \geq 1$  statisztikát szokás használni.

### 2.6.4. A $\chi^2$ -próba

A  $\chi^2$ -próbák az alábbi tételt használják fel.

**2.3. tétel.** *Legyen  $A_1, A_2, \dots, A_r$  egy teljes eseményrendszer ( $r \geq 3$ ), legyen  $p_i = \mathbb{P}(A_i) > 0, i = 1, \dots, r$ . Ismételjük a kísérletet  $n$ -szer egymástól függetlenül. Jelölje  $X_i$  az  $A_i$  esemény bekövetkezésének számát. Belátható, hogy ekkor a*

$$\sum_{j=1}^r \frac{(X_j - np_j)^2}{np_j}$$

*valószínűségi változó eloszlása  $n \rightarrow \infty$  esetén  $\chi_{r-1}^2$  eloszláshoz konvergál.*

A  $\chi^2$  eloszlás kvantiliseit függvény-táblázatokban megtalálhatjuk.

A  $\chi^2$ -próba legfontosabb alkalmazási területei az (1.) illeszkedés-, (2.) függetlenség- és (3.) homogenitásvizsgálat. Témánkhoz a függetlenség-vizsgálat tartozik hozzá, így a továbbiakban ezt részletezzük. A  $\chi^2$ -próba iránt érdeklődőknek a [65] magyar nyelvű irodalmat ajánljuk.

---

<sup>1</sup>Gondolkozzunk el azon, hogy mi történne, ha  $\alpha$ -nak nagyon kis értéket választanánk!

### 2.6.5. Függetlenségvizsgálat

Legyen  $A_1, A_2, \dots, A_r$  és  $B_1, B_2, \dots, B_s$  két teljes eseményrendszer. Végezzünk  $n$  kísérletet. Nullhipotézisünk az, hogy az eseményrendszerek függetlenek.

$$H_0 : \mathbb{P}(A_i, B_j) = \mathbb{P}(A_i)\mathbb{P}(B_j), \quad i = 1, \dots, r \quad j = 1, \dots, s$$

Ha az események valószínűségei adottak, akkor tiszta illeszkedés vizsgálati feladatról beszélünk, ahol

$$H_0 : \mathbb{P}(A_i \cap B_j) = p_i q_j$$

hiszen  $p_i, q_j$  értékek adottak. Jelölje  $k_{ij}$  az  $A_i \cap B_j$  esemény bekövetkezésének számát. Ekkor ki kell számítanunk a

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(k_{ij} - np_i q_j)^2}{np_i q_j}$$

ún. próbastatisztika értéket. Jobban megvizsgálva  $\chi^2$ -et láthatjuk, hogy az egy  $\sum \frac{(\text{megfigyelt érték} - \text{várt érték})^2}{\text{várt érték}}$  jellegű kifejezés. Amennyiben  $\chi^2$  kicsi, akkor a megfigyelt értékek közel vannak azokhoz, amit  $H_0$  fennállása esetén vártunk, tehát a nullhipotézist elfogadjuk.

Hogy pontosan mit jelent az, hogy „kicsi”, azt a 2.3-as tétel alapján  $\chi_{rs-1}^2$  és az  $\alpha$  paraméter határozza meg. Táblázatból keressük ki, hogy a  $\chi_{rs-1}^2$  eloszlás hol veszi fel az  $1 - \alpha$  értéket. Amennyiben ez nagyobb a fent kiszámított  $\chi^2$  értéknél, akkor a nullhipotézist elfogadjuk, ellenkező esetben elvetjük.

A gyakorlatban sokkal többször fordul elő az az eset, amikor az események valószínűségeit nem ismerjük. Ekkor a valószínűségeket az események relatív gyakoriságával becsüljük meg. Jelöljük az  $A_i$  esemény gyakoriságát  $k_{i\cdot}$ -vel, tehát  $k_{i\cdot} = \sum_{j=1}^s k_{ij}$  és hasonlóan  $B_j$  esemény gyakoriságát  $k_{\cdot j}$ -vel.  $\chi^2$ -próbák során az adatok szemléltetésének gyakran használt eszköze az ún. kontingencia-táblázat. Ez egy többdimenziós táblázat, amely celláiban a megfelelő esemény bekövetkezésének száma található. Egy ilyen 2-dimenziós kontingencia-táblázatot láthatunk a következő ábrán.

	$B_1$	$B_2$	$\dots$	$B_s$	$\Sigma$
$A_1$	$k_{11}$	$k_{12}$	$\dots$	$k_{1s}$	$k_{1\cdot}$
$A_2$	$k_{21}$	$k_{22}$	$\dots$	$k_{2s}$	$k_{2\cdot}$
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$A_r$	$k_{r1}$	$k_{r2}$	$\dots$	$k_{rs}$	$k_{r\cdot}$
$\Sigma$	$k_{\cdot 1}$	$k_{\cdot 2}$	$\dots$	$k_{\cdot s}$	$n$

Az  $A_i \cap B_j$  megfigyelt értéke  $k_{ij}$ , várt értéke  $H_0$  esetén  $n \cdot \frac{k_{i\cdot}}{n} \cdot \frac{k_{\cdot j}}{n}$ . Ezek alapján  $\chi^2$  értéke:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{(k_{ij} - \frac{k_{i\cdot} \cdot k_{\cdot j}}{n})^2}{\frac{k_{i\cdot} \cdot k_{\cdot j}}{n}}$$

Mivel a függetlenség fennállása esetén  $r - 1$  darab  $p_i$ -t és  $s - 1$  darab  $q_j$  valószínűséget kell megbecsülni, így a fenti  $H_0$  fennállása esetén  $\chi^2_{rs-1-(r+s-2)} = \chi^2_{(r-1)(s-1)}$  eloszlású.

A  $\chi^2$  eloszlás közelítése csak abban az esetben pontos, ha a  $k_{ij}$  értékek nagyok. Persze nincs pontos szabály arra nézve, hogy mennyire kell nagyoknak lennie. Azt szokták mondani, hogy a kontingencia táblázat elemeinek 90%-a nagyobb legyen ötnél.

### 2.6.6. Student t-próba

## 2.7. Algoritmus-elmélet

Terjedelmi okok miatt csak felsorolni tudjuk azokat az algoritmusokat, amelyeket az olvasónak ismernie kell. Ezek pedig: lineáris-, bináris keresés, mélységi, szélességi bejárás, Kruskal algoritmus a minimális súlyú feszítőfa meghatározásához stb. Emellett feltételezzük, hogy az olvasó tisztában van az NP-teljesség (vagy általánosabban a bonyolultság) elméletének alapjaival.

## 2.8. Adatstruktúrák

Feltételezzük, hogy az olvasó tisztában van a lista (vektor) és a tömb fogalmával. Az adatbányászatban további közkezdvelt adatstruktúrái az ún. *szófa* (trie), vagy más néven prefix-fa (prefix-tree), a piros-fekete fa, illetve a hash-tábla.

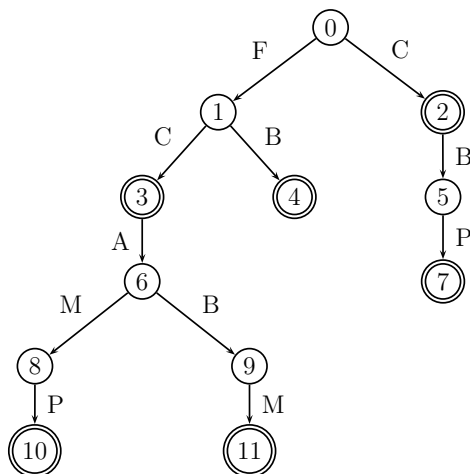
### 2.8.1. Szófák

A szófát eredetileg szótár szavainak tárolásánál alkalmazták, annak érdekében, hogy gyorsan el lehessen dönteni, hogy egy adott szó szerepel-e a szótárban [33], [44]. A szavak az abc felett értelmezett sorozatok, így általánosan azt mondhatjuk, hogy egy szófa egy adott véges elemhalmaz feletti sorozatok tárolására és gyors visszakeresésére alkalmas adatstruktúra. A szófa angol neve (trie, amit úgy ejtünk, mint a try szót) a visszakeresés angol fordításából származik (retrieval). A továbbiakban az alaphalmazt  $\mathcal{J}$ -vel, az alaphalmaz felett értelmezett, adott sorozatok halmazát szótárnak hívjuk. A 2.1 ábrán egy szófát láthatunk, mely az  $C$ ,  $FC$ ,  $FB$ ,  $CBP$ ,  $FCAMP$ ,  $FCABM$  sorozatokat tárolja.

A szófa egy (lefelé) irányított gyökeres címkézett fa. Egy  $d$ -edik szintű pontból csak  $d+1$ -edik szintű pontba mutathat él. Néha a hatékonyság kedvéért minden pontból a pont szülőjére is mutat él. A gyökeret 0. szintűnek tekintjük. A címkék az  $\mathcal{J}$ -nek egy-egy elemei. Minden pont egy elemsorozatot reprezentál, amely a gyökérből ebbe a pontba vezető éleken található elemekből áll. Akkor tartalmazza a szófa az  $S$  sorozatot, ha van olyan pont, amely az  $S$ -t reprezentálja.

Ha egy sorozatot tartalmaz egy szófa, akkor annak tetszőleges prefixét is tartalmazza. A prefix azonban nem biztos, hogy eleme a szótárnak. Ezt a problémát kétféleképpen lehet kiküszöbölni. Egyrészt megkülönböztetünk *elfogadó* és *nem elfogadó* pontokat. Egy sorozatot akkor tartalmazza a szófa, ha van olyan elfogadó állapot, amely a sorozatot reprezentálja. Másrészt bevezethetünk egy speciális elemet, amit minden sorozat végére illesztünk, továbbá sorozatot csak levél reprezentálhat.

A szófának két implementációját különböztetjük meg attól függően, hogy milyen technikát alkalmazunk az élek tárolására. Az ún. *táblázatos implementációban* (tabular implementation)



2.1. ábra. Példa szófára

[44] minden ponthoz egy rögzített hosszúságú, mutatókat tartalmazó vektort veszünk fel. Az  $i$ -edik mutató mutat az  $i$ -edik elemhez tartozó él végpontjára. Ha a pontnak nincs ilyen címkéjű éle, akkor a mutató értéke NULL. A vektor hossza az  $J$  elemszámával egyezik meg.

A *láncolt listás implementációban* [33] az éleket egy láncolt listában tároljuk. A lista elemei élcímke, gyermekmutató párok. A láncolt lista következő elemére mutató mutatókat megspórolhatjuk, ha egy vektort alkalmazunk, aminek hossza megegyezik a pont éleinek számával, és elemei szintén címke, mutató párok. Ez azért is jó megoldás, mert egy lépéssel tudunk tetszőleges indexű elemre lépni (a címke, mutató pár memóriaszükségletének ismeretében), és nem kell a mutatókon keresztül egyesével lépegetnünk.

Szófák esetében a legfontosabb elemi művelet annak eldöntése, hogy egy adott pontnak van-e adott címkéjű éle, és ha van, akkor ez hova mutat. Táblázatos implementációnál ezt a feladatot egy lépésben megoldhatjuk a megfelelő indexű elem megvizsgálásával. Láncolt listás, illetve változó hosszúságú vektor esetén a megoldás lassabb művelet. A vektor minden párját ellenőriznünk kell, hogy a pár címkéje megegyezik-e az adott címkével. A hatékonyságot növelhetjük, ha a párokat címkék szerint rendezve tároljuk, és bináris keresést végzünk.

Érdekes összehasonlítani a két vektoros implementációban a pontok memóriaigényét. Amennyiben a mutatók, és a címkék is 4 bájtot foglalnak, akkor a táblázatos implementációban egy pont memóriaigénye (a vektor fejléc memóriaigényétől eltekintve)  $|J| \cdot 4$  bájt, a listás implementációé  $n \cdot 2 \cdot 4$  bájt, ahol  $n$  az adott pontból induló élek száma, amire igaz, hogy  $0 \leq n \leq |J|$ . Ha a szófa pontjai olyanok, hogy kevés élük van, akkor a listás implementációnak lesz kevesebb memóriára szüksége, sok élnél azonban táblázatos implementáció a jobb megoldás. A két technikát ötvözhetjük akár egy adott szófán belül is [123], [144]: ha a pont éleinek száma meghalad egy korlátot (általában  $J/2$ -t), akkor táblázatos implementációt használunk, ellenkező esetben maradunk a listás megoldásnál.

Megemlítünk két szófa leszármazottat. Ezek a *nyesett szófák* (pruned trie) és a *PATRICIA fák*. Mindkét fa abban különbözik az eredeti szófától, hogy kiiktatják az olyan utakat a fából, amelyekben nincsen elágazás. A nyesett fánál ezt kizárólag levélhez vezető utakkal teszik, PATRICIA fáknál ez a korlátozás nem áll fenn.



## Patrícia-fák

Egy irányított utat láncnak hívunk, ha minden pontjának csak egy gyereke van. A Patrícia-fa a szófaból származtatható úgy, hogy a szófa nem bővíthető láncait egy-egy éllé vonjuk össze. Az új él a lánc utolsó pontjába mutat, címkéje a lánc éleinek címkéiből álló sorozat. Ha a láncösszevonást csak a levélben végződő láncokra hajtjuk végre, akkor ún Patrícia\* fát kapunk.

Ha a szófa sok láncot tartalmaz, akkor a Patrícia-fa sokkal hatékonyabb. Ellenkező esetben viszont több memóriát használ, mivel a címkéket vektorokban tároljuk, ami egyetlen elem tárolása esetén nem célravezető a nagy többletköltség miatt.

### 2.8.2. Piros-fekete fák

A piros-fekete (RB-tree vagy symmetric binary B-trees) fák a kiegyensúlyozott bináris fák (balanced binary tree) egy típusa. Minden csúcsnak színe van, hagyományosan piros vagy fekete. Speciális forgatásokat használó beszúrás művelet biztosítja, hogy bármely a gyökérből levélbe vezető út hossza ne legyen nagyobb, mint a legrövidebb ilyen út hosszának kétszerese. Egy piros-fekete fa a következő tulajdonságokkal rendelkezik:

- I. Minden csúcsnak a színe piros vagy fekete.
- II. Minden levél színe fekete.
- III. Minden piros csúcsnak mindkét fia fekete.
- IV. Bármely két, azonos csúcsból induló, levélig vezető úton ugyanannyi fekete csúcs van.

A fentiekből következik, hogy bármely  $n$  belső csúccsal rendelkező piros-fekete fa magassága legfeljebb  $2 \lg(n+1)$ . A bizonyítás és a fa építésének menete megtalálható az irodalomban (pl. [79]).

### 2.8.3. Hash-tábla

A hash-tábla magyar elnevezése hasító-tábla ( $\odot$ ), de mi ezt a szót nem fogjuk használni. A hash-tábla elemek gyors elhelyezésére és visszakeresésére használt adatstruktúra. A táblázatnak cellái vannak, amibe elemeket helyezhetünk. Minden cellának van egy címe (vagy indexe). A hash-táblás tárolásban központi szerepet tölt be az elemeken értelmezett ún. *hash-függvény*, ami megadja az elem *hash-értékét*. Egy elemet arra a címre helyezünk be a hash-táblában, amelyet a hash-értéke megad. Előfordulhat, hogy különböző elemekhez a hash-függvény ugyanazokat a hash-értéket rendeli. Ezt *ütközésnek* hívjuk. A hash-táblákról önmagában fejezeteket lehet írni, ennyi bevezető azonban elég ahhoz, hogy megértsük a jegyzet további részeit.

## 2.9. Számítógép-architektúrák

Sok kutató alkalmazza a külső táras modellt az algoritmusának hatékonyságának vizsgálatokor. Mára az óriási memóriaméreteknek köszönhetően a legtöbb adatbázis elfér a memóriában, valamilyen szűrt formában. Ilyen esetekben az elemzéshez használt modell leegyszerűsödik az egyszerűbb közvetlen hozzáférésű (RAM-) modellre (amely Neumann-modell

[139] néven is ismert, mivel a magyar születésű Neumann János javasolta először ezt az architektúrát). A programokat olyan modern processzorokon futtatják, amely sokkal kifinomultabb a RAM-modellnél. A modell túlzott egyszerűsítése ahhoz vezet, hogy az elemzéseknek semmi köze nincs a valósághoz. Az új modell új elvárásokat támaszt az algoritmusokkal szemben. Ezekről egy kiváló áttekintés olvasható a [94] tanulmányban. A modern processzorok legfontosabb sajátossága a többszintű memória és a csővezetékes (pipeline-) feldolgozás.

### 2.9.1. Többszintű memória, adatlokalitás

A memória nem egyetlen nagy blokk, sokkal inkább különböző méretű, késleltetésű memóriákból álló hierarchikus rendszer. Minél nagyobb a memória mérete, annál több idő kell a hozzáféréshez. A hierarchia elemei, méret szerint növekvő sorrendben a következők: regiszterek, pár kilobájtos elsőszintű gyorsítótár, pár megabájtos másodsztintű gyorsítótár, esetleges harmadsztintű gyorsítótár, rendszermemória és merevlemez. Az adatot a rendszermemóriából a másodsztintű gyorsítótárba, a másodsztintűből az elsőszintű gyorsítótárba blokkonként másolhatjuk. A blokkméret egy Pentium 4-es processzor esetén 128 bájt.

A blokkonkénti feldolgozás más megvilágításba helyezi az algoritmusok vizsgálatát: egyetlen bit eléréséhez és beolvasásához egy lassú memóriából ugyanannyi idő kell, mint a bitet tartalmazó teljes blokk eléréséhez és beolvasásához. Másik adat elérése ugyanebben a blokkban viszont nem igényli már a hozzáférést a lassú memóriához. Így rendkívül fontos követelménnyé válik az adatlokalitás, azaz hogy az adatok, amelyeket időben egymáshoz közel dolgozunk fel, a memóriában is közel legyenek egymáshoz.

Az adatot feldolgozásakor be kell hozni a regiszterekbe. Előfordulhat, hogy már eleve ott van, mert az előző műveletekhez szükség volt rá. A korlátozott számú regiszterek miatt azonban sokkal valószínűbb, hogy az egyik gyorsítótárban vagy a rendszermemóriában helyezkedik el. Sőt, az is lehet, hogy a merevlemezen található, ha az algoritmus memóriaigénye annyira nagy, hogy lapozásra van szükség. Ha a másodsztintű gyorsítótárban vagy a rendszermemóriában helyezkedik el a kívánt adat, akkor az adathozzáférés ún. cache miss-t okoz. Amíg ez az adat bekerül a regiszterekbe, a processzor végrehajthat más műveleteket (ezer alapművelet, például összeadás elvégzésére képes ez idő alatt), ennek ellenére a teljesítménye messze elmaradhat ilyenkor a maximálistól. Tehát az adatstruktúra, algoritmus pár tervezésekor törekednünk kell a minél jobb adatlokalitásra a cache miss-ek elkerülése érdekében.

### 2.9.2. Csővezetékes feldolgozás, elágazás-előrejelzés

A programozók által használt műveleteket a fordító mikroutasítások sorozatára bontja. A műveleteket nem külön-külön, egymás után hajtjuk végre, hanem párhuzamosan dolgozzuk fel őket, csővezeték használatával. Sajnos azonban az adatfüggőség és a feltételes ugrások sokat rontanak a párhuzamos feldolgozás hatékonyságán. Adatfüggőségről akkor beszélünk, ha egy utasítás egy előző utasítás eredményétől függ. Elágazás-előrejelzésnél megjósoljuk a feltétel kimenetét, és betöltjük a csővezetékbe az ennek megfelelő utasításokat. Ha a jóslás hamisnak bizonyul, akkor a csővezeték ki kell üríteni, és be kell tölteni a helyes utasításokat. Ezt a problémát gyakran kiküszöbölhetjük különböző technikák alkalmazásával, (mint például kódátrendezéssel) amelyet automatikusan elvégez a fordító. Számításiigényes algoritmus tervezésekor azonban nekünk kell ügyelnünk az adatfüggetlenségre és az elágazás-előrejelzésre.

A csővezetékes feldolgozás lehetővé teszi, hogy egy órajel alatt több műveletet is elvégezzünk. A fent említett problémák miatt azonban a processzor átlagos teljesítménye messze nem éri el az optimumot. A felesleges feltételek ronthatják a hatékonyságot. Az elágazás-előrejelzés intelligens olyan szempontból, hogy ha egy feltétel kimenete sohasem változik, akkor a processzor ezt figyelembe veszi és a későbbiekben ennek megfelelően jósol. Így a mindig igaz (vagy hamis) feltételek nem befolyásolják a hatékonyságot.

## 3. fejezet

# Előfeldolgozás, hasonlósági függvények

Ebben a fejezetben ismertetjük, hogy milyen elterjedt mértékek vannak elemek közötti hasonlóságra majd rátérünk a legfontosabb előfeldolgozási műveletekre. Mindenek előtt azt kell tisztáznunk, hogy milyen típusú attribútumok léteznek matematikus szemmel.

### 3.1. Attribútum típusok

Jelöljük az  $A$  attribútum két értékét  $a$ -val és  $a'$ -vel.

- I. A *kategória típusú (nominal)* attribútumnál az attribútum értékei között csak azonosságot tudunk vizsgálni. Tehát csak azt tudom mondani, hogy  $a = a'$  vagy azt, hogy  $a \neq a'$ . A kategória típusú attribútum egy speciális esete a *bináris attribútum*, ahol az attribútum csak két értéket vehet fel. A kategória típusú attribútumokat az irodalom néha *felsorolás* (enumerated) vagy *diszkrét* típusnak is hívja. Másodlagos jelentésük miatt a tanulmányban ezeket az elnevezéseket nem használjuk. Például a felsorolás típus említésénél a legtöbb informatikusnak a C++, java, C#-beli felsorolás típusú változó jut eszébe, amelyek mindig egyértelmű megfeleltetésben állnak egy egész számmal.
- II. A *sorrend típusú (ordinal)* attribútumoknál az értékeket sorba tudjuk rendezni, azaz az attribútum értéken teljes rendezést tudunk megadni. Ha tehát  $a \neq a'$ , akkor még azt is tudjuk, hogy  $a > a'$  és  $a < a'$  közül melyik igaz.
- III. Ha az eddigiek mellett meg tudunk adni egy  $+$  függvényt, amivel az elemek csoportot alkotnak, akkor *intervallum típusú (interval scale)* attribútumról beszélünk.
- IV. Ha egy intervallum típusú attribútumnál meg lehet adni zérus értéket, vagy pontosabban az attribútum elemei gyűrűt alkotnak, akkor az attribútum *arány skálájú (ratio scale)*. Az arány skálájú attribútumot gyakran fogjuk *valós* attribútumnak hívni, hiszen a gyakorlati esetek többségében az arány skálájú attribútumok megadásához valós számokat használunk. Azonban ne felejtjük el az arány skálájú attribútum eredeti definícióját, illetve azt, hogy az arány skálájú attribútumok nem feltétlenül valós számokat tartalmaznak.

Például egy ügyfeleket leíró adatbázisban vannak bináris (pl.: büntetett előéletű-e), kategorikus (pl.: vallás, családi állapot) és intervallum (pl.: dátum) típusú attribútumok is.

Furcsa mód nem mindig triviális, hogy egy attribútum milyen típusú. Például az időjárás jellemzésére használt **napsütéses**, **borús**, **esős** értékekre mondhatjuk, hogy ez egy kategória attribútum elemei. Ugyanakkor érezzük, hogy a **borús** valahol a **napsütéses** és az **esős** között helyezkedik el, így inkább sorrend típusúnak mondanánk az attribútumot.

Az intervallum típusú attribútumok megadására is számokat használunk, amelyeknél értelme van a különbség számításának, de a hányados képzésnek nincs túl sok. Tulajdonképpen azt, hogy egy attribútum esetében mikor beszélünk intervallum és mikor arány skálájú típusról az dönti el, hogy egyértelmű-e a zérus pont definiálása. Gondoljuk meg, hogy például az évszámoknál hány fajta nullát ismerünk. Hasonló a helyzet a hőmérséklet esetében (Fahrenheit kontra Celsius).

---

**Weka 3.5.7** A *weka* saját fájlformátumát *Arff*-nak nevezik.

Az *Arff* formátum (*Attribute-Relation File Format*) egy olyan ASCII szöveges fájl formátum, mely azonos attribútummal rendelkező rekordok tárolására alkalmas. Két részből áll: fejléc (*header*) és adat (*data*). A fejléc tartalmazza az attribútumokat és azok típusát. A *weka*-ban használt adattípusok a következők: kategória (*nominal*), szám (*numeric*), karakterlánc (*string*) és dátum (*date*). Kategória típusú attribútumoknál fel kell sorolnunk az attribútum lehetséges értékeit. A sorrend fontos lehet bizonyos előfeldolgozási szűrőknél. A dátum típusú attribútumoknál megadhatjuk a dátum formátumát is. A *Data* részben minden sorban egy rekord szerepel, amelynek attribútumértékei vesszővel vannak elválasztva. A hiányzó értékeket a *?* jelzi. A *%* jellel kezdődő sorok a megjegyzéseket jelölik.

Ha az adatbázisban sok nulla érték szerepel (az gyakori elemhalmazok és az asszociációs szabályok kinyerésénél általában ez a helyzet) akkor a *sparse arff* formátumot célszerű használni. Ennél a formátumnál a *data* részben attribútum sorszám, attribútum érték párok vesszővel elválasztott sorozata áll. A nulla értékeket nem rögzítjük.

A `weka.filters.unsupervised.attribute.MergeTwoValues` szűrő összevonja egy kategória típusú attribútum két értékét. Ha az eredeti attribútum *k* különböző értéket vehet fel, akkor a szűrő alkalmazása után már csak  $(k-1)$ -et.

A `weka.filters.unsupervised.attribute.ChangeDateFormat` szűrő egy dátum formátumú attribútumom formátumát átalakít egy másik formátumba. Így egy részletes dátumformátumból (például év, hónap, nap, óra, perc, másodperc) részinformációt (például óra, perc) nyerhetünk ki.

A `weka.filters.unsupervised.attribute.NominalToBinary` minden kategória típusú attribútumot átvált bináris attribútummá. Minden olyan *A* attribútumot, amely *k* különböző értéket vehet fel ( $k > 2$ ), helyettesítünk *k* darab bináris attribútummal. Ha egy elem *A* attribútumának értéke az *i*-edik attribútum érték volt, akkor csak *i*-edik új attribútum értéke lesz egy, a többi pedig nulla.

A `weka.filters.unsupervised.attribute.NumericToNominal` szűrő a szám típusú attribútumokból kategória típusúakat állít elő. Ezt egyszerűen úgy végzi, hogy minden egyes számot kategória típusú attribútum egy értékeként kezel, és hozzáadja az attribútum



*értékhalmozáshoz.*

---

Szinte minden adatbányász/statisztikai program megadja minden intervallum típusú attribútumnak a legfontosabb statisztikáit. Ezek a

- középértékre vonatkozó adatok: mintaátlag, medián, módusz,
- szóródásra vonatkozó adatok: empirikus szórásnégyzet, minimum, maximum, terjedelem (max és min érték közötti különbség)
- eloszlásra vonatkozó adatok: empirikus kvantilisek, ferdeség, lapultság.

A *ferdeség* egy eloszlás szimmetriáját próbálja megadni. Ha a ferdeség nulla, akkor az eloszlás szimmetrikus (például normális eloszlásoknál), ellenkező esetben a várható értéktől balra (negatív ferdeség esetében) vagy jobbra „nyúlik el”. A ferdeségnek több mutatóját definiálták; ezek közül a legelterjedtebb a  $\gamma_1 = \frac{D^3[X]}{(D^2[X])^{3/2}}$ , de szokás még a  $\beta_1 = \sqrt{\gamma_1}$ -et is használni.

Szintén nem az alapfogalmak közé tartozik a *lapultság* fogalma, ami egy eloszlás csúcsosságát adja meg. A lapultságnak is több elfogadott definíciója létezik. Legelterjedtebb a  $\beta_2 = \frac{D^4[X]}{(D^2[X])^2}$  (kurtosis proper), és a  $\gamma_2 = \beta_2 - 3$  (kurtosis excess) értékek. A normális eloszlás  $\beta_2$  lapultsági értéke három, a normálisnál laposabbaké háromnál kisebb. A ferdeséget és a lapultságot annak eldöntésénél szokták használni, hogy egy adott minta származhat-e normális eloszlásból.

Kategória típusú attribútum esetén általában grafikusan ábrázolják az eloszlásokat vagy gyakoriságokat. A legjellemzőbb ábrázolási módok a kördiagrammok és a hisztogramok.

## 3.2. Hasonlósági mértékek

Az adatbányászban gyakran szükségünk lesz arra, hogy attribútumokkal leírt elemek között hasonlóságot definiáljunk. Természetesen elvárjuk, hogy ha minél inkább több azonos érték szerepel az attribútumaik között annál hasonlóbba legyenek az elemek. A gyakorlatban hasonlósági mérték helyett *különbözőségi* mértékkel dolgozunk, amely a hasonlóság inverze (minél hasonlóbba, annál kevésbé különbözők). Elvárjuk, hogy két elem különbözőségét ( $d(x, y)$ ) ki lehessen fejezni egy pozitív valós számmal, továbbá egy elem önmagától ne különbözzön, szimmetrikus legyen ( $d(x, y) = d(y, x)$ ), és teljesüljön a háromszög egyenlőtlenség ( $d(x, y) \leq d(x, z) + d(y, z)$ ). Tehát a különbözőség metrika legyen. Két elem különbözősége helyett gyakran mondunk majd két elem *távolságát*.

A következőkben sorra vesszük, hogyan definiáljuk a távolságot különböző típusú attribútumok esetében, és azt, hogy miként lehet egyes attribútumok fontosságát (súlyát) megnövelni.

### 3.2.1. Bináris attribútum

Egy bináris attribútum olyan kategória típusú attribútum, amely két értéket vehet fel (pl.: 0 és 1). Hogyan határozzuk meg  $x$  és  $y$  elemek hasonlóságát, ha azok  $m$  darab bináris attribútummal vannak leírva? Készítsük el a következő összefoglaló táblázatot.

	1	0	$\Sigma$
1	q	r	q+r
0	s	t	s+t
$\Sigma$	q+s	r+t	m

Például az 1-es sor 0-ás oszlopához tartozó érték azt jelenti, hogy  $r$  darab olyan attribútum van, amelyek az  $x$  elemnél 1-et,  $y$ -nál 0-át vesznek fel.

Ez alapján definiálhatjuk az ún. invariáns és variáns hasonlóságot. Az invariáns hasonlóságot olyan eseményeknél használjuk, amikor a bináris attribútum két értéke ugyanolyan fontos (szimmetrikus attribútum), tehát mindegy, hogy melyiket kódoljuk 0-val, illetve 1-essel. Ilyen attribútum például egy ember neme. Azért kapta ez a hasonlóság az invariáns jelzöt, mert nem változik az értéke, ha valaki máshogy kódolja az attribútumokat (tehát kódolás invariáns). A legegyszerűbb invariáns hasonlóság az eltérő attribútumok relatív száma:

$$d(x, y) = \frac{r+s}{m}.$$

Aszimmetrikus attribútum esetében a két lehetséges érték nem egyenrangú. Ilyen attribútum lehet például egy orvosi vizsgálat eredménye. Nagyobb súlya van annak a ténynek, hogy valaki fertőzött beteg, mint annak, hogy nem az. A konvencióknak megfelelően 1-essel kódoljuk a lényeges (általában ritka) kimenetet. A legegyszerűbb variáns hasonlósági mérték a Jaccard-koefficiens komplementere:

$$d(x, y) = 1 - \frac{q}{m-t} = \frac{r+s}{m-t},$$

ahol nem tulajdonítunk jelentőséget a nem jelentős kimenetek egyezésének.

Amennyiben szimmetrikus és aszimmetrikus értékek is szerepelnek a bináris attribútumok között, akkor azokat vegyes attribútumként kell kezelni (lásd a 3.2.5-os részt).

### 3.2.2. Kategória típusú attribútum

Általános esetben a kategória típusú attribútum nem csak kettő, hanem véges sok különböző értéket vehet fel. Ilyen attribútum például az ember szeme színe, családi állapota, vallása stb. A legegyszerűbb hasonlóság a nemegyezések relatív száma:

$$d(x, y) = \frac{u}{m},$$

ahol  $m$  a kategória típusú attribútumok száma,  $u$  pedig azt adja meg, hogy ezek közül mennyi nem egyezett. Természetesen a kategória típusú attribútumok sem feltétlenül szimmetrikusak, mert lehet, hogy az alapértelmezett értékek egyezése nem igazán fontos. A Jaccard-koefficiens komplementerét kategória típusú attribútumokra is felírhatjuk.

### 3.2.3. Sorrend típusú attribútum

Sorrend típusú attribútum például az iskolai végzettség: 8 általános, befejezett középiskola, érettségi, főiskolai diploma, egyetemi diploma, doktori cím. Vannak arány skálájú attribútumok, amelyeket inkább sorrend típusú attribútumnak kezelünk. Például a Forma 1-es versenyeken sem az egyes körök futási ideje számít, hanem az, hogy ki lett az első, második ...

A sorrend típusú attribútumokat általában egész számokkal helyettesítik – tipikusan 1 és  $M$  közötti egész számokkal. Ha több sorrend típusú attribútumunk van, amelyek a fontos állapotok számában eltérnek, akkor célszerű mindegyiket a  $[0,1]$  intervallumba képezni az  $\frac{x-1}{M-1}$  művelettel. Így mindegyik egyenlő súllyal szerepel majd a végső hasonlósági mértékben. Ezután alkalmazhatjuk valamelyik intervallum típusú hasonlóságot.

### 3.2.4. Intervallum típusú attribútum

Az intervallum típusú attribútumokat általában valós számok írják le. Ilyen attribútumra példa egy ember súlya, magassága, egy ország éves átlaghőmérséklete stb. Tekinthezünk úgy egy elemre, mint egy pontra az  $m$ -dimenziós vektortérben. Az elemek közötti különbözőséget a vektoraik különbségének normájával (hosszával) definiáljuk ( $d(\vec{x}, \vec{y}) = \|\vec{x} - \vec{y}\|$ ). Legtermészetesebb talán az Euklideszi-norma, de alkalmazhatjuk a Manhattan-normát is. Mindkét mérték a Minkowski-norma speciális esete.

**Euklideszi-norma:**  $L_2(\vec{z}) = \sqrt{|z_1|^2 + |z_2|^2 + \dots + |z_m|^2}$

**Manhattan-norma:**  $L_1(\vec{z}) = |z_1| + |z_2| + \dots + |z_m|$

**Minkowski-norma:**  $L_p(\vec{z}) = (|z_1|^p + |z_2|^p + \dots + |z_m|^p)^{1/p}$

A  $p = \infty$  esetén két vektor távolsága megegyezik a koordinátáinak a legnagyobb eltéréssel ( $L_\infty(\vec{z}) = \max_i \{|z_i|\}$ ).

Habár az elemek leírásában már csak számok szerepelnek, a háttérben megbújó mértékegységeknek nagy szerepük van. Gondoljuk meg, ha méter helyett milliméterben számolunk, akkor sokkal nagyobb értékek fognak szerepelni az elemek leírásában, és így a különbségek is megnőnek. A nagy értékészletű attribútumoknak nagyobb hatásuk van a hasonlóság értékére, mint a kis értékészletűeknek. Jogos tehát az egyes attribútumok normalizálása, azaz transzformáljuk őket pl. a  $[0,1]$  intervallumba, majd ezen transzformált attribútumok alapján számítsuk a távolságokat (3.3.6 rész).

Gyakran előfordul, hogy a különbözőség megállapításánál bizonyos attribútumokra nagyobb súlyt szeretnénk helyezni. Például két ember összehasonlításánál a hajszínek nagyobb szerepe van, mint annak, hogy melyik lábujja a legnagyobb. Ha figyelembe vesszük az attribútumok súlyait, akkor például az Euklideszi-távolság így módosul:

$$d(x, y) = \sqrt{w_1|x_1 - y_1|^2 + w_2|x_2 - y_2|^2 + \dots + w_m|x_m - y_m|^2},$$

ahol  $w_i$ -vel jelöltük  $i$ -edik attribútum súlyát és legyen  $\sum_{i=1}^m w_i = 1$ .

Előfordulhat, hogy olyan attribútummal van dolgunk, amely értékeit nemlineáris léptékben ábrázoljuk (nemlineáris növekedésű attribútumnak szokás hívni ezeket). Például a baktérium populációk növekedését vagy algoritmusok futási idejét exponenciális skálán érdemes ábrázolni. Az ilyen attribútumoknál nem célszerű közvetlenül intervallum alapú hasonlóságot alkalmazni, mert ez óriási különbözőségeket eredményez azokon a helyeken, ahol kis különbözőséget várunk.

„Az ideális korkülönbség férj és feleség között hat év. Egy svéd kutatás szerint ilyen esetben van maximális lehetőség az utódok születésére.”  
 Forrás: [http://hvg.hu/egeszseg/20070913\\_idealis\\_korkulonbseg.aspx](http://hvg.hu/egeszseg/20070913_idealis_korkulonbseg.aspx)



Két megközelítés között szokás választani. Egyrészt használhatjuk az intervallum alapú hasonlóságot, de nem az attribútum eredeti értékén, hanem annak logaritmusán. Másrészt diszkrétizálhatjuk az értékeket, és vehetjük csak a sorrendet a hasonlóság alapjául.

### 3.2.5. Vegyes attribútumok

Az előző részekben azt tekintettük át, hogyan definiáljuk a hasonlóságot két elem között adott típusú attribútumok esetén. Mit tegyünk akkor, ha egy objektum leírásánál vegyesen adottak a különböző típusú – intervallum, bináris, kategória stb. – attribútumok? Csoportosítsuk az egyes attribútumokat típusuk szerint, és határozzuk meg a két elem hasonlóságát minden csoportra nézve. A kapott hasonlóságokat képezzük a  $[0,1]$  intervallumba. Minden attribútumnak feleltessünk meg egy dimenziót a térben, így két elem hasonlóságához hozzárendelhetünk egy vektort a vektortérben. A hasonlóság értékét feleltessük meg a vektor hosszának.

Ennek a megközelítésnek a hátránya, hogy ha például egyetlen kategória típusú attribútum van, akkor az ugyanolyan súllyal fog szerepelni, mint akár tíz bináris attribútum összesen. Célszerű ezért az egyes attribútumtípusok által szolgáltatott értékeket súlyozni a hozzájuk tartozó attribútumok számával.

### 3.2.6. Speciális esetek

Egyre több olyan alkalmazás kerül elő, ahol a fent definiált általános hasonlóságok nem ragadják meg jól két elem különbözőségét. A teljesség igénye nélkül bemutatunk két olyan esetet, amikor speciális távolságfüggvényre van szükség.

#### Elemsorozatok hasonlósága

Elemsorozaton egy véges halmazból vett elemek sorozatát értjük. Például a magyar nyelven értelmezett szavak elemsorozatok. Nézzük az  $S = \langle abcde \rangle$  sorozatot. Legtöbben azt mondanánk, hogy a  $\langle bcdxye \rangle$  sorozat jobban hasonlít  $S$ -re, mint az  $\langle xxxddd \rangle$  sorozat. Nem ezt kapnánk, ha a pozíciókban megegyező elemek relatív számával definiálnánk a hasonlóságot.

Egy elterjedt mérték az elemsorozatok hasonlóságára az ún. *szerkesztési távolság*. Két sorozatnak kicsi a szerkesztési távolsága, ha az egyik sorozatból kevés elem törlésével ill. beszúrásával megkaphatjuk a másikat. Pontosabban, két sorozat szerkesztési távolsága adja meg, hogy legkevesebb hány beszúrás és törlés művelettel kaphatjuk meg az egyik sorozatból a másikat. A szerkesztési távolság alapján csoportosíthatunk dokumentumokat, weboldalakat, DNS sorozatokat, vagy kereshetünk illegális másolatokat.

#### Bezárt szög alapú hasonlóság

Vannak alkalmazások, ahol nem a vektorok különbségének a hossza a lényeges, hanem a vektorok által bezárt szög. Például dokumentumok hasonlóságával kapcsolatban számos okfejtést olvashatunk, hogy miért jobb szögekkel dolgozni, mint a távolságokkal. Emlékeztetőül a koszinusz-mérték pontos képlete:

$$d(x, y) = \arccos \frac{\vec{x}^T \vec{y}}{\|\vec{x}\| \cdot \|\vec{y}\|}.$$

### 3.3. Előfeldolgozás

---

**Weka 3.5.7** Az előfeldolgozási módszereket az *explorer preprocess* fülén keresztül érhetjük el. Itt adhatjuk meg a bemeneti adatot tároló fájlt (*Open file...*), url (*Open URL...*) vagy adatbázis nevét (*Open DB...*). Az *Edit...* gombra klikkelve könnyen olvasható formában megjelenik az adat, amelyet közvetlen módosíthatunk is. Ha valamelyi oszlop fejlécére kattintunk, akkor az adatokat az oszlop szerint rendezve láthatjuk.

Minden előfeldolgozási eljárást két csoportba soroljuk. A *supervised* (felügyelt) módszereknél meg kell adni egy osztályattribútumot, az *unsupervised* (felügyelet nélküli) módszereknél minden attribútum egyforma. Ezen csoportokon belül megkülönböztetünk *attribute* és *instance* eljárásokat attól függően, hogy attribútumokra (oszlopokra) vagy elemekre/objektumokra (sorokra) vonatkoznak.

Minden szűrőnél (még a felügyelet nélkülieknél is) az *ignoreClass* bináris paraméterrel állíthatjuk be, hogy a szűrés során figyelembe vegye-e az osztályattribútumot. Alapértelmezés szerint az osztályattribútum az utolsó attribútum.

---

#### 3.3.1. Hiányzó értékek kezelése

Az adatbányászati algoritmusok csak olyan elemeket tudnak kezelni, amelyeknek minden attribútuma adott. A való élet adatbázisainál ez nem mindig áll fenn, könnyen lehet, hogy bizonyos cellákat nem töltött ki az adatot begépelő személy. Sok oka lehet a hiányos mezőknek. Például az orvos bizonyos teszteken nem végzett el a páciensen, mert nem találta szükségesnek.

Persze törölhetjük azokat az elemeket, amelyek tartalmazznak hiányzó attribútumokat, de lehet, hogy ekkor annyira lecsökken az adatbázis mérete, hogy az alkalmatlan lesz az elemzésre, vagy legalábbis adott konfidencia mellett keveset tudunk mondani az összefüggésekről. A hiányzó értékeket tartalmazó elemek hasznos információt tartalmazhatnak, ne dobjuk el őket, ha nem muszáj.

A hiányzó cellákat fel kell töltenünk valamilyen értékkel, vagy a hiányzást mint külön attribútumértéket kell kezelnünk. Ez utóbbit teszi például a C4.5 nevű döntési fák előállító módszer [113] is.

Sokféleképpen helyettesíthetjük a hiányzó értékeket. Ha a hiányzó attribútum kategória típusú, akkor vehetünk egy alapértelmezett értéket, vagy az attribútum leggyakrabban előforduló értékét (módusz). Létrehozhatunk egy elem helyett sok új teljes elemet úgy, hogy a hiányzó attribútum helyére az összes lehetséges értéket beírjuk. Intervallum attribútumok esetén szokás az átlagot vagy a mediánt választani.

---

**Weka 3.5.7** A *weka.filters.unsupervised.attribute.ReplaceMissingValues* szűrő a hiányzó értékek helyettesítésére szolgál. Kategória típusú attribútumoknál a leggyakrabban előforduló értékkel (módusz), szám típusúaknál pedig az átlaggal helyettesít.



Ha osztályozási feladattal van dolgunk, akkor a fenti értékek számításánál szorítkozhatunk csak az adott osztályba tartozó elemekre. Sőt ezt a gondolatot vihetjük tovább és az értékek számításánál tekinthetjük csak azokat az elemeket (ha vannak ilyenek), amelyek attribútumainak értékei megegyeznek a hiányzó értéket tartalmazó elem ismert attribútumainak értékeivel. Itt érdemes gondosan eljárni és csak a fontos attribútumokat vizsgálni (gondoljuk meg, ha a vizsgálatnál nem zárjuk ki az azonosító attribútumot, akkor egyetlen elemet sem fogunk figyelembe venni).

### 3.3.2. Attribútum transzformációk

#### Új attribútumok létrehozása

Előfordulhat, hogy egy attribútumérték jóslásánál bizonyos attribútumok függvénye játszhat szerepet. Például rendelkezésünkre állhatnak az emberek magassága és a tömege, a betegség jóslásánál azonban a testtömeg index játszhat szerepet. Persze elvárhatjuk, hogy ezt a függvényt az osztályozó automatikusan felismerje elvégre – mint azt látni fogjuk – az osztályozás maga egy függvény approximáció. Az előzetes ismeretek, apriori tudás bevitelével azonban szinte mindig javul az osztályozás minősége. Ne várjunk csodát az osztályozótól, amikor tudunk, segítsünk neki.

---

**Weka 3.5.7** A `weka.filters.unsupervised.attribute.Add` szűrő egy új attribútumot hoz létre. Minden elem ezen attribútuma üres (hiányzó) lesz. Kategória típusú attribútum létrehozásához meg kell adnunk a lehetséges értékeket.

A `weka.filters.unsupervised.attribute.AddExpression` szűrővel új attribútumot származtathatunk meglévő attribútumokból. Az meglévő attribútumokra, mint `a1`, `a2`, ... hivatkozhatunk. A felhasználható operátorok a következők: összeadás, kivonás, szorzás, osztás, hatványozás, logaritmus, exponenciális, szinus, coszinus, tangens, egészrész képzés és a kerekítés.

A `weka.filters.unsupervised.attribute.AddID` szűrő egy azonosító attribútumot ad az adathalmazhoz. Minden elem (sor) azonosítója egyedi lesz.

A `weka.filters.unsupervised.attribute.Copy` egy meghatározott attribútumhalmazt duplikál. Ezt a szűrőt általában olyan más szűrőkkel együtt szokás használni, amelyek felülírják az adatokat. Ebben az esetben lehetővé válik az eredeti attribútum megőrzése az új mellett.

A `weka.filters.unsupervised.attribute.FirstOrder` szűrő egy  $k$  elemből álló szám típusú attribútum intervallumból készít egy  $(k-1)$ -eleműt, az egymást követő tagok különbségének képzésével. Például az  $1,2,1$  sorozatból  $1,-1$  sorozatot készít.

A `weka.filters.unsupervised.attribute.MathExpression` végrehajt egy megadott függvényt a kiválasztott típusú attribútumokon. A függvényt az `expression` paraméterrel adjuk meg ('A' betűvel lehet az attribútumra hivatkozni). A `MIN`, `MAX`, `MEAN`, `SD` változók



az attribútum minimumát, maximumát, átlagát és szórását jelölik. A támogatott műveletek listája az alábbi: +, -, \*, /, pow, log, abs, cos, exp, sqrt, tan, sin, ceil, floor, rint, (, ), A, COUNT, SUM, SUMSQUARED, ifelse.

`weka.filters.unsupervised.attribute.NumericTransform` szűrő a szám típusú attribútumokon végrehajt egy eljárást. A `className` paraméterrel adható meg az osztály, amely a felhasználni kívánt eljárást tartalmazza. A `methodName` opció segítségével adjuk meg a metódus nevét.

---

### Attribútumok törlése

Az adatbányász algoritmustól elvárjuk, hogy a lényegtelen attribútumokat ne vegye figyelembe. Szokták mondani, hogy a döntési fák nagy előnye, hogy a döntését csak a lényeges attribútumok alapján hozza meg. Ez azt sugallja, hogy nyugodtan összekapcsolhatjuk az adattáblákat és létrehozhatunk egy sok attribútumot tartalmazó táblát, a csodamódszerek majd figyelmen kívül hagyják a lényegtelen attribútumokat. Sajnos ez csak elméletben van így, a felesleges attribútumok által okozott zaj ugyanis rontja a módszerek teljesítményét. Erre a problémára majd a döntési fáknál visszatérünk.

Ha tehetjük segítsünk az adatbányász módszereken és töröljük azokat az attribútumokat (például egyedi azonosító), amelyek nem fontosak az elemzés céljából. Minden adatbányász eszköz kínál erre támogatást.

---

**Weka 3.5.7** A `weka.filters.unsupervised.attribute.Remove` törli az általunk megadott attribútumokat. Használjuk a `weka.filters.unsupervised.attribute.RemoveType` szűrőt, ha az összes, adott típusú attribútumot törölni kívánjuk.

`weka.filters.unsupervised.attribute.RemoveUseless` szűrő a haszontalan attribútumokat törli. Ezek egyáltalán nem vagy nagyon sokat változnak. A haszontalan attribútumok nem játszhatnak szerepet semmilyen adatbányászati módszerben. Minden konstans értékű attribútum haszontalan, de a másik végtel is igaz. Ha egy attribútum túl sok különböző értéket vesz fel, akkor az is haszontalan. A haszontalanság megítélésénél a különböző értékek számának arányát az összes sorhoz képest a `maximumVariancePercentageAllowed` paraméterrel adhatjuk meg.

---

### 3.3.3. Hibás bejegyzések, a zaj eltávolítása

---

**Weka 3.5.7** A `weka.filters.unsupervised.attribute.InterquartileRange` szűrő a külön pontokat és az extrém értékeket



deríti fel. Jelöljük  $Q_1$ ,  $Q_3$ -mal a 25% és 75% tartozó kvantiliseket, legyen  $IQR = Q_3 - Q_1$ , továbbá  $OF$  és  $EVF$  a felhasználó által megadott két érték (*Outlier Factor* és *Extreme Value Factor*). Extrémnek nevezünk egy értéket, ha az nagyobb, mint  $Q_3 + EVF * IQR$ , vagy kisebb, mint  $Q_1 - EVF * IQR$ . Különc pontok közé soroljuk azokat az értékeket, amelyen nem extrémek és nem esnek a  $[Q_1 - OF * IQR, Q_3 + OF * IQR]$  intervallumba sem. Ha az `outputOffsetMultiplier` paramétert igazra állítjuk, akkor a szűrő új attribútumot hoz létre, amelynek értéke a  $(A - median) / IQR$  lesz ( $A$  az attribútum érték jelöli).

A `weka.filters.unsupervised.instance.RemoveWithValues` szűrővel azokat az elemeket törölhetjük az adathalmazból, amelyek adott attribútuma adott értéket vesz fel.

A `weka.filters.unsupervised.attribute.NumericCleaner` szűrő a `maxThreshold` paraméternél nagyobb értékeket `maxDefault` értékkel, a `minThreshold` paraméternél kisebbeket `minDefault` értékkel és a `closeTo` paraméterhez közeli (`closeToTolerance`) értékeket `closeToDefault` értékkel helyettesíti.

A `weka.filters.unsupervised.instance.RemoveMisclassified` lefuttat egy osztályozó módszert, majd törli a rosszul osztályozott elemeket.

### 3.3.4. Adatok elrontása, összezagválása

Miért akarnánk elrontani az adathalmazt? Több okunk is lehet rá. Például vizsgálni szeretnénk, hogy egy adott módszer mennyire érzékeny a zajra. Az is lehet, hogy egy cég publikussá teszi bizonyos adatait, de először azt kicsit átalakítja/lerontja úgy, hogy az adatelemzés technikailag kivitelezhető legyen, de a konkurencia ne tudjon hasznos információhoz jutni. Általában a kutatóknak átadott adathalmazoknál a kutatók nem ismerik az egyes attribútumok eredeti jelentését.

**Weka 3.5.7** A `weka.filters.unsupervised.attribute.AddNoise` osztály az elemek adott részének megváltoztatja adott attribútumának értékét.

A `weka.filters.unsupervised.attribute.Obfuscate` szűrő megváltoztatja az attribútumok nevét és átnevezi az attribútumértékeket.

### 3.3.5. Diszkretizálás

A diszkretizálás/kvantálás során szám típusú attribútumot kategória típusúvá alakítjuk. Az attribútum értékészletét intervallumokra/csoportokra osztjuk és minden intervallumhoz egy kategóriát rendelünk. A diszkretizálás során nyilván információt veszítünk viszont segíthetünk az adatbányász algoritmuson. Számos módszer létezik diszkretizációra.



Kialakíthatunk egyenő szélességű vagy egyenő gyakoriságú intervallumokat. Az egyenlő gyakoriságú intervallumoknál minden intervallumba ugyanannyi adatpont esik.

---

**Weka 3.5.7** *A fenti két módszert a `weka.filters.unsupervised.attribute.Discretize` szűrőn keresztül érhetjük el. A `useEqualFrequency` paraméterrel adhatjuk meg, hogy a két lehetőség közül melyiket választjuk.*

---



PKI (Proportional k-Interval Discretization) diszkrétizációs módszerként hivatkoznak arra az esetre, amikor egyenlő gyakoriságú intervallumokat alakítunk ki és az intervallumok száma az adatpontok négyzetgyökével egyezik meg [143].

---

**Weka 3.5.7** *A PKI módszert a `weka.filters.unsupervised.attribute.PKIDiscretize` osztály implementálja.*

---



## 1R módszer

Az 1R tulajdonképpen egy egyszerű osztályozó módszer, amely tartalmaz egy diszkrétizációs eljárást. Egy példán keresztül szemléltetjük az algoritmust. A diszkrétizálandó attribútum a hőmérsékletet adja meg Fahrenheitban mérve. A tanítómintában az egyes hőmérsékletekhez a következő osztályértékek tartoznak (az attribútumértékeket nagyság szerint növekvően sorba kell rendezni):

64	65	68	69	70	71	72	72	75	75	80	81	83	85
1	0	1	1	1	0	0	1	1	1	0	1	1	0

Egy lehetséges csoportosítás szerint induljuk el a legkisebb értékektől és akkor zárjuk le az aktuális intervallumot, ha változik az osztály. A példában nyolc csoportot hoznánk létre:

64	65	68	69	70	71	72	72	75	75	80	81	83	85
1	0	1	1	1	0	0	1	1	1	0	1	1	0
1	0	1	1	1	0	0	1	1	0	0	1	0	0

A határokat a felezőpontokban megválasztva a következő határokat hoznánk létre: 64.5, 66.5, 70.5, 72, 77.5, 80.5, 84. A felosztás persze nem egyértelmű, hiszen ugyanahhoz a ponthoz tarthatnak különböző osztályok is. Erre példa a 72. Ha van egy osztály, amely a leggyakrabban fordul elő a kérdéses tanítópontok között, akkor azt az osztályt rendeljük a ponthoz. Ellenkező esetben a leggyakoribb osztályok közül azt, amelyik a legkevesebb csoportot/felosztást adja.

A túl sok kicsi intervallum létrehozásának elkerülése végett célszerű megadni egy minimális elemszám küszöböt, legalább ennyi elemet kell tartalmaznia minden csoportnak, kivéve az utolsót. Ha ez a minimum érték három, akkor a következő csoportokat hozzuk létre:

64	65	68	69	70	71	72	72	75	75	80	81	83	85
1	0	1	1	1	0	0	1	1	1	0	1	1	0
		1					1				0	v	1

Amikor a szomszédos csoportokban megegyezik a legtöbbször előforduló osztályérték, akkor a két csoport közötti határt eltörölhetjük. Ez alapján csak két intervallumot fogunk előállítani, a határvonalat a 77.5 adja. Az utolsó csoporthoz önkényesen rendeltük a 0-ás osztályértéket. Ha nem így teszünk, akkor egyáltalán nem jelölünk ki határt és minden pont egy intervallumba tartozik.

Lássuk, hogy különböző felosztás kaphatunk, attól függően, hogy a sor melyik végétől kezdjük a módszert.

### Entrópia alapú diszkrétizálás

---

**Weka 3.5.7** `weka.filters.supervised.attribute.Discretize`

---



### 3.3.6. Normalizálás

Normalizáláson azt értjük, hogy az attribútum elemeit egy másik intervallum elemeivel helyettesítjük úgy, hogy a helyettesített értékek eloszlása megegyezzen az eredeti értékek eloszlásával. Tegyük fel, hogy az  $A$  attribútum  $a_1, a_2, \dots, a_l$  értékeket vesz fel. Az  $a_j$ ,  $j = 1, \dots, l$  érték normáját  $a'_j$ -vel jelöljük. Normalizálásra két módszer terjedt el.

**Min-max normalizálás:** Itt egy sima lineáris transzformációt végzünk:  $a'_j = \frac{a_j - \min_A}{\max_A - \min_A}$ , ahol  $\min_A$  ( $\max_A$ ) jelöli az  $A$  attribútum legkisebb (legnagyobb) értékét. Minden elem a  $[0,1]$  intervallumba fog esni.

**Standard normalizálás (z-score normalization):**  $a'_j = \frac{a_j - \bar{A}}{\sigma_A}$ , ahol  $\bar{A}$  az  $A$  attribútum átlaga,  $\sigma_A$  pedig a szórása. A hagyományos szórás ( $\sqrt{\frac{\sum_{i=1}^l (a_i - \bar{A})^2}{l}}$ ) helyett az abszolút szórást is használni szokták ( $\frac{\sum_{i=1}^l |a_i - \bar{A}|}{l}$ ). Ennek előnye, hogy csökkenti az átlagtól távol eső pontok (különcök, outlier-ek) hatását.

---

**Weka 3.5.7** A két normalizáló eljárást a `weka.filters.unsupervised.attribute.Normalize` és a `weka.filters.unsupervised.attribute.Standardize` szűrők implementálják. Itt kell megemlítenünk a `weka.filters.unsupervised.attribute.Center` osztályt, amely csak annyit tesz, hogy minden értékből kivonja az átlagot ( $a'_j = a_j - \bar{A}$ ).

---



### 3.3.7. Mintavételezés

Az adatbányászati algoritmusok általában erőforrás-igényesek. Ha a bemeneti adathalmaznak csak egy kis szeletét, kis mintáját dolgozzuk fel, akkor hamarabb kapunk eredményt. A

mintavételezés következménye, hogy az így kapott eredmény nem biztos, hogy pontos, azaz lehet, hogy nem azt az eredményt kapjuk, mint amikor a teljes adathalmazt dolgozzuk fel. Vannak esetek, amikor a pontos eredménynél fontosabb a gyors adatfeldolgozás. Ilyen esetekben nagyon hasznos egy olyan mintaméret meghatározása, aminél az algoritmus gyors, és a hibázás valószínűsége kicsi.

A hiba mértékéről csak abban az esetben tudunk bővebben nyilatkozni, ha tudjuk, milyen jellegű összefüggéseket nyerünk ki. Most azt a speciális esetet nézzük meg, amikor elemek előfordulásának valószínűségét akarjuk közelíteni a relatív gyakoriságukkal. Gyakori minták, asszociációs szabályok,  $\chi^2$  alapú függetlenségvizsgálatnál ez az eset áll fenn.

Tegyük fel, hogy elemek halmazából egy tetszőleges  $x$  elem előfordulásának valószínűsége  $p$  és  $m$  megfigyelés/minta áll rendelkezésünkre. A mintavételezés hibázik, amennyiben  $x$  relatív gyakorisága eltér  $p$ -től, pontosabban a mintavételezés hibája:

$$\text{hiba}(m) = \mathbb{P}\left(\left|\text{rel. gyakoriság}(x) - p\right| \geq \epsilon\right).$$

Jelölje  $X_i$  azt a valószínűségi változót, amely 1, ha  $x$ -et választottuk egy  $i$ -edik húzásnál, különben 0, és legyen  $Y = \sum_{i=1}^m X_i$ . Mivel a húzások egymástól függetlenek, az  $Y$  eloszlása  $m, p$  paraméterű binomiális eloszlást követ. Ezt felhasználva:

$$\begin{aligned} \text{hiba}(m) &= \mathbb{P}\left(\left|\frac{Y}{m} - p\right| \geq \epsilon\right) = \mathbb{P}\left(\left|Y - m \cdot p\right| \geq m \cdot \epsilon\right) \\ &= \mathbb{P}\left(\left|Y - \mathbb{E}[Y]\right| \geq m \cdot \epsilon\right) \\ &= \mathbb{P}\left(Y \geq m \cdot (\mathbb{E}[X] + \epsilon)\right) + \mathbb{P}\left(Y \leq m \cdot (\mathbb{E}[X] - \epsilon)\right) \end{aligned}$$

A második egyenlőségnél kihasználtuk, hogy a binomiális eloszlás várható értéke  $m \cdot p$ . Tetszőleges eloszlás esetén a várható értékétől való eltérés valószínűségére több ismert korlát is létezik [128]. A Csernov-korlát (amely a Hoeffding korlát egy speciális esete) a következőket adja:

$$\mathbb{P}\left(Y \geq m \cdot (\mathbb{E}[X] + \epsilon)\right) \leq e^{-2\epsilon^2 m}$$

és

$$\mathbb{P}\left(Y \leq m \cdot (\mathbb{E}[X] - \epsilon)\right) \leq e^{-2\epsilon^2 m}$$

amiből megkapjuk, hogy:

$$\text{hiba}(m) \leq 2 \cdot e^{-2\epsilon^2 m}.$$

Amennyiben a hibakorlátot  $\delta$ -val jelölöm, akkor az alábbiak kell igaznak lennie, hogy

$$m \geq \frac{1}{2\epsilon^2} \ln \frac{2}{\delta}.$$

„Az Elevit hatékonyságát igazoló klinikai vizsgálatok közel tízezer magyar kismama bevonásával végezték. A vizsgálatok során az Elevit szedésével kilencvenkét százalékkal csökkent az idegrendszeri fejlődési rendellenességek előfordulása.” Forrás: Baba Patika X. évfolyam 10. szám, 44. oldal, 2007. október:



$\epsilon$	$\delta$	$ \mathcal{M} $
0.05	0.01	1060
0.01	0.01	27000
0.01	0.001	38000
0.01	0.0001	50000
0.001	0.01	2700000
0.001	0.001	3800000
0.001	0.0001	5000000

3.1. táblázat. A minimális minta mérete rögzített  $\epsilon, \delta$  mellett

Ha például azt szeretnénk, hogy a mintavételezés során tetszőleges elem minta, – illetve előfordulásának valószínűsége – 0.01-nál nagyobb eltérés valószínűsége kisebb legyen 1%-nál, akkor a minta mérete legalább 27000 kell legyen. A 3.1 táblázatban adott eltérés- és valószínűségkorlátokhoz tartozó minimális mintaméret található.

Gyanús, hogy a végső képletben nem szerepel  $p$ . Ez nem Csernov hibája, hanem a mienk, túl gyenge korlátot használtunk, olyat, amelyik nem vette figyelembe az  $X$  eloszlását. A Csernov-Hoeffding korlát feltételezi, hogy  $X$  binomiális eloszlású és 0,1 értékeket vehet fel.

$$\text{hiba}(m) \leq e^{-D(p+\epsilon||p)m} + e^{-D(p-\epsilon||p)m},$$

ahol  $D$  a Kullback-Leibler divergenciafüggvényt jelöli

$$D(a||b) = a \log \frac{a}{b} + (1-a) \log \frac{1-a}{1-b}.$$

A Csernov korlátot megkapjuk, ha észrevesszük, hogy  $D(p+\epsilon||p) \geq 2\epsilon^2$ .

Minek vacakolunk mi mindenféle korláttal amikor ismerjük  $Y$  sűrűségfüggvényét, így tetszőleges intervallumra meg tudjuk mondani az előfordulás valószínűségét:

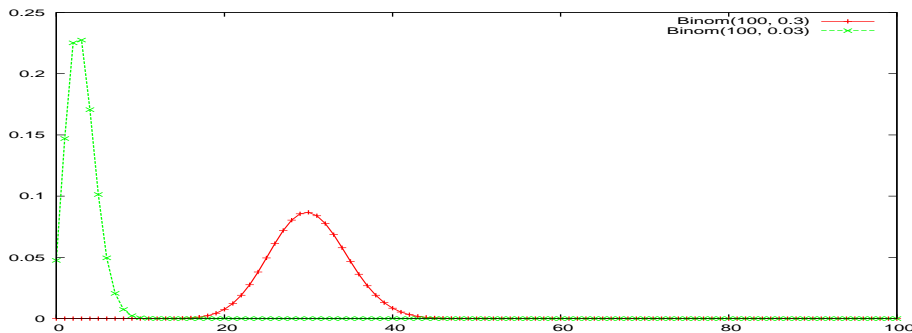
$$\begin{aligned} \mathbb{P}\left(\left|Y - m \cdot p\right| \geq m \cdot \epsilon\right) &= 1 - \sum_{i=\max\{\lfloor mp - m\epsilon \rfloor, 0\}}^{\min\{\lfloor mp + m\epsilon \rfloor, m\}} \binom{m}{i} p^i (1-p)^{m-i} \\ &= 1 + F(\max\{\lfloor mp - m\epsilon \rfloor, 0\}, m, p) - F(\min\{\lfloor mp + m\epsilon \rfloor, m\} - 1, m, p), \end{aligned}$$

ahol  $F(x, m, p)$ -vel az  $(m, p)$  paraméterű binomiális eloszlás eloszlásfüggvényét jelöljük. Sajnos a fenti képlet alapján nem tudunk szép zárt képletet adni a minta méretének alsó korlátja és az  $\epsilon, \delta$  páros közötti kapcsolatra.

Mit gondolunk? Rögzített  $m$  és  $\epsilon$  esetén kis vagy nagy  $p$  esetén lesz kicsi a hiba (mivel a binomiális eloszlás szimmetrikus, ezért szorítkozzunk  $p \leq 0.5$  esetekre)? A bevezető példa azt sugallja, hogy minél kisebb a  $p$ , annál nagyobb mintát kell venni. Ez sajnos nem így van.

Amennyiben  $p \leq \epsilon$ , akkor a  $mp - m\epsilon \leq 0$  és így a hiba  $1 - F(\lfloor mp + m\epsilon \rfloor, m, p)$ -re egyszerűsödik. Ez viszont nullához tart, amennyiben  $p \rightarrow 0$ , hiszen

$$1 - F(\lfloor mp + m\epsilon \rfloor, m, p) \leq 1 - F(\lfloor m\epsilon \rfloor, m, p) = \mathbb{P}(Y \geq \lfloor m\epsilon \rfloor) \leq \frac{mp}{\lfloor m\epsilon \rfloor}.$$

3.1. ábra. Különböző  $p$  paraméterű binomiális eloszlások

Az utolsó egyenlőtlenségénél a Markov egyenlőtlenséget használtuk fel. A 0 határértéket megkaphattuk volna úgy is, ha a Hoeffding-korlát határértékét számítjuk ki  $p \rightarrow 0$  esetén. Az eredmény ellentmond elvárásainknak, hiszen eszerint kis valószínűségeket kisebb mintával tudunk jól közelíteni.

Na, és mi van  $p \geq \epsilon$  esetén? Továbbra is igaz, hogy a  $p$  növelésével növekszik a hiba? A válasz igenlő. Ezt az állítást csak szemléltetni fogjuk. Vessünk egy pillantást a 3.1 ábrára, amelyen két, különböző  $p$  paraméterű binomiális eloszlást láthatunk.

Két dolgot vehetünk észre. A kisebb  $p$ -hez tartozó maximális valószínűség nagyobb. A nagy valószínűségek a várható érték kisebb környezetében találhatóak. Az észrevételeink általánosan is igazak. A második észrevétel például a szórással van kapcsolatban. A kisebb  $p$  paraméterű eloszlás szórása kisebb. Legyen a két paraméter  $p$  és  $q$  és legyen  $p < q < 0.5$ . Ekkor

$$\begin{aligned} mp(1-p) = \sigma_p^2 &< \sigma_q^2 = mq(1-q) \\ p - p^2 &< q - q^2 \\ 0 &< (q-p)(1-p-q) \end{aligned}$$

A kisebb valószínűségeknél a várható érték szűkebb környezetében vannak a nagy valószínűségek, ezért a várható érték  $\pm \epsilon m$  környezetén kívüli pontok valószínűségének összege kisebb, azaz a hiba kisebb!

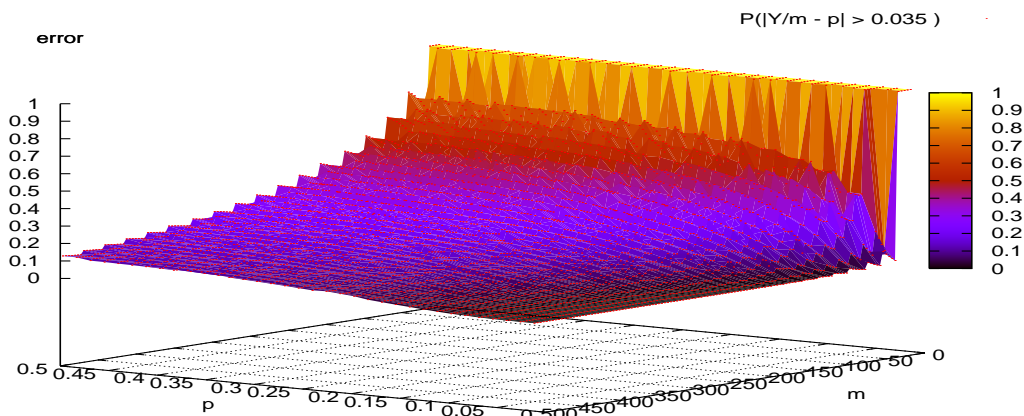
A következő ábrákon az érvelést támasztjuk alá. A 3.2 ábrán a hibát ábrázoljuk a minta mérete és a valószínűség függvényében rögzített  $\epsilon$  mellett. Látjuk, hogy ha növekszik  $p$  (vagy csökken  $m$ ), akkor csökken a hiba valószínűsége.

A 3.3 ábrán megint a mintavételezés hibáját ábrázoltuk, de most az  $\epsilon$  (0.035) mellett a minta mérete (200) is rögzítve van. Itt még jobban látszik, hogy ahogy csökken  $p$  úgy csökken a hiba is.

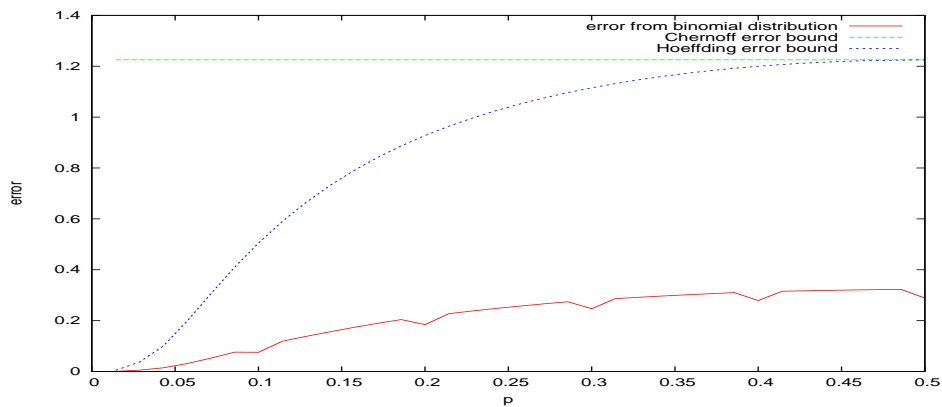
A 3.2 táblázatban a binomiális eloszlásból számol hibát és a Hoeffding-korlátot láthatjuk néhány  $p$  valószínűségre. Nyilvánvaló, hogy a Hoeffding-korlát használhatóbb, mint a Csernov-korlát és jól mutatja, hogy a  $p$  csökkenésével a hiba is csökken, ugyanakkor a tényleges valószínűségek elég távol vannak a felső korláttól.

Ha ezeknél a paramétereknél a Csernov-korlátot alkalmazzuk, akkor azt kapjuk, hogy a hiba kisebb 1.2-nél. Mivel a hibát egy valószínűséggel definiáltuk ez elég semmitmondó korlát.

Idézzük fel a kiinduló kérdést: Mit gondolunk? Rögzített  $m$  és  $\epsilon$  esetén kis vagy nagy  $p$  esetén lesz kicsi a hiba? Hát, nem mondhatjuk, hogy a várt választ kaptuk. Az ember valamiért



3.2. ábra. A mintavételezés hibája a minta méretének és az előfordulás valószínűségének függvényében



3.3. ábra. A mintavételezés hibája és a hibára adott felső korlátok az előfordulás valószínűségének függvényében ( $m = 200, \epsilon = 0.035$ )

$p$	$\mathbb{P}\left(\left Y - m \cdot p\right  \geq m \cdot \epsilon\right)$	Hoeffding
0.02	0.00078	0.01420
0.04	0.00728	0.08386
0.06	0.02431	0.21903
0.1	0.07547	0.50479
0.2	0.18433	0.92763
0.4	0.27896	1.19989

3.2. táblázat. A mintavételezés hibája és a hibára adott Hoeffding korlát néhány előfordulás valószínűsége  $m = 200$  és  $\epsilon = 0.035$  esetén

öszönösen ragaszkodik ahhoz a válaszhoz, hogy kisebb valószínűség mellett nagyobb lesz a hiba. Elemzésünk azonban pont az ellenkezőjét adta. Meg kell békélnünk ezzel, vagy tehetünk valamit a zavaró válasz ellen?

Térjünk vissza a hiba definíciójához:  $\text{hiba}(m) = \mathbb{P}\left(\left|\text{rel. gyakoriság}(x) - p\right| \geq \epsilon\right)$ , azaz hibát követünk el, ha a relatív gyakoriság és a tényleges valószínűség közötti különbség nagyobb egy adott konstansnál, amelyet  $\epsilon$ -nal jelöltünk. A relatív gyakoriságnak a valószínűség egy rögzített szélességű környezetében kell lennie.

Szerencsés az, hogy a hibát a relatív gyakoriság és a valószínűség különbségével mérjük? Ez alapján például ugyanakkora hibát követünk el, ha  $p = 0.8$  esetén a relatív gyakoriság 0.81 és ha  $p = 0.01$  esetén a relatív gyakoriság nulla, azaz az esemény nem következett be egyszer sem. Az embernek az az érzése van, hogy az első esetben kisebbet hibáztunk.

A fenti érvelés alapján célszerűbb a hibát a valószínűség és a relatív gyakoriság hányadosával mérni. Jobban érdekel minket az, hogy hány százalékkal nagyobb vagy kisebb a relatív gyakoriság a valószínűségnél, mint az abszolút különbség. Ha elfogadjuk ezt az érvelést, akkor a hibát a következőképpen definiáljuk:

$$\begin{aligned} \text{hiba}(m) &= \mathbb{P}\left(\text{rel. gyakoriság}(x)/p \geq 1 + \epsilon\right) + \mathbb{P}\left(\text{rel. gyakoriság}(x)/p \leq \frac{1}{1 + \epsilon}\right) \\ &= 1 - \mathbb{P}\left(\frac{1}{1 + \epsilon} < \text{rel. gyakoriság}(x)/p < 1 + \epsilon\right) \\ &= 1 + F(\lfloor mp/(1 + \epsilon) \rfloor, m, p) - F(\min\{\lceil mp(1 + \epsilon) \rceil, m\} - 1, m, p), \end{aligned}$$

ahol  $\epsilon > 0$  valós szám.

Felső korlát ismét létezik [52]

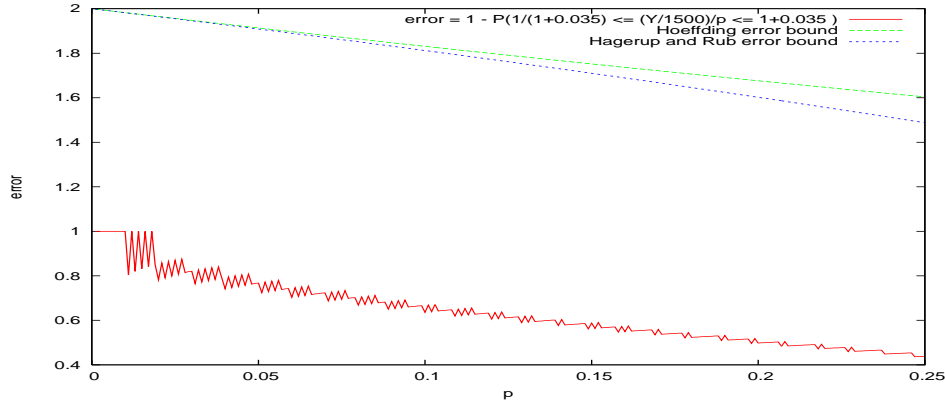
$$\mathbb{P}\left(Y/mp \geq 1 + \epsilon\right) \leq \left(\frac{e^\epsilon}{(1 + \epsilon)^{(1 + \epsilon)}}\right)^{mp},$$

továbbá

$$\mathbb{P}\left(Y/mp \leq 1 - \epsilon'\right) \leq e^{-mp\epsilon'^2/2},$$

amelyből  $\epsilon' = \epsilon/(1 + \epsilon)$  helyettesítéssel kapjuk, hogy

$$\mathbb{P}\left(Y/mp \leq 1 - \frac{\epsilon}{1 + \epsilon} = \frac{1}{1 + \epsilon}\right) \leq e^{-mp\epsilon^2/(2(1 + \epsilon)^2)},$$



3.4. ábra. A mintavételezés hibája az előfordulás valószínűségének függvényében relatív hibamérés esetében és a hibára adott felső korlát

amiből kapjuk, hogy

$$\text{hiba}(m) \leq \left( \frac{e^\epsilon}{(1+\epsilon)(1+\epsilon)} \right)^{mp} + e^{-mp\epsilon^2/(2(1+\epsilon)^2)}.$$

Hagerup és Rüb a fentieknél élesebb korlátot adott [52]. Bevezetve a  $Y \succeq k \iff Y \geq k$ , amennyiben  $k \geq mp$  és  $Y \succeq k \iff Y \leq k$  különben, fennáll, hogy

$$\mathbb{P}(Y \succeq k) \leq \left( \frac{mp}{k} \right)^k \left( \frac{m-mp}{m-k} \right)^{m-k}.$$

A hiba új definíciójánál már igaz lesz – nagyvonalakban, – hogy minél kisebb az előfordulás valószínűsége, annál nagyobb lesz a hiba, tehát annál nagyobb mintát kell vennünk. Ezt támasztja alá a 3.4 ábra is.

Az ábra mutatja, hogy tényleg csak nagyvonalakban igaz, hogy kisebb  $p$ -knél nagyobb a hiba. Szigorúan véve ugyanis ez nem igaz. Ennek oka, hogy a binomiális eloszlás diszkrét eloszlás és ezért ahogy csökkentjük a  $p$ -t és úgy tolódik nem hibát jelentő intervallum a nulla pont felé és fordulhat elő az, hogy egy újabb pont bekerül az intervallumba. Például  $\epsilon = 0.035$  és  $m = 1500$  esetében a  $[pm/(1+\epsilon), pm(1+\epsilon)]$  intervallumba nem esik egész érték  $p = 0.007$  esetében (hiszen a nem hibát jelentő intervallum  $[10.1, 10.9]$ ), míg  $p = 0.006$  esetén igen (ekkor a vizsgált intervallum  $[8.7, 9.3]$ ).

Ha  $p$  tart nullához, akkor a hiba egyhez tart. Amennyiben a  $p$  kisebb  $1/m(1+\epsilon)$ , akkor a  $(\frac{mp}{1+\epsilon}, mp(1+\epsilon))$  intervallumba nem eshet egész érték, ezért az  $X$  előfordulásától függetlenül a hiba értéke egy lesz.

Az adatbányász cikkekben mintavételezés esetén a Csernov-korlátos megközelítéssel támasztják alá, hogy az általuk használt minta miért elég nagy. Most már tudjuk, hogy ez az elemzés meglehetősen elnagyolt. Egyrészt a hiba definíciója sem túl jó, másrészt a Csernov-korlát alkalmazása sem ad pontos eredményt.

Jobb megoldás a hibát a valószínűség és a relatív gyakoriság hányadosából származtatni és Csernov-korlát helyett a binomiális eloszlást használni. Mivel a végeredmény nem egy zárt képlet lesz, ezért a hiba vagy a szükséges mintaméret kiszámítása bonyolultabb.

A binomiális eloszlás sem a legpontosabb eredményt adja. Az elemzés során ugyanis feltételeztük, hogy az esemény bekövetkezésének valószínűsége ismert. A valóságban a mintát egy nagy alaphalmazból vesszük. Például a népszavazást megelőző közvélemény-kutatásokban a mintát a felnőtt lakosságból vesszük, amely egy véges halmaz. Ha úgy tesszük fel a kérdést, hogy egy  $M$  alaphalmazból mekkora  $m$  mintát kell vennünk, hogy a mintában az  $x$  relatív gyakorisága kis mértékben térjen el az  $x$   $M$ -beli relatív gyakoriságától, akkor a binomiális eloszlás helyett hipergeometrikus eloszlást kell használnunk.

---

**Weka 3.5.7** A `weka.filters.unsupervised.instance.RemovePercentage` szűrő az elemek egy adott százalékát törli.

A `weka.filters.unsupervised.instance.RemoveFolds` megkeveri az adatbázist, majd egyenlő méretű részekre osztja és megtartja az egyik részt. A szűrőt kereszt-validációhoz szokták használni (lásd 6.10 rész). Amennyiben azt szeretnék, hogy az osztályok eloszlása megegyezzen minden részben (rétegzettség fogalma – lásd 6.10 rész), akkor használjuk a `weka.filters.unsupervised.instance.StratifiedRemoveFolds` szűrőt.

A `weka.filters.unsupervised.instance.Resample` szűrő egy véletlenszerű részhalmazát képi az adathalmaznak. A `sampleSizePercent` opcióval százalékosan fejezhetjük ki az részhalmaz méretét az eredeti adathalmazhoz képest. Megadjuk, hogy a mintavételezéshez visszatevéses, vagy visszatevés nélküli módszert használjon. Az eredeti adathalmaznak el kell férnie a memóriában. A szűrő felügyelt változata (`weka.filters.unsupervised.instance.Resample`) abban különbözik a felügyelet nélküli változattól, hogy befolyásolhatjuk a mintában az osztály eloszlását. Ha a `biasToUniformClass` értéke 0, akkor nem változik az osztály eloszlása, ha viszont 1, akkor minden osztály ugyanannyiszor fog előfordulni.

A `weka.filters.unsupervised.instance.ReservoirSample` Vitter "R" algoritmusát felhasználva mintavételez.

A `weka.filters.supervised.instance.SpreadSubsample` mintavételező szűrő olyan részhalmazt fog előállítani, amelyben a leggyakoribb és a legritkább osztályok előfordulásának hányadosa kisebb, mint egy előre megadott konstans (`distributionSpread`).

---

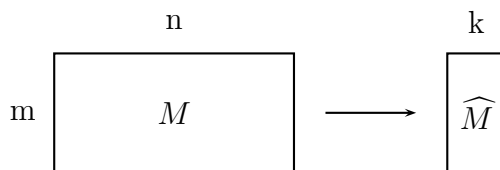
### 3.3.8. Dimenziócsökkentés

Az adatbányászati alkalmazásokban az adathalmaz mérete általában nagy. Felmerül a kérdés, hogy lehet-e ezt a nagy adathalmazt egy kisebb méretűvel helyettesíteni úgy, hogy a kisebb adathalmaz valamilyen szempont szerint hűen reprezentálja a nagy adathalmazt. Természetesen az adatbányászati feladattól függ az, hogy mit jelent pontosan a hű reprezentáció.

Ebben a részben dimenzió-csökkentésről lesz szó, melynek során az objektumok sok attribútummal való leírását szeretnénk helyettesíteni kevesebb attribútumot használó leírással.



*Hasonlóságtartó* dimenzió-csökkentésről fogunk beszélni, ami azt jelenti, hogy tudunk adni egy olyan hasonlósági definíciót az új leírásban, ami jó becslése az eredeti hasonlóságnak.



Az eredeti adathalmazt az  $m \times n$ -es  $M$  mátrixszal jellemezzük, az új leírást pedig az  $m \times k$ -s  $\widehat{M}$  mátrixszal. Az  $n$  nagyon nagy lehet (az interneten együtt előforduló szópárok keresésénél például  $10^9$  körüli volt az értéke), ami azt jelenti, hogy az adatbázis nem biztos, hogy elfér a memóriában. Ezt a problémát szeretnénk megkerülni azzal, hogy az  $M$ -et az  $\widehat{M}$  mátrixszal helyettesítjük úgy, hogy  $k \ll n$  annyira, hogy  $\widehat{M}$  elférjen a memóriában. Ezáltal lehetővé válik olyan algoritmusok futtatása, amelyek feltételezik, hogy az adatokat leíró mátrix a gyors elérési memóriában található.

Két speciális feladatot tárgyalunk. Az elsőben az attribútumok valós számok és két objektum különbözőségén (hasonlóság inverze) az Euklideszi távolságukat értjük. A második esetben az attribútumok csak binárisak lehetnek, és két objektum hasonlóságát a Jaccard-koefficiens (lásd 3.2.1 rész) adja meg.

A dimenziócsökkentés során csak a legfontosabb (esetleg származtatott) dimenziókat tartjuk meg, azokat, amelyekről úgy gondoljuk, hogy a legnagyobb szerepet játszanak két objektum hasonlóságának megállapításánál. A többi attribútumot elhagyjuk, ezért a dimenziócsökkentés zajsűrésnek is tekinthető.

### Szinguláris felbontás

A szinguláris felbontás az elméleti szempontból egyik legtöbbet vizsgált, klasszikus lineáris algebrai eszközöket használó dimenzió-csökkentési eljárás<sup>1</sup>. Ennek alkalmazása után nyert  $\widehat{M}$  mátrix soraiból jól közelíthető az euklideszi távolság, illetve az attribútumok vektoraiból számított skaláris szorzattal mért hasonlóság. Utóbbi megegyezik a koszinusz mértékkel, ha a mátrix sorai normáltak. Ebben a szakaszban néhány jelölés és alapvető fogalom után definiáljuk a szinguláris felbontást, igazoljuk a felbontás létezését, majd megmutatjuk, hogy miként használható a felbontás dimenzió-csökkentésre. Megjegyezzük, hogy a szakasz nem mutat a gyakorlatban numerikus szempontból jól alkalmazható módszert a felbontás kiszámítására. Kiseb adathalmaz esetén általános lineáris algebrai programcsomag (Matlab, Octave, Maple) használata javasolt, míg nagyobb adatbázisoknál az adatok sajátosságát kihasználó szinguláris felbontó program (SVDPack) használata ajánlott.

Egy  $U \in \mathbb{R}^{n \times n}$  mátrixot *ortogonálisnak* nevezünk, ha oszlopai ortogonális rendszert alkotnak, azaz  $U^T U = I_n$ , ahol  $I_n$  az  $n \times n$  méretű egységmátrixot, és  $U^T$  az  $U$  transzponáltját jelöli. Másképpen mondva  $U$  invertálható és  $U^{-1}$  inverzére  $U^{-1} = U^T$  teljesül. Mátrix ortogonalitásának szemléletes tárgyalásához szükségünk lesz a vektorok hosszának általánosítására, a

<sup>1</sup>A statisztikusok a szinguláris felbontást főkomponens analízisnek (angolul: principal component analysis) hívják

$$M_{m \times n} = \overbrace{\begin{pmatrix} | & & | \\ u_1 & \dots & u_m \\ | & & | \end{pmatrix}}^{U_{m \times m}} \cdot \overbrace{\begin{pmatrix} \sigma_1 & & & & \\ & \ddots & & & \\ & & \sigma_r & & \\ & & & 0 & \\ & & & & \ddots \\ & & & & & 0 \end{pmatrix}}^{\Sigma_{m \times n}} \cdot \overbrace{\begin{pmatrix} - & v_1^T & - \\ & \vdots & \\ - & v_n^T & - \end{pmatrix}}^{V_{n \times n}^T}$$

3.5. ábra. A szinguláris felbontás sematikus vázlatja.

norma fogalmára. Egy  $v \in \mathbb{R}^n$  vektor  $\|v\|_2$ -vel jelölt 2-normáját a  $\|v\|_2 = \sqrt{\sum_i v_i^2}$  egyenlőséggel definiáljuk. Egyszerűen látható, hogy  $\|v\|_2^2 = v^T v$  teljesül. A 2-norma általánosítása a tetszőleges  $M \in \mathbb{R}^{m \times n}$  mátrix esetén értelmezett  $\|M\|_F$  Frobenius-norma, amelynek definíciója  $\|M\|_F = \sqrt{\sum_{i=1}^m \sum_{j=1}^n M_{i,j}^2}$ .

Visszatérve az ortogonalitás szemléletes jelentésére, egy ortogonális mátrix által reprezentált lineáris transzformációra úgy gondolhatunk, mint egy forgatásra, amely a vektorok hosszát nem változtatja. A szemlélet alapja, hogy tetszőleges  $U \in \mathbb{R}^{n \times n}$  ortogonális mátrix és  $x \in \mathbb{R}^n$  vektor esetén

$$\|Ux\|_2 = \|x\|_2$$

teljesül. Az azonosság az alábbi elemi lépésekből következik:  $\|Ux\|_2^2 = (Ux)^T (Ux) = x^T (U^T U)x = x^T x = \|x\|_2^2$ . Hasonlóan belátható, hogy tetszőleges  $X \in \mathbb{R}^{m \times n}$  mátrix esetén és  $U \in \mathbb{R}^{m \times m}$  illetve  $V \in \mathbb{R}^{n \times n}$  ortogonális mátrixok esetén igaz, hogy

$$\|UXV^T\|_F = \|X\|_F.$$

A rövid bevezető után rátérünk a szinguláris felbontás definíciójára. Egy nem szükségszerűen négyzetes  $M \in \mathbb{R}^{m \times n}$  mátrix *szinguláris érték felbontásán* (singular value decomposition, SVD) az olyan

$$M = U\Sigma V^T$$

szorzattá bontást értjük, ahol  $U \in \mathbb{R}^{m \times m}$ ,  $V \in \mathbb{R}^{n \times n}$  ortogonális mátrixok, továbbá a  $\Sigma$  mátrix  $M$ -mel megegyező méretű és a bal felső sarokból 45°-ban lefele elhelyezkedő  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  pozitív számokat csupa 0 követ és a többi elem szintén 0. A  $\sigma_i$  számokat *szinguláris értékeknek* nevezzük, és a  $\sigma_i = 0$  választással terjesztjük ki az  $i > r$  esetre. A felbontásból látható, hogy  $\text{rang}(M) = \text{rang}(\Sigma) = r$ . Az  $U$  és a  $V$  oszlopait *bal-, illetve jobboldali szinguláris vektoroknak* mondjuk. A jelölések áttekintése a 3.5. ábrán látható.

**3.1. tétel.** *Tetszőleges  $M \in \mathbb{R}^{m \times n}$  mátrixnak létezik szinguláris érték felbontása, azaz léteznek  $U \in \mathbb{R}^{m \times m}$ ,  $V \in \mathbb{R}^{n \times n}$  ortogonális mátrixok, melyekkel*

$$M = U\Sigma V^T,$$

ahol

$$\Sigma \in \mathbb{R}^{m \times n}, \quad \Sigma = \begin{pmatrix} \Sigma_+ & 0 \\ 0 & 0 \end{pmatrix},$$



továbbá  $\Sigma_+$  egy  $r \times r$  méretű diagonális mátrix, amelynek főátlójában a  $\sigma_1 \geq \sigma_2 \geq \dots \geq \sigma_r > 0$  számok helyezkednek el sorrendben.

*Bizonyítás:* Az  $M^T M$  mátrix szimmetrikus, ezért ortogonális transzformációval diagonalizálható és sajátértékei valósak. Továbbá pozitív szemidefinit, mert tetszőleges  $x \in \mathbb{R}^{n \times n}$  vektor esetén  $x^T M^T M x = (Mx)^T (Mx) = \|Mx\|_2^2 \geq 0$ , ezért a sajátértékek nem negatívak. A sajátértékek legyenek  $\sigma_1^2 \geq \sigma_2^2 \geq \dots \geq \sigma_r^2 > 0$ . Az ezekhez tartozó sajátvektorokból alkotott ortogonális mátrixot jelölje  $V$ , ekkor

$$V^T M^T M V = \begin{pmatrix} \Sigma_+^2 & 0 \\ 0 & 0 \end{pmatrix}.$$

A mátrixot két részre osztva  $V = (V_r \ V_2)$ , ahol  $V_r \in \mathbb{R}^{n \times r}$  a pozitív sajátértékhez tartozó sajátvektorokat tartalmazza. Vagyis

$$V_r^T M^T M V_r = \Sigma_+^2.$$

Vezessük be az

$$U_r = M V_r \Sigma_+^{-1}$$

jelölést, ekkor

$$M = U_r \Sigma_+ V_r^T.$$

Az  $U_r$  vektorai ortogonális vektorrendszert alkotnak, ezt tetszőlegesen kiegészítve  $U = (U_r \ U_2)$  ortogonális mátrixszá

$$M = U \begin{pmatrix} \Sigma_+ & 0 \\ 0 & 0 \end{pmatrix} V^T.$$

■

Most megmutatjuk, hogy szinguláris felbontás segítségével hogyan lehet dimenziócsökkentést végrehajtani. Emlékeztetünk rá, hogy az  $M$  mátrix  $n$ -dimenziós sorvektorai objektumokat jellemeznek. Dimenzió-csökkentéskor az  $n$  attribútumot szeretnénk  $k < n$  dimenziójú vektorokkal jellemezni úgy, hogy közben az objektumok euklideszi távolsága vagy skaláris szorzattal mért hasonlósága csak kis mértékben változzon. A mátrixszorzás elemi tulajdonsága, hogy a szinguláris felbontás az alábbi formában is írható.

$$M = U \Sigma V^T = \sum_{i=1}^r \sigma_i u_i v_i^T,$$

ahol  $u_i v_i^T$  a bal- illetve a jobboldali szinguláris vektorokból képzett diádszorzat, azaz egy oszlop- és egy sorvektor szorzataként felírt  $m \times n$  méretű 1-rangú mátrix. Látható, hogy az  $u_i v_i^T$  diádok monoton csökkenő  $\sigma_i$  súllyal szerepelnek az összegben. Innen adódik az ötlet, hogy  $k < r$  esetén csak az első  $k$  legnagyobb súlyú diád összegével közelítsük az  $M$  mátrixot. Azaz

$$M_k = \sum_{i=1}^k \sigma_i u_i v_i^T = U_k \Sigma_k V_k^T,$$

ahol  $U_k = (u_1 \ u_2 \ \dots \ u_k)$  és  $V_k = (v_1 \ v_2 \ \dots \ v_k)$ , valamit  $\Sigma_k$  egy  $k \times k$  méretű diagonális mátrix, melynek főátlójában a  $\sigma_1, \sigma_2, \dots, \sigma_k$  értékek vannak. Könnyen látható, hogy  $M_k$  sorai egy  $k$ -dimenziós altérben helyezkednek el, hiszen  $\text{rang}(M_k) = \text{rang}(\Sigma_k) = k$ . Sokkal mélyebb eredmény a következő, melynek bizonyítását mellőzzük.

**3.2. tétel.** Legyen  $M$  egy legalább  $k$  rangú mátrix és legyen  $M_k$  a fenti módon számított közelítése. Ha a közelítés hibáját Frobenius-normával mérjük, akkor a  $k$ -rangú mátrixok közül az  $M_k$  mátrix a lehető legjobban közelíti  $M$ -et, azaz

$$\|M - M_k\|_F = \min_{N: \text{rang}(N)=k} \|M - N\|_F.$$

Továbbá a közelítés hibája a  $\sigma_i$  szinguláris értékekkel kifejezhető:

$$\|M - M_k\|_F = \sqrt{\sum_{i=k+1}^r \sigma_i^2}.$$

A közelítés relatív pontosságán a hibanégyzet egytől vett különbségét értjük, azaz

$$\frac{\sum_{i=1}^k \sigma_i^2}{\sum_{i=1}^r \sigma_i^2}. \quad (3.2)$$

Az  $M_k$  mátrix sorai az  $M$ -éhez hasonlóan  $n$  méretűek, de most már egy  $k$ -dimenziós altérnek az elemei. Ennek az altérnek egy bázisát alkotják a  $V_k^T$  sorai, és az

$$M' = U_k \Sigma_k$$

mátrix  $k$ -dimenziós sorvektorai e bázisban fejezik ki az  $M_k$  sorait. Tehát a dimenzió-csökkentés eredménye, hogy az  $M$  mátrix  $n$ -dimenziós sorait a vetítés után az  $M'$  mátrix  $k$ -dimenziós soraival közelítjük. A  $V_k^T$  sorainak ortogonalitásából könnyen belátható, hogy az  $M_k$ , illetve az  $M'$  soraiból számított euklideszi távolságok és skaláris szorzatok is megegyeznek. Tehát a közelítés alatt torzítás kizárólag az  $M$ -ből  $M_k$ -ba történő vetítés során történik, melynek mértéke a 3.2. tétel alapján felülről becsülhető.

---

**Weka 3.5.7** *A SVD-t a `weka.attributeSelection.LatentSemanticAnalysis`*

osztályon keresztül érhetjük el. Amennyiben a **rank** értéke egynél nagyobb, akkor a **rank** a  $k$ -t adja meg (figyelembe vett szinguláris értékek számát). Ellenkező esetben a közelítés relatív pontosságát definiálhatjuk (lásd a 3.2 képlet). Ha a **normalize** paraméternek igaz értéket adunk, akkor a weka az SVD elvégzése előtt az attribútumokat normalizálni fogja. Célszerű a normalizálást elvégezni, ugyanis az SVD során a hibát a Frobeniusz normával számítjuk, amely attribútumértékek különbségének négyzetével számol, így a nagy értékekkel rendelkező attribútumok nagy jelentőséget kapnak.

---

## Főkomponens analízis

---

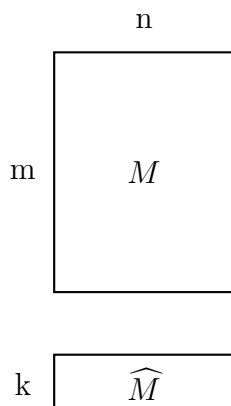
**Weka 3.5.7** *A főkomponens analízist a `weka.filters.unsupervised.attribute.PrincipalComponents` szűrő hajtja végre.*

---



### Minhash alapú lenyomat

A következőkben az adathalmaz sok oszlopot és még több sort tartalmaz. Célunk a sorok számának csökkentése. A feladatot a következő ábra szemlélteti.



Az  $M$  mátrix bináris és két oszlop (vektor) hasonlóságát a Jaccard-koefficiens adja meg:

$$d_{i,j} = \frac{\|m^i \cap m^j\|}{\|m^i \cup m^j\|} = \frac{(m^i)^T m^j}{\|m^i\|^2 + \|m^j\|^2 - (m^i)^T m^j},$$

hiszen az  $m^i(m^j)^T$  bináris vektorok esetében az azonos pozíciókban lévő 1-esek számát adja meg,  $\|m^i\|^2$  pedig a vektor egyeseinek számát. Feltételezzük, hogy a bináris vektorok ritkák azaz, ha  $r$ -el jelöljük a sorokban az 1-esek átlagos számát, akkor  $r \ll n$ .

Az  $\widehat{M}$  mátrixot az  $M$  lenyomatmátrixának fogjuk hívni. A lenyomatmátrixnak nem kell binárisnak lennie, de azt természetesen most is elvárjuk, hogy a memóriaigénye jóval kevesebb legyen, mint az  $M$  memóriaigénye. További kikötés, hogy az adatok sorfolytonosan vannak tárolva, azaz először kiolvashatjuk az első sort, majd a másodikat, és így tovább.

Ez a helyzet áll fel hasonló weboldalak kiszűrésénél, koppintások, kalózmásolatok felderítésénél, hasonló tulajdonságú felhasználók keresésénél stb. Továbbá ezt a módszert alkalmazhatjuk, amikor hasonló eladású termékpárokat keresünk. Amennyiben a termékeket kis tételben értékesítik, akkor az asszociációs szabályokat kinyerő technikák (lásd 5 fejezet) nem alkalmazhatóak.

Gondolkozzunk el azon, hogy működik-e az alábbi algoritmus. Válasszunk ki néhány sort véletlenszerűen és tekintsük ezeket lenyomatoknak. Két lenyomat hasonlóságának várható értéke meg fog egyezni az oszlopaik hasonlóságával. Ez alapján azt mondhatnánk, hogy a sorok egy véletlenszerűen választott halmaza jó lenyomat.

A fentiek ellenére ez az egyszerű módszer nagyon rossz eredményt adna. Ennek oka az, hogy a mátrixunk nagyon ritka ( $r \ll n$ ), tehát egy oszlopban a legtöbb elem 0, így nagy valószínűséggel a legtöbb lenyomat is csupa 0 elemből állna.

A minhash alapú lenyomat egy elemét a következőképpen állítjuk elő. Véletlenszerűen permutáljuk meg a sorokat, majd válasszuk az  $j$ -edik oszlopok hash értékének ( $h$ ) azt a legkisebb sorindexet, ahol 1-es szerepel a  $j$ -edik oszlopban. A véletlen permutáció természetesen csak elméleti megközelítés, diszken található nagy adatbázis esetén túl lassú művelet. Ehelyett sorsoljunk ki minden sorhoz egy véletlen hash értéket. Amennyiben feltehetjük, hogy a

mátrix sorainak száma  $2^{16}$ -nál kisebb, akkor a születésnap paradoxon<sup>2</sup> alapján válasszunk 32 bit szélességű egyenletes eloszlású véletlen számot. Az algoritmus tényleges implementálása során tehát egyesével olvassuk a sorokat, véletlen számot generálunk, és minden oszlopnak folyamatosan frissítjük azt a változóját, ami megadja a legkisebb, 1-est tartalmazó sorindexet.

Mivel egy lenyomatnak  $k$  darab eleme van, ezért minden oszlophoz  $k$  darab véletlen számot állítunk elő, és  $k$  darab hash értéket tároló változót tartunk karban. Vegyük észre, hogy a lenyomat előállításához egyszer megyünk végig a mátrixon.

Két lenyomat hasonlóságát a páronként egyező lenyomatok számának  $k$ -hoz vett aránya adja meg, azaz

$$\widehat{d}_{ij} = \frac{|\{\ell : \widehat{M}_{i,\ell} = \widehat{M}_{j,\ell}\}|}{k},$$

ahol  $\widehat{M}_{i,\ell}$  az  $\widehat{M}$  mátrix  $i$ -edik oszlopának  $\ell$ -edik elemét jelöli.

Be fogjuk bizonyítani, hogy  $\widehat{d}_{ij}$  jó becslése  $d_{ij}$ -nek abban az értelemben, hogy ha  $i$  és  $j$  oszlopok nagyon hasonlóak, akkor azok lenyomatai is nagy valószínűséggel hasonlóak. Ehhez a következő észrevételt használjuk fel.

**3.3. észrevétel.** *Tetszőleges  $(i, j)$  oszloppárra igaz, hogy*

$$\mathbb{P}[\widehat{M}_{i,\ell} = \widehat{M}_{j,\ell}] = d_{ij}.$$

*Bizonyítás:* Csak akkor lehet a két lenyomat azonos, ha a legalább az egyik oszlopban az 1-est tartalmazó indexek közül olyan sor kapta a legkisebb véletlen számot, amelynél mindkét oszlopban 1-es szerepel. Ennek valószínűsége éppen  $d_{ij}$ , amennyiben a permutáció egyenletesen szórja szét az egyeseket. ■

És most a hasonlóság megőrzésével kapcsolatos állítás:

**3.4. tétel.** *Legyenek  $0 < \delta < 1$ , és  $\epsilon > 0$  valós számok. Amennyiben  $k > -\frac{\ln \delta/2}{2\epsilon^2}$ , akkor  $\delta$ -nál kisebb a valószínűsége annak, hogy a lenyomat és az eredeti hasonlóság különbsége  $\epsilon$ -nál nagyobb.*

*Bizonyítás:* Tekintsük az  $i, j$  oszlopokat. Definiáljuk  $X_l$  valószínűségi változót, ami 1  $\widehat{M}_{i,\ell} = \widehat{M}_{j,\ell}$  esetén, különben 0. Legyen  $Y = X_1 + \dots + X_k$ .

$X_l$  binomiális eloszlású és az előzőekben kimondott észrevétel miatt  $E[X_l] = p = \mathbb{P}(\widehat{M}_{i,\ell} = \widehat{M}_{j,\ell}) = d_{ij}$ . A lenyomatok hasonlóságának definíciójából adódik, hogy  $\widehat{d}_{ij} = \frac{Y}{k}$ . Írjuk fel  $Y$ -re 2.5.2-es tételét:

$$\mathbb{P}(|Y - E[Y]| > k\epsilon) \leq 2e^{-2\epsilon^2 k},$$

amiből adódik, hogy

$$\mathbb{P}(|\widehat{d}_{ij} - d_{ij}| > \epsilon) \leq 2e^{-2\epsilon^2 k}.$$

---

<sup>2</sup>A születésnap paradoxonnal kapcsolatos kérdés a következő: „Mekkora a valószínűsége annak az eseménynek, hogy emberek egy véletlenszerűen választott  $r$  fős csoportjában van legalább két személy, akik egy napon ünneplik a születésnapjukat?”. Elemi kombinatorikus úton a válasz meghatározható:  $p_r = 1 - \frac{\binom{365}{r} \cdot r!}{365^r} \approx 1 - \exp\left(-\frac{r^2}{3 \cdot 365}\right)$ . A feladat következménye az az állítás, miszerint  $2^n$  elemnek  $2^{2n}$  elemű halmazból kell egyenletes eloszlás szerint véletlenszerűen egyesével kulcsot sorsolni, hogy kicsi ( $\exp(-3) < 0.05$ ) legyen annak valószínűsége, hogy két elem ugyanazt a kulcsot kapja.

## 4. fejezet

# Gyakori elemhalmazok

A gyakori elemhalmazok kinyerése az adatbányászat eltulajdoníthatatlan területe. A feladat vásárlói szokások kinyerésénél merült fel részfeladatként. A nagy profitot elsősorban a gyakran együtt vásárolt termékek, termékhalmozok jelentik, így ezek kinyerése jelentette az első lépést a feladat megoldásánál.

Egyes alkalmazásokban a gyakori részstruktúrák, gyakori minták meghatározásánál elemhalmazok helyett sorozatok, gyökeres fák, címkézett gráfok vagy bool-formulákat kerestek. A 10 fejezetben bemutatjuk a gyakori minták bányászatának absztrakt modelljét, majd egyesével vesszük a különböző típusú mintákat és megvizsgáljuk, hogy milyen technikák alkalmazhatók.

A gyakori elemhalmazok bányászata nagyon népszerű kutatási terület. A publikált algoritmusokkal könyveket lehetne megtölteni. Ebben a jegyzetben csak a leghíresebb algoritmusokat és ötleteket ismertetjük.

### 4.1. A gyakori elemhalmaz fogalma

Legyen  $\mathcal{I} = \{i_1, i_2, \dots, i_m\}$  elemek halmaza és  $\mathcal{T} = \langle t_1, \dots, t_n \rangle$  az  $\mathcal{I}$  hatványhalmaza felett értelmezett sorozat, azaz  $t_j \subseteq \mathcal{I}$ . A  $\mathcal{T}$  sorozatot *bemeneti sorozatnak* hívjuk, amelynek  $t_j$  elemei a *tranzakciók*. Az  $I \subseteq \mathcal{I}$  elemhalmaz *fedése* megegyezik azon tranzakciók sorozatával, amelyeknek részhalmaza az  $I$ . Az  $I$  elemhalmaz *támogatottsága* a fedésének elemszámával egyezik meg (jelölésben  $\text{supp}(I)$ ). Az  $I$  *gyakori*, amennyiben támogatottsága nem kisebb egy előre megadott konstansnál, amelyet hagyományosan *min-supp*-pal jelölünk, és *támogatottsági küszöbnek* hívunk. A gyakori elemhalmazok keresése során adott egy  $\mathcal{I}$  elemhalmaz,  $\mathcal{T}$  bemeneti sorozat, *min-supp* támogatottsági küszöb, feladatunk meghatározni a gyakori elemhalmazokat és azok támogatottságát. Az egyszerűség kedvéért a halmazt jelölő kapcsos zárójeleket (sőt az elemek határoló vesszőt) gyakran elhagyjuk, tehát például az  $\langle \{A, C, D\}, \{B, C, E\}, \{A, B, C, E\}, \{B, E\}, \{A, B, C, E\} \rangle$  sorozatot  $\langle ACD, BCE, ABCE, BE, ABCE \rangle$  formában írjuk.

Az általánosság megsértése nélkül feltehetjük, hogy az  $\mathcal{J}$  elemein tudunk egy rendezést definiálni, és a minták illetve a tranzakciók elemeit minden esetben nagyság szerint növekvő sorrendben tároljuk. Ezen rendezés szerinti lexikografikusan tudjuk rendezni az azonos méretű halmazokat.

A keresési teret úgy képzelhetjük el, mint egy irányított gráfot, amelynek csúcsai az elemhalmazok, és az  $I_1$ -ből él indul  $I_2$ -be, amennyiben  $I_1 \subset I_2$ , és  $|I_1| + 1 = |I_2|$ . A keresési tér bejárásán mindig ezen gráf egy részének bejárását fogjuk érteni. Tehát például a keresési tér szélességi

bejárása ezen gráf szélességi bejárását jelenti.

Elterjedt, hogy a támogatottság helyett *gyakoriságot*, a támogatottsági küszöb helyett *gyakorisági küszöböt* használnak, melyeket  $freq(I)$ -vel, illetve  $min\_freq$ -kel jelölnék. Az  $I$  elemhalmaz gyakoriságán a  $supp(I)/|T|$  hányadost értjük.

A gyakorlatban előforduló adatbázisokban nem ritka, hogy az elemek száma  $10^5 - 10^6$ , a tranzakcióké pedig  $10^9 - 10^{10}$ . Elméletileg már az eredmény kiírása is az  $T$  elemszámában exponenciális lehet, hiszen előfordulhat, hogy  $T$  minden részhalmaza gyakori. A gyakorlatban a maximális méretű gyakori elemhalmaz mérete  $|T|$ -nél jóval kisebb (legfeljebb 20-30). Ezen kívül minden tranzakció viszonylag kicsi, azaz  $|t_j| \ll |T|$ . A keresési tér tehát nagy, ami azt jelenti, hogy az egyszerű nyers erő módszerek (határozzuk meg minden elemhalmaz támogatottságát, majd válogassuk ki a gyakoriakat) elfogadhatatlanul lassan futnának.

A későbbiekben gyakran használjuk majd tranzakciók esetén a „szűrt” jelzőt. Egy tranzakció szűrt tranzakcióját úgy kaphatjuk meg, ha töröljük belőle a ritka elemeket. A szűrt tranzakciók minden információt tartalmaznak a gyakori elemhalmazok kinyeréséhez, ezért a legtöbb algoritmus első lépése a gyakori elemek meghatározása, majd a szűrt tranzakciók előállítás. Ezután az eredeti adatbázist nem használják többé.

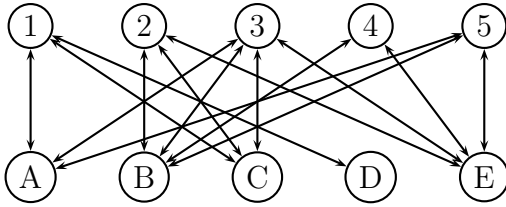
A bemenetet illetően három adattárolási módot szoktak elkülöníteni. *Horizontális adatbázisról* beszélünk, ha a tranzakciókat azonosítóval látjuk el, és minden azonosítóhoz tároljuk a tranzakcióban található elemeket. *Vertikális adatbázisnál* minden elemhez tároljuk az elemet tartalmazó tranzakciók azonosítóit (sorszámát). A vertikális tárolás nagy előnye, hogy gyorsan megkaphatjuk egy elemhalmaz fedését (az elemekhez tartozó kosarak metszetét kell képezni), amiből közvetlen adódik a támogatottság. Mind a horizontális, mind a vertikális ábrázolási módnál használhatunk az elemek vagy tranzakciók felsorolása helyett rögzített szélességű bitvektorokat. Az  $i$ -edik elem (tranzakció) meglétét az  $i$ -edik pozícióban szereplő 1-es jelzi.

tranzakció	elemhalmaz	elem	tranzakcióhalmaz	tranzakció	elem
1	C	A	2	1	C
2	A,B,C	B	2	2	A
		C	1,2	2	B
				2	C

4.1. táblázat. Horizontális-, vertikális- és relációs tárolási mód

Tudjuk, hogy egy tranzakcióban változó számú elem lehet (és fordítva: egy elem változó számú tranzakcióban szerepelhet). A legtöbb mai adatbázis *relációs táblák* formájában van elmentve, amelyekben csak rögzített számú attribútum szerepelhet. A valóságban ezért a tranzakciók két attribútummal rendelkező relációs tábla formájában találhatók, ahol az első attribútum a tranzakciót, a második pedig az elemet adja meg (pontosabban a tranzakciók és az elemek azonosítóit). A három tárolási módszerre mutatnak példát a 4.1 táblázatok.

A bemenetet elemhalmazok sorozataként definiáltuk. Ábrázoljuk ezt, mint  $G = (I, T, R)$  irányítatlan, páros gráf, vagy mint  $B$  bináris mátrix. Ha a  $t$  tranzakció tartalmazza az  $i$  elemet, akkor és csak akkor az  $(i, t)$  eleme  $R$ -nek. Vagy mátrix esetén a  $t$  sorának  $i$  eleme 1 (különben 0). A  $\langle ACD, BCE, ABCE, BE, ABCE \rangle$  bemenethez tartozó gráf és bináris mátrix a 4.1 és a 4.2 ábrán látható.



4.1. ábra. Gráfes ábrázolási mód

	A	B	C	D	E
1	1		1	1	
2		1	1		1
3	1	1	1		1
4		1			1
5	1	1	1		1

4.2. ábra. Bináris mátrixos ábrázolási mód

A bemeneti adatot szokták a *sűrű* (dense) illetve a *ritka* (sparse) jelzővel illetni, amellyel a bináris mátrixban található 1-esek számára utalnak. Vásárlói kosarakat ábrázoló mátrix tipikusan ritka, ugyanis a kosarakban általában jóval kevesebb termék van (50-100), mint az összes termék száma (10 000-100 000).

A tranzakciók száma általában nagy, de a mai tárolókapacitások mellett, még egészen nagy adatbázisok is elérnek a memóriában. Gondoljuk meg például, hogy egy  $10^7$  tranzakciót tartalmazó adatbázis csak 120 MB helyet kíván, amennyiben a tranzakciók átlagos mérete 6 elem. Csak extrém nagy adathalmazok esetén nem alkalmazhatók azok az algoritmusok, amelyek feltételezik, hogy a bemenet (vagy a szűrt tranzakciók) elérnek a memóriában.

Mielőtt bemutatjuk az APRIORI módszert elemhalmazok esetén, gondolkozzunk el azon, vajon működne-e az alábbi egyszerű algoritmus a gyakorlatban. Olvassuk be a háttértárolóból az adatbázis első blokkját, és vizsgáljuk meg az első tranzakciót. Ennek a  $t_1$  tranzakciónak az összes részhalmazát tároljuk el a memóriában és mindegyikhez rendeljük egy számlálót 1 kezdeti értékkel. Az  $I$  elemhalmazhoz rendelt számláló fogja tárolni  $I$  támogatottságát. Az első tranzakció feldolgozása után vizsgáljuk meg sorban a többit: a  $t_i$  tranzakció minden részelemhalmazának számlálóját növeljük eggyel, vagy vegyük fel a memóriába egy új számlálót, amennyiben az eddig feldolgozott tranzakcióban még nem fordult elő. Az adatbázis teljes végigolvasása után az összes – valahol előforduló – elemhalmaz támogatottsága rendelkezésünkre áll, amiből könnyen megkaphatjuk a gyakoriakat.

Látható, hogy ennél az egyszerű algoritmusnál IO szempontjából gyorsabbat nem lehet találni, mert az adatbázis egyszeri végigolvasása mindenképpen szükséges a támogatottság meghatározásához és ennél az algoritmusnál elég is. A gyakorlatban mégsem használják ezt a gyors és egyszerű algoritmust. Ennek oka, hogy az életben előforduló adatbázisokban nem ritka, hogy valamelyik tranzakció sok elemet tartalmaz. Egy átlagos supermarketben mindennapos, hogy valaki 60 különböző elemet vásárol. Ekkor csak a számlálók mintegy 16 ezer TB-ot foglalnának a memóriából, amennyiben a számlálók 4 byte-osak. A számlálókat mindenképpen a memóriában szeretnénk tartani, hogy elkerüljük a folyamatos swappelést, hiszen egy új tranzakció vizsgálatánál nem tudjuk előre, hogy melyik számlálót kell növelni.

Abban az esetben, ha biztosan tudjuk, hogy a tranzakciók egyike sem tartalmaz sok elemet, vagy az adatbázis bináris értékeket tartalmazó mátrix formájában adott, ahol az oszlopok (attribútumok) száma kicsi, akkor a fenti algoritmus hatékonyan használható.

A fenti algoritmus kis módosítását javasolták [9]-ben. Egyrészt csak olyan elemhalmazokat vizsgáltak, amelyek mérete nem halad meg egy előre megadott korlátot, másrészt a vizsgált elemhalmazokat és számlálóikat – a gyors visszakeresés érdekében – szófában tárolták. A módszernek két súlyos hátránya van: nem teljes (az algoritmus nem találja meg azokat az elemhalmazokat, amelyek mérete nagyobb az előre megadott küszöbnél), továbbá túlságosan nagy a memóriaigénye (sok lehet a hamis jelölt).

Amennyiben az adatbázisunk kicsi, akkor még a fenti egyszerű algoritmusokat sem kell leprogramoznunk, mert egy teljesen szabványos adatbázis-lekérdező nyelv segítségével megkaphatjuk a gyakori elemhalmazokat. Az alábbi SQL parancs a gyakori elem párokat adja eredményül.

```
SELECT I.elem,J.elem, COUNT(I.tranzakció)
FROM tranzakciók I, tranzakciók J
WHERE I.tranzakció=J.tranzakció AND I.elem<J.elem
GROUP BY I.elem, J.elem
HAVING COUNT(I.tranzakció) >= min_supp
```

4.3. ábra. SQL utasítás gyakori elem párok kinyeréséhez

Látnunk kell, hogy a fenti parancs az összekapcsolás (FROM mezőben két tábla) művelet miatt nem fog működni, ha az adatbázis mérete túl nagy.

A következőkben bemutatjuk a három leghíresebb gyakori elemhalmazokat kinyerő (GYEK) algoritmust. Mindhárman az üres mintából indulnak ki. Az algoritmusok egy adott fázisában *jelöltnek* hívjuk azokat az elemhalmazokat, amelyek támogatottságát meg akarjuk határozni. Az algoritmus akkor *teljes*, ha minden gyakori elemhalmazt megtalál és *helyes*, ha csak a gyakoriakat találja meg.

Mindhárom algoritmus három lépést ismétel. Először jelölteket állítanak elő, majd meghatározzák a jelöltek támogatottságát, végül kiválogatják a jelöltek közül a gyakoriakat. Természetesen az egyes algoritmusok különböző módon járják be a keresési teret (az összes lehetséges elemhalmazt), állítják elő a jelölteket, és különböző módon határozzák meg a támogatottságokat.

Az általánosság megsértése nélkül feltehetjük, hogy az  $\mathcal{I}$  elemein tudunk definiálni egy teljes rendezést, és a jelöltek, illetve a tranzakciók elemeit ezen rendezés szerint tároljuk. Más szóval az elemhalmazokat sorozatokká alakítjuk. Egy sorozat  $\ell$ -elemű *prefixén* a sorozat első  $\ell$  eleméből képzett részsorozatát értjük. A példákban majd, amennyiben a rendezésre nem térünk ki külön, az ábécé szerinti sorrendet használjuk. A GYEK algoritmusok általában érzékenyek a használt rendezésre. Ezért minden algoritmusnál megvizsgáljuk, hogy milyen rendezést célszerű használni annak érdekében, hogy a futási idő, vagy a memóriaszükséglet a lehető legkisebb legyen.

A jelölt-előállítás *ismétlés nélküli*, amennyiben nem állítja elő ugyanazt a jelöltet többféle módon. Ez a hatékonyság miatt fontos, ugyanis ismétléses jelölt-előállítás esetében minden jelölt előállítása után ellenőrizni kellene, hogy nem állítottuk-e elő már korábban. Ha ezt nem tesszük, akkor feleslegesen kötünk le erőforrásokat a támogatottság ismételt meghatározásánál. Mindhárom ismertetett algoritmusban a jelöltek előállítása ismétlés nélküli lesz, amit a rendezéssel tudunk garantálni.

Az algoritmusok pszeudokódjaiban *GY*-vel jelöljük a gyakori elemhalmazok halmazát, *J*-vel a jelöltekét és *j*.számláló-val a *j* jelölt számlálóját. Az olvashatóbb kódok érdekében feltesszük, hogy minden számláló kezdeti értéke nulla, és az olyan halmazok, amelyeknek nem adunk kezdeti értéket (például *GY*), nem tartalmaznak kezdetben egyetlen elemet sem.



## 4.2. Az Apriori algoritmus

Az APRIORI algoritmus az egyik legelső GYEK algoritmus. Szélességi bejárást valósít meg, ami azt jelenti, hogy a legkisebb mintából (ami az üres halmaz) kiindulva szintenként halad előre a nagyobb méretű gyakori elemhalmazok meghatározásához. A következő szinten (iterációban) az eggyel nagyobb méretű elemhalmazokkal foglalkozik. Az iterációk száma legfeljebb eggyel több, mint a legnagyobb gyakori elemhalmaz mérete.

A jelöltek definiálásánál a következő egyszerű tényt használja fel: *Gyakori elemhalmaz minden részhalmaza gyakori*. Az állítást indirekten nézve elmondhatjuk, hogy egy elemhalmaz biztosan nem gyakori, ha van ritka részhalmaza. Ennek alapján ne legyen jelölt azon elemhalmaz, amelynek van ritka részhalmaza. Az APRIORI algoritmus ezért építkezik letről. Egy adott iterációban csak olyan jelöltet veszünk fel, amelynek összes valódi részalmazáról tudjuk, hogy gyakori. Az algoritmus onnan kapta a nevét, hogy az  $\ell$ -elemű jelölteket a bemeneti sorozat  $\ell$ -edik átolvasásának megkezdése előtt (a priori) állítja elő. Az  $\ell$ -elemű jelöltek halmazát  $J_\ell$ -l, az  $\ell$ -elemű gyakori elemhalmazokat pedig  $GY_\ell$ -l jelöljük.

---

### Algorithm 1 APRIORI

---

**Require:**  $\mathcal{T}$ : tranzakciók sorozata,  
 $min\_supp$ : támogatottsági küszöb,  
 $\ell \leftarrow 0$   
 $J_\ell \leftarrow \{\emptyset\}$   
**while**  $|J_\ell| \neq 0$  **do**  
    támogatottság\_meghatározás( $\mathcal{T}, J_\ell$ );  
     $GY_\ell \leftarrow$  gyakoriak\_kiválogatása( $J_\ell, min\_supp$ );  
     $J_{\ell+1} \leftarrow$  jelölt\_előállítás( $GY_\ell$ );  
     $\ell \leftarrow \ell + 1$ ;  
**end while**  
return  $GY$

---

A kezdeti értékek beállítása után egy ciklus következik, amely akkor ér véget, ha nincsen egyetlen  $\ell$ -elemű jelölt sem. A cikluson belül először meghatározzuk a jelöltek támogatottságát. Ehhez egyesével vesszük a tranzakciókat, és azon jelöltek számlálóját növeljük eggyel, amelyeket tartalmaz a vizsgált tranzakció. Ha rendelkezésre állnak a támogatottságok, akkor a jelöltek közül kiválogathatjuk a gyakoriakat.

---

**Weka 3.5.7** *Gyakori elemhalmazokat úgy nyerhetünk ki, ha asszociációs szabályokat keresünk APRIORI algoritmussal. Ehhez az Associate fülre kell kattelnünk, majd a Apriori kell választanunk, mint Associator. Alapból a módszer csak asszociációs szabályokat nyer ki, de ha az outputItemSets paramétert igazra állítjuk, akkor megkaphatjuk a gyakori elemhalmazokat is. A módszer fő hátránya, hogy asszociációs szabályokat keres, nem pedig gyakori elemhalmazokat ezért nem tudjuk elérni, hogy az adott min\_supp-nál nagyobb támogatottságú elemhalmazokat adja meg. A weka.associations.Apriori osztályról bővebben az asszociációs szabályok fejezetben írunk a 106. oldalon.*

---



### 4.2.1. Jelöltek előállítás

A JELÖLT-ELŐÁLLÍTÁS függvény az  $\ell$ -elemű gyakori elemhalmazokból  $(\ell + 1)$ -elemű jelölteket állít elő. Azok és csak azok az elemhalmazok lesznek jelöltek, amelyek minden részhalmaza gyakori.

A jelöltek előállításán olyan  $\ell$ -elemű, gyakori  $I_1, I_2$  elemhalmaz párokat keresünk, amelyekre igaz, hogy

- $I_1$  lexikografikusan megelőzi  $I_2$ -t,
- $I_1$ -ből a legnagyobb elem törlésével ugyanazt az elemhalmazt kapjuk, mintha az  $I_2$ -ből törölnénk a legnagyobb elemet.

Ha a feltételeknek megfelelő párt találunk, akkor képezzük a pár unióját, majd ellenőrizzük, hogy a kapott elemhalmaznak minden valódi részhalmaza gyakori-e. A támogatottság anti-monotonitása miatt szükségtelen az összes valódi részhalmazt megvizsgálni; ha mind az  $\ell + 1$  darab  $\ell$ -elemű részhalmaz gyakori, akkor az összes valódi részhalmaz is gyakori. Az  $I_1, I_2$  halmazokat a jelölt *generátorainak* szokás hívni.

**4.1. példa.** Legyenek a 3-elemű gyakori elemhalmazok a következők:  $GY_3 = \{ABC, ABD, ACD, ACE, BCD\}$ . Az  $ABC$  és  $ABD$  elemhalmazok megfelelnek a feltételnek, ezért képezzük az uniójukat. Mivel  $ABCD$  minden háromelemű részhalmaza a  $GY_3$ -nak is eleme, az  $ABCD$  jelölt lesz. Az  $ACD, ACE$  pár is megfelel a két feltételnek, de uniójuknak van olyan részhalmaza ( $ADE$ ), amely nem gyakori. Az APRIORI a következő iterációban tehát már csak egyetlen jelölt támogatottságát határozza meg.

A fenti módszer csak akkor alkalmazható, ha  $\ell > 0$ . Az egyelemű jelöltek előállítását egyszerű: minden egyelemű halmaz jelölt, amennyiben az üres elemhalmaz gyakori ( $|T| \geq \min\_supp$ ). Ez összhangban áll azzal, hogy akkor lehet egy elemhalmaz jelölt, ha minden részhalmaza gyakori.

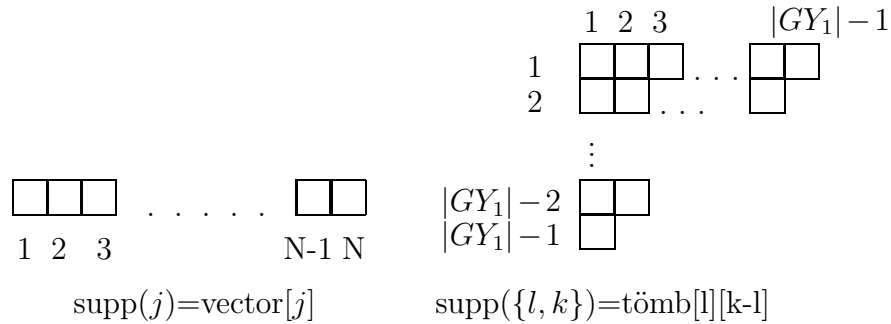
### 4.2.2. Jelöltek támogatottságának meghatározása

A jelöltek előfordulásait össze kell számolni. Ehhez egyesével vizsgáljuk a kosarakat, és azon jelöltek számlálóját növeljük eggyel, amelyeket tartalmaz a kosár.

#### 1- és 2-elemű jelöltek támogatottsága

Könnyű dolgunk van, amennyiben a jelöltek mérete 1 vagy 2. A feladatot megoldhatjuk egy olyan lista, illetve féltömb segítségével, amelyekben a számlálót tároljuk. Az elemek támogatottságának meghatározásánál a lista  $j$ -edik eleme tárolja a  $j$ -edik elem számlálóját. A tranzakciók feldolgozásánál végigmegyünk a tranzakció elemein és növeljük a megfelelő cellákban található számlálót.

Az első végigolvasás után kiválogathatjuk a gyakori elemeket. A továbbiakban már csak ezekkel az elemekkel dolgozunk, így új sorszámokat adhatunk nekik a  $[1..|GY_1|]$  intervallumból (emlékeztetőül  $GY_j$ -vel jelöljük a  $j$ -elemű gyakori mintákat). Az  $l$  és  $k$ -edik elemekből álló pár



4.4. ábra. Adatstruktúrák az 1- és 2-elemű jelöltek támogatottságának meghatározásához.

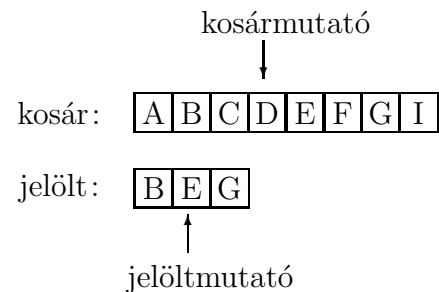
támogatottságát a tömb  $l$ -edik sorának  $(k-l)$ -edik eleme tárolja (az általánosság megsértése nélkül feltehetjük, hogy  $l < k$ ).

Ha egy számláló 4 byte-ot foglal, akkor a tömb helyigénye nagyjából  $4 \cdot \binom{|GY_1|}{2}$  byte. Azon elempárokhoz tartozó tömbelem értéke, amelyek sosem fordulnak elő együtt, 0 lesz. Helyet takaríthatunk meg, hogy ha csak akkor vesszük fel egy jelöltpár számlálóját, ha a párt legalább egy tranzakció tartalmazza [101]. A párok támogatottságának meghatározása kevesebb memóriát fog igényelni, de ezzel együtt lassabb is lesz.

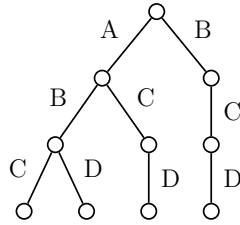
### Nagyobb elemhalmazok támogatottsága

Vizsgáljuk meg részletesebben az 5. sort. Adott egy tranzakció és  $\ell$ -méretű jelöltek egy halmaza. Feladatunk meghatározni azon jelölteket, amelyek a tranzakció részhalmazai. Megoldhatjuk ezt egyszerűen úgy, hogy a jelölteket egyesével vesszük, és eldöntjük, hogy tartalmazza-e őket a tranzakció. Rendezett halmazban rendezett részhalmaz keresése elemi feladat.

Vegyünk fel két mutatót, amelyek a kosár, illetve a jelölt elemein fognak végighaladni. Kezdetben mutasson mindkét mutató az elemhalmazok első elemeire. Amennyiben a két mutató által mutatott elemek megegyeznek, akkor léptessük mindkét mutatót a következő elemre. Ha a tranzakcióban található elem kisebb sorszámú, akkor csak a kosár mutatóját léptessük, ellenkező esetben pedig álljunk meg; ekkor a kosár biztosan nem tartalmazza a jelöltet. Ha a jelölt utolsó eleme is megegyezik a kosár valamelyik elemével, akkor a kosár tartalmazza a jelöltet.



Ennek az egyszerű módszernek a hátránya, hogy sok jelölt esetén lassú, hiszen annyiszor kell a tranzakció elemein végighaladni, amennyi a jelöltek száma. A gyorsabb működés érdekében a jelölteket szófában vagy hash-fában (hash-tree) célszerű tárolni. A szófát szokás prefix-fának vagy lexikografikus fának is hívni [2]. Az eredeti APRIORI implementációban hash-fát alkalmaztak, azonban tesztek bizonyítják, hogy a szófa gyorsabb működést eredményez, mint a hash-fa. A hash-fa szófával való helyettesítéséről már a [96]-ban írtak, ahol a szófát alkalmazó APRIORI algoritmust SEAR-nek nevezték el. A továbbiakban a szófában való keresést ismertetjük (a szófák felépítéséről és típusairól az alapfogalmak 2.8.1 részében már szóltunk).

4.5. ábra. Az  $ABC$ ,  $ABD$ ,  $ACD$ ,  $BCD$  jelölteket tároló szófa.

A szófa éleinek címkei elemek lesznek. Minden csúcs egy elemhalmazt reprezentál, amelynek elemei a gyökérből a csúcsig vezető út éleinek címkeivel egyeznek meg. Feltehetjük, hogy az egy csúcsból induló élek, továbbá az egy úton található élek címkeként rendezve vannak (pl. legnagyobb elem az első helyen). A jelöltek számlálóját a jelölteket reprezentáló levélhez rendeljük. A 4.5. ábrán egy szófat láthatunk.

A  $t$  tranzakcióban az  $\ell$ -elemű jelölteket úgy találjuk meg, hogy a jelölteket leíró fa gyökérből kiindulva, rekurzív módon bejárunk bizonyos részfákat. Ha egy  $d$  szintű belső csúcsba a tranzakció  $j$ -edik elemén keresztül jutunk, akkor azon élein keresztül lépünk egyvel mélyebb szintre, amelyeknek címkeje megegyezik a tranzakció  $j'$ -edik elemével, ahol  $j < j' \leq |t| - \ell + d$  (ugyanis  $\ell - d$  elemre még szükség van ahhoz, hogy levélbe érjünk). Ha ily módon eljutunk egy  $\ell$  szintű csúcsba, az azt jelenti, hogy a csúcs által reprezentált elemhalmazt tartalmazza  $t$ , így ennek a levélnek a számlálóját kell növelnünk egyvel.

A szófat prefix fának is szokták hívni, ami arra utal, hogy a közös prefixeket csak egyszer tárolja. Ettől lesz gyorsabb a szófas támogatottság-meghatározás a naiv módszernél. A közös prefixeket összevonjuk, és csak egyszer foglalkozunk velük.

A szófa nagy előnye a gyors támogatottság-meghatározás mellett, hogy a jelölt-előállítást is támogatja. Tudjuk, hogy két gyakori elemhalmaz akkor lesz generátor, ha a legnagyobb sorszámú elemük elhagyásával ugyanazt az elemhalmazt kapjuk, vagy más szavakkal, a két gyakori elemhalmaz  $\ell - 1$  hosszú prefixei megegyeznek. A támogatottság-meghatározásában használt szófat felhasználhatjuk a következő iterációs lépés jelöltjeinek az előállítására, hiszen a szófa tárolja a jelölt-előállításhoz szükséges gyakori elemhalmazokat.

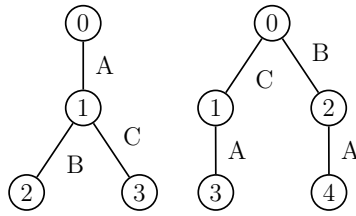
Az egész algoritmus alatt tehát egyetlen szófat tartunk karban, amely az algoritmus kezdetekor csak egy csúcsból áll (ez reprezentálja az üres halmazt). A támogatottság-meghatározás után töröljük azon leveleket, amelyek számlálójuk kisebb  $min\_supp$ -nál. Az iterációs lépés végére kialakuló szófa alapján előállítjuk a jelölteket, amely során a szófa új, egyvel mélyebb szinten lévő levelekkel bővül. A jelölt-előállítás során arra is lehetőségünk van, hogy az előző iterációban gyakran talált elemhalmazokat és azok számlálóját kiírjuk (a kimenetre vagy a háttértárolóra).

## A rendezés hatása az Apriori algoritmusra

Amennyiben a szófa hatékony adatstruktúra sorozatok tárolására, és gyors visszakeresésére, akkor ugyanez mondható el elemhalmazok esetére is. Ha tehát elemhalmazok adóttak és az a feladat, hogy gyorsan megállapítsuk, hogy egy elemhalmaz szerepel-e a megadottak között, akkor elég definiálnunk az elemeken egy teljes rendezést, ami alapján a halmazokat sorozatokká alakíthatjuk.

Különböző rendezések különböző sorozatokat állítanak elő, amelyek különböző szófatokat eredményeznek. Erre mutat példát a következő ábra, ahol két olyan szófat láthatunk, amelyek

a  $AB, AC$  elemhalmazokat tárolják. Az első szófa az ABC szerint csökkenő sorrendet használ ( $C \prec B \prec A$ ), míg a második ennek ellenkezőjét.



4.6. ábra. Példa: különböző rendezést használó szófák

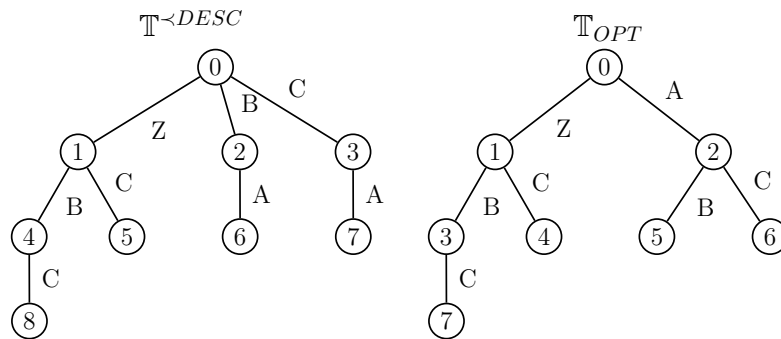
Egy szófa memóriaigénye arányos a szófa pontjainak számával, így jogos az az igény, hogy azt a teljes rendezést válasszuk, amely a legkevesebb pontú, azaz minimális méretű szófát adja. Ez az ún. minimális szófa előállításának feladata. Sajnos ez egy nehéz feladat.

**4.2. tétel.** *A minimális szófa probléma NP-nehéz.*

Eredetileg a feladatot  $n$ -esekre bizonyították, de ebből következik, hogy halmazokra is érvényes. Legyen ugyanis az alaphalmaz  $J$ . Ekkor minden halmazt felfoghatunk, mint egy  $|J|$  hosszú bináris értékeket tartalmazó vektort.

A fenti példát szemlélve az embernek az az érzése támad, hogy az a rendezés adja a legkevesebb csúcsú szófát, amelyben az elemek a halmazokban való előfordulások számának arányában csökkenő sorba vannak rendezve. Ugyanis a gyakori elemek fognak a halmazok kapott sorozatok elejére kerülni, és ezek az elemek, mivel gyakoriak sok sorozat elején lesznek megtalálhatók. A szófa a közös prefixeket csak egyszer tárolja, így akkor lesz a szófa mérete várhatóan a legkisebb, ha minél több sorozatnak van közös prefixe. Az előző ábra is ezt sugallta.

Sajnos a fenti módon kapott szófa nem feltétlenül adja a legkevesebb pontot tartalmazó szófát. Ezt a legegyszerűbben egy ellenpéldával tudjuk bizonyítani. Legyenek a halmazaink a következők:  $AB, AC, CZ, BCZ, BZ, Z$ . A  $Z$  elem gyakorisága 4, a  $B, C$ -é 3 és az  $A$  elemé 2. Ha felrajzoljuk ezen gyakoriságok alapján kapott rendezés ( $Z > B > C > A$ , de a  $C, B$  elemek sorrende tetszőleges lehet) szerinti szófát, akkor a bal oldali szófát kapjuk. Ha az  $A$  és  $B, C$  elemek sorrendjét felcseréljük, akkor eggyel kevesebb pontot tartalmazó szófát kapunk (jobb oldal).



4.7. ábra. Ellenpélda arra, hogy az előfordulás szerinti csökkenő sorrend adja a minimális méretű szófát

Tapasztalatok alapján gyakoriság szerint csökkenő rendezés kisebb szófát eredményez, mint a gyakoriság szerint növekvő rendezés, vagy más véletlenszerűen megválasztott rendezések.

Ennek ellenére olyan szófat célszerű alkalmazni, amelyben az elemeken értelmezett rendezés a gyakoriság szerint növekvő sorrendnek felel meg. Ennek ugyanis két előnye van. Egyrészt a szófa pontjai kisebbek lesznek (kevesebb él indul ki belőlük), de ami még fontosabb, hogy a ritka elemek lesznek közel a gyökérhez. A ritka elemekkel kevesebb kosárbeli elem fog egyezni, ezáltal a szófa kisebb részét járjuk be a támogatott jelöltek meghatározása során.

A továbbiakban bemutatunk néhány gyorsítási ötletet, amelynek segítségével nagymértékben lecsökkenthető a szófa alapú APRIORI algoritmus futási ideje és memóriaigénye.

### 4.2.3. Ritka jelöltek törlése

Ritka jelöltek törléséhez és a gyakori jelöltek kiírásához be kell járnunk a szófat és amikor levélbe érünk, akkor össze kell hasonlítani a levél számlálóját a támogatottsági küszöbvel. Ha a számláló nagyobb, akkor az eredményfájlba írjuk a levél által reprezentált halmazt és a számlálót. Ellenkező esetben töröljük a levelet.

A bejárást megtakaríthatjuk, ha a fenti műveletet a jelöltek előállításával együtt tesszük meg. A jelöltek előállításánál is be kell járni a szófat. A testvér levelek közül töröljük a ritka jelölteket, majd a megmaradtakból generáljunk új, eggyel nagyobb méretű jelölteket.

### 4.2.4. Zsácutca nyesés

Szükségtelen tárolni azon csúcsoakat, amelyekből az összes elérhető levelet töröltük. Ezek ugyanis lassítják a támogatottságok meghatározását (miközben szerepet nem játszanak benne) és feleslegesen foglalják a memóriát.

Nem mindegy azonban, hogy mikor távolítjuk el a zsácutcákat. Ha például a  $AB$ ,  $AC$ ,  $BC$  két-elemű jelölt lett gyakori, akkor a  $BC$  levélből (de az  $AC$ -ből sem) nem fogunk új levelet felvenni, azaz a  $BC$  levél zsácutca lesz. Ezt a levelet azonban nem törölhetjük az  $ABC$  új jelölt felvétele előtt, hiszen a  $BC$  halmaz az  $ABC$ -nek valódi részhalmaza, így szükséges, hogy szerepeljen a fában.

Könnyű belátni, hogy tetszőleges  $I$  halmaz nem generátor, eggyel kisebb méretű, valódi részhalmaza lexikografikus rendezés szerint  $I$  után következik. Ezért, ha preorder bejárást használunk a jelöltek előállítása során, akkor egy levelet azonnal törölhetünk, ha belőle nem tudtunk új levelet felvenni. Garantált, hogy egyetlen részhalmazt sem töröltünk még, hiszen a valódi, nem generátor részhalmozokat csak később fogjuk meglátogatni a preorder bejárás szerint.

### 4.2.5. A bemenet tárolása

Amikor megvizsgálunk egy kosarat annak érdekében, hogy eldöntsük, mely jelölteket tartalmazza, akkor az operációs rendszer a háttértárolóból bemásolja a tranzakciót a memóriába. Ha van elég hely a memóriában, akkor a tranzakció ott is marad, és amikor ismét szükség van rá, nem kell lassú IO műveletet végeznünk. A bemenetet tehát szükségtelen explicit eltárolnunk a memóriában, hiszen az operációs rendszer ezt megteszi helyettünk. Sőt, ha a program eltárolja a bemeneti adatot (például egy listában), akkor a valóságban duplán lesz eltárolva.

A bemenet tárolásának vannak előnyei is. Például összegyűjthetjük az azonos tranzakciókat és ahelyett, hogy többször hajtánánk végre ugyanazon a tranzakción a támogatott jelöltek meghatározását, ezt egyszer tesszük meg. Szükségtelen az eredeti tranzakciókat tárolni. Az

első végigolvasás után rendelkezésre állnak a gyakori elemek. A ritka elemek úgysem játszanak szerepet, ezért elég a tranzakcióknak csak a gyakori elemeit tárolni. Ennek további előnye, hogy sokkal több azonos „szűrt” tranzakció lehet, ezáltal tovább csökken a támogatott jelölteket kereső eljárás meghívásának száma. Ráadásul az  $\ell$ -edik végigolvasás során törölhetjük azokat a szűrt tranzakciókat, amelyek nem tartalmazzak egyetlen  $\ell$ -elemű jelöltet sem.

A szűrt tranzakciókat célszerű olyan adatstruktúrában tárolni, amit gyorsan fel lehet építeni (azaz gyorsan tudjuk beszúrni a szűrt tranzakciókat) és gyorsan végig tudunk menni a beszűrt elemeken. Alkalmazhatunk erre a célra egy szófát, de tesztek azt mutatják, hogy egy piros-fekete fa (kiegyensúlyozott bináris fa), amelynek csúcsaiban egy-egy szűrt tranzakció található, még jobb megoldás, mert jóval kisebb a memóriaigénye.

#### 4.2.6. Tranzakciók szűrése

A feldolgozás során a tranzakciókat módosíthatjuk/törölhetjük annak érdekében, hogy az APRIORI még hatékonyabb legyen. A tranzakció szűrése alatt a tranzakció olyan elemeinek törlését értjük, amelyek nem játszanak szerepet az algoritmus kimenetének előállításában. A nem fontos elemek lassítják az algoritmust, gondoljunk itt a támogatottság meghatározásának módjára. A szófa egy belső csomópontjánál meg kell határoznunk a közös elemeket az élek címkéinek és a tranzakció elemeinek halmazában. Minél több elem van a tranzakcióban, annál tovább tart ez a művelet.

Szűrésnek tekinthetjük az első iteráció után végrehajtott lépést:

**1. szűrő ötlet.** *Minden tranzakcióból töröljük a ritka elemeket.*

Egyszerű szűrő ötletek a következők:

**2. szűrő ötlet.** *Az  $\ell$ -edik iterációban a  $t$  tranzakció feldolgozása után töröljük a  $t$ -t, amennyiben a  $t$  elemeinek száma nem nagyobb, mint  $\ell$ . Nyilvánvaló, hogy ez a tranzakció nem tartalmaz olyan elemhalmazt, amely a későbbi iterációban lesz jelölt.*

**3. szűrő ötlet.** *Töröljük a tranzakciót, amennyiben az nem tartalmaz jelöltet.*

Ennek az ötletnek a javított változata:

**4. szűrő ötlet.** *Töröljük a tranzakció azon elemeit, amelyek nem elemei egyetlen olyan jelöltnek sem, amelyet tartalmaz a tranzakció.*

Amennyiben az így keletkezett tranzakció mérete  $\ell$ , akkor töröljük teljesen a tranzakciót. Például, ha a háromelemű jelöltek halmaza  $\{ABC, ABD, BCD, FGH\}$  és  $t = ABCDH$ , akkor a  $H$  elemet törölhetjük a tranzakcióból.  $t' = ABCGH$  esetében a teljes tranzakciót töröljük.

Az előző szűrőötletet tovább szigoríthatjuk. Mi kell ahhoz, hogy egy elem eleme legyen majd egy olyan  $\ell + 1$ -elemű  $j$  jelöltnek a következő iterációban, amelyet tartalmaz az aktuális jelölt. Szükséges feltétel, hogy a  $j$  minden  $\ell$ -elemű részhalmazát tartalmazza a tranzakció. A  $j$  egy eleme pontosan  $\ell$  darab részhalmaznak az eleme. Ez alapján:

*A vérnyomás és a nemzetek boldogsága közötti összefüggésekre mutatott rá egy amerikai kutató: „Az amerikai kutatók szerint az eredmények egészen egyszerűek, a boldog országokból származó emberek - svédok, dánok, britek és hollandok - kevesebbet szenvednek a magas vérnyomástól, mint a németek, vagy a portugálok, akik az európai boldogságskála végén találhatók.”*  
 Forrás: <http://www.karpatinfo.net/article38511.html>

**5. szűrő ötlet.** *Töröljük a tranzakció azon elemeit, amelyek nem elemei  $\ell$  darab olyan jelöltnek, amelyet tartalmaz a tranzakció.*

Természetesen most is igaz, hogy ez után a szűrés után alkalmazzuk a második szűrő ötletet, ha ez lehetséges. A fenti példában a  $t'' = ABCFGH$  tranzakciót ez a szűrés teljes egészében törli.

#### 4.2.7. Equisupport nyesés

Az egyenlő támogatottságú elemhalmazok alapján történő, ún. equisupport nyesés talán a legelterjedtebb trükk a gyakori elemhalmazok kinyerésének meggyorsítására. A nyesés a 4.3 tulajdonság egy következményét használja ki. A támogatottság meghatározásánál kihagyhatjuk azokat a halmazokat, amelyeknek van olyan  $\ell$ -elemű valódi részhalmazuk, amelyek támogatottsága egyenlő egy  $(\ell-1)$ -elemű részhalmazukéval.

**4.3. tulajdonság.** *Legyen  $X \subset Y \subseteq \mathcal{J}$ . Ha  $\text{supp}(X) = \text{supp}(Y)$ , akkor  $\text{supp}(X \cup Z) = \text{supp}(Y \cup Z)$  teljesül minden  $Z \subseteq \mathcal{J}$ -re.*

Ez az állítás minden  $Z \subseteq \mathcal{J}$  elemhalmazra igaz, de nekünk elég lesz csak a  $Z \subseteq \mathcal{J} \setminus Y$  halmazokra koncentrálnunk.

Az equisupport nyesés és a zárt elemhalmazok közötti összefüggés egyértelmű. Az  $X$  elemhalmaz nem zárt, és lezártja  $Y$ , amennyiben  $X \subset Y$ ,  $\text{supp}(X) = \text{supp}(Y)$ , továbbá nem létezik olyan elemhalmaz, amelynek  $Y$  valódi részhalmaza, és támogatottsága megegyezik  $Y$  támogatottságával. Egy  $X$  elemhalmaz akkor, és csak akkor lehet egy egzakt (100% bizonyosságú) asszociációs szabály feltétel része, ha  $X$  nem zárt elemhalmaz. Az  $X$  elemhalmaz *kulcs minta* [12], ha nincs vele egyenlő támogatottságú valódi részhalmaza.

Ha az  $Y$  jelöltnek a támogatottsága megegyezik az  $X$ -el jelölt prefixe támogatottságával, akkor felesleges az  $Y$ -t tartalmazó  $Y \cup Z$  halmazokat mint új jelölteket előállítani, a 4.3 tulajdonság alapján ezek támogatottsága  $X \cup Z$  részhalmazukból közvetlenül számítható [48].

Az alulról építkező algoritmusoknál (APRIORI, ECLAT, FP-GROWTH, stb.) a prefixek támogatottsága mindig elérhető, így a *prefix equisupport nyesést* (az  $X$  az  $Y$  prefixe és  $|X| + 1 = |Y|$ ) bármikor alkalmazhatjuk. A prefix equisupport nyesés a következőképpen működik: miután kiszámoltuk a  $P$  elemhalmaz gyerekeinek támogatottságát, a ritka elemek elhagyásakor ellenőrizzük, hogy a támogatottságuk egyenlő-e a szülő támogatottságával, azaz  $\text{supp}(P)$ -vel. Az ezt teljesítő elemeket nem kell figyelembe vennünk mint generátorokat a következő jelölt-előállítás során. Ezen jelölteket töröljük és az utolsó elemeiket egy halmazban tároljuk el, amit equisupport halmaznak hívunk és  $P$ -hez rendeljük. Vegyük észre, hogy az elemhalmazháló prefix bejárásnak köszönhetően a jelölt-előállítás során az  $X \setminus Y \prec z$  minden  $z \in Z$ , ahol  $\prec$  az elemhalmaz bejárásánál használt rendezés.

Amikor kiírjuk a  $GY$  gyakori elemhalmazt, vele együtt kiírjuk minden  $E' \subseteq E$  halmazokkal vett unióját is, ahol  $E$  a  $GY$  prefixeinek equisupport halmazainak uniója.

**4.4. példa.** *Legyenek a kételemű,  $A$  prefixű gyakori elemhalmazok a következők:  $\{AB, AC, AD\}$  és  $\text{supp}(A) = \text{supp}(AB) = \text{supp}(AC) = 4$  továbbá  $\text{supp}(AD) = 3$ . A többi  $A$  prefixű jelölt előállításához egyedül az  $AD$  elemhalmazt kell figyelembe vennünk. Azonban egy*



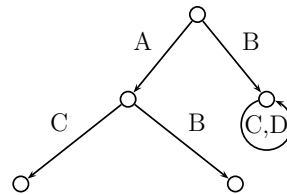
jelölt létrehozásához mind az APRIORI, az ECLAT- és az FP-GROWTH algoritmusnál legalább két elemhalmaz szükséges, így itt véget is ér az  $A$  prefixű halmazok feldolgozása. Az  $AD$  és  $A$  elemhalmazok kiírásakor  $BC$  minden részalmazát is hozzájuk kell venni, így végül az  $AD$ ,  $ABD$ ,  $ACD$ ,  $ABCD$ , valamint az  $A$ ,  $AB$ ,  $AC$ ,  $ABC$  halmazok kerülnek kiírásra; az előbbiek támogatottsága 3, utóbbiaké 4 lesz.

Ha az adatbázis csak zárt elemhalmazokat tartalmaz, akkor nem tudjuk ezt a nyesést alkalmazni, a támogatottságok egyenlőségének vizsgálata viszont lelassítja az algoritmust. A tapasztalat azonban azt mutatja, hogy az ellenőrzés gyors (például az APRIORI algoritmusnál nem kell újra bejárni a szófát), és nem okoz cache miss-t. A kevés nemzárt elemhalmazt tartalmazó adatbázisoknál elenyésző a futásiidő növekedése. Az equisupport nyesés ezért biztonságos gyorsítási trükknek tekinthető.

A fenti leírásban nem használtuk ki az APRIORI algoritmus sajátosságait, csak azt, hogy az algoritmus alulról építkező és az elemhalmaz bejárás során definiálva van egy rendezés és így a prefix is. A továbbiakban jobban a részletekbe mélyedünk és megnézzük, hogy mit kell tennünk az APRIORI algoritmusban, ha a prefix equisupport nyesést kívánjuk alkalmazni.

Az APRIORI algoritmus szófás megközelítése esetén minden csúcshoz egy listát kell hozzávennünk, mely az equisupport halmaz elemeit tartalmazza. A ritkának bizonyuló jelöltek eltávolításakor ellenőrizzük, hogy a levél támogatottsága megegyezik-e prefixének támogatottságával. Ha igen, a levelet törölhetjük a szófából, és az éle címkéjét hozzáírjuk a szülő equisupport halmazához. Minden  $i$  elem egy equisupport halmazban tekinthető egy  $i$  címkéjű hurokélnak. A hurokéleket nem kell figyelembe venni a támogatottság meghatározásakor, de a jelölt-előállításnál igen.

**4.5. példa.** Legyenek  $AB$ ,  $AC$ ,  $BC$ ,  $BD$  a gyakori párok.  $\text{supp}(AB) \neq \text{supp}(A) \neq \text{supp}(AC)$  és  $\text{supp}(B) = \text{supp}(BC) = \text{supp}(BD)$ . A 4.8 ábra a szófa ritka jelöltek eltávolítása utáni állapotát mutatja. Vegyük észre, hogy ha a hurokéleket figyelmen kívül hagytuk volna a jelöltgenerálás során, akkor az  $ABC$  elemhalmazt nem állítottuk volna elő mint jelölt, holott minden részalmazza gyakori.



4.8. ábra. Példa: equisupport levelek eltávolítása

Ez a példa az equisupport nyesés és a zsákutcanyesés közti összefüggésre is felhívja a figyelmet. Láttuk, hogy a  $B$  csomópont nem vezet 2 mélységű levélbe, így a zsákutcanyesés törölte volna ezt a csúcst, és nem lett volna jelölt az  $ABC$  elemhalmaz. Újra kell értelmeznünk a csomópontok mélységét a zsákutcanyesésnél azért, hogy ne töröljön olyan leveleket, amikre szükség lehet a jelölt-előállítás során. Az  $X$  elemhalmaz támogatottsága megegyezik az  $X$  olyan bővítésének támogatottságával, ahol a hozzáadott elem az  $X$  valamely prefixéhez tartozó equisupport halmaz egy eleme. Így amikor figyelembe vesszük az  $X$  csomópont mélységét a zsákutcanyesés során, hozzá kell adnunk  $X$  aktuális mélységéhez a gyökérből az  $X$ -be vezető

pontok equisupport halmazainak összméretét. Például a 4.8 ábrán látható szófán a  $B$  mélysége 1 helyett 3.

A szófa hatékony megvalósításának részleteit és további gyorsítási ötleteket a [17, 20, 41] írásokban találhatunk. Egy olyan programcsomag, amely szófa alapú APRIORI implementációt tartalmaz (továbbá hatékony ECLAT és Fp-growth implementációt) és kutatási célokra szabadon letölthető a

`http://www.cs.bme.hu/~bodon/en/fim_env`

oldalról.

#### 4.2.8. Borgelt-féle támogatottság-meghatározás

Ha a tranzakciókat szófában vagy Patrícia-fában tároljuk, akkor egy másik technikát is használhatunk a támogatottságok meghatározására [19, 20]. Ezt a módszert alkalmazza Christian Borgelt a világhírű APRIORI implementációja utolsó változataiban.

Az a megfigyelés áll az ötlet mögött, hogy két tranzakció a közös prefixig ugyanazt a programfutást eredményezi a támogatottság meghatározásakor (ugyanazt a szófarészt járjuk be). Ha szófában tároljuk a tranzakciókat, akkor rendelkezésre áll minden szükséges információ a közös prefixekről. Megoldható, hogy ugyanazokat a prefixeket csak egyszer dolgozzuk fel, és ne annyiszor, ahányszor előfordulnak.

A tranzakciófába minden csomóponthoz egy számlálót rendelünk. Az  $I$  elemhalmaz számlálója azoknak a tranzakcióknak a számát tárolja, amelyek prefixe  $I$ . Ebből a szempontból ez a megoldás eltér a bemenet tárolásánál bemutatott (lásd 4.2.5-es rész) szófa alapú megoldástól (és inkább egy olyan FP-fára hasonlít, amelyből elhagytuk a keresztéleket és a fejléc táblát, lásd 84 oldal). A tranzakció szófánál és a jelölt szófánál használt rendezésnek meg kell egyeznie. Ez hátrány, mivel az egyes szófákhoz más-más rendezés lenne optimális.

Sajnos a [19]-ben nincsen részletesen kidolgozva az algoritmus, de vélhetően a következőképp működik: Párhuzamosan bejárjuk a jelölt- és a tranzakció szófát duplán rekurzív módon. Két mutatót használunk, melyek kezdetben az egyes gyökerekre mutatnak. Ezután végigmegyünk mindkét csúcs élein. Ha a tranzakciószófa aktuális címkéje kisebb vagy egyenlő a másik címkénél, akkor rekurzívan továbblépünk a tranzakciószófában a gyerekcsomópontra (az szófa aktuális csomópontmutatója nem változik). Amennyiben a két címke egyenlő, a rekurziót azokkal a gyerekekkel folytatjuk, amelyekre a mutatók által mutatott élek mutatnak. A 76 oldalon található pszeudó-kód a Borgelt-féle támogatottság-meghatározás egy tovább optimalizált változatát adja meg.

A fenti megoldásnak hátránya, hogy sok olyan utat jár be a jelölt szófában, amelyet az eredeti támogatottság meghatározó módszer nem tenne, mert nem vezet levélbe. A módszer nem veszi figyelembe, hogy a tranzakciónak csak egy részét kell kiértékelnünk. Megoldhatjuk a problémát, ha hozzárendelünk egy számlálót a tranzakció szófa minden pontjához. A számláló adja meg a pontból kiinduló leghosszabb út hosszát. A támogatottság meghatározása során nem vesszük figyelembe azokat a csomópontokat, melyek számlálója kisebb, mint  $\ell - 1$ , ahol  $\ell$  azon lépések számát adja, amelyeket meg kell még tenni a jelölt szófa aktuális pontjából, hogy levélbe jussunk. Az algoritmus gyorsítható, ha a tranzakció szűrésének ötletét (lásd 4.2.6-ös rész) is alkalmazzuk. További részletek tudhatunk meg a [19] tanulmányból.

**Algorithm 2** BORGELT\_SUPPCOUNT**Require:**  $n_c$ : a szófa aktuális csomópontja, $n_t$ : a tranzakciófa aktuális csomópontja, $\ell$ : az  $n_c$ -ből levélbe vezető út hossza, $i$ : az  $n_c$  legkisebb olyan élének indexe, amely címkéje nagyobb, mint az  $n_t$ -be vezető él címkéje**if**  $\ell = 0$  **then** $n_c$ .számláló  $\leftarrow n_c$ .számláló +  $n_t$ .számláló**else****for**  $j = 0$  to  $n_t$ .élszám  $- 1$  **do****while**  $i < n_c$ .élszám AND  $n_c$ .él[ $i$ ].címké  $< n_t$ .él[ $j$ ].címké **do** $i \leftarrow i + 1$ **end while****if**  $i < n_c$ .élszám AND  $n_c$ .él[ $i$ ].címké  $\geq n_t$ .él[ $j$ ].címké **then**BORGELT\_SUPPCOUNT( $n_c$ ,  $n_t$ .él[ $j$ ].gyermek,  $\ell$ ,  $i$ )**if**  $n_c$ .él[ $i$ ].címké =  $n_t$ .él[ $j$ ].címké **then**BORGELT\_SUPPCOUNT( $n_c$ .él[ $i$ ].gyermek,  $n_t$ .él[ $j$ ].gyermek,  $\ell - 1$ , 0) $i \leftarrow i + 1$ **end if****else**

break

**end if****end for****end if****4.2.9. Futási idő és memóriaigény**

A GYEK feladat megadásakor elmondtuk, hogy már az eredmény kiírása – ami a futási időnek a része – az  $|I|$ -ben exponenciális lehet. A memóriaigényről is hasonló mondható el. Az  $(\ell + 1)$ -elemű jelöltek előállításához szükségünk van az összes  $\ell$ -elemű jelöltre, amelyek száma akár  $\binom{|I|}{\lfloor |I|/2 \rfloor}$  is lehet. Ezek a felső korlátok élesek is, hiszen  $\min\_supp = 0$ -nál minden elemhalmaz gyakori.

Az algoritmus indítása előtt tehát nem sokat tudunk mondani a futási időről. A futás során, azonban egyre több információt gyűjtünk, így felmerül a kérdés, hogy ezt fel tudjuk-e használni az algoritmus maradék futási idejének jóslására. Például, ha a gyakori elemek száma négy, akkor tudjuk, hogy a legnagyobb gyakori elemhalmaz mérete legfeljebb négy (azaz még legfeljebb háromszor olvassuk végig az adatbázist), az összes jelölt maximális száma pedig  $\binom{4}{2} + \binom{4}{3} + \binom{4}{4} = 11$ . A következőkben megvizsgáljuk, hogy mit tudunk elmondani a jelöltek számáról és a maximális jelöltek méretéről, ha adottak az  $\ell$ -elemű gyakori elemhalmazok ( $GY_\ell$ ).

A következő rész fontos fogalma a kanonikus reprezentáció lesz.

**4.6. lemma.** *Adott  $n$  és  $\ell$  pozitív egészek esetében a következő felírás egyértelmű:*

$$n = \binom{m_\ell}{\ell} + \binom{m_{\ell-1}}{\ell-1} + \cdots + \binom{m_r}{r},$$

ahol  $r \geq 1$ ,  $m_\ell > m_{\ell-1} > \cdots > m_r$  és  $m_j \geq j$  minden  $j = r, r+1, \dots, \ell$  számra.

Ezt a reprezentációt hívják  $\ell$ -kanonikus reprezentációnak. Meghatározása nagyon egyszerű:  $m_\ell$ -nek ki kell elégítenie a  $\binom{m_\ell}{\ell} \leq n < \binom{m_\ell+1}{\ell}$  feltételt,  $m_{\ell-1}$ -nek a  $\binom{m_{\ell-1}}{\ell-1} \leq n - \binom{m_\ell}{\ell} < \binom{m_{\ell-1}+1}{\ell-1}$  feltételt, és így tovább, amíg  $n - \binom{m_\ell}{\ell} - \binom{m_{\ell-1}}{\ell-1} - \dots - \binom{m_r}{r}$  nulla nem lesz.

Legyen  $\mathcal{I} = \{i_1, i_2, \dots, i_m\}$  elemek halmaza és  $GY_\ell$  egy olyan  $\mathcal{I}$  feletti halmazcsalád<sup>1</sup>, amelynek minden eleme  $\ell$ -elemű. Az  $\ell$ -nél nagyobb méretű  $I \subseteq \mathcal{I}$  halmaz *fedí* a  $GY_\ell$ -et, ha  $I$  minden  $\ell$ -elemű részhalmaza eleme  $GY_\ell$ -nek. Az összes lehetséges  $(\ell+p)$ -méretű  $GY_\ell$ -et fedő halmazokból alkotott halmazcsaládot  $J_{\ell+p}(GY_\ell)$ -lél jelöljük. Nem véletlen, hogy ezt a halmazt ugyanúgy jelöltük, mint az APRIORI algoritmus jelöltjeit, ugyanis az  $(\ell+p)$ -méretű jelöltek ezen halmazcsaládnak az elemei, és ha az algoritmus során minden jelölt gyakori, akkor az  $(\ell+p)$ -méretű jelöltek halmaza megegyezik  $J_{\ell+p}(GY_\ell)$ -lél.

A következő tétel megadja, hogy adott  $GY_\ell$  esetén legfeljebb mennyi lehet a  $J_{\ell+p}(GY_\ell)$  elemeinek száma.

**4.7. tétel.** *Ha*

$$|GY_\ell| = \binom{m_\ell}{\ell} + \binom{m_{\ell-1}}{\ell-1} + \dots + \binom{m_r}{r}$$

$\ell$ -kanonikus reprezentáció, akkor

$$|J_{\ell+p}(GY_\ell)| \leq \binom{m_\ell}{\ell+p} + \binom{m_{\ell-1}}{\ell-1+p} + \dots + \binom{m_s}{s+p},$$

ahol  $s$  a legkisebb olyan egész, amelyre  $m_s < s+p$ . Ha nincs ilyen egész szám, akkor a jobb oldal 0.

A fenti tétel a Kruskal–Katona tétel következménye, ezért a tételben szereplő felső korlátot a továbbiakban  $KK_\ell^{\ell+p}(|GY_\ell|)$ -el jelöljük.

**4.8. tétel.** *A 4.7. tételben szereplő felső korlát éles, azaz adott  $n, \ell, p$  számokhoz mindig létezik  $GY_\ell$ , amelyre  $|GY_\ell| = n$ , és  $|J_{\ell+p}(GY_\ell)| = KK_\ell^{\ell+p}(|GY_\ell|)$ .*

A kanonikus reprezentáció segítségével egyszerű éles felső becslést tudunk adni a legnagyobb jelölt méretére (jelölésben  $\maxsize(GY_\ell)$ ) is. Tudjuk, hogy  $|GY_\ell| < \binom{m_\ell+1}{\ell}$ , ami azt jelenti, hogy nem létezhet olyan jelölt, amelynek mérete nagyobb  $m_\ell$ -nél.

**4.9. következmény.** *Amennyiben a  $|GY_\ell|$  számnak az  $\ell$ -kanonikus reprezentációjában szereplő első tag  $\binom{m_\ell}{\ell}$ , akkor  $\maxsize(GY_\ell) \leq m_\ell$ .*

Az  $m_\ell$  számot a továbbiakban  $\mu_\ell(|GY_\ell|)$ -el jelöljük. Ez az érték azt is megmondja, hogy mekkora jelölméretnél válik nullává a felső korlát, azaz:

**4.10. következmény.**  $\mu_\ell(|GY_\ell|) = \ell + \min\{p \mid KK_\ell^{\ell+p}(|GY_\ell|) = 0\} - 1$

A maradék futási idő jóslására a következő állítás nyújt segítséget.

**4.11. következmény.** *Az összes lehetséges  $\ell$ -nél nagyobb méretű jelölt száma legfeljebb*

$$KK_\ell^{\text{összes}}(|GY_\ell|) = \sum_{p=1}^{\mu_\ell(|GY_\ell|)} KK_\ell^{\ell+p}(|GY_\ell|).$$

<sup>1</sup>A  $H$ -t az  $\mathcal{I}$  feletti halmazcsaládnak nevezzük, amennyiben  $H \subseteq 2^{\mathcal{I}}$ .

A fenti korlátok szépek és egyszerűek, mivel csak két paramétert használnak: az  $\ell$  aktuális méretet és az  $\ell$ -elemű gyakori elemhalmazok számát ( $|GY_\ell|$ ). Ennél jóval többet tudunk. Nem csak a gyakori elemhalmazok számát ismerjük, hanem már pontosan meghatároztuk őket magukat is! Az új információ segítségével számos esetben jobb felső korlátot adhatunk. Például, ha a  $GY_\ell$ -ben csak páronként diszjunkt elemhalmazok vannak, akkor nem állítunk elő jelölteket. A 4.7. tételben szereplő felső korlát azonban jóval nagyobb lehet nullánál. A következőkben bemutatjuk, hogyan lehet a meglévő felső korlátot az  $\ell$  méretű gyakori elemhalmazok *struktúrájára* rekurzívan alkalmazni. Ehhez feltesszük, hogy egy teljes rendezést tudunk definiálni az  $\mathcal{I}$  elemein, ami alapján tetszőleges elemhalmaznak meg tudjuk határozni a legkisebb elemét. Vezessük be a következő két jelölést:

$$GY_\ell^i = \{I - \{i\} \mid I \in GY_\ell, i = \min I\},$$

A  $GY_\ell^i$  halmazt úgy kapjuk  $GY_\ell$ -ből, hogy vesszük azon halmazokat, amelyek legkisebb eleme  $i$ , majd töröljük ezekből az  $i$  elemet.

Ezek után definiálhatjuk a következő rekurzív függvényt tetszőleges  $p > 0$ -ra:

$$KK_{\ell,p}^*(GY_\ell) = \begin{cases} \binom{|GY_\ell|}{p+1} & , \text{ ha } \ell = 1 \\ \min\{KK_{\ell}^{\ell+p}(|GY_\ell|), \sum_{i \in \mathcal{I}} KK_{\ell-1,p}^*(GY_\ell^i)\} & , \text{ ha } \ell > 1. \end{cases}$$

A definícióból következi, hogy  $KK_{\ell,p}^*(GY_\ell) \leq KK_{\ell}^{\ell+p}(|GY_\ell|)$ , továbbá

**4.12. tétel.**  $|J_{\ell+p}(GY_\ell)| \leq KK_{\ell,p}^*(GY_\ell)$ .

*Bizonyítás:* A bizonyítás teljes indukción alapul, az  $\ell = 1$  eset triviális. Tulajdonképpen csak azt kell belátni, hogy

$$|J_{\ell+p}(GY_\ell)| \leq \sum_{i \in \mathcal{I}} KK_{\ell-1,p}^*(GY_\ell^i)$$

Az egyszerűség kedvéért vezessük be a következő jelölést:  $H \cup i = \{I \cup \{i\} \mid I \in H\}$ , ahol  $H$  egy  $\mathcal{I}$  feletti halmazcsalád. Vegyük észre, hogy  $GY_\ell = \sum_{i \in \mathcal{I}} GY_\ell^i \cup i$  és  $GY_\ell^i \cap GY_\ell^j = \emptyset$  minden  $i \neq j$  elempárra. Azaz a  $GY_\ell$  halmazcsalád egy partícióját képeztük.

Amennyiben  $I \in J_{\ell+p}(GY_\ell)$ , és  $I$ -nek legkisebb eleme  $i$ , akkor  $I \setminus \{i\} \in J_{\ell-1+p}(GY_\ell^i)$ , hiszen  $I \setminus \{i\}$  minden  $(\ell - 1)$ -elemű részhalmaza  $GY_\ell^i$ -beli. Ebből következik, hogy

$$J_{\ell+p}(GY_\ell) \subseteq \bigcup_{i \in \mathcal{I}} J_{\ell-1+p}(GY_\ell^i) \cup i.$$

Abból, hogy az  $GY_\ell^i$  halmazcsaládok páronként diszjunktak következik, hogy  $J_{\ell-1+p}(GY_\ell^i) \cup i$

„Angol kutatók állítják, apukáddal való kapcsolatod befolyásolja, milyen férfiakat találsz vonzónak. Szerintük az egészséges apa-lánya kapcsolatot ápoló lányok inkább az úgynevezett alfa-hím típusú férfiakhoz vonzódnak: a veszélyes kinézetű, széles állú, dús szemöldökű típusokhoz, míg azok a lányok, akiknek kevésbé pozitív a kapcsolatuk a családfenntartóval, azok inkább a finom vonású, már szinte nőies kinézetű férfiakat részesítik előnyben.” Forrás: <http://shape.proweb.hu/main.php?rovat=6&cikk=507>

is páronként diszjunkt halmazcsaládok. Ebből következik az állítás, hiszen:

$$\begin{aligned}
|J_{\ell+p}(GY_\ell)| &\leq \left| \bigcup_{i \in \mathcal{I}} J_{\ell-1+p}(GY_\ell^i) \cup i \right| \\
&= \sum_{i \in \mathcal{I}} |J_{\ell-1+p}(GY_\ell^i) \cup i| \\
&= \sum_{i \in \mathcal{I}} |J_{\ell-1+p}(GY_\ell^i)| \\
&\leq \sum_{i \in \mathcal{I}} KK_{\ell-1,p}^*(GY_\ell^i),
\end{aligned}$$

ahol az utolsó egyenlőtlenségél az indukciós feltevést használtuk. ■

A páronként diszjunkt halmazok esete jó példa arra, hogy a minimum kifejezésben szereplő második tag kisebb lehet az elsőnél. Előfordulhat azonban az ellenkező eset is. Például legyen  $GY_2 = \{AB, AC\}$ . Könnyű ellenőrizni, hogy  $KK_2^3(|GY_2|) = 0$ , ugyanakkor a második tagban szereplő összeg 1-et ad. Nem tudhatjuk, hogy melyik érték a kisebb, így jogos a két érték minimumát venni.

Javíthatjuk a legnagyobb jelölt méretére, illetve az összes jelölt számára vonatkozó felső korlátokon is. Legyen  $\mu_\ell^*(GY_\ell) = \ell + \min\{p | KK_{\ell+p}^*(GY_\ell) = 0\} - 1$  és

$$KK_{\text{összes}}^*(GY_\ell) = \sum_{p=1}^{\mu_\ell^*(GY_\ell)} KK_{\ell+p}^*(GY_\ell).$$

**4.13. következmény.**  $\text{maxsize}(GY_\ell) \leq \mu_\ell^*(GY_\ell) \leq \mu_\ell(|GY_\ell|)$ .

**4.14. következmény.** *Az összes lehetséges  $\ell$ -nél nagyobb méretű jelölt száma legfeljebb  $KK_{\text{összes}}^*(GY_\ell)$  lehet, és  $KK_{\text{összes}}^*(GY_\ell) \leq KK_\ell^{\text{összes}}(|GY_\ell|)$ .*

A  $KK^*$  érték függ a rendezéstől. Például a  $KK_{2,1}^*(\{AB, AC\})$  értéke 1, amennyiben a rendezés szerinti legkisebb elem  $A$ , és 0 bármely más esetben. Elméletileg meghatározhatjuk az összes rendezés szerinti felső korlátot, és kiválaszthatjuk azt, amelyik a legkisebb értéket adja. Ez a megoldás azonban túl sok időbe telne. A szófa által használt rendezés szerinti felső korlátot viszonylag könnyen meghatározhatjuk. Ehhez azt kell látnunk, hogy a gyökér  $i$  címkéjű éléhez tartozó részfa levelei reprezentálják a  $GY_\ell^i$  elemeket. A szófa egyetlen bejárásával egy egyszerű rekurzív módszer segítségével minden csúcshoz kiszámíthatjuk a  $KK_{\ell-d,p}^*(GY_{\ell-d}^I)$  és  $KK_{\ell-d}^{\ell-d+p}(|GY_{\ell-d}^I|)$  értékeket, ahol  $d$  a csúcs mélységét jelöli,  $GY_{\ell-d}^I$  pedig az adott csúcshoz tartozó részfa által reprezentált elemhalmazokat. A gyökérhez kiszámított két érték adja meg a  $KK$  és  $KK^*$  korlátokat.

Ha a maradék futási idő becslésére kívánjuk használni a fenti felső korlátot, akkor tudnunk kell, hogy a jelöltek támogatottságának meghatározása függ az APRIORI algoritmusban felhasznált adatstruktúrától. Szófa esetében például egy jelölt előfordulásának meghatározásához el kell jutnunk a jelöltet reprezentáló levélhez, ami a jelölt méretével arányos lépésszámú műveletet igényel. A maradék futási idő pontosabb felső becsléséhez a  $KK_{\ell+p}^*(GY_\ell)$  értékeket súlyozni kell  $(\ell+p)$ -vel.

### 4.3. Az Eclat algoritmus

Az ECLAT az üres mintából indulva egy rekurzív, mélységi jellegű bejárást valósít meg. A rekurzió mélysége legfeljebb eggyel több, mint a legnagyobb gyakori elemhalmaz mérete. Az APRIORI-val szemben mindig egyetlen jelöltet állít elő, majd ennek azonnal meghatározza a támogatottságát. Az  $(\ell + 1)$ -elemű,  $P$  prefixű jelöltek, ahol  $|P| = \ell - 1$  az  $\ell$ -elemű,  $P$  prefixű gyakori elemhalmazokból állítja elő egyszerű páronkénti unióképzéssel.

Az algoritmus központi fogalma az ún. TID-halmaz. Egy elemhalmaz *TID-halmazának* (Transaction IDentifier) elemei azon bemeneti sorozatok azonosítói (sorszámjai), amelyek tartalmazzák az adott elemhalmazt. Más szóval egy TID-halmaz a vertikális adatbázis egy megfelelő sora. Például  $\langle AD, AC, ABCD, B, AD, ABD, D \rangle$  bemenet esetén az  $\{A, C\}$  elemhalmaz TID-halmaz  $\{1, 2\}$ , amennyiben egy tranzakció azonosítója megegyezik a bemeneti sorozatban elfoglalt helyével, és a helyek számozását nullától kezdjük.

A TID-halmaz két fontos tulajdonsággal bír:

- I. Az  $I$  elemhalmaz TID-halmazának mérete megadja az  $I$  támogatottságát.
- II. Egy jelölt TID-halmazát megkaphatjuk a generátorainak TID-halmazából egy egyszerű metszetképzéssel.

Az ECLAT pszeudokódja a 80 oldalol található.

---

#### Algorithm 3 ECLAT

---

**Require:**  $\mathcal{T}$ : tranzakciók sorozata,  
 $min\_supp$ : támogatottsági küszöb,  
 támogatottság\_meghatározás( $\mathcal{T}, J_1$ );  
 $GY_1 \leftarrow$  gyakoriak\_kiválogatása( $J_1, min\_supp$ );  
**for**  $i \leftarrow 1$  to  $|\mathcal{T}|$  **do**  
   **for all**  $j \in t_i \cap GY_1$  **do**  
      $j.TID \leftarrow j.TID \cup \{i\}$   
   **end for**  
**end for**  
 return  $GY_1 \cup$  ECLAT-SEGÉD( $\emptyset, GY_1, min\_supp$ )

---

Először meghatározzuk a gyakori elemeket, majd felépítjük a gyakori elemek TID-halmazait. A későbbiekben nem használjuk a bemenetet, csak a TID-halmazokat. Az algoritmus lényege a ECLAT-SEGÉD rekurziós eljárás. Jelöljük a  $P$  prefixű,  $P$ -nél eggyel nagyobb méretű gyakori elemhalmazokból alkotott halmazcsaládot  $GY^P$ -vel. Nyilvánvaló, hogy  $GY^0 = GY_1$ .

Az ECLAT jelölt-előállítás megegyezik az APRIORI jelölt-előállításával, azzal a különbséggel, hogy nem ellenőrizzük az unióképzéssel kapott halmaznak minden részalmazára, hogy gyakori-e (a mélységi bejárást miatt ez az információ nem is áll rendelkezésünkre). Látható, hogy az ECLAT abban is különbözik az APRIORI-tól, hogy egy jelölt előállítása után azonnal meghatározza a támogatottságát, mielőtt újabb jelöltet állítana elő. Nézzünk egy példát a keresési tér bejárására.

**4.15. példa.** Legyen  $\mathcal{T} = \langle ACDE, ACG, AFGM, DM \rangle$  és  $min\_supp = 2$ . Első lépésben meghatározzuk a gyakori elemeket:  $A, C, D, G, M$ , ami nem más, mint  $GY^0$ . Ezután előállítjuk és

**Algorithm 4** ECLAT-SEGÉD**Require:**  $P$ : prefix elemhalmaz. $GY^P$ :  $P$  prefixű,  $P$ -nél eggyel nagyobb méretű gyakori elemhalmazokból alkotott halmazcsalád, $min\_supp$ : támogatottsági küszöb,**for all**  $gy \in GY^P$  **do****for all**  $gy' \in GY^P$ ,  $gy \prec gy'$  **do** $j \leftarrow gy \cup gy'$  $j.TID \leftarrow gy.TID \cap gy'.TID$ **if**  $|j.TID| \geq min\_supp$  **then** $GY^{gy} \leftarrow GY^{gy} \cup \{j\}$ **end if****end for****if**  $|GY^{gy}| \geq 2$  **then** $GY \leftarrow GY \cup GY^{gy} \cup \text{ECLAT-SEGÉD}(gy, GY^{gy}, min\_supp)$ **else** $GY \leftarrow GY \cup GY^{gy}$ **end if****end for**return  $GY$ 

azonnal meg is határozzuk az  $(A, C)$ ,  $(A, D)$ ,  $(A, G)$ ,  $(A, M)$  párok unióját. Ezek közül csak az  $AC$ ,  $AG$  halmazok gyakoriak. A következő rekurziós lépésben ennek a két halmaznak vesszük az unióját, állítjuk elő a  $TID$ -halmazát, amely alapján kiderül, hogy az  $ACG$  ritka, és a rekurzió ezen ága véget ér. Ezután a  $C$  elemnek vesszük az unióját a sorban utána következő elemekkel egyesével és így tovább.

Látnunk kell, hogy az ECLAT legalább annyi jelöltet állít elő, mint az APRIORI. A mélységi bejárás miatt ugyanis egy jelölt előállításánál nem áll rendelkezésünkre az összes részhalmaz. Az előző példa esetében például az  $\{A, C, G\}$  támogatottságát hamarabb vizsgálja, mint a  $\{C, G\}$  halmazét, holott ez utóbbi akár ritka is lehet. Ebben a tekintetben tehát az ECLAT rosszabb az APRIORI-nál, ugyanis több lesz a ritka jelölt.

Az ECLAT igazi ereje a jelöltek támogatottságának meghatározásában van. A jelöltek  $TID$ -halmazainak előállítása egy rendkívül egyszerű és nagyon gyors művelet lesz. Emellett ahogy haladunk egyre mélyebbre a mélységi bejárás során, úgy csökken a  $TID$ -halmazok mérete, és ezzel a támogatottság meghatározásának ideje is. Ezzel szemben az APRIORI-nál ahogy haladunk az egyre nagyobb méretű jelöltek felé, úgy nő a szófa mélysége, és lesz egyre lassabb minden egyes jelölt támogatottságának meghatározása (persze a zsákutca nyesés segít ezen egy kicsit).

A keresési tér bejárása függ a prefix definíciójától, amit az elemeken definiált rendezés határoz meg. Melyek lesznek azok a jelöltek, amelyek az APRIORI-ban nem lennének jelöltek (tehát biztosan ritkák), illetve várhatóan melyik az a rendezés, amely a legkevesebb ilyen tulajdonságú halmazt adja? Ha egy elemhalmaz jelölt az ECLAT algoritmusban, de az APRIORI-ban nem, akkor van olyan részhalmaza, amely ritka. Amennyiben feltételezzük, hogy az elemek függetlenek, akkor azon részhalmaz előfordulásának lesz legkisebb a valószínűsége (és ezzel



együtt az esélye annak, hogy ritka), amely a leggyakoribb elemet nem tartalmazza. A jelölt prefixe generátor, tehát gyakori, így akkor lesz a legnagyobb esélye annak, hogy minden részhalmoz gyakori, ha a prefix a leggyakoribb elemet nem tartalmazza. Az ECLAT algoritmusnál a legkevesebb ritka jelöltet és így a legjobb futási időt tehát a gyakoriság szerint növekvő rendezéstől várhatjuk.

**4.16. példa.** *Ennek a gondolatmenetnek az illusztrálására nézzük a következő példát. Legyenek gyakori halmazok a következők:  $A, B, C, D, AB, AC, BC, AD, ABC$ , továbbá  $\text{supp}(D) \prec \prec \text{supp}(C) \prec \text{supp}(B) \prec \text{supp}(A)$ . Amennyiben az ECLAT algoritmus a gyakoriság szerint csökkenő sorrendet használja, akkor az előállítás sorrendjében a következő halmazok lesznek jelöltek:  $A, B, C, D, AB, AC, AD, ABC, ABD, ACD, BC, BD, CD$ . Ugyanez gyakoriság szerint növekvő sorrendnél  $D, C, B, A, DC, DB, DA, CB, CA, CBA, BA$ . Az utóbbi esetben tehát négy ritka jelölt helyett ( $ABD, ACD, BD, CD$ ) csak kettő lesz ( $CD, BD$ ). Megjegyezzük, hogy ez a két elemhalmaz az APRIORI esetében is jelölt lesz. A gyakoriság szerint csökkenő esetben egyszer állítunk elő olyan háromelemű jelöltet, amelynek van olyan kételemű részhalmoz, amelyet nem vizsgáltunk. Ez a jelölt a  $CBA$  és a nem megvizsgált részhalmoz a  $BA$ . Mivel a részhalmoz éppen a leggyakoribb elemeket tárolja, ezért van nagy esélye annak, hogy gyakori (főleg ha hozzá vesszük, hogy a jelölt két generátora,  $CB$  és  $CA$  is gyakori).*

Javíthatunk az algoritmus hatékonyságán, ha nem a jelöltek TID-listáit tároljuk, hanem a jelölt és prefixe TID-listájának különbségét. A prefix támogatottságából és a TID listák különbségéből a támogatottság egyértelműen megadható. A különbségi listák akár nagyobbak is lehetnek az eredeti TID-listáknál (például, ha a  $I$  támogatottsága kicsi, de a prefixének támogatottsága nagy), így a legjobb megoldást a két technika ötvözése adhatja (például 4-nél kisebb elemszámnál TID lista, utána különbségi listák) [147]. A különbségi listát használó algoritmusok nagy fölénnyel verik a többi algoritmust, amennyiben a bemenet sűrű, és nagy méretű gyakori minták is vannak.

### 4.3.1. kdc

### 4.3.2. lcm

## 4.4. Az FP-growth algoritmus

Az FP-GROWTH algoritmus<sup>2</sup>[56] egy mélységi jellegű, rekurzív algoritmus, a keresési tér bejárása tekintetében megegyezik az ECLAT-tal. A támogatottságok meghatározását az egyelemű gyakori halmazok meghatározásával, majd a bemenet *szűrésével* és *vetítésével* valósítja meg rekurzív módon. A bemenet szűrése azt jelenti, hogy az egyes tranzakciókból töröljük a bennük előforduló ritka elemeket. A  $\mathcal{T}$  elemhalmaz  $P$  elemhalmazra vetítését (jelölésben  $\mathcal{T}|P$ ) pedig úgy kapjuk, hogy vesszük a  $P$ -t tartalmazó tranzakciókat, majd töröljük belőlük a  $P$ -t. Például  $\langle ACD, BCE, ABCE, BE, ABCE \rangle | B = \langle CE, ACE, E, ACE \rangle$ . Az algoritmus pszeudokódja a következőkben olvasható.

<sup>2</sup>Az FP a Frequent Pattern rövidítése, ami miatt az algoritmust *mintanövelő* algoritmusnak is hívják. Ez az elnevezés azonban félrevezető, ugyanis szinte az összes GYEK algoritmus mintanövelő abban az értelemben, hogy egy új jelölt a generátorainak egyelemű bővítése, vagy más szóval növelése. Az FP-GROWTH sajátja nem a jelöltek előállítása, hanem a jelöltek támogatottság-meghatározásának módja.

**Algorithm 5** FP-GROWTH

---

**Require:**  $\mathcal{T}$ : tranzakciók sorozata,  
 $min\_supp$ : támogatottsági küszöb,  
 FP-GROWTH-SEGÉD( $\mathcal{T}, min\_supp, \emptyset$ )

---

A segédeljárás harmadik paramétere ( $P$ ) egy prefix elemhalmaz, az első paraméter pedig az eredeti bemenet  $P$ -re vetítése. Az eredeti bemenet  $\emptyset$ -ra vetítése megegyezik önmagával.

**Algorithm 6** FP-GROWTH-SEGÉD

---

**Require:**  $\mathcal{T}$ : vetített bemenet,  
 $min\_supp$ : támogatottsági küszöb,  
 $P$ : prefix elemhalmaz,  
 támogatottság\_meghatározás( $\mathcal{T}, J_1$ );  
 $GY_1 \leftarrow$  gyakoriak\_kiválogatása( $J_1, min\_supp$ );  
 $T^* \leftarrow$  SZŰRÉS( $T, GY_1$ )  
**for all**  $gy \in GY_1$  **do**  
    $T^*|gy \leftarrow$  VETÍTÉS( $T^*, gy$ )  
    $GY \leftarrow GY \cup \{P \cup \{gy\}\}$  FP-GROWTH-SEGÉD( $T^*|gy, min\_supp, P \cup \{gy\}$ )  
    $T^* \leftarrow$  TÖRLÉS( $T^*, gy$ )  
**end for**  
 return  $GY$

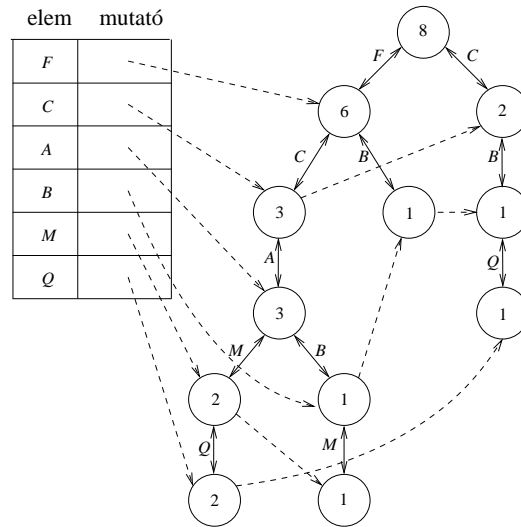
---

Egy rekurziós lépés három fő lépésből áll. Először meghatározzuk azon elemek támogatottságát, amelyek előfordulnak valamelyik tranzakcióban. Ezekből kiválasztjuk a gyakoriakat. Ezután minden  $gy$  gyakori elemet egyesével vesszünk. Meghatározzuk a  $gy$ -hez tartozó vetített bemenetet, majd meghívjuk az algoritmust rekurzívan a  $\mathcal{T}|gy$  bemenetre. Törölnünk kell a  $gy$  elemet a  $\mathcal{T}^*$ -beli tranzakciók elemei közül annak érdekében, hogy egy jelöltet csak egyszer állítsunk elő.

A jelöltek előállításának tekintetében az FP-GROWTH algoritmus a legegyszerűbb. Ha az  $I$  elemhalmaz gyakori, akkor a következő rekurziós szinten azon  $I \cup j$  halmazok lesznek a jelöltek, ahol  $j$  az  $I$ -re vetített bemenetben előforduló elem és  $I \cup j$  nem volt jelölt korábban. Tulajdonképpen az FP-GROWTH a nagy elemszámú jelöltek támogatottságának meghatározását visszavezeti három egyszerű műveletre: egyelemű gyakori elemhalmazok kiválogatása, szűrés és vetített bemenet előállítás.

A szűrés után egyesével vesszük a gyakori elemeket. Ezt valamilyen rendezés szerint kell tennünk és ez a rendezés határozza meg, hogy milyen sorban járjuk be a keresési teret, milyen vetített bemeneteket állítunk elő és mely elemhalmazok lesznek a hamis jelöltek. Az ECLAT-nál elmondottak itt is élnek; várhatóan abban az esetben lesz a hamis jelöltek száma minimális, amennyiben a prefixben a legritkább elemek vannak, azaz a 9. sorban gyakoriság szerint növekvő sorban vesszük az elemeket.

Az FP-GROWTH algoritmus szerves része az *FP-fa*, amelyben a szűrt bemenetet tároljuk. Az FP-fa segítségével könnyen előállíthatjuk a vetített bemeneteket, azokban könnyen meghatározhatjuk az elemek támogatottságát, amiből előállíthatjuk a vetített, majd szűrt bemenetet. Ezt a vetített és szűrt bemenetet szintén egy FP-fában tároljuk, amelyet *vetített FP-fának* hívunk.



4.9. ábra. Az  $\langle ACFMQ, ABCFM, BF, BCQ, ACFMQ, C, F, F \rangle$  szűrt bemenetet tároló FP-fa.

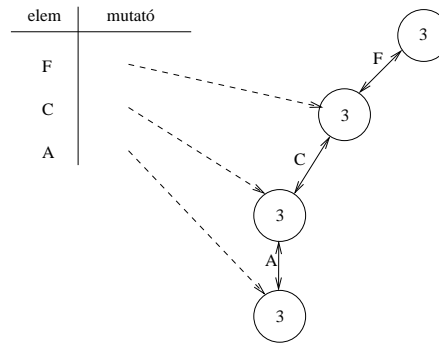
Az FP-fa egy keresztélelkel és egy fejléc táblával kibővített szófa. Az élek címkei gyakori elemek. Az egyszerűbb leírás kedvéért egy (nemgyökér) csúcs címkején a csúcsba mutató él címkejét értjük. Minden csúcs egy elemhalmazt reprezentál, amelynek elemei a gyökérből a csúcsig vezető út csúcsainak címkeivel egyeznek meg. Minden csúcsához egy számlálót rendelünk. Ez a számláló adja meg, hogy a csúcs által reprezentált halmaz mennyi bemeneti (vagy vetített) elemhalmaznak a prefixe. Az azonos címkejű csúcsok láncolt listaszerűen össze vannak kötve keresztirányú élekkel. A lánc legelső elemére mutat a fejléctáblának az adott eleméhez tartozó mutatója.

**4.17. példa.** *Tegyük fel, hogy bemenetként a  $\langle ACDFMQ, ABCFMO, BFO, BCKSQ, ACFMQ, CS, DFJ, FHI \rangle$  sorozat van adva, és  $\min\_supp = 3$ . A gyakori elemek: A, B, C, F, M, Q, amelyek támogatottsága rendre 3, 3, 5, 6, 3, 3. Ekkor a szűrt bemenetet ( $\langle ACFMQ, ABCFM, BF, BCQ, ACFMQ, C, F, F \rangle$ ) reprezentáló FP-fa, amely gyakoriság szerint csökkenő sorrendet ( $Q \prec M \prec B \prec A \prec C \prec F$ ) használ, a 4.9. ábrán látható*

Egy FP-fát hasonló módon építünk fel, mint egy szófát. Különbség, hogy egy  $I$  elemhalmaz beszúrásánál nem csak az  $I$ -t reprezentáló levélnek a számlálóját növeljük eggyel, hanem minden olyan csúcsot, amelyet érintünk a beszúrás során (hiszen ezen csúcsokat reprezentáló halmazok az  $I$  prefixei). A keresztirányú éleket és a fejléctáblát is egyszerűen megkaphatjuk. Legyen a fejléctábla mutatóinak kezdeti értéke NIL. Amikor beszúrunk egy új,  $i$  címkejű csúcsot, akkor két dolgot kell tennünk. Az új csúcs keresztél mutatója felveszi a fejléctábla  $i$ -hez tartozó bejegyzését, majd ezt a bejegyzést az új csúcs címére cseréljük. Ezzel tulajdonképpen olyan láncot készítünk, amelyben a csúcsok a beszúrási idejük szerint csökkenően vannak rendezve (az először beszűrt elem van leghátul) és a lista a fejléctáblában kezdődik.

A fejléc mutatókból kiindulva és a keresztéleket követve megkaphatjuk a vetített bemenetet és meghatározhatjuk a vetített bemenetben gyakori elemeket. Az adott tranzakciók előfordulása megegyezik a keresztélek által mutatott pontok számlálójával. Ezek alapján a vetített bemenetet szűrhetjük és belőle egy újabb FP-fát építhetünk fel. Ezt a fát vetített FP-fának hívjuk. A

következő ábrán az  $M$  elemhez tartozó vetített és szűrt bemenet FP-fáját láthatjuk (amelyet a  $Q$  elem feldolgozása után kapunk).



4.10. ábra. példa: vetített FP-fa

Az FP-fa mérete – hasonlóan a szófa méretéhez – függ az elemeken definiált rendezéstől. Az FP-GROWTH algoritmus akkor lesz hatékony, ha a fa elfér a memóriában, ezért fontos lenne azt a rendezést használni, ami várhatóan a legkisebb fát eredményezi. Az APRIORI esetében már elmondtuk, hogy az a heurisztika, amely az elemek gyakoriság szerint csökkenő rendezését használja, általában kis méretű fát eredményez.

Egyszerű lesz a vetített bemenet előállítás és a szűrt bemenetből egy elem törlése, amennyiben a legritkább gyakori elemet ( $gy_r$ ) vesszük először. Ez összhangban áll azzal, hogy a pszeudokód 9. sorában az elemeket gyakoriság szerint növekvő sorrendben vesszük. A  $gy_r$  csak levél címkéje lehet. Mivel a fából törölni fogjuk a  $gy_r$  címkéjű csúcsokat a rekurziós művelet után (13. sor), a következő elem is csak levél címkéje lesz.

Nézzük most meg, hogy amennyiben a szűrt bemenet egy FP-fában van tárolva, akkor hogyan kaphatjuk meg a  $gy_r$  elemre vett vetítésben az elemek támogatottságát. A fejléctábla  $gy_r$  eleméhez tartozó mutatóból kiindulva a keresztélek alkotta láncban pontosan azok a csúcsok vannak, amelyek  $gy_r$ -t tartalmazó bemeneti elemet reprezentálnak. Az egyes elemhalmazok előfordulását a  $gy_r$  címkéjű csúcsokhoz rendelt számláló adja meg, az elemeket pedig a gyökérig felsétálva kaphatjuk. A lista utolsó csúcsának feldolgozása után rendelkezésünkre állnak a  $gy_r$  elemhez tartozó vetített bemenetben valahol előforduló elemek támogatottságai, amely alapján kiválogathatjuk a vetített bemenetben gyakori elemeket.

Ugyanilyen bejárással kaphatjuk meg a vetített, majd szűrt bemenetet tartalmazó FP-fát. A fejléctáblából kiindulva végigmegyünk a láncolt lista elemein. A csúcs által reprezentált elemhalmazból töröljük a ritka elemeket, majd a kapott elemhalmazt beszúrjuk az új FP-fába. A kis memóriaigény érdekében a gyakoriság szerint csökkenő sorrendet használjuk. Ezt a sorrendet a vetített bemenet alapján állítjuk fel (lévén az új fa a vetített és szűrt bemenetet fogja tárolni), ami különbözhet az eredeti FP-fában alkalmazott rendezéstől.

**4.18. példa.** Folytassuk az előző példát és állítsuk

„A nyugalom megóv az UV sugáraktól Amerikai kutatók szerint a stressz és az UV-sugárzás együttesen tudnak csak igazán veszélyesek lenni. Ez az eredmény azt az ismert tényt erősíti meg, hogy a krónikus stressz lecsökkenti a bőr védekező képességét. Ha tehát nem idegeskedünk, nem kell félnünk a napsugaraktól.”  
 Forrás: [http://www.habostorta.hu/hab/tomy/tudomany/200507/a\\_nyugalom\\_megov\\_az\\_uvsugaraktol?print=1](http://www.habostorta.hu/hab/tomy/tudomany/200507/a_nyugalom_megov_az_uvsugaraktol?print=1)

*elő a legritkább gyakori elemhez ( $Q$ ) tartozó vetített és szűrt bemenetet. A fejléctábla  $Q$  eleméhez tartozó mutatóból kiindulva mindössze két csúcsot látogatunk meg, ami azt jelenti, hogy a vetített bemenet két különböző elemhalmazt tartalmaz: az FCAM-et kétszer, a CB-t egyszer. Ez alapján a vetített bemenetben egyetlen gyakori elem van,  $C$ . Ez a rekurziós ág nem folytatódik, hanem visszatér a QC gyakori elemhalmazzal. Az FP-fából törölhetjük a fejléctábla  $Q$  bejegyzéséhez tartozó mutatóból, keresztirányú élek segítségével elérhető csúcsokat. A következő vizsgált elem az  $M$ . Az  $M$  vetített bemenetében három gyakori elem van, és a vetített szűrt bemenet az FCA elemhalmazt tartalmazza háromszor. Ezt a vetített, szűrt bemenetet egy egyetlen útból álló FP-fa fogja reprezentálni. A többi FP-fa ugyanilyen egyszerűen megkapható.*

Hatékonysági szempontból rendkívül fontos, hogy a rekurziót ne folytassuk, ha a vizsgált FP-fa egyetlen útból áll. A rekurzió helyett képezzük inkább az út által reprezentált elemhalmaz minden részhalmazát. A részhalmaz támogatottságát annak a csúcsnak a számlálója adja meg, amely a legmélyebben van a részhalmazt meghatározó csúcsok között.

#### 4.4.1. Az FP-growth\* algoritmus

2003 novemberében megszervezték az első gyakori elemhalmaz-kinyerő algoritmusok versenyét [49]. Bárki benevezhetett egy általa készített programot. Ezeket központilag tesztelték különböző adatbázisokon, különböző támogatottsági küszöbökkel. Nem volt olyan implementáció, amely minden esetben a legjobban szerepelt, de ki lehet emelni néhány olyat, amelyek szinte mindig az elsők között végeztek. A szervezők végül annak adták a fődíjat (egy sört és egy pelenkát!), aki az FP-GROWTH\* algoritmust [50] küldte be.

Az FP-GROWTH\* algoritmus az FP-GROWTH módosítása. Előnye, hogy gyorsabban állítja elő a vetített fát, amiért viszont memóriával fizet. Nézzük meg, hogy pontosan mi történik egy rekurziós lépésben. Először ellenőrizzük, hogy a fa egyetlen útból áll-e. Ha nem, akkor a legritkább elemből kiindulva előállítjuk a vetített fákat, és rekurzívan meghívjuk az algoritmust. A vetített fában első lépésként meg kell határozni a vetített bemenetben szereplő elemek támogatottságát, második lépésként pedig előállítjuk a vetített FP-fát. Ez tulajdonképpen az aktuális fa adott elemhez tartozó ágainak kétszeri bejárását jelenti. Az első bejárást lehet meggyorsítani egy segéd tömb használatával.

Az FP-fa építésénél töltsünk fel egy, kezdetben 0 értékeket tartalmazó tömböt is. Amikor beszurunk egy  $t$  (akár vetített) tranzakciót az (akár vetített) FP-fába, növeljük eggyel a tömb  $(i, j)$ -edik celláját, amennyiben az  $i$  és  $j$  elemei  $t$ -nek. A fa felépítése után rendelkezésünkre áll egy tömb, ami tartalmazza az elempárok előfordulását. Ha ezek után egy vetített fát akarunk készíteni, akkor szükségtelen időt töltenünk az első lépéssel, hiszen a tömb megfelelő sorából közvetlen megkaphatjuk a támogatottságokat. Összességében az első lépés gyorsabb (nem kell a fában bolyonganunk, csak a tömb elemeit kiolvasni), a második lassabb (a tömböt is fel kell tölteni), a memóriefogyasztás pedig nagyobb (a tömb méretével).

#### 4.4.2. Patricia

### 4.5. Elavult technikák

1993 óta kétszáz körüli cikk jelent meg gyakori elemhalmazokat kinyerő algoritmusok témájában. Legtöbb cikk egy új gyorsítási trükköt vagy egy új módszert mutatott be és a szerzői azt állítják, hogy az ő módszerük a legjobb. Ez nyilvánvalóan képtelenség. A rengetek módszer miatt kialakult káoszt néhány kutató megelégette és megrendezték 2003-ban és 2004-ben a gyakori elemhalmazokat kinyerő algoritmusok versenyét. Sebesség tekintetében az ECLAT és az FP-GROWTH különböző módosításai voltak a legjobbak, a memória terén pedig az APRIORI volt kiemelkedő.

A továbbiakban felsoroljuk azokat az algoritmusokat, amelyeket mai napig megtalálhatunk különböző adatbányászati tankönyvekben, a legtöbb kutató által ismertek, de a versenyek során egyik sem bizonyította be, hogy beférne az elit algoritmusok körébe. A tanulmány korábbi verzióiban ezek az elavult módszerek leírásai is ott szerepeltek az APRIORI, ECLAT és FP-GROWTH mellett, mára azonban csak a következő listában kapnak helyet: SETM [60], APRIORI-TID [5], APRIORI-HYBRID [5], DHP [103], DIC, Patrícia, Tree projection, DF-apriori [109],

### 4.6. Mintavételező algoritmus elemzése

Az egyszerű mintavételező algoritmust bemutatunk a 10.5.4 részben. Itt azt vizsgáljuk, hogy mekkora mintát célszerű venni annak érdekében, hogy az algoritmus minden gyakori elemhalmazt megtaláljon.

#### 4.6.1. Mintavétel nagysága

Mintavételezésen alapuló eljárásoknál a minta mérete központi kérdés. Ha a minta túl kicsi, akkor a mintából nyert információ távol állhat a teljes adatbázisban található globális „helyzettől”. Mivel főlegesen nagy minta lassú algoritmusokat eredményez, ezért fontos egy kicsi, de már pontos képet adó mintaméret meghatározása. A 3.3.7 részben megadtuk, hogy mekkora mintát kell választani, ha azt akarjuk, hogy a relatív gyakoriságok megegyezzenek az előfordulások valószínűségével. Használjuk most is a A 3.3.7 részben bevezetett elnevezéseket és jelöléseket.

Nézzük, hogy mennyivel kell csökkenteni a gyakorisági küszöböt ( $min\_freq'$ ) ahhoz, hogy kicsi legyen annak valószínűsége, hogy tetszőleges gyakori elem mintához tartozó gyakorisága kisebb a csökkentett küszöbnél, tehát:

$$\mathbb{P}(\text{gyakoriság}(x, m) < min\_freq') = \mathbb{P}\left(\frac{Y}{m} < min\_freq'\right)$$

egy adott küszöbnél ( $\delta'$ ) kisebb kell legyen és tudjuk, hogy

$$p > min\_freq$$

A fenti egyenletre alkalmazva a Hoeffding-korlátot azt kapjuk, hogy

$$\mathbb{P}\left(\frac{Y}{m} < min\_freq'\right) =$$

$$\begin{aligned} & \mathbb{P}\left(\frac{Y}{m} - p < \min\_freq' - p\right) < \\ & \mathbb{P}\left(\frac{Y}{m} - p < \min\_freq' - \min\_freq\right) \\ & \leq e^{-2(\min\_freq' - \min\_freq)^2 m} \end{aligned}$$

tehát ahhoz, hogy a hibázás valószínűsége kisebb legyen  $\delta'$ -nél teljesülnie kell, hogy

$$\min\_freq' < \min\_freq - \sqrt{\frac{1}{2m} \ln \frac{1}{\delta'}}$$

A 4.2 táblázat azt mutatja, hogy rögzített hibakorlát mellett ( $\delta' = 0.001$ ) adott mintamérethez mennyi legyen a csökkentett küszöb.

min_freq (%)	Minta mérete			
	20000	40000	60000	80000
0.25	0.13	0.17	0.18	0.19
0.50	0.34	0.38	0.40	0.41
0.75	0.55	0.61	0.63	0.65
1.00	0.77	0.83	0.86	0.88
1.50	1.22	1.30	1.33	1.35
2.00	1.67	1.77	1.81	1.84

4.2. táblázat. A küszöb csökkentése adott mintaméretekre rögzített  $\delta = 0.001$  mellett

## 4.7. Elemhalmazok Galois lezárja

Egy minta zárt, ha nincs vele egyező támogatottságú bővebb minta. Esetünkben ez azt jelenti, hogy ha egy elemhalmaz nem zárt, akkor pontosan azokban a bemeneti elemekben fordul elő, amelyekben a lezártja. Ha például az  $A$  elem lezártja az  $AB$  halmaz, akkor tudjuk, hogy az  $A$  halmaz soha nem fordul elő a bemeneti elemekben a  $B$  elem nélkül.

Ebben a részben a lezárt további tulajdonságait fogjuk megismerni. Azért illetjük a lezártat a Galois jelzővel, mert teljesülni fog a lezáras operátorra a Galois elméletből jól ismert 3 tulajdonság. Mielőtt erre rátérünk nézzük meg, hogy az elemhalmazokat tartalmazó mintakörnyezet egyértelmű-e a zártságra nézve.

**4.19. lemma.** *Az elemhalmazokat tartalmazó mintakörnyezet a zártságra nézve egyértelmű.*

*Bizonyítás:* Indirekt tegyük fel, hogy az  $I$  elemhalmaznak létezik két lezártja, azaz létezik  $I', I''$  különböző elemhalmazok, amelyekre a minimalitás mellett teljesülnek a  $I \subset I', I \subset I'', |I'| = |I''|, \text{supp}(I') = \text{supp}(I'')$  feltételek. Ez azt jelenti, ahogy azon tranzakciók, amelyek  $I$ -t tartalmaznak, tartalmazzák az  $I' \setminus I$  és az  $I'' \setminus I$  halmazokat is. De ebből következik, hogy ezek a tranzakciók  $I' \cup I''$  is tartalmazzák, azaz  $I' \cup I''$  is lezártja  $I$ -nek, így sem  $I'$  sem  $I''$  nem lehet minimális.

■

A fentiek miatt a gyakori zárt elemhalmazokból és azok támogatottságaiból egyértelműen meg tudjuk határozni a gyakori elemhalmazokat és azok támogatottságát. A gyakori zárt minták tehát a zárt minták egy veszteségmentes tömörítése, érdemes csak ezeket meghatározni és eltárolni [104–106, 150].

#### 4.7.1. A zárt elemhalmazok fogalma

Az  $I$  elemhalmaz zárt, amennyiben nincs nála bővebb halmaz, amelynek támogatottsága megegyezik  $I$  támogatottságával. Jelöljük  $cover$ -rel azt a függvényt, amely egy elemhalmazhoz az azt tartalmazó tranzakciók halmazát adja meg.

A zárt elemhalmazokra adhatunk egy másik definíciót is. Vezessük, be a  $cover'$  függvényt:

**4.20. definíció.** Legyen  $\mathcal{T} = \langle t_1, \dots, t_n \rangle$  tranzakciók sorozata, amelynek minden eleme az  $\mathcal{J}$ -nek egy részhalmaza. Definiáljuk a  $cover' : 2^{\mathcal{N}} \rightarrow 2^{\mathcal{J}}$  függvényt a következőképpen

$$cover'(T) = \{i \in \mathcal{J} \mid \forall j \in T, i \in cover(t_j)\} = \bigcap_{t \in T} cover(t)$$

Tehát  $cover'(T)$  megadja azon közös elemeket, amelyeket minden olyan tranzakció tartalmaz, amelynek sorszámja  $T$ -beli.

A  $(cover, cover')$  függvényt az  $\mathcal{T}$  és  $\mathcal{J}$  hatványhalmazai közötti Galois-kapcsolatnak hívjuk. Legyen a példaadatbázisunk a következő:  $\langle ACD, BCE, ABCE, BE, ABC E \rangle$ . Ekkor:  $cover(\{A, C\}) = \{1, 3, 5\}$ ,  $cover(\emptyset) = \{1, 2, 3, 4, 5\}$ ,  $cover'(\{1, 2, 3\}) = \{C\}$ ,  $cover'(\{1, 4\}) = \emptyset$ .

Az alábbi tulajdonságok igazak tetszőleges  $t, t_1, t_2 \subseteq \mathcal{T}$  és  $I, I_1, I_2 \subseteq \mathcal{J}$  halmazokra:

- (1)  $I_1 \subseteq I_2 \Rightarrow cover(I_1) \supseteq cover(I_2)$     (1')  $T_1 \subseteq T_2 \Rightarrow cover'(T_1) \supseteq cover'(T_2)$   
 (2)  $T \subseteq cover(I) \iff I \subseteq cover'(T)$

**4.21. definíció.** A  $h = cover' \circ cover$  (vagy  $h' = cover \circ cover'$ ) operátort Galois-lezárás operátornak hívjuk.

Belátható, hogy tetszőleges halmaznak a lezártja tartalmazza magát a halmazt, továbbá a Galois-lezárás operátora idempotens és monoton, tehát

$$\begin{array}{ll} (I) I \subseteq h(I) & (I') T \subseteq h'(T) \\ (II) h(h(I)) = h(I) & (II') h'(h'(T)) = h'(T) \\ (III) I_1 \subseteq I_2 \Rightarrow h(I_1) \subseteq h(I_2) & (III') T_1 \subseteq T_2 \Rightarrow h'(T_1) \subseteq h'(T_2) \end{array}$$

**4.22. definíció (zárt elemhalmaz).**  $I$  elemhalmaz zárt, amennyiben  $I = h(I)$ .

Tetszőleges elemhalmazt  $(I)$  tartalmazó minimális elemszámú zárt elemhalmazt a lezárás operátor alkalmazásával kaphatunk meg; ez éppen  $h(I)$  lesz. A példaadatbázisban található zárt elemhalmazok alábbiak:



zárt elemhalmazok
$\{\emptyset\}, \{C\}, \{B,E\},$ $\{B,C,E\}, \{A,C\}, \{A,B,C,E\},$ $\{A,C,D\}, \{A,B,C,D,E\}$

Adósok vagyunk még annak bizonyításával, hogy a két definíció ekvivalens, azaz, ha  $h(C) = C$ , akkor  $C$ -nél nincs bővebb halmaz, amely támogatottsága megegyezne  $C$  támogatottságával, illetve fordítva. A két állítás közvetlen adódik a következő tételből.

**4.23. tétel.** *Minden elem támogatottsága megegyezik lezártjának támogatottságával, tehát*

$$\text{supp}(I) = \text{supp}(h(I))$$

*Bizonyítás:* A lezárás (1) tulajdonsága miatt  $\text{supp}(I) \geq \text{supp}(h(I))$ . Ugyanakkor

$$\text{supp}(h(I)) = |\text{cover}(h(I))| = |\text{cover}(\text{cover}'(\text{cover}(I)))| = |h'(\text{cover}(I))| \leq \text{supp}(I)$$

a (III') miatt, amiből következik az egyenlőség. ■

A 10.4.2 részben bemutatjuk, hogy a gyakori mintákból hogyan választhatjuk ki a zártakat, illetve az APRIOR-CLOSE algoritmust, ami már eleve csak a gyakori zárt mintákat állítja elő. Az APRIOR-CLOSE algoritmusnál léteznek gyorsabb algoritmusok (CHARM [149], CLOSET [108], CLOSET+ [140], MAFIA [24]), ezek ismertetésétől eltekintünk.

## 4.8. Kényszerek kezelése

Ebben a részben azt a speciális feladatot nézzük meg, hogy miként lehet csökkenteni a bemenetet, ha az anti-monoton kényszerek mellett monoton kényszereket is megadunk. Már az általános mintakeresésnél megtárgyaltuk, hogy tetszőleges anti-monoton kényszer könnyűszerrel beépíthető az APRIORI algoritmusba. Most azt nézzük meg, hogy a monoton kényszerek hogyan alkalmazhatók a bemeneti tér csökkentésére.

Adott egy bemeneti sorozat, minimális támogatottsági küszöb és monoton kényszerek  $\mathcal{C}$  halmaza. Feladat a bemenet csökkentése oly módon, hogy bármely teljes algoritmus a csökkentett bemeneten is teljes legyen.

### 4.8.1. ExAnte

Az ExAnte [80] algoritmus kétféle lépést ismétél egészen addig, amíg ez valamilyen változást jelent. Az első lépés azon tranzakciók törlése, amelyek nem adnak igaz értéket minden  $\mathcal{C}$ -beli kényszeren. Az ilyen tranzakciók csak olyan minták támogatottságát növelik, amelyek úgysem elégítik ki a kényszereket (ez következik a kényszerek monoton tulajdonságából). A második lépésben a bemenet elemei közül töröljük a ritkákat, hiszen azok úgysem játszanak szerepet a támogatottság meghatározásánál.

Látunk kell, hogy az első lépésbeli törlés új ritka elemekhez vezethet, ami csökkenti bizonyos tranzakciók méretét, ami viszont ahhoz vezethet, hogy ezek újabb kényszereket fognak sérteni. Jogos tehát, hogy a két módszert felváltva futtassuk addig, amíg van valami változás. Az algoritmus a bemenet csökkentése mellett előállítja azon gyakori elemeket, amelyekre minden kényszer teljesül. Gyakori elemhalmaz csak ezekből az elemekből épülhetnek fel.

Nézzünk egy példát. Az adatbázisban 8 elem és 9 tranzakció van. Legyen  $min\_supp = 4$ . Minden elemnek van egy ára. Az egyetlen kényszer ( $\sum(i.ár) > 44$ ) szerint a halmazban található termékek árának összege 44-nél nagyobb legyen. A következő két táblázat adja meg az adatokat.

termék	ár	TID	tranzakció	ár összeg
A	5	1	<i>B, C, D, G</i>	58
B	8	2	<i>A, B, D, E</i>	63
C	14	3	<i>B, C, D, G, H</i>	70
D	30	4	<i>A, E, G</i>	31
E	20	5	<i>C, D, F, G</i>	65
F	15	6	<i>A, B, C, D, E</i>	77
G	6	7	<i>A, B, D, F, G, H</i>	76
H	12	8	<i>B, C, D</i>	52
		9	<i>B, E, F, G</i>	49

Az első végigolvasás során meghatározzuk az elemek támogatottságát azon tranzakciókban, amelyek kielégítik a kényszert (a 4-es kivételével mindegyik). Ezután töröljük a ritka elemeket (*A, E, F, H*). Ismét végigmegyünk az adatbázison, de most már ezeket az elemeket nem nézzük, aminek következtében újabb tranzakciók esnek ki (2,7,9). A kiesett tranzakciók miatt csökkennek a támogatottságok, így újabb elem lesz ritka (*G*). Ezt így folytatjuk, amíg van változás. A 4. végigolvasás után azt kapjuk, hogy csak az 1,3,6,8 tranzakciókat és a *B, C, D* elemeket kell figyelembe venni.

## 4.9. Többszörös támogatottsági küszöb

Az univerzális támogatottsági küszöbnek vannak előnyei és hátrányai. Előnye, hogy felhasználhatjuk azt a tényt, hogy gyakori minta minden részmintája gyakori, ami alapján hatékony algoritmusokat adhatunk. Hátránya, hogy a ritkán előforduló, de mégis fontos mintákat csak akkor tudjuk kinyerni, ha a támogatottsági küszöböt alacsonyra állítjuk. Ez viszont rengeteg gyakori mintához fog vezetni, ha egyáltalán le tud futni az algoritmus.

Különböző támogatottsági küszöbök (vagy másként támogatottsági küszöb függvényének) megadásával ez a probléma elkerülhető: a nem lényeges mintáknak legyen nagy a küszöbük, a lényegesebbeknek legyen alacsony.

Egyedi támogatottsági küszöbök bevezetésével azonban felborul eddigi kényelmes világunk, amelyet az biztosított, hogy nem lehet egy minta gyakori, ha van ritka részmintája. A részminták támogatottsági küszöbe ugyanis nagyobb lehet, így hiába nagyobb a támogatottsága, ettől még lehet ritka. A következőkben bemutatjuk a legelső és legegyszerűbb támogatottsági küszöb függvényt, majd bemutatjuk az MSApriori algoritmust, amely ezt hatékonyan kezeli.

### 4.9.1. MSApriori algoritmus

Kézzel megadni a  $2^J$  minden elemének támogatottsági küszöbét fáradságos, sőt nagy  $|J|$  esetén kivitelezhetetlen feladat. Az MSApriori algoritmusnál csak az egyelemű elemhalmazok támogatottsági küszöbét lehet megadni. Jelöljük az  $i$  elem küszöbét  $MIS(i)$ -vel. Az  $I$  elemhalmaz támogatottsági küszöbe legyen a legkisebb támogatottsági küszöbvel rendelkező elemének támogatottsági küszöbe ( $MIS(I) = \min_{i \in I} \{MIS(i)\}$ ). Akkor gyakori az  $I$  halmaz, ha támogatottsága nagyobb vagy egyenlő  $MIS(I)$ -nél.

A definícióból következik, hogy tényleg nem mondhatjuk, hogy gyakori minta minden részmintája gyakori. Például az  $ABC$  elemhalmaz  $BC$  részalmazának nagyobb lehet MIS értéke. Ha a feladat megoldására az APRIORI algoritmust használjuk úgy, hogy csak a gyakori elemhalmazok kiválasztásának módját módosítjuk ( $min\_supp$  cseréje  $MIS(I)$ -re), akkor nem garantált, hogy jó megoldást kapunk. Ha például a  $BC$  ritka, akkor az  $ABC$  halmaz nem lenne a jelöltek között annak ellenére, hogy akár gyakori is lehet.

Szerencsére a probléma könnyen orvosolható. Csak azt kell észrevennünk, hogy mi okozhatja a hibát. Az általánosság megsértése nélkül feltehetjük, hogy az elemek MIS értékük alapján növekvő sorba vannak rendezve. A MIS definíciójából következik, hogy tetszőleges  $\ell$ -elemű  $I = \{i_1, \dots, i_\ell\}$  halmaz  $\ell - 1$  darab  $(\ell - 1)$ -elemű részalmazának MIS értéke megegyezik  $I$  MIS értékével, ami  $MIS(i_1)$ . Ezeknek a részalmazoknak tehát gyakornak kell lenniük, hiszen a támogatottság monotonitása most is fennáll. Az egyetlen részalmaz, amely lehet ritka, az  $I$  legelső elemét nem tartalmazó részalmaz. Ezt a részalmazt tehát ne vizsgáljuk a jelölt előállítás második lépése során. Kivétel ez alól azon eset, amikor a második elem MIS értéke megegyezik az első elem MIS értékével, mert ekkor még ennek a részalmaznak is gyakornak kell lennie.

Amennyiben  $\ell > 2$ , akkor biztos, hogy a generátorok egyike sem egyezik meg a legkisebb elemet nem tartalmazó részalmazmal ( $\ell > 2$  esetében ugyanis a generátorok  $(\ell - 2)$ -elemű prefixei megegyeznek, amelyek biztos, hogy tartalmazzák a jelölt első elemét). Ez pedig garantálja, hogy az algoritmus teljes, amennyiben az összes gyakori elempárt megtaláltuk. Nézzük meg most az egy- és kételemű jelöltek esetét.

Gyakori elemek meghatározásánál a szokásos eljárást követjük: minden elem jelölt. Elempárok esetében azonban nem állíthatjuk, hogy egy pár akkor jelölt, ha mindkét eleme gyakori. Például az  $AB$  pár lehet gyakori akkor is, ha az  $A$  ritka. Ha ugyanis  $B$ -nek MIS értéke kisebb  $A$ -nak MIS értékénél, akkor az  $AB$ -nek a MIS értéke megegyezik  $B$ -nek a MIS értékével, így  $AB$  lehet gyakori. Szerencsére szükségtelen az összes elemet figyelembe venni. Ha például az  $A$  elem ritka és az  $A$  MIS értéke a legkisebb, akkor a támogatottság monotonitásából következik, hogy az  $A$ -t tartalmazó halmazok ritkák. Ha tehát MIS érték szerint növekvően vannak rendezve az elemek, akkor a legkisebből kiindulva keressük meg az első gyakori elemet. Az összes utána következőt figyelembe kell venni a jelöltpárok előállításánál akkor is, ha valamelyik ritka.

„Vakságot okoz a nyakkendő. A kutatás szerint a szorosan megkötött nyakkendő csökkentheti a nyaki véna hatékonyságát, ezáltal a szem vérellátását, és hályog kialakulásához, legsúlyosabb esetben pedig részleges vagy teljes vaksághoz vezethet. Még veszélyesebb a helyzet a vékony nyakú emberek esetében, mert az ő vénájuk érzékenyebb – mutatnak rá az orvosok.” Forrás: <http://pvg.uw.hu/cikk/nyakkendo.html>

## 5. fejezet

# Asszociációs szabályok

A gyakori elemhalmazokat felhasználhatjuk arra, hogy gyakori elemhalmazokra vonatkozó szabályokat nyerjünk ki belőlük. Az  $I_1 \rightarrow I_2$  asszociációs szabály azt állítja, hogy azon bemeneti elemek, amelyek tartalmazzák  $I_1$ -et, tartalmazzák általában  $I_2$ -t is. Például a pelenkát vásárlók sört is szoktak venni.

Mi az értelme ezeknek a szabályoknak? Például az, hogy supermarket extra profithoz juthat az alábbi módon: Ha  $I_1 \rightarrow I_2$  szabály igaz, akkor óriási hírverés közepette csökkentjük  $I_1$  termékek árát (mondjuk 15%-kal). Emellett diszkréten emeljük meg  $I_2$  termék árát (mondjuk 30%-kal) úgy, hogy az  $I_1$  árcsökkenéséből származó profitcsökkenés kisebb legyen, mint az  $I_2$  áremeléséből származó profitnövekedés. Az  $I_1$  és  $I_2$  termékek eladásai együtt mozognak, tehát az  $I_2$  termék eladása is nőni fog. Amit veszünk a réven, azt megnyerjük a vámon: összességében a profitunk nőni fog, és a leárazás reklámnak is jó volt.

Korunkra jellemző olcsó internetes üzletek is ilyen szabályok alapján dolgoznak. Tudják milyen terméket vásárolnak együtt. Sokszor az együtt vásárlást elő is írják azzal, hogy nem adják el önmagában az olcsó árucikket, csak akkor, ha megveszi az ügyfél a drága kiegészítőt is.

Az ilyen szabályokból nyert információt használhatják emellett áruházak terméktérképének kialakításához is. Cél a termékek olyan elrendezése, hogy a vevők elhaladjanak az őket érdekelhető termékek előtt. Gondoljuk meg, hogyan lehet kiaknázni e célból egy asszociációs szabályt.

Elemhalmazok sorozatát ábrázolhatjuk bináris értékeket tartalmazó táblával is. Ekkor az asszociációs szabályok attribútumok közötti összefüggést mutatnak: ha az  $I_1$  attribútumok értékei 1-es, akkor nagy valószínűséggel az  $I_2$  attribútumok értéke is az. A valószínűség értékét a szabály *bizonyossága* adja meg. Csak olyan szabályok lesznek érdekesek, amelyek bizonyossága magas. Például a házasságban élők 85%-ának van gyermekük.

Az asszociációs szabályok felhasználási területe egyre bővül. A piaci stratégia meghatározásán túl egyre fontosabb szerepet játszik a döntéstámogatás és pénzügyi előrejelzések területén is.

Nézzük most az asszociációs szabály pontos definícióját.

## 5.1. Az asszociációs szabály fogalma

Használjuk a 4.1 részben bevezetett definíciókat és jelöléseket (elemhalmaz, kosár, támogatottság, fedés, gyakori elemhalmaz stb.).

**5.1. definíció (asszociációs szabály).** Legyen  $\mathcal{T}$  az  $\mathcal{I}$  hatványhalmaza felett értelmezett sorozat. Az  $R: I_1 \xrightarrow{c,s} I_2$  kifejezést  $c$  bizonyosságú,  $s$  támogatottságú asszociációs szabálynak nevezzük, ha  $I_1, I_2$  diszjunkt elemhalmazok, és

$$c = \frac{\text{supp}_{\mathcal{T}}(I_1 \cup I_2)}{\text{supp}_{\mathcal{T}}(I_1)},$$

$$s = \text{supp}_{\mathcal{T}}(I_1 \cup I_2)$$

A szabály bal oldalát feltétel résznek, a jobb oldalát pedig következmény résznek nevezzük.

Az  $R: I_1 \rightarrow I_2$  szabály bizonyosságára gyakran  $\text{conf}(R)$ -ként hivatkozunk.

Feladat egy adott kosársorozatban azon asszociációs szabályok megtalálása, amelyek gyakoriak (támogatottságuk legalább  $\text{min\_supp}$ ), és bizonyosságuk egy előre megadott korlát felett van. Jelöljük ezt a bizonyossági korlátot  $\text{min\_conf}$ -fal. A feltételt kielégítő szabályokat *érvényes asszociációs szabályoknak* hívjuk, az 1 bizonyossággal rendelkezőket pedig *egzakt asszociációs szabálynak*.

„Felmérések igazolják, hogy azok a legboldogabb párok, akik nemcsak hétköznapi problémájukat osztják meg egymással, de mernék a titkos álmaikról is beszélni.”  
Forrás: Wellness 2007. októberi szám 106. oldal

**5.2. definíció (érvényes asszociációs szabály).**  $\mathcal{T}$  kosarak sorozatában,  $\text{min\_supp}$  támogatottsági és  $\text{min\_conf}$  bizonyossági küszöb mellett az  $I_1 \xrightarrow{c,s} I_2$  asszociációs szabály *érvényes*, amennyiben  $I_1 \cup I_2$  gyakori elemhalmaz, és  $c \geq \text{min\_conf}$

A fenti feladatot két lépésben oldjuk meg. Először előállítjuk a gyakori elemhalmazokat, majd ezekből az érvényes asszociációs szabályokat. Az első lépésről szól a 4. fejezet, nézzük most a második lépést.

Minden  $I$  gyakori termékalmazt bontunk fel két diszjunkt nem üres részre ( $I = I_1 \cup I_2$ ), majd ellenőrizzük, hogy teljesül-e a  $\frac{\text{supp}(I)}{\text{supp}(I_1)} \geq \text{min\_conf}$  feltétel. Amennyiben igen, akkor a  $I_1 \rightarrow I_2$  egy érvényes asszociációs szabály. A támogatottság anti-monoton tulajdonságát felhasználhatjuk annak érdekében, hogy ne végezzünk túl sok felesleges kettéosztást.

**5.3. észrevétel.** Amennyiben  $I_1, I$  gyakori elemhalmazok a  $\mathcal{T}$  bemeneti sorozatban, és  $I_1 \subset I$ , illetve  $I_1 \rightarrow I \setminus I_1$  nem érvényes asszociációs szabály, akkor  $I'_1 \rightarrow I \setminus I'_1$  sem érvényes semmilyen  $I'_1 \subset I_1$ -re.

*Bizonyítás:* Az  $I_1 \xrightarrow{c,s} I \setminus I_1$  nem érvényes szabály, tehát  $c = \frac{\text{supp}(I_1 \cup (I \setminus I_1))}{\text{supp}(I_1)} = \frac{\text{supp}(I)}{\text{supp}(I_1)} < \text{min\_conf}$ . Mivel a támogatottság anti-monoton, ezért  $\text{supp}(I'_1) \geq \text{supp}(I_1)$ , amiből  $\frac{1}{\text{supp}(I'_1)} \leq \frac{1}{\text{supp}(I_1)}$ , és ebből, ha  $c'$ -vel jelöljük az  $I'_1 \rightarrow I \setminus I'_1$  szabály bizonyosságát, akkor

$$c' = \frac{\text{supp}(I)}{\text{supp}(I'_1)} \leq \frac{\text{supp}(I)}{\text{supp}(I_1)} < \text{min\_conf}$$

tehát  $I'_1 \rightarrow I \setminus I'_1$  sem érvényes asszociációs szabály.



**Weka 3.5.7** Az asszociációs szabályokkal kapcsolatos osztályokat az Explorer Associate fülén keresztül érhetjük el.



### 5.1.1. Maximális következményű asszociációs szabály

A maximális méretű gyakori mintákból az összes gyakori mintát meghatározhatjuk. Ez abból következik, hogy gyakori minta minden részmintája gyakori. Asszociáció szabályoknál is vannak olyanok, amelyekből más szabályok levezethetők. Nézzünk két egyszerű levezetési szabályt. Tegyük fel, hogy  $I_1 \rightarrow I_2$  érvényes asszociációs szabály, ekkor

- $I_1 \rightarrow I'_2$  is érvényes, minden  $I'_2 \subseteq I_2$ -re.
- $I_1 \cup i \rightarrow I_2 \setminus \{i\}$  is érvényes minden  $i \in I_2$ -re. Ezek szerint a következményrészről tetszőleges elemet áttehetünk a feltételrészbe.

Mindkét állítás a támogatottság anti-monoton tulajdonságából közvetlenül adódik.

Ezek szerint minden asszociációs szabály levezethető a maximális következményrészrel rendelkező asszociációs szabályokból. Persze a levezethetőség nem a legjobb szó, ugyanis a szabályok paramétereire nem tudunk következtetni.

### 5.1.2. Egzakt asszociációs szabályok bázisa

A 100%-os bizonyossággal rendelkező asszociációs szabályokat *egzakt asszociációs szabályoknak* hívjuk. Az egzakt asszociációs szabályokra érvényes tranzitivitás is, tehát  $I_1 \rightarrow I_2$  és  $I_2 \rightarrow I_3$ -ból következik, hogy  $I_1 \rightarrow I_3$ . Matematikus beállítottságú emberek agyában azonnal felmerül, hogy van-e az egzakt asszociációs szabályoknak egy minimális bázisa, amelyből minden egzakt asszociációs szabály levezethető. Ehhez a bázishoz a *pszeudo-zárt elemhalmazokon* keresztül jutunk.

**5.4. definíció.** Az  $I$  elemhalmaz  $\mathcal{T}$ -re nézve zárt, amennyiben nem létezik olyan  $I'$  minta, amelynek  $I$  valódi részhalmaza, és  $I'$  támogatottsága megegyezik  $I$  támogatottságával ( $\text{supp}(I') = \text{supp}(I)$ ).

**5.5. definíció.**  $I \subseteq \mathcal{J}$  pszeudo-zárt elemhalmaz, ha nem zárt, és minden pszeudo-zárt  $I' \subset I$  elemhalmazra fennáll, hogy lezártja valódi része  $I$ -nek.

Az üres halmaz pszeudo-zárt, amennyiben az nem zárt. Zárt elemhalmaz fogalmához

A pszeudo-zárt elemhalmazok segítségével tudunk egy olyan szabálybázist megadni, amelyekből az összes egzakt asszociációs szabály megkapható.

„Pici péniszt okozhat a parfüm”  
 Forrás: <http://www.ma.hu/page/cikk/aj/0/166581/1>

**5.6. definíció.** Legyen  $FP$  a pszeudo-zárt elemhalmazok halmaza  $\mathcal{T}$ -ben. Ekkor a Duquenne–Guigues-bázist a következőképpen definiáljuk:

$$DG = \{r : I_1 \rightarrow h(I_1) \setminus I_1 \mid I_1 \in FP, I_1 \neq \emptyset\},$$

ahol az  $I$  lezártját  $h(I)$ -vel jelöltük.

**5.7. tétel.** A Duquenne–Guigues-bázisból az összes egzakt szabály levezethető és a bázis minimális elemszámú, tehát az egzakt szabályoknak nincsen olyan kisebb elemszámú halmaza, amelyből az összes egzakt asszociációs szabály levezethető.

A Duquenne–Guigues-bázis meghatározásához a pszeudo-zárt elemhalmazokra van szükség, amelyek a nem zárt gyakori elemhalmazokból kerülnek ki. A pszeudo-zártság eldöntéséhez a definícióból indulunk ki: amennyiben  $I$  nem zárt gyakori termékhalmoznak létezik olyan részhalmozza, amely lezártja tartalmazza  $I$ -t, akkor  $I$  nem pszeudo-zárt elemhalmaz. Ellenkező esetben az. Jelöljük az  $i$ -elemű gyakori, illetve gyakori zárt halmazokat  $GY_i$  és  $ZGY_i$ -vel.

Az algoritmus menete a következő: Vegyük fel az üres halmazt a pszeudo-zártak közé, amennyiben az nem zárt. Ezután vizsgáljuk  $GY_1 \setminus ZGY_1$ ,  $GY_2 \setminus ZGY_2$ ,  $\dots$ ,  $GY_m \setminus ZGY_m$  halmazokat. Az  $I \in GY_i \setminus ZGY_i$  pszeudo-zártságának eldöntéséhez, az összes eddig megtalált kisebb elemszámú pszeudo-zárt elemhalmazra ellenőrizzük, hogy részhalmozza-e  $I$ -nek és ha igen akkor lezártja tartalmazza-e  $I$ -et. Amennyiben tehát létezik olyan  $I' \in FP_j$  ( $j < i$ ), amire fennáll, hogy  $I' \subset I$  és  $I \subseteq h(I')$ , akkor  $I$  nem pszeudo-zárt, ellenkező esetben igen. Ekkor  $I$  lezártja az  $I$ -t tartalmazó legkisebb zárt halmaz.

## 5.2. Érdekességi mutatók

Az asszociációs szabályok gyakorlati alkalmazása során az alábbi három súlyos probléma jelentkezett:

- I. Az asszociációs szabályok száma túl nagy. Ha magasra állítjuk a két küszöbszámot, akkor kevés szabály lesz érvényes, azonban ekkor számos – amúgy érdekes – szabály rejtve marad. Ellenkező esetben azonban rengeteg szabály jön létre, amelyek közül kézzel kiválogatni a fontosakat szinte lehetetlen feladat.
- II. Az asszociációs szabályok félrevezetőek lehetnek. Mivel az adatbányászat fontos stratégiai döntéseknek adhat alapot, félrevezető szabály rossz stratégiát eredményezhet. Fejtsük ki ezt egy kicsit bővebben. Egy asszociációs szabályra szoktak úgy tekinteni (helytelenül!!! lásd 5.6 rész), mint egy valószínűségi okozatiság viszonyra: adott termékhalmoz megvásárlása nagy valószínűséggel másik termékhalmoz megvásárlását „okozza”. Az okozatiság valószínűségét a szabály bizonyossága adja meg. Csak ennek az értékét vizsgálni azonban nem elég!

Képzeljünk el egy büfét, ahol az alábbiak teljesülnek. Az emberek egyharmada hamburgert vesz, egyharmada hot-dogot, egyharmada hamburgert és hot-dogot egyszerre. Azok és csak azok vesznek majonézt, akik hamburgert esznek. Ezek szerint a „kosarak” 66% tartalmaz hot-dogot és 50%-uk hot-dogot és majonézt is. Emiatt a hot-dog  $\rightarrow$  majonézt

érvényes asszociációs szabály lehet. Felhasználva az asszociációs szabályok bevezetésénél bemutatott trükköt, a hot-dogért felelős részleg vezetője (©) úgy dönt, hogy a nagyobb értékesítés reményében csökkenti a hot-dog árát és növeli a majonézét. A várakozásokkal ellentétben a profit csökkenni fog! Miért? Azért, mert a hamburger fogyasztók a hot-dog kedvező ára miatt inkább hot-dogot vesznek, aminek valójában semmi köze a majonézhez, azaz annak eladása nem fog nőni. Következtetésünk az, hogy egy asszociációs szabály *nem jelent okozatiságot*.

A példa jól szemlélteti, hogy a bizonyosság nem a legtökéletesebb mutató az összefüggések méréséhez. Gondoljunk arra, hogy egy szabály bizonyossága a következményrész feltételes valószínűségét próbálja becsülni, tehát  $I_1 \xrightarrow{c,s} I_2$  esetén  $c = p(I_2|I_1) = \frac{p(I_1, I_2)}{p(I_1)}$ . Amennyiben  $p(I_2|I_1)$  megegyezik  $p(I_2)$ -vel, akkor a szabály nem hordoz semmi többlet- hasznos információt (kivéve azt, hogy  $I_2$  az  $I_1$ -et tartalmazó kosarakban is ugyanolyan gyakori, mint általában. De ilyen szabály rengeteg van!).

- III. A legtöbb szabály nem érdekes. Pontosabban a szabályok nagy része bizonyos más szabályoknak semmitmondó speciális esetei, apró módosításai. Szükség lenne valahogy a szabályokat fontosságuk alapján sorba rendezni, vagy minden szabályhoz egy érdekességi mutatót rendelni.

A második problémára a függetlenségi mutató bevezetése lesz a megoldás. A harmadik problémának is köze van a függetlenséghez. Érdekes szabályt, ha „felhígítunk” egy kicsit független elemekkel, akkor még kaphatunk érdekes szabályt. A felhígított szabály azonban egy extra feltételt tartalmaz így feleslegesen speciálist. Többet ér egy általános szabály, mint sok speciális szabály felsorolása.

## 5.3. Szabályok függetlensége

Az összefüggőség mérésére számos mutatószámot vezettek be a kutatók.

### 5.3.1. lift érték

Egy szabály nem érdekes, ha a feltétel és a következményrészek függetlenek egymástól. Valószínűségi számításbeli ismereteinket felidézve: az  $X$  és az  $Y$  események függetlenek egymástól, ha  $P(X, Y) = P(X)P(Y)$ , azaz ha a  $\frac{P(X, Y)}{P(X)P(Y)}$  hányados értéke 1. Minél jobban eltér a hányados egytől, annál inkább összefüggők az események. Ez alapján egy szabály *lift* értékét, amely a függetlenséget szándékozik megragadni, a következőképpen definiáljuk:

$$\text{lift}(I \rightarrow I') = \frac{\text{freq}(I \cup I')}{\text{freq}(I) \cdot \text{freq}(I')},$$

ahol *freq* a gyakoriságot jelöli. Csendben feltételeztük, hogy a valószínűséget a relatív gyakorisággal közelíthetjük.

Ha ezek után egy adatbázisból a rejtett összefüggéseket asszociációs szabályok formájában akarjuk kinyerni, akkor a támogatottsági és bizonyossági küszöb mellett függetlenségi küszöböt (*min\_lift*) is megadhatunk. Például, ha  $\text{min\_lift} = 1.3$ , akkor azok a szabályok érdekesek, amelyekre  $\text{lift}(R) \geq 1.3$  vagy  $\text{lift}(R) \leq \frac{1}{1.3}$ .



Gyakori termékhalmból alkotott asszociációs szabály lift értékének meghatározásához minden adat rendelkezésünkre áll, így könnyedén megkaphatjuk az értékét.

A lift érték előnye, hogy könnyű értelmezni, még a matematika iránt kevésbé fogékonyak is megértik. Írjuk át a lift definícióját a következő alakra:  $\text{lift}(I \rightarrow I') = \frac{\text{freq}(I \cup I')}{\frac{\text{freq}(I) \cdot \text{freq}(I')}{\text{freq}(I \cup I')}} = \frac{\text{freq}(I \cup I')}{\text{freq}(I) \cdot \text{freq}(I')}$ . Ez az  $I'$  feltételes relatív gyakoriságának és az  $I$  relatív gyakoriságának a hányadosa. Ha például vásárlói szokások elemzésénél a sör  $\rightarrow$  pelenka szabály lift értéke 2, akkor a sört vásárlók körében a pelenkát vásárlók aránya dupla annyi, mint úgy általában a pelenkát vásárlók aránya.

### 5.3.2. Empirikus kovariancia, empirikus korreláció

A lift érték bevezetésénél használt logika alapján mondhatnánk, hogy két esemény akkor független, ha a  $P(X, Y)$  és a  $P(X)P(Y)$  szorzat különbsége 0. Minél jobban eltér a különbség nullától, annál nagyobb az összefüggés  $X$  és  $Y$  között. Legyen tehát a függetlenségi mutatónk

$$\text{cov}(I \rightarrow I') = \text{freq}(I \cup I') - \text{freq}(I) \cdot \text{freq}(I').$$

Relatív gyakoriságváltozás helyett abszolút gyakoriságváltozást használunk. De mi köze mindezen a címben említett empirikus kovarianciához? Egyáltalán, mi az az empirikus kovariancia?!?

Az  $X$  és  $Y$  valószínűségi változók kovarianciája  $\text{cov}(X, Y) = E[(X - \mu)(Y - \nu)] = E[X \cdot Y] - \mu \cdot \nu$ , ahol  $\mu$  és  $\nu$  az  $X$  és  $Y$  várható értékét jelöli. Könnyű belátni, hogy a kovariancia nulla, amennyiben  $X$  és  $Y$  függetlenek. Ha sűrűségfüggvényeket nem ismerjük, hanem csak megfigyelések  $(x_i, y_i)$ -k állnak rendelkezésünkre, akkor empirikus kovarianciáról beszélünk, amelynek definíciója:  $\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})(y_j - \bar{y})$ , ahol  $\bar{x}$  és  $\bar{y}$  a mintaátlagokat jelölik.

Az  $I$  és  $I'$  valószínűségi változók jelölhetik két termék megvételét. Az asszociációs szabályoknál bevezetett jelöléseket használva a mintaátlaga ekkor a gyakorisággal egyezik meg az  $i_j$  pedig 1, amennyiben a  $j$ -edik kosár tartalmazza az  $i$  terméket. Ekkor

$$\begin{aligned} \text{cov}(I \rightarrow I') &= \frac{1}{n} \sum_{j=1}^n (i_j - \text{freq}(I))(i'_j - \text{freq}(I')) \\ &= \frac{1}{n} \left( \sum_{j=1}^n i_j i'_j - \text{freq}(I) \sum_{j=1}^n i'_j - \text{freq}(I') \sum_{j=1}^n i_j + n \text{freq}(I) \text{freq}(I') \right) \\ &= \text{freq}(I \cup I') - \text{freq}(I) \text{freq}(I') - \text{freq}(I) \text{freq}(I') + \text{freq}(I) \text{freq}(I') \\ &= \text{freq}(I \cup I') - \text{freq}(I) \text{freq}(I'). \end{aligned}$$

A kovariancia normalizálásából adódik a korreláció:  $\text{corr}(X, Y) = \frac{\text{cov}(X, Y)}{\sigma_X \sigma_Y}$ . A korreláció értéke mindig -1 és 1 közé esik. Számítsuk ki egy asszociációs szabály empirikus korrelációját. Mivel egynek és nullának a négyzete egy és nulla, ezért  $\sigma_X^2 = E[X^2] - E^2[X] = E[X] - E^2[X]$ . Ebből

$$\begin{aligned} \text{corr}(I \rightarrow I') &= \frac{\text{Cov}(I \rightarrow I')}{\sigma_I \sigma_{I'}} = \frac{\text{freq}(I \cup I') - \text{freq}(I) \text{freq}(I')}{\sqrt{E[I](1 - E[I])} \cdot \sqrt{E[I'](1 - E[I'])}} \\ &= \frac{\text{freq}(I \cup I') - \text{freq}(I) \text{freq}(I')}{\sqrt{\text{freq}(I) \text{freq}(\bar{I}) \text{freq}(I') \text{freq}(\bar{I}')}}. \end{aligned}$$

„Ausztrál kutatók állítása szerint a sok stressz elhízáshoz vezet.”

Forrás: [http://www.hirtv.hu/eletmod/?article\\_hid=165457](http://www.hirtv.hu/eletmod/?article_hid=165457)

### 5.3.3. A $\chi^2$ -statisztika

Valójában a lift mutató nem ragadja meg kellőképpen a két esemény ( $X$  és  $Y$  előfordulása) statisztikai függetlenségét. Tudjuk, hogy az  $X$ ,  $Y$  események függetlenek, ha  $P(X)P(Y) = P(X, Y)$ , amelyet átírhatunk  $1 = P(Y|X)/P(X)$  alakra. A jobb oldal annyiban tér el a függetlenségi mutatótól, hogy abban a valószínűségek helyén relatív gyakoriságok szerepelnek. Pusztán a relatív gyakoriságok hányadosa nem elég jó mérték a függetlenség mérésére. Nézzünk például a következő két esetet. Első esetben négy tranzakció van,  $\text{supp}(I) = 2$ ,  $c = 0.5$ , amiből lift 1 adódik. A másodikban a tranzakciók száma négyezer,  $\text{supp}(I) = 1992$ ,  $c = 0.504$ , amiből lift 1.012 következik. Ha csak a függetlenségi mutatókat ismernénk, akkor azt a téves következtetést vonhatnánk le, hogy az első esetben a két esemény függetlenebb, mint a második esetben. Holott érezzük, hogy az első esetben olyan kevés a tranzakció, hogy abból nem tudunk függetlenségre vonatkozó következtetéseket levonni. Minél több tranzakció alapján állítjuk, hogy két elemhalmaz előfordulása összefüggésben van, annál jobban kizárjuk ezen állításunk véletlenségének (esetlegességének) esélyét.

A függetlenség mérésére a statisztikusok által alkalmazott eszköz az ún.  $\chi^2$  próbastatisztika. Az  $A_1, A_2, \dots, A_r$  és  $B_1, B_2, \dots, B_s$  két teljes eseményrendszer  $\chi^2$  próbastatisztikáját az alábbi képlet adja meg:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^s \frac{\left(k_{ij} - \frac{k_i \cdot k_j}{n}\right)^2}{\frac{k_i \cdot k_j}{n}}$$

ahol  $k_{ij}$  az  $A_i \cap B_j$  esemény,  $k_i = \sum_{j=1}^s k_{ij}$  az  $A_i$  esemény és  $k_j = \sum_{i=1}^r k_{ij}$  a  $B_j$  esemény bekövetkezésének számát jelöli. Minél kisebb a próbastatisztika, annál inkább függetlenek az események. A jelölést megjegyzését segítő kétszer kettős kontingenciátáblát a következő ábra mutatja.

	$X$	nem $X$	$\Sigma$
$Y$	$k_{1,1}$	$k_{1,2}$	$k_1$
nem $Y$	$k_{2,1}$	$k_{2,2}$	$k_2$
$\Sigma$	$k_1$	$k_2$	$n$

A mi esetünkben az egyik eseményrendszer az  $I$  elemhalmaz a másik az  $I'$  elemhalmaz előfordulásához tartozik, és mindkét eseményrendszernek két eseménye van<sup>1</sup> (előfordul az elemhalmaz az adott tranzakcióban, vagy sem). A következő táblázat mutatja, hogy a  $\chi^2$  próbastatisztika kiszámításához szükséges értékek közül melyek állnak rendelkezésünkre támogatottság formájában.

	$I$	nem $I$	$\Sigma$
$I'$	$\text{supp}(I \cup I')$		$\text{supp}(I')$
nem $I'$			
$\Sigma$	$\text{supp}(I)$		$ \mathcal{T} $

<sup>1</sup>Amennyiben mindkét eseményrendszer két eseményből áll, akkor az eredeti képletet módosítani szokás a Yates-féle korrekciós együtthatóval, azaz  $\chi^2 = \sum_{i=1}^2 \sum_{j=1}^2 \left( \left| k_{ij} - \frac{k_i \cdot k_j}{n} \right| - \frac{1}{2} \right)^2 / \frac{k_i \cdot k_j}{n}$ .

A hiányzó értékeket a táblázat ismert értékei alapján könnyen pótolni, hiszen például  $k_{2,1} = \text{supp}(I) - \text{supp}(I \cup I')$ .

A  $\chi^2$  próbastatisztika helyett használhatjuk mutatószámunk a próba  $p$ -értékét. A  $p$ -érték megegyezik azzal a legnagyobb próbaszinttel, amely mellett a hipotézisünket (függetlenség) elfogadjuk.

A  $\chi^2$ -próba közelítésen alapul ezért akkor működik jól, ha a kontingencia táblázat elemei nagyok. Kétszer kettes táblázat esetében az ökölszabály az, hogy mind a négy elem nagyobb legyen 10-nél.

Mielőtt teljes elégedettségben hátradölnénk a karosszékünkben, mert találtunk egy tudományosan megalapozott módszert, olvassuk el a következőket.

**5.8. állítás.** *Kétszer kettes kontingenciatáblák esetében a  $\chi^2$  próbastatisztika értéke megegyezik az empirikus korreláció négyzetének  $n$ -szeresével, ahol  $n$ -nel a minták számát jelöljük.*

*Bizonyítás:* Írjuk fel a  $\chi^2$  próbastatisztika értékét kétszer kettes kontingenciatáblák esetére:

$$\begin{aligned} \chi^2 &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\left(k_{ij} - \frac{k_{i \cdot} k_{\cdot j}}{n}\right)^2}{\frac{k_{i \cdot} k_{\cdot j}}{n}} = \sum_{i=1}^2 \sum_{j=1}^2 \frac{\left(k_{ij} - \frac{(k_{i1} + k_{i2})(k_{1j} + k_{2j})}{k_{11} + k_{12} + k_{21} + k_{22}}\right)^2}{\frac{(k_{i1} + k_{i2})(k_{1j} + k_{2j})}{n}} \\ &= \sum_{i=1}^2 \sum_{j=1}^2 \frac{\frac{(k_{11}k_{22} - k_{12}k_{21})^2}{n^2}}{\frac{(k_{i1} + k_{i2})(k_{1j} + k_{2j})}{n}} = \frac{(k_{11}k_{22} - k_{12}k_{21})^2}{n} \sum_{i=1}^2 \sum_{j=1}^2 \frac{1}{(k_{i1} + k_{i2})(k_{1j} + k_{2j})} \\ &= \frac{(k_{11}k_{22} - k_{12}k_{21})^2}{n} \cdot \left( \frac{1}{k_{11} + k_{12}} \left( \frac{1}{k_{11} + k_{13}} + \frac{1}{k_{12} + k_{22}} \right) + \frac{1}{k_{21} + k_{22}} \left( \frac{1}{k_{11} + k_{13}} + \frac{1}{k_{12} + k_{22}} \right) \right) \\ &= \frac{(k_{11}k_{22} - k_{12}k_{21})^2}{n} \cdot \left( \left( \frac{1}{k_{11} + k_{12}} + \frac{1}{k_{21} + k_{22}} \right) \left( \frac{1}{k_{11} + k_{13}} + \frac{1}{k_{12} + k_{22}} \right) \right) \\ &= \frac{(k_{11}k_{22} - k_{12}k_{21})^2}{n} \cdot \left( \frac{k_{11} + k_{12} + k_{21} + k_{22}}{(k_{11} + k_{12})(k_{21} + k_{22})} \cdot \frac{k_{11} + k_{12} + k_{21} + k_{22}}{(k_{11} + k_{21})(k_{12} + k_{22})} \right) \\ &= \frac{n(k_{11}k_{22} - k_{12}k_{21})^2}{k_{1 \cdot} k_{2 \cdot} k_{1 \cdot} k_{2 \cdot}} = \frac{n^3 \left(k_{11} - \frac{k_{1 \cdot} k_{\cdot 1}}{n}\right)^2}{k_{1 \cdot} k_{2 \cdot} k_{1 \cdot} k_{2 \cdot}} = \frac{n(f_{11} - f_{1 \cdot} f_{\cdot 1})^2}{f_{1 \cdot} f_{2 \cdot} f_{\cdot 1} f_{\cdot 2}}, \end{aligned}$$

ahol  $f_{ij} = k_{ij}/n$ . A bizonyítás során többször felhasználtuk, hogy  $n = k_{11} + k_{12} + k_{21} + k_{22}$ . ■

Ha a  $\chi^2$  próbastatisztika csak egy megbonyolított korreláció, amely pedig egy normalizált kovariancia, a kovariancia pedig a lift érték „testvére”, akkor most miért is mond többet a  $\chi$ -próbastatisztika a lift értéknél?

Egyrészt, az eredményként egy eloszlásfüggvényt kapunk, nem csak egy számot. Ez olyan, mint amikor megkérdezzük az útvonaltervező programtól, hogy mennyi időbe fog telni, hogy eljussunk  $A$  pontból  $B$ -be. Egy kezdetleges program egy konkrét számot adna eredményül. A valóságban azonban a helyes válasz egy eloszlásfüggvény, amelynek meghatározhatjuk például a várható értékét és a szórását. A szórás, amely a bizonytalanságra utal, szintén fontos paraméter.

Másrésről, mert figyelembe veszi az adatbázis méretét. Nem nekünk kell meghatároznunk egy jó lift értéket, amely adatbázisonként más lesz, hanem csak a próba szintjét kell megadnunk és máris szűrhetjük ki azokat a szabályokat, amelyek feltétel- és következményrésze között nincs szignifikáns kapcsolat. Olyan, mintha a szűrésre használt küszöböt is automatikusan állítanánk elő.

### 5.3.4. A binomiális próba

A  $\chi^2$ -próba és az ebből adódó  $p$ -érték nem használható, ha a  $2 \times 2$ -es kontingenciátáblázat valamely eleme kisebb, mint 10. Hagyjuk a statisztikát és térjünk vissza az elemi valószínűségszámításhoz.

Induljunk ki abból, hogy az  $I$  és az  $I'$  termékek függetlenek egymástól, azaz  $P(I, I') = P(I)P(I')$ . Legyen  $Z_j = I_j \cdot I'_j$ , azaz  $Z_j = 1$ , amennyiben a  $j$ -edik kosárban előfordul az  $I$  és az  $I'$  termék is. A  $Z = \sum_{j=1}^n Z_j$  binomiális eloszlású valószínűségi változó  $n$  és  $P(I, I')$  paraméterekkel. A  $P(I, I')$  valószínűséget a  $freq(I)freq(I')$  értékkel közelítjük.

Azt kell eldöntenünk, hogy a megfigyeléseink  $(z_1, \dots, z_n)$  ellentmondanak-e a kiindulási feltételből kapott következtetésnek. Jelöljük a próba szintjét  $1 - \alpha$ -val és legyen  $z = \sum_{j=1}^n z_j$ . Határozzuk meg azt a legszűkebb  $[l, u]$  intervallumot, amelyre igaz, hogy  $\sum_{k=l}^u P(Z = k) \leq 1 - \alpha$ . Amennyiben  $z$  a  $[l, u]$  intervallumba esik, akkor  $X$  és  $Y$  (tehát az  $I$  és  $I'$  termékhalmozok) függetlenek egymástól.

Ha ezt a megközelítést használjuk egy asszociációs szabály függetlenségének megadására, akkor legyen a függetlenségi mutató a szabály  $p$ -értéke. Határozzuk meg azt az  $[l', u']$  intervallumot, amelynek minden  $k$  elemére igaz, hogy  $P(Z=k) > P(Z=z)$ . A  $p$ -érték ekkor  $\sum_{k=l'}^{u'} P(Z=k)$ .

### 5.3.5. Fisher-féle egzakt próba

A binomiális próba a  $P(I, I')$  valószínűséget a  $freq(I)freq(I')$  relatív gyakoriságával közelíti. A közelítés pontatlansághoz vezet. Gondoljuk meg, hogy a binomiális eloszlás nemnulla valószínűséget fog rendelni az  $n$ -nél kisebb,  $\min\{supp(I), supp(I')\}$ -nél nagyobb értékekhez. Azonban ezeknek a valószínűségeknek nullának kellene lenniük. Nem fordulhat az elő, hogy az  $I$ -nél nagyobb,  $I$ -t részalmazként tartalmazó halmaznak  $supp(I)$ -nél nagyobb legyen a támogatottsága. Hasonló mondható el az  $n - supp(I) - supp(I')$  értékekre, amennyiben  $n - supp(I) - supp(I') > 0$ . A Fisher-féle egzakt próba a közelítés helyett a pontos valószínűségeket használja.

Tegyük fel, hogy a kontingenciátáblázat ún. marginális értékei  $(k_{1.}, k_{2.}, k_{.1}, k_{.2})$  és így a minták száma is adva vannak. Ez az asszociációs szabályoknál azt jelenti, hogy a kosarak száma,  $supp(I) = k_{1.}$  és  $supp(I') = k_{.1}$  rögzítettek. A kérdés a következő: Ha tudjuk, hogy a  $k_{1.}$  darab  $I$  termék és a  $k_{.1}$  darab  $I'$  termék egyenletes eloszlás szerint véletlenszerűen van szétszórva az  $n$  kosárban, akkor mennyi az esélye annak, hogy az  $I'$ -t tartalmazó kosarakból  $X$  darabban lesz  $I$ . Elvonatkoztatva a részletektől ez ugyanaz a kérdés, mint amelyet a hipergeometrikus eloszlás bemutatásakor tettünk fel (lásd a 2.5.1 rész). Ezek szerint

$$P(X, n, k_{1.}, k_{.1}) = \frac{\binom{k_{1.}}{X} \binom{k_{2.}}{k_{.1}-X}}{\binom{n}{k_{.1}}}.$$

Ez a valószínűség már önmagában egy jó mutatószám. Minél nagyobb az értéke, annál függetlenebbek az  $I$  és az  $I'$  termékek. Ha a  $\chi^2$  statisztikához hasonló  $p$ -értéket szeretnénk

kapni, akkor ki kell számolni az összes olyan  $X'$ -re a  $P(X', n, k_1, k_1)$  valószínűséget, amelyre  $P(X', n, k_1, k_1) \leq P(X, n, k_1, k_1)$ . Ezeket az  $X'$  értékeket hívjuk *extrémebb*, azaz kisebb valószínűségű értékeknek. A  $p$ -érték ezen extrém értékhez rendelt valószínűségek összege Formálisan:

$$p_{\text{Fisher}}(I \rightarrow I') = \sum_{X': P(X', n, \text{supp}(I), \text{supp}(I')) \leq P(\text{supp}(I \cup I'), n, \text{supp}(I), \text{supp}(I'))} \mathbb{P}(X', n, \text{supp}(I), \text{supp}(I'))$$

A Fisher-próbát nem csak kis értékeknél használhatjuk, tulajdonképpen függetlenség eldöntésére ez a módszer mindig a legjobb eredményt adja. Hátránya, hogy nagy  $n, k_1, k_1$  értékeknél nehéz a valószínűségeket kiszámítani. Így jutunk el a  $\chi^2$  próbához. Amennyiben  $k_1 \ll N$ , akkor a hipergeometrikus eloszlást közelíthetjük az  $k_1, k_1/n$  paraméterű binomiális eloszlással. A binomiális eloszlást pedig a normális eloszlással közelíthetjük. Standard normális eloszlású valószínűségi változók négyzetének összege pedig olyan valószínűségi változót ad, amelynek eloszlása a  $\chi^2$  eloszlás. Tyű, a mindenit, de szép ez az egész!

### Értékinvariancia

Egy függetlenségi mutatót értékinvariánsnak hívunk, amennyiben a kontingencia-táblázat tetszőleges sorait vagy oszlopait felcserélve ugyanazt a kimenetet ( $p$ -értéket) kapjuk. Bináris esetre gondolva ez azt jelenti, hogy  $X$  és  $Y$  függetlensége esetén,  $X$  és  $\bar{Y}$  (továbbá  $\bar{X}, Y$  és  $\bar{X}, \bar{Y}$ ) is az. Ha például megállapítjuk, hogy a tejkvásárlás és kenyérvásárlás függetlenek egymástól, akkor tejkvásárlás, nem kenyérvásárlás is függetlenek.

Könnyű belátni, hogy a Fisher-féle egzakt próba és a  $\chi^2$  próba megfelel a fenti elvárásnak, de a binomiális próba nem. A Fisher-féle egzakt próbához csak azt kell meggondolnunk, hogy

$$P(X, n, k_1, k_1) = \frac{\binom{k_1}{X} \binom{k_2}{k_1 - X}}{\binom{n}{k_1}} = \frac{\binom{n - k_1}{k_1 - X} \binom{n - k_2}{X}}{\binom{n}{k_1}} = P(k_1 - X, n, n - k_1, k_1),$$

tehát attól, hogy a két sort (vagy a két oszlopot) felcseréljük még ugyanazt a hipergeometrikus eloszlást kapjuk. A  $\chi^2$  próbára vonatkozó állítás közvetlen adódik a  $\chi^2$  statisztika definíciójából.

A binomiális próba esetét egy példával vizsgáljuk. Induljunk ki a bal oldali kontingenciátáblából majd cseréljük fel a két sorát.

	$X$	nem $X$	$\Sigma$
$Y$	2	0	2
nem $Y$	0	1	1
$\Sigma$	2	1	3

	$X$	nem $X$	$\Sigma$
$Y$	0	1	1
nem $Y$	2	0	2
$\Sigma$	2	1	3

A bal oldali kontingenciátáblához  $(3, 4/9)$  paraméterű binomiális eloszlás tartozik. A kettőhöz nagyobb valószínűség tartozik, mint a nullához és a háromhoz, ezért a  $p$ -érték  $1 - 3 \cdot \frac{4}{9} \cdot \frac{5^2}{9^2} = 0.588$ . A jobb oldali kontingenciátábla binomiális eloszlásához tartozó valószínűség  $2/9$ . A legnagyobb valószínűséget a  $(3, 2/9)$  paraméterű binomiális eloszlás nullánál veszi fel a maximumát ezért a  $p$ -érték egy.

## Érdekesség

Most, hogy tudjuk hogyan kell függetlenséget meghatározni, feltehetjük azt a kérdést, hogy legalább hány megfigyelésnek kell rendelkezésünkre állnia ahhoz, hogy összefüggést állapítsunk meg.

Adott  $1-\alpha$  próbaszint mellett csak akkor tudunk összefüggést megállapítani (függetlenséget elutasítani), ha az elfogadási tartományon kívül van olyan pont, amelyet felvehet azoknak a megfigyeléseknek a száma, amelyre mindkét vizsgált tulajdonság fenáll. Az elfogadási tartományba a legnagyobb valószínűséggel rendelkező pontok esnek. Amennyiben a legkisebb valószínűségű pont valószínűsége kisebb  $\alpha$ -nál, akkor ez a pont nem esik az elfogadási tartományba. Kétoldali próbánál két legkisebb valószínűségi pont is lehet, így ezen valószínűségek összege kell  $\alpha$ -nál kisebbnek lennie. Ha  $n$  páratlan, akkor csak egy legkisebb valószínűségi pont lehet, éljünk ezért ezzel a feltétellel. Az általánosság megsértése nélkül feltehetjük, hogy  $k_1 \leq k_{.1}$  és a hipergeometrikus eloszlás módusza ( $\lfloor \frac{(k_1+1)(k_{.1}+1)}{n+2} \rfloor$ ) nem nagyobb, mint az értelmezése tartomány ( $[\max(0, n - k_1 - k_{.1}), \min(k_1, k_{.1})]$ ) felezőpontja. A legkisebb valószínűségi pont ekkor a  $k_1$ , amelynek valószínűsége

$$\begin{aligned} P(k_1, n, k_1, k_{.1}) &= \frac{\binom{k_1}{k_1} \binom{k_{.1}}{k_1 - k_1}}{\binom{n}{k_1}} = \frac{(n-k_1)!}{(k_1-k_1)!(n-k_1)!} \\ &= \frac{(n-k_1)(n-k_1-1) \cdots (k_1-k_1+1)}{n(n-1) \cdots (k_1+1)} \\ &= \prod_{i=0}^{n-k_1+1} \left(1 - \frac{k_1}{n-i}\right) \end{aligned}$$

A fenti valószínűség rögzített  $n$  esetén akkor lesz a legnagyobb, ha  $k_1$  minél nagyobb, tehát  $k_1 = k_{.1}$ . Ekkor viszont

$$P(k_1, n, k_1, k_{.1}) = \frac{1}{\binom{n}{k_1}},$$

amely  $k_1 = \lfloor n/2 \rfloor$ -nél és  $k_1 = \lceil n/2 \rceil$  veszi fel a minimumát. Az 5.1 táblázat második oszlopa megadja a legkisebb valószínűséget néhány  $n$ -re Ezek szerint 97%-os bizonyossággal már

$n$	$P(\lfloor n/2 \rfloor, n, \lfloor n/2 \rfloor, \lfloor n/2 \rfloor)$	$p$ -érték	
		binom	$\chi^2$
3	33.33%	29.76 %	8.32%
5	10.00%	18.34 %	2.53%
7	2.85%	12.11 %	0.81%
9	0.79%	8.24 %	0.27%
11	0.21%	5.70 %	0.09%
13	0.06%	4.00 %	0.03%
15	0.02%	2.82 %	0.01%

5.1. táblázat.  $p$ -értékek extrém kontingencai-táblázat esetén

hét megfigyelésből összefüggőséget állapíthatunk meg. Ehhez a legextrémebb eseménynek kell

bekövetkeznie, nevezetesen, 7 megfigyelésből háromra teljesül egy tulajdonság ( $X$ ) és csak erre a három megfigyelésre egy másik tulajdonság ( $Y$ ) is teljesül. Tehát a kontingenciatáblázat:

	$X$	nem $X$	$\Sigma$
$Y$	3	0	3
nem $Y$	0	4	4
$\Sigma$	3	4	7

A  $P(\lfloor n/2 \rfloor, n, \lfloor n/2 \rfloor, \lfloor n/2 \rfloor)$  érték egyben annak a tesztnek a  $p$ -próbája, amelyben a megfigyelések száma  $n$  és  $k_{11} = k_{1.} = k_{.1} = \lfloor n/2 \rfloor$ . Ha a próba szintje ennél az értéknél nagyobb, akkor elutasítjuk a függetlenségre tett feltételt, ellenkező esetben elfogadjuk. A függetlenség eldöntésére használhatnánk más próbát is. Az 5.1 táblázat harmadik és negyedik oszlopa a megfigyelés  $p$ -értékét adja meg binomiális és  $\chi^2$  próba esetén. Láthatjuk, hogy a binomiális próba jóval nagyobb  $p$ -értékeket ad ugyanarra a megfigyelésre, azaz a binomiális próba „függetlenség felé húz”. Például  $n = 11$  és  $\alpha = 5\%$  estén a Fisher próba elutasítja a függetlenséget a binomiális próba pedig elfogadja azt.

Ha megszorítkozunk olyan kontingenciatáblákra, amelyeknél  $k_{1,1} = k_{1.} - 1$ , tehát nem a legextrémebb eset következik be, akkor a Fisher-féle próba  $p$ -értékei a következőképpen alakulnak:

$n$	$p$ -érték					
	$k_{1,1} = k_{1.} - 1$			$k_{1,1} = k_{1.} - 2$		
	fisher	binom	$\chi^2$	fisher	binom	$\chi^2$
5	100%	58.17%	70.9%	40%	100%	13.6%
7	48.57%	61.95%	27.0%	100%	100%	65.9%
9	20.62%	39.33%	9.89%	100%	69.4%	76.4%
11	8.00%	25.45%	3.56%	56.71%	70.75%	37.6%
13	2.91%	16.75%	1.27%	28.61%	49.42%	16.9%
15	1.01%	11.19%	0.45%	13.19%	34.30%	7.21%
17	0.35%	7.59%	0.16%	5.67%	23.74%	2.95%
19	0.11%	5.21%	0.05%	2.30%	16.44%	1.17%

A 97%-os bizonyosság megtartásához most már 13 megfigyelés kell (binomiális próba szerint 21,  $\chi^2$ -próba szerint is 13). Ha hat-hat megfigyelésnél teljesül az  $X$ ,  $Y$  tulajdonságok, akkor abban az esetben állapítunk meg összefüggést, ha az  $X$ ,  $Y$  tulajdonsággal együtt rendelkező megfigyelések száma 0, 5 vagy 6.

### 5.3.6. További mutatószámok

A lift,  $\chi$ -statisztika, vagy  $p$ -érték mellett még számos elterjedt mutatószám létezik függetlenség mérésére. A teljesség igénye nélkül felsorolunk néhányat

név	jelölés	képlet	megjegyzés
empirikus kovariancia	$\phi$	$freq(I \cup I') - freq(I)freq(I')$	Az általános képlet átírásából adódik, felhasználva, hogy $\bar{I} = freq(I)$ és $\sum_{j=1}^n I_j = supp(I)$
empirikus korreláció		$\frac{freq(I \cup I') - freq(I)freq(I')}{\sqrt{freq(I)freq(\bar{I})}\sqrt{freq(I')freq(\bar{I}')}}}$	Az általános képlet átírásából adódik, a fentiek mellett felhasználva, hogy $I_j^2 = I_j$ .
esélyhányados	$\alpha$	$\frac{freq(I \cup I') \cdot freq(\bar{I}, I')}{freq(I, \bar{I}') \cdot freq(\bar{I}, I)}$	odds ratio, cross-product ratio
Yule féle Q érték	Q	$\frac{\alpha - 1}{\alpha + 1}$	
Yule féle Y érték	Y	$\frac{\sqrt{\alpha - 1}}{\sqrt{\alpha + 1}}$	<i>measure of colligation</i>
conviction	V	$\frac{freq(I)freq(\bar{I}')}{freq(I, I')}$	az $I \rightarrow I'$ implikáció logikai megfelelője alapján definiálják.
conviction*	V*	$\max\{V(I, I'), V(I', I)\}$	
Jaccard-koefficiens	$\varsigma$	$\frac{freq(I \cup I')}{freq(I) + freq(I') - freq(I \cup I')}$	
koszinusz mérték	cos	$\arccos\left(\frac{freq(I \cup I')}{\sqrt{freq(I)freq(I')}}\right)$	
normált kölcsönös entrópia	$H^n$	$\frac{H(I' I)}{H(I)}$	

### 5.3.7. Asszociációs szabályok rangsora

Az asszociációs szabályok kinyerésének feladatában adott bemeneti sorozat és küszönszámok mellett célunk volt meghatározni az asszociációs szabályokat. Ennyi. Aztán mindenki kezdjen a szabályokkal, amit akar.

A gyakorlatban általában sok érvényes asszociációs szabályt találunk, hasznos lenne őket sorba rendezni. Ha a három paraméterhez (támogatottság/gyakoriság, bizonyosság, függetlenségi mutató) tudnánk súlyt rendelni fontosságuk szerint, akkor az alapján sorrendet tudnánk felállítani. A marketinges a támogatottságot részesítené előnyben a statisztikus a függetlenségi mutatót. Elvégre kit érdekel a két termékhalmaz támogatottsága, ha a két termékhalmaz független egymástól.

Függetlenség kifejezésére több mutatószámot adtunk meg: lift érték, empirikus kovariancia, empirikus korreláció,  $\chi^2$ -statisztika, p-érték. Ráadásul  $\chi^2$ -statisztika helyett használhatunk hipergeometrikus (vagy binomiális) eloszlás alapján definiált p-értéket is. Matematikusokban azonban felmerül a kérdés, hogy ugyanazt a sorrendet adják-e az egyes függetlenségi mutatók.

A  $\chi^2$ -statisztika és az ebből származtatott p-érték ugyanazt a sorrendet fogja adni, hiszen a p-érték a  $\chi^2$ -statisztika szigorúan monoton függvénye. A  $\chi^2$ -statisztika és az empirikus korreláció között teremt szigorúan monoton kapcsolatot az 5.8 állítás.

Az empirikus korreláció és az empirikus kovariancia adhat különböző sorrendet. A korreláció a kovariancia normált változata. Ha két asszociációs szabály közül az elsőnek nagyobb a kovarianciája, attól még lehet kisebb a korrelációja, amennyiben az első szabályhoz tartozó két



binomiális valószínűségi változó szórásának szorzata, mint a második szabályhoz tartozó két változó szórásának szorzata.

A lift érték és az empirikus kovariancia között nincs monoton kapcsolat, azaz a két mutató alapján különböző sorrendet kaphatunk. Ehhez csak azt kell meggondolnunk, hogy  $a, b, c, d$  nulla és egy közötti számokra sem igaz, hogy

$$\frac{a}{b} < \frac{c}{d} \not\Rightarrow a - b < c - d.$$




---

**Weka 3.5.7** *Asszociációs szabályokat a `weka.associations`-  
Apriori osztály segítségével nyerhetünk ki. Az osztály nem a klasszikus asszociációs szabály kinyerésének feladatát oldja meg – adott `min_supp`, `min_conf`, `min_lift` mellett határozzuk meg az érvényes asszociációs szabályokat – hanem csak a legjobb `numRules` darab szabályt adja meg, ahol `numRules` a felhasználó által megadott paraméter. Ehhez a `min_supp` értéket egy kiindulási értékről (`upperBoundMinSupport` paraméter) mindig  $\delta$  értékkel csökkenti és ellenőrzi, hogy van-e legalább `numRules` darab érvényes szabály. Ha van, akkor kiírja a legjobb `numRules` szabályt, ha nincs, akkor tovább csökkenti `min_supp`-ot. A minimális támogatottsági küszöböt nem csökkenti annyira, hogy az kisebb legyen a `lowerBoundMinSupport` paraméternél.*



A `metricType` paraméterrel adhatja meg a felhasználó, hogy mi alapján rangsorolja az asszociációs szabályokat a `weka`. Az empirikus kovarianciát a `Leverage` jelöli. Javasoljuk, hogy a `Conviction` mutatót sose használjuk; ez tulajdonképpen csak egy elbaltázott függetlenségi mutató.

Lehetőségünk van egy osztályattribútumot kijelölni a `car` paraméter igazra állításával és a `classIndex` megadásával. Ekkor csak olyan szabályokat fog a `weka` előállítani, amelyek következményrészében csak az osztályattribútum szerepel.

---

## 5.4. Általánosság, specialitás

A lift mutató gyengéje, hogy ha találunk egy érdekes szabályt, akkor „az mögé elbújva” sok érdektelen szabály átmegy a szűrésen, azaz érdekesnek bizonyul. Szemléltetésképpen nézzünk egy példát. Legyen az  $I_1 \rightarrow I_2$  érvényes és érdekes asszociációs szabály, továbbá  $I_3$  egy olyan gyakori termékhalmoz, amely független  $I_1$  és  $I_2$ -től ( $supp(I_1 \cup I_3) = freq(I_1) \cdot freq(I_3)$ ,  $freq(I_2 \cup I_3) = freq(I_2) \cdot freq(I_3)$ ) és támogatottsága olyan nagy, hogy még a  $supp(I_1 \cup I_2 \cup I_3) \geq min\_supp$  egyenlőtlenség is fennáll. Könnyű belátni, hogy ekkor az  $I_1 I_3 \rightarrow I_2$  is érvényes és

érdekes asszociációs szabályok, hiszen

$$\begin{aligned} \text{lift}(I_1 I_3 \rightarrow I_2) &= \frac{\text{supp}(I_1 \cup I_2 \cup I_3)}{\text{supp}(I_1 \cup I_3) \text{supp}(I_2)} = \frac{\text{supp}(I_1 \cup I_2) \text{supp}(I_3)}{\text{supp}(I_1) \text{supp}(I_2) \text{supp}(I_3)} = \\ &= \text{lift}(I_1 \rightarrow I_2) \geq \text{min\_lift}, \end{aligned}$$

$$\text{conf}(I_1 I_3 \rightarrow I_2) = \frac{\text{supp}(I_1 \cup I_2 \cup I_3)}{\text{supp}(I_1 \cup I_3)} = \frac{\text{supp}(I_1 \cup I_2) \text{supp}(I_3)}{\text{supp}(I_1) \text{supp}(I_3)} \geq \text{min\_conf}$$

Ezek alapján, egy adatbázisból kinyert érdekes asszociációs szabályok között a többség haszontalan, amennyiben sok a nagy támogatottságú, más termékektől független termék. Ha a valóságban  $n$  darab érdekes szabályunk van, de az adatbázis tartalmaz  $c$  darab a fenti tulajdonsággal rendelkező gyakori elemet, akkor az érdekességi mutató alapú szűrésen  $n2^c$  szabály fog átcsúszni a fenti módon.

A fenti problémát kiküszöbölhetjük, ha a feltételrész minden elemét megnézzük függetlenül a feltételrész többi elemének uniójától. Ha független, akkor dobjuk ki az elemet, csak feleslegesen bonyolítja életünket. Sőt, az egész szabályt kidobhatjuk. Az eredményként kapott szabály ugyanis ott kell legyen az érvényes szabályok között, hiszen a független elem törlése esetén a függetlenségi mutató és a bizonyosság nem változik a támogatottság pedig nő.

## 5.5. Asszociációs szabályok általánosítása

Számos általánosítást találtak ki a kutatók az asszociációs szabályoknak. Ebben a részben ezekből szemezgetünk.

### 5.5.1. Hierarchikus asszociációs szabályok

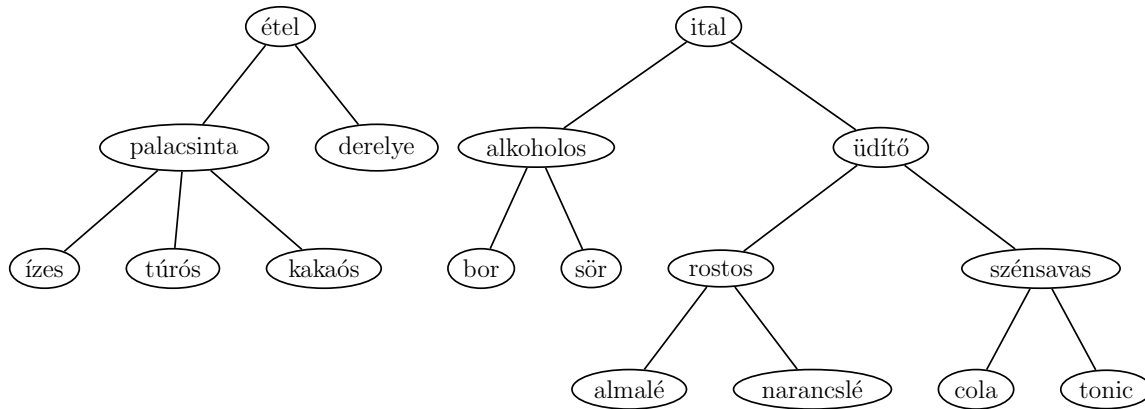
A hierarchikus asszociációs szabályok kinyerése a gyakorlatban tényleg előkerülő elvárásként jelentkezett [43, 45, 53, 125, 129, 135]. Vásárlási szokások elemzése közben a marketingesek új igénnyel álltak elő. Olyan szabályokat is ki szerettek volna nyerni, amelyek termékkategóriák között mondanak ki összefüggéseket. Például a sört vásárlók 70%-ban valami chips félét is vesznek. Lehet, hogy egyetlen sör és chips közötti asszociációs szabályt nem nyerünk ki, amennyiben sokfajta sör és chips létezik, ugyanis ezen termékek között a támogatottság „elaprózódik”. Például a sör  $\rightarrow$  chips támogatottsága lehet 5000, de ha ötféle sör létezik, akkor termék szinten könnyen lehet, hogy mindegyik, sört tartalmazó, asszociációs szabály támogatottsága 1500 alatt lesz és nem lesz érvényes.

Egy üzletnek a kategória szintű asszociációs szabályok legalább annyira fontosak lehetnek, mint a termékeken értelmezett szabályok (pl.: akciót hirdetünk: '17"-os monitorok óriási árengedményekkel', miközben más számítástechnikai alkatrészek – például monitorvezérlő kártya – árait megemeljük).

Ahhoz, hogy kategóriák is szerepelhessenek asszociációs szabályokban, ismernünk kell az elemek kategóriákba, a kategóriák alkategóriákba sorolását, azaz ismernünk kell az elemek

„ A rendszeres alkoholfogyasztás növeli az emlékezőtehetséget, különösen a nőknél – jelentette be a Sunday Telegraph a londoni University College kutatói információira hivatkozva.” Forrás: <http://www.origo.hu/nagyvilag/20040805azalkohol.html>

*taxonómiáját*, közgazdász nyelven szólva az elemek *nomenklatúráját*. A termék-taxonómia nem más, mint egy gyökeres címkézett fa, vagy fák sorozata. A fa leveleiben található az egyes termékek, a belső csomópontokban pedig a kategóriák. Egy képzeletbeli büfé termék-taxonómiája az alábbi ábrán látható.



5.1. ábra. Példa: képzeletbeli büfé termék-taxonómiája

Ha a kategóriák halmazát  $\hat{\mathcal{J}}$ -vel jelöljük, akkor a bemenet továbbra is az  $\mathcal{J}$  felett értelmezett sorozat, a mintatér elemei azonban  $\mathcal{J} \cup \hat{\mathcal{J}}$  részhalmazai lesznek. Azt mondjuk, hogy az  $I$  kosár tartalmazza  $I'$  elemhalmazt, ha minden  $i \in I'$ -re vagy  $i \in I$ , vagy  $\exists i' \in I$ , hogy  $i \in \text{ős}(i')$ <sup>2</sup>. Tehát egy kosár tartalmaz egy elemhalmazt, ha annak minden elemét, vagy annak leszármazottját tartalmazza. Nyilvánvaló, hogy ha a taxonómia egyetlen fenyőből áll, akkor a gyökérben található kategóriát minden nem üres kosár tartalmazza.

Hasonlóan módosítanunk kell az asszociációs szabályok definícióját, hiszen a 94. oldalon található definíció szerint minden  $X \xrightarrow{100\%,s} \hat{X}$  szabály érvényes lenne, ha  $\hat{X} \subseteq \text{ős}(X)$ , és  $X$  gyakori termékhalmoz.

**5.9. definíció (hierarchikus asszociációs szabály).** *Adott a termékek taxonómiája. A benne található termékeket és kategóriákat reprezentáló levelek, illetve belső csomópontok halmazát jelöljük  $\mathcal{J}$ -vel.  $I_1 \xrightarrow{c,s} I_2$ -t hierarchikus asszociációs szabálynak nevezzük, ha  $I_1, I_2$  diszjunkt részhalmazai  $\mathcal{J}$ -nek, továbbá egyetlen  $i \in I_2$  sem őse egyetlen  $i' \in I_1$ -nek sem.*

A támogatottság ( $s$ ), és bizonyosság ( $c$ ) definíciója megegyezik az 5.1. részben megadottal.

Hierarchikus asszociációs szabályok kinyerése csöppnyit sem bonyolultabb a hagyományos asszociációs szabályok kinyerésénél. Amikor a gyakori elemhalmazokat nyerjük ki (pl.: az APRIORI módszerrel), akkor képzeletben töltsük fel a kosarakat a kosarakban található elemek őseivel. Természetesen nem kell valóban előállítani egy olyan adatbázist, ami a feltöltött kosarakat tartalmazza, elég akkor előállítani ezt a kosarat, amikor a tartalmát vizsgáljuk.

Ha nem akarunk kinyerni olyan asszociációs szabályokat, amelyben bárhogy elosztva egy elem és őse is szerepel, akkor szükségtelen az is, hogy ilyen elemhalmazokkal foglalkozzunk. Ne állítsunk elő olyan jelöltet, amely ilyen tulajdonságú [129].

<sup>2</sup>Gyökeres gráfoknál definiálhatjuk a szülő, gyermek, ős, leszármazott fogalmakat. Ezt az alapfogalmak gráfelmélet részében megtettük.

A fentitől különböző megközelítést javasoltak a [45, 53]-ben. Az algoritmus azt az észrevételt használja ki, hogy ha egy tetszőleges kategória ritka, akkor annak minden leszármazottja is ritka. Éppen ezért, az adatbázis első végigolvasása során csak a fenyők gyökerében (első szinten) található kategóriák lesznek a jelöltek. A második végigolvasásnál a gyakorinak talált elemek gyerekei, a harmadik végigolvasásnál pedig a második olvasásból kikerült gyakori elemek gyerekei, és így tovább. Akkor nincs szükség további olvasásra, ha vagy egyetlen elem sem lett gyakori, vagy a jelöltek között csak levélelemek voltak.

A gyakori elempárok meghatározásához először ismét csak a gyökerekben található kategóriákat vizsgáljuk, természetesen csak azokat, amelyeknek mindkét eleme gyakori. A következő lépésben a pár egyik tagjának a második szinten kell lennie, és hasonlóan: az  $i$ -edik végigolvasásnál a jelöltpárosok egyik tagja  $i$ -edik szintbeli.

A fenti eljárást könnyű általánosítani gyakori elemhármassok és nagyobb méretű gyakori termékhalmazok megtalálására. A leállási feltétel hasonló az APRIORI algoritmuséhoz: ha a jelöltek közül senki sem gyakori, akkor minden gyakori hierarchikus termékhalmazt megtaláltunk. A továbbiakban az algoritmust nem tárgyaljuk, részletek és futási eredmények találhatóak [53] cikkben.

### Hierarchikus szabály „érdekessége”

Kategóriák bevezetésével az érvényes asszociációk száma nagymértékben nőhet. Ennek oka az, hogy a kategóriák támogatottsága mindig nagyobb, mint a bennük szereplő termékeké, így sokszor szerepelnek majd gyakori termékhalmazokban, amelyekből az érvényes asszociációs szabályokat kinyerjük. A szabályok között sok semmitmondó is lesz, amelyek csökkentik az áttekinthetőséget, és a tényleg fontos szabályok megtalálását. Egy ilyen semmitmondó szabályt az alábbi példa szemléltet:

Egy élelmiszerüzletben háromféle tejet lehet kapni: zsírszegényt, félzsírosat, és normált. Az emberek egynegyede zsírszegény tejet iszik. Hierarchikus szabályok kinyerése során többek között az alábbi két érvényes szabály is szerepel (a szabály harmadik paramétere a lift értéket adja):

$$\begin{aligned} \text{tej} &\xrightarrow{80\%,4.8\%,2} \text{zabpehely} \\ \text{zsírszegény tej} &\xrightarrow{80\%,1.2\%,2} \text{zabpehely} \end{aligned}$$

Látható, hogy a második szabály kevésbé általános, mint az első és nem hordoz semmi többletinformációt. Jogos tehát az a kijelentés, hogy egy szabály nem érdekes, ha annak bizonyossága és támogatottsága nem tér el a nála általánosabb szabály paramétere alapján becsült értékektől. A pontos definíciók magadásával nem terheljük az olvasót.

### 5.5.2. Kategória asszociációs szabályok

Az asszociációs szabályok kinyerésénél a bemenet elemhalmazok sorozataként van megadva (plussz néhány küszöbszám). Ábrázolhatjuk a bemenetet, mint bináris mátrix, amelynek az

„A kék szemű emberek hatékonyabban képesek tanulni, és jobban teljesítenek a vizsgákon, mint sötét szeműek - állapították meg amerikai kutatók.”  
 Forrás: <http://inforadio.hu/hir/tudomany/hir-143110>

$i$ -edik sor  $j$ -edik eleme egy, ha szerepel az  $i$ -edik tranzakcióban a  $j$ -edik elem, különben nulla. Tetszőleges bináris relációs táblát is választhatunk bemenetként, ekkor például  $nem = férfi \rightarrow tájékozódási\ készség = jó$  jellegű szabályokat nyerünk ki.

Könnyen kaphatunk kategória asszociációs szabályt a meglévő módszereinkkel. Minden olyan  $A$  attribútumot, amely  $k$  különböző értéket vehet fel ( $k > 2$ ), helyettesítsünk  $k$  darab bináris attribútummal. Ha egy elem  $A$  attribútumának értéke az  $i$ -edik attribútumérték volt, akkor csak  $i$ -edik új attribútum értéke legyen egy, a többi pedig nulla. Az így kapott bináris táblán már futtathatjuk a kedvenc asszociációs szabályokat kinyerő algoritmusunkat.




---

**Weka 3.5.7** *A weka.associations.Apriori algoritmus a fenti konverziót automatikusan elvégzi, ha kategória típusú attribútummal találkozik. Ha kézzel szeretnénk mindezt megtenni, akkor használhatjuk a weka.filters.unsupervised.attribute.NominalToBinary szűrőt.*

---



## 5.6. A korreláció nem jelent implikációt

A támogatottság és a fontosabb érdekességi mutatók ( $\chi^2$ -próbat statisztika,  $p$ -érték) szimmetrikus függvények, nem veszik figyelembe, hogy melyik termékhalmoz szerepel a szabály feltételrészében és melyik a szabály következményrészében. A bizonyosság az egyetlen aszimmetrikus függvény, amely meghatározza a szabály irányát. Az asszociációs szabályokban a nyilat használjuk az irány kijelölésére. Ez nagyon rossz döntés volt és rengeteg hamis következtetésnek adott alapot.

Ha megvizsgáljuk az asszociációs szabályok három paraméterét, akkor rájöhethetünk, hogy egyik paraméter sem jelent okozatiságot. A függetlenségi paraméter csak azt mondja meg, hogy a feltételrész nem független a következményrészétől. Okozatiságról szó sincs. Biztosan csak azt állíthatjuk, hogy nincs okozatisági viszony olyan jelenségek között, amelyek között korreláció sem áll fenn (azaz függetlenek). A korreláció és az okozatiság összekeverése nagyon gyakori hiba, amelyre a latin *cum hoc ergo propter hoc* (magyarul: vele, ezért miatta) kifejezéssel hivatkoznak.

Ha  $A$  és  $B$  között korreláció van, akkor lehet, hogy  $A$  okozza  $B$ -t, de lehet, hogy másféle kapcsolat áll fenn köztük. Az is lehet, hogy

- I.  $B$  okozza  $A$ -t.
- II. egy harmadik  $C$  jelenség okozza  $A$ -t és  $B$ -t is. Az okozatisági viszonyok ennél jóval bonyolultabb lehetnek.
- III. lehet, hogy a korrelációt véletlenek különös együttállása okozza. Emlékezzünk, hogy a statisztikai tesztek sosem mondanak teljes bizonyossággal semmit. Az elsőfajú hiba adja meg annak esélyét, hogy mi azt állítjuk, hogy két esemény között összefüggés áll fenn, holott azok függetlenek egymástól.
- IV.  $A$  és  $B$  egymást is okozhatják kölcsönösen megerősítő módon.

Nézzünk néhány példát.

- *Az egy cipőben alvás erős összefüggésben áll a fejfájással ébredéssel. Tehát a cipőben alvás fejfájást okoz.* Nyilvánvalóan hibás ez a következtetés. Sokkal kézenfekvőbb az a magyarázat, hogy az ittas állapot okozza mindkét dolgot.
- A következő állítás egy magyar kereskedelmi rádióban hangzott el. Forrásnak amerikai kutatókat neveztek meg. *A magassarkú cipő skizofréniát okoz.* Az állítás nyilván teljes blödség és csak azért hangzott el, hogy felkeltse a hallgatók figyelmét.
- Az alábbi állítás viszont a Nature magazinban hangzott el 1993-ban. *Valószínűbb, hogy rövidlátók lesznek azok a gyerekek, akik égő lámpa mellett alszanak.* Későbbi kutatások kimutatták, hogy valójában a szülők rövidlátása és a gyerekek rövidlátása között van összefüggés továbbá a rövidlátó szülők hajlamosabbak a lámpát égve hagyni, mint úgy általában a szülők.

Ha vásárlói kosarak elemzéséhez kanyarodunk vissza, akkor ezek szerint  $I \rightarrow I'$  nem az jelenti, hogy az  $I$  termék az  $I'$  termék megvásárlását okozza. Ha mind  $I$ , mind  $I'$  megvételét egy harmadik  $I''$  terméknek köszönhetjük, akkor csak pénzt veszítenénk, ha az  $I$  termék árát csökkentenék a  $I'$ -ét pedig növelnék. Az  $I$  eladásának növekedése ugyanis nem okozza az  $I'$  eladását, tehát nem nyernénk vissza az  $I'$ -vel az  $I$  árcsökkenéséből származó profitkiesést.

A valóságban a termékek csoportokat alkotnak, amelyekben a termékek eladása kölcsönösen megerősítik egymást. Például, ha veszünk egy fényképezőgépet, akkor sokan memóriakártyát és tokot is vesznek. Ha okozati kapcsolatok csak a fényképező  $\rightarrow$  memóriakártya és a fényképező  $\rightarrow$  tok lennének, akkor matematikailag a fényképező és a memóriakártya eladásának nem kéne változnia, ha a tok árusítását megszüntetnénk. Legtöbbünknek azonban igenis számít, hogy egy helyen lehet megvásárolni mindhárom terméket, ezért az eladások igenis csökkennének. A fényképezőgép, memóriakártya és tok termékhalmaz egy olyan halmaz, amelynek elemei kölcsönösen megerősítik egymás eladását.

## 5.7. Asszociációs szabályok és az osztályozás

A következő részben az osztályozással és kicsit a regresszióval fogunk foglalkozni. Mik a hasonlóságok és mik a különbségek az asszociációs szabályok kinyerése és az osztályozás között? Mindkét feladatban attribútumok közötti összefüggéseket tárunk fel.

Az asszociációs szabályok előnye, hogy tetszőleges két attribútumhalmaz között található összefüggést. Ezzel szemben osztályozásnál kijelölünk egy attribútumot és csak azt vizsgáljuk, hogy ezt az attribútumot hogyan határozzák meg a többi attribútumok. Asszociációs szabályok jellemző alkalmazási területe a vásárlási szokások elemzése, ahol minden termékösszefüggés érdekes lehet.

Asszociációs szabályoknál bináris attribútumokkal dolgozik. Ha a feltéte részben szereplő attribútumok értéke egy, akkor a következményrészben szereplő attribútum is egy lesz. Ha a feltételrész értéke nulla, akkor nem tudunk semmilyen megállapítást tenni a következményrészre vonatkozóan. Osztályozásnál ilyen nincs, ha tudjuk a magyarázó attribútumok értékét, akkor tudjuk a magyarázandóét is. Az attribútumtípusokra annyi megkötés van, hogy a magyarázandó attribútum kategória típusú legyen (regressziónál numerikus).

Más az egyes területek tudományos cikkeinek témája is. Az asszociációs szabályokról szóló cikkek nagy része gyakori elemhalmazok kinyeréséről szól. A fő cél az, hogy minél gyorsabb algoritmust adjunk erre az adott feladatra. A feladat értelmét nem vonják kétségbe

(sem azt, hogy tényleg szükség van-e olyan gyors algoritmusokra, amelyek gigabájt méretű adatokat tudnak feldolgozni másodpercek alatt és gigabájt méretű kimenetet generálnak). A cikkekben algoritmikus és adatstruktúrális megoldásokat mutatnak be, implementációs és párhuzamosíthatósági kérdéseket vizsgálnak, nem ritkán egy módszer elemzésénél a hardver tulajdonságait is számításba veszik.

Ezzel szemben osztályozásnál az osztályozás pontosságának javítása a fő cél, a hatékonyságbeli kérdések csak másodlagosak. Az osztályozás kutatói általában jóval komolyabb statisztikai tudással rendelkeznek.

## 6. fejezet

# Osztályozás és regresszió

### 6.1. Bevezetés

Ismeretlen, előre nem megfigyelhető változók, attribútumok értékének előrejelzése más ismert, megfigyelhető változók, attribútumok ismeretében régóta aktív kutatás tárgyát képezi. A kérdés gyakorlati jelentőségét nehéz lenne túlértékelni. Ebben a fejezetben vázlatosan ismertetjük, hogy miként alkalmazhatók a statisztika és gépi tanulás területén kifejlesztett módszerek az adatbányászatban.

A megnevezések tisztázása érdekében előrebocsátjuk, hogy a tanulmányban akkor beszélünk regresszióról vagy előrejelzésről (predikcióról), ha a magyarázott változót intervallum skálán mérjük. Amennyiben a magyarázott változó diszkrét értékészletű, nominális vagy ordinális skálán mért, akkor osztályozásról vagy klasszifikációról (csoportba sorolásról) beszélünk. Fogalmaink szerinti előrejelzést és klasszifikációt a statisztikai irodalom általában regresszió-számítás, valamint diszkriminancia elemzés és klasszifikáció néven illeti. A gépi tanulás területén az eljárásokat összefoglalóan felügyelt tanulásnak (supervised learning) nevezik.

Az adatbányászatban leggyakrabban alkalmazott előrejelző és klasszifikáló módszerek a következők:

- I. Legközelebbi szomszéd módszerek
- II. Lineáris és logisztikus regresszió
- III. Mesterséges neurális hálózatok
- IV. Döntési szabályok, sorozatok és fák
- V. Naiv Bayes klasszifikáció és Bayes hálózatok
- VI. SVM
- VII. Metaalgoritmusok (boosting, bagging, randomization, stb. )

Mindegyik eljárásról elmondható, hogy (legalább) két lépcsőben működik. Először az ún. tanító adatbázison felépítjük a modellt, majd később azt alkalmazzuk olyan új adatokra, amelyeken a magyarázott változó értéke nem ismert, de ismerni szeretnénk. Amikor előrejelző, vagy klasszifikáló módszert választunk a következő tulajdonságait célszerű figyelembe venni:



- Előrejelzés teljesítménye: Milyen értékes információt ad számunkra a modell a nem megfigyelhető magyarázó változóról (lásd 6.2 szakasz)?
- Gyorsaság: A modell előállításának és használatának időigénye.
- Robusztusság: Érzékeny-e a modell hiányzó, vagy outlier adatokra.
- Skálázhatóság: Használható-e a modell nagyon nagy adathalmazokra is?
- Értelmezhetőség: Kinyerhetünk-e az emberek számára értelmezhető tudást a modell belső szerkezetéből?
- Skála-invariancia: A klaszterezés lehetlenség-elméletét adaptálva (lásd 7.1 rész) skála-invariánsnak hívunk egy osztályozó eljárást, ha a módszer kimenete nem változik abban az esetben, ha tetszőleges intervallum típusú magyarázó változó helyett annak  $\alpha > 0$ -szorosát vesszük.

Az adatbányász közösség leginkább a korábban is ismert előrejelző és klasszifikáló eljárások skálázhatóságának továbbfejlesztésében ért el eredményeket. Különösen a döntési fák területén fejlesztettek ki olyan algoritmusokat, amelyek akár milliós esetszámú tanuló adatbázis esetén is alkalmazhatók.

„ A ritkábban borotválkozók korábban halnak.” Forrás: <http://gondola.hu/cikkek/31731>

A fejezet hátralévő részében először a klasszifikálók és előrejelzők teljesítményének értékelésével foglalkozunk, majd az eljárásokat ismertetjük. A hagyományos statisztikai módszerek (diszkriminancia analízis, lásd. pl.: [65] ismertetésétől eltekintünk, helyettük inkább az „egzotikusabbakra” koncentrálunk: a döntési fák, a mesterséges neuronhálózatok, a Bayeshálózatok, és négy további eljárás főbb jellemzőit mutatjuk be [68], [55], [46] és [93] írások alapján.




---

**Weka 3.5.7** A wekában az osztályozó módszereket a *Classify* fülön keresztül érjük el.

---



## 6.2. Az osztályozás és a regresszió feladata

Az osztályozás és regresszió során  $n$ -esekkel (angolul *tuple*) fogunk dolgozni, amelyeket objektumoknak/elemeknek hívunk. Adott lesz objektumok sorozata (vagy zsákja), amelyet tanító mintáknak, tanító pontoknak, tanító halmaznak (habár a halmaz szó használata itt helytelen, hiszen ugyanaz az objektum többször is előfordulhat) nevezünk. A tanítópontok számát  $m$ -mel vagy  $|T|$ -val fogjuk jelölni. Valójában tanításra a tanító pontok egy részét használjuk. A többi pont szerepe a tesztelés lesz.

Az  $n$ -es  $j$ -edik elemét  $j$ -edik *attribútumnak* hívjuk. Egy attribútumra névvel is hivatkozhatunk (pl. kor, magasság, szélesség attribútumok), nem csak sorszámmal. Minden attribútumnak saját értékészlete van. Az  $A$  attribútumváltozón olyan változót értünk, amely az  $A$  értékészletéből vehet fel értékeket.

Általános módon egy klasszifikáló vagy előrejelző módszer teljesítményét várható hasznosságával mérhetjük. Legyen a magyarázandó attribútumváltozó  $Y$ , a magyarázó attribútumváltozó(k) pedig  $X$ , eljárásunkat jelöljük  $f$ -fel (Az  $f$  az  $X$  értékészletéről az  $Y$  értékészletére képez). Ekkor célunk  $E[U(Y, f(X))]$  maximalizálása, ahol  $U(y, \hat{y})$  jelöli az előrejelzett  $\hat{y}$  hasznosságát, miközben a valódi érték  $y$ . Bináris  $Y$  esetén *bináris osztályozásról* beszélünk.

A feladatot fordítva,  $\mathbb{E}[L(Y, f(X))]$  minimalizálásaként is megfogalmazhatjuk, ahol  $L$  az  $U$  inverze, egy veszteséget mérő függvény. A  $\mathbb{E}[L(Y, f(X))]$  értéket *várható osztályozási hibának* (*expected prediction error*) nevezzük és *VOH*-val jelöljük. Mivel a várható érték változóiban additív és a konstanssal való eltolás nem változtat az optimalizáláson, ezért  $L(y, \hat{y})=0$  feltehető. A hibát a gyakorlatban egy távolságfüggvénnyel definiálják (lásd 3.2 rész).

### 6.2.1. Az elméleti regressziós görbe

Regresszió esetén a két legelterjedtebb megoldás a hiba mérésére a négyzetes hiba  $L(y, \hat{y}) = (y - \hat{y})^2$  és az abszolút hiba  $L(y, \hat{y}) = |y - \hat{y}|$  alkalmazása. Fejtsük ki a várható értéket:

$$\begin{aligned} VOH(f) &= \mathbb{E} \left[ (Y - f(X))^2 \right] \\ &= \int (y - f(x))^2 \mathbb{P}(dx, dy) \end{aligned}$$

A legkisebb hiba akkor adódik, ha

$$f(x) = \mathbb{E}[Y|X = x], \quad (6.1)$$

ugyanis

$$\begin{aligned} \mathbb{E} \left[ (Y - f(X))^2 \right] &= \mathbb{E} \left[ (Y - \mathbb{E}[Y|X] + \mathbb{E}[Y|X] - f(X))^2 \right] \\ &= \mathbb{E} \left[ (Y - \mathbb{E}[Y|X])^2 \right] + \mathbb{E} \left[ (\mathbb{E}[Y|X] - f(X))^2 \right] \geq \mathbb{E} \left[ (Y - \mathbb{E}[Y|X])^2 \right], \end{aligned}$$

mert

$$\begin{aligned} \mathbb{E} \left[ (Y - \mathbb{E}[Y|X]) (\mathbb{E}[Y|X] - f(X)) \right] &= \mathbb{E} \mathbb{E} \left[ (Y - \mathbb{E}[Y|X]) (\mathbb{E}[Y|X] - f(X)) | X \right] \\ &= \mathbb{E} \left[ (\mathbb{E}[Y|X] - f(X)) \mathbb{E}[Y - \mathbb{E}[Y|X] | X] \right] \\ &= \mathbb{E} \left[ (\mathbb{E}[Y|X] - f(X)) (\mathbb{E}[Y|X] - \mathbb{E}[Y|X]) \right] = 0 \end{aligned}$$

A második egyenlőségnél felhasználtuk, hogy  $\mathbb{E}(V) = \mathbb{E}\mathbb{E}(V|W)$ , a harmadik egyenlőségnél felcseréltük a szorzat két tagját és felhasználtuk, hogy a  $\mathbb{E}[Y|X] - f(X)$  független  $Y$ -től, ezért a várható érték elé mozgatható. Végezetül ismét a  $\mathbb{E}(V) = \mathbb{E}\mathbb{E}(V|W)$  trükköt használtuk,  $V = \mathbb{E}[Y|X]$  és  $W = X$  helyettesítéssel.

Az  $f(x) = \mathbb{E}[Y|X = x]$  függvényt *elméleti regressziós görbének* nevezik.

Ha a hiba mérésénél a négyzetösszeg helyett ( $L_2$  norma) az különbségösszeget használjuk ( $L_1$  norma), akkor az elméleti regressziós görbe:

$$f(x) = \text{median}(Y|X = x). \quad (6.2)$$

### 6.2.2. Maximum likelihood osztályozó

Osztályozás esetén négyzetes hibáról nincs értelme beszélnünk. Hibafüggvény helyett,  $k$  osztály esetén, egy  $c \times c$  méretű hibamatrixot ( $L$ ) adhatunk meg, amely  $i$ -edik sorának  $j$ -edik eleme ( $L[i, j]$ ) megadja a hiba mértékét, ha  $i$ -edik osztály helyett a  $j$ -edik osztályt jelezzük előre. A mátrix fődiagonálisában nulla értékek szerepelnek.

A várható osztályozási hiba

$$VOH(f) = \mathbb{E} [L[Y, f(X)]],$$

amelyből

$$f(x) = \operatorname{argmin}_{y_\ell \in Y} \sum_{i=1}^c L(y_i, y_\ell) \mathbb{P}(y_i | X = x)$$

A legismertebb veszteség mátrix a nulla-egy mátrix, amelyben a fődiagonálison kívül minden elem egy. A fenti kifejezés a következőre egyszerűsödik:

$$f(x) = \operatorname{argmin}_{y_l \in Y} [1 - \mathbb{P}(y_l | X = x)],$$

vagy egyszerűen:

$$f(x) = y_k, \text{ amennyiben } \mathbb{P}(y_k | X = x) = \max_{y_l \in Y} \mathbb{P}(y_l | X = x).$$

A fenti osztályozó a Bayes vagy maximum likelihood osztályozó, amely azt állítja, hogy a  $\mathbb{P}(Y|X)$  feltételes valószínűség szerinti legnagyobb valószínűségű osztály lesz az osztályozó kimenete adott megfigyelés esetén.

Ha a várható értéket meghatározó valódi eloszlásokat ismernénk, akkor megtalálható a legjobb előrejelző / klasszifikáló. Például (azonos kovarianciájú) többdimenziós normális eloszlásokat feltételezve egyszerű kvadratikus (lineáris) döntési szabályokat kapunk [133], [65]. Az eloszlás paramétereit általában még akkor is becsülnünk kell, ha feltételezhető / feltételezünk egy adott típusú eloszlás.

Az adatbányászat területén a normalitás nem reális feltevés (gondoljunk a sok nominális változóra). Az adatbányászati módszerek nem élnek feltevésekkel az eloszlással kapcsolatban.

Ugyanakkor a módszerek összetettségük folytán – ha hagyjuk őket – képesek nem csak a tanító adatbázis szabályszerűségeit, hanem a mintaadatokban lévő egyedi hibákat és torzításokat is megtanulni (ami kifejezetten káros). Így általában pusztán a tanító adatbázis segítségével nem megalapozott a várható haszon / költség nagyságát megbecsülni. Mennyire jó egy osztályozó módszer, amely 100% pontosságot ad a tanító mintákon, de 0%-ot a tesztelő adathalmazon?

A túlzott modellbonyolultság elkerülésére pl.: a regressziószámítás területén modellszelekciós kritériumok (módosított  $R^2$ , Akaike Schwartz, stb.), illetve heurisztikus eljárások (stepwise regresszió) állnak rendelkezésre. Az osztályozó módszer kiértékeléséről, illetve osztályozók összehasonlításáról a 6.10 részben írunk bővebben. Most lássuk a legismertebb osztályozó módszereket.

### 6.3. $k$ -legközelebbi szomszéd módszere

A  $k$ -legközelebbi szomszéd módszere egy „lusta” klasszifikáló eljárás, amely nem épít modellt. Alapelgondolása, hogy a hasonló attribútumú objektumok hasonló tulajdonságokkal

bírnak. A hasonlóságot (igazából a különbözőséget (lásd 3.2. rész)) a klaszterelemzésnél is használt távolságfüggvénnyel mérjük. A tanuló adatbázist eltároljuk és amikor egy ismeretlen objektumot kell klasszifikálnunk, akkor megkeressük a távolságfüggvény szerinti  $k$  darab legközelebbi pontot, és az objektumot abba a kategóriába soroljuk, amely a legtöbbször előfordul (leggyakoribb) a  $k$  szomszéd között (többségi szavazás). A módszer egyfajta lokális sűrűségfüggvény becslő eljárásnak is tekinthető. Regresszió esetén a szomszédok osztályértékeinek átlaga lesz a kimenet.

A módszer regresszióra is használható. Ekkor a többségi szavazás helyett a szomszédok osztályértékének átlagaként szokás megadni a jóslást.

Idézzük fel az optimális előrejelzőre tett megállapításunkat (lásd 6.1 egyenlőség), regresszió esetén:

$$f(x) = \mathbb{E}[Y|X = x],$$

ahol a tetszőleges pontban az optimális osztályozó értéke megegyezik a feltételes várható értékkel. Osztályozás esetén pedig

$$\hat{f}(x) = y_l, \text{ amennyiben } \mathbb{P}(y_l|X = x) = \max_{y_i \in Y} \mathbb{P}(y_i|X = x).$$

A  $k$ -legközelebbi szomszéd a következő regressziós függvényt adja tetszőleges  $x$  pontra

$$\hat{f}(x) = \frac{\sum y_i}{k},$$

ahol  $x_i \in N_k(x)$ , osztályozás esetén pedig:

$$\hat{f}(x) = y_k, \text{ amennyiben } \text{freq}(y_k|x_i \in N_k(x)) = \max_{y_\ell \in Y} \text{freq}(y_\ell|x_i \in N_k(x)),$$

ahol  $N_k(x)$  az  $x$  pont  $k$ -legközelebbi szomszédját,  $Ave$  az átlagot,  $\text{freq}$  pedig a gyakoriságot jelöli.

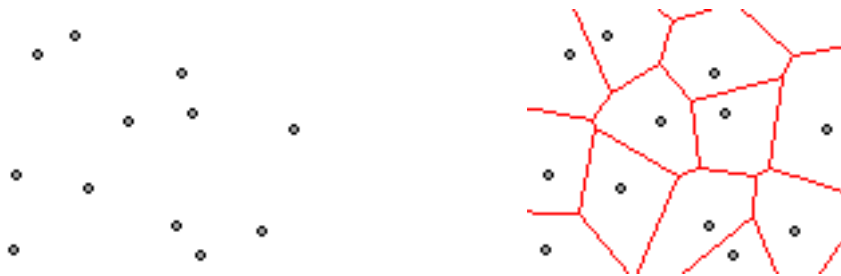
Az  $\hat{f}(x)$  tulajdonképpen az  $f(x)$  közelítése. A közelítés két okból következik:

- I. regresszió esetén a várható érték helyett a mintaátlagot használtuk, osztályozás esetén pedig a valószínűség helyett a relatív gyakoriságot,
- II. az  $x$  pontban vett feltétel helyett az  $x$  környezetét vettük.

Sok tanítópont esetében tetszőleges ponthoz közel lesznek a szomszédai, továbbá az átlag egyre stabilabb lesz, amennyiben  $k$  egyre nagyobb. Be lehet látni, hogy  $P(X, Y)$ -ra tett enyhe feltételek mellett  $\hat{f}(x) \rightarrow \mathbb{E}[Y|X = x]$ , amennyiben  $m, k \rightarrow \infty$  és  $N/k \rightarrow 0$ . Ezek szerint a  $k$ -legközelebbi szomszéd egy univerzális approximátor, nem is érdemes további osztályzókkal foglalkoznunk?!!

Legtöbb esetben nem áll rendelkezésünkre sok tanítópont. Ha fel tudunk tenni az osztályozásra valamilyen struktúrális feltételt (pl. linearitás), akkor ezt kihasználva pontosabb modellt építhetünk, mint azt a  $k$ -legközelebbi szomszéd adna. További probléma, hogy magas dimenziószám mellett (tehát sok attribútum esetén) a konvergencia lassú.

A legközelebbi szomszéd módszer ábrázolásánál ( $k = 1$  esetén) kedvelt eszköz a Voronoi diagramm. A területet felosztjuk tartományokra úgy, hogy minden tartományba egy tanító



6.1. ábra. Tanítópontok a síkon (bal oldali ábra) és a Voronoi tartományok (jobb oldali ábra)

pont essen és igaz legyen, hogy a tartományon belüli bármely pont a tanítópontok közül a tartomány tanítópontjához van a legközelebb. Egy ilyen felosztást láthatunk a 6.1 ábrán <sup>1</sup>.

Az osztályozáshoz természetesen nem kell meghatározni a tartományokat és megnézni, hogy az osztályozandó pont melyik tartományba tartozik. Egyszerűen nézzük végig a tanítópontokat és válasszuk ki a leginkább hasonlót.

### 6.3.1. Dimenzióátok - Curse of dimensionality

A legközelebbi szomszéd módszer egy univerzális approximátor, tetszőleges osztályozó függvényt képes reprodukálni, csak elég tanítópont kell hozzá. A módszert lokális approximátornak is szokás hívni, mert tetszőleges pont osztályértékét a (lokális) környezetének tanítóértékeinek átlagával helyettesíti. A módszer jól működik alacsony dimenzióknál, de magas dimenzióknál csődöt mond. Erre, mint dimenzióátok szoktak hivatkozni és számos szemléltetése, interpretációja létezik. A legkönnyebben érthető az alábbi.

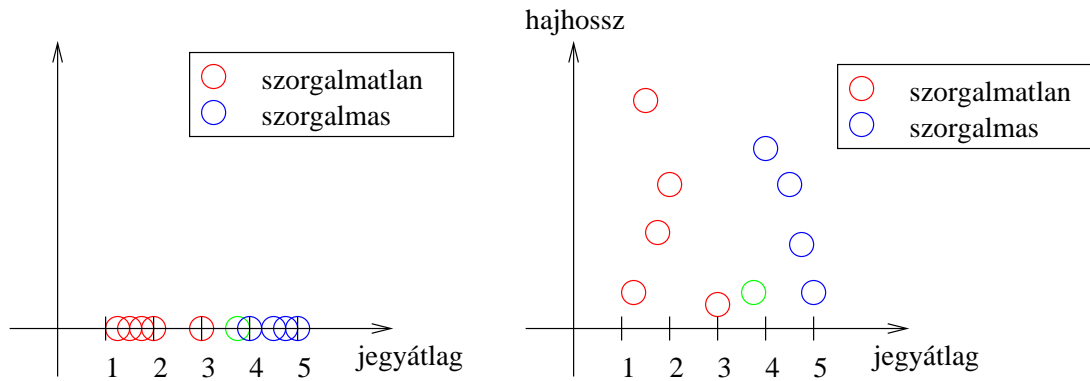
Tetszőleges pont környezetében elég tanítópontnak kell lenni. Ideális esetben tetszőleges  $x$  pont környezetén az  $x$ -től legfeljebb  $\epsilon$  távolságra lévő pontokat értjük. Ez egydimenziós esetben egy  $2\epsilon$  hosszú szakaszt, kétdimenziós esetben  $\epsilon$  sugarú kört, háromdimenziós esetben  $\epsilon$  sugarú gömböt jelent. Ha azt szeretnénk, hogy a keresési térben a tanítópontok sűrűsége rögzített legyen, akkor a tanítópontok számának exponenciálisan kell nőnie a dimenzió növelésével. A gyakorlatban a tanítópontok adottak, ami általában behatárolja a dimenziók és így a figyelembe vehető attribútumok számát.

Ez nem jelenti azt, hogy magas dimenziókban nem lehet osztályozó függvényt találni, csak megkötést kell tennünk az osztályozó függvény típusára vonatkozóan. Például, ha úgy gondoljuk, hogy az osztályozó egy hipersíkkal leírható, akkor a dimenziók számának növelésével csak lineárisan növekszik a szükséges tanítópontok száma (hiszen kétdimenziós esetben két pont határoz meg egy egyenest, három dimenzióknál három pont határoz meg egy síkot, stb.).

### 6.3.2. A legközelebbi szomszéd érzékenysége

A legközelebbi szomszéd módszer hátránya, hogy érzékeny a független attribútumokra. Lássunk egy példát. Feladatunk, hogy egy döntési modellt adjunk a szorgalmas diákokra. Az egyik attribútum a görgetett tanulmányi átlag a másik a hajhossz. A 6.2 ábra mutatja a tanító

<sup>1</sup>A szemléltető ábrát a [http://www.manifold.net/doc/7x/transform\\_voronoi\\_operators.htm](http://www.manifold.net/doc/7x/transform_voronoi_operators.htm) oldalról töltöttük le.



6.2. ábra. Független attribútumok hatása a legközelebbi szomszéd osztályozásra

pontokat, cél a zölddel jelölt tanuló osztályozása. Ha csak a jegyátlagot tekintjük, akkor a szorgalmasak közé soroljuk. Ha a távolság megállapításánál a hajhossz is figyelembe vesszük, akkor egy olyan hallgató lesz hozzá a legközelebb, akiről tudjuk, hogy szorgalmatlan. Sőt, ha euklideszi távolságot használunk és a független attribútum értékei jóval nagyobbak a függő attribútum értékeinél, akkor a független attribútum „elnyomja” a függő attribútumot.

Számos megoldást javasolnak a független attribútum által okozott hiba kiküszöbölésére. (1.) Ha tehetjük használjunk több tanító pontot, (2.) kérdezzük meg az alkalmazási terület szakértőjét, hogy a távolság meghatározásánál mely attribútumokat vegyük számításba, vagy (3.) alkalmazzunk statisztikai tesztet a függetlenség megállapítására. Amennyiben nincs sok attribútumunk, akkor meghatározhatjuk az osztályozás pontosságát az összes attribútum részhalmaz esetén majd kiválaszthatjuk a legjobbat.

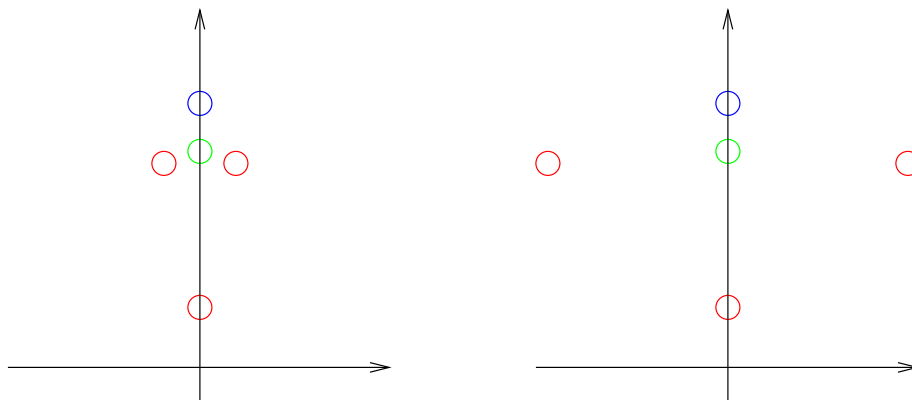
Sok attribútum esetén az összes attribútumhalmaz kipróbálása túl sok időt/erőforrást kíván. Egy (4.) mohó, bővítő eljárás egyesével bővítené a tesztelendő attribútumhalmazt úgy, hogy az a legjobb osztályozást adja. Ha az osztályozás minősége nem javul, akkor befejeznénk a bővítést. Ez a módszer kiselejtezné az  $X_1$  és  $X_2$  bináris attribútumokat annál az osztályozásnál, amelyben a magyarázandó attribútum értéke  $X_1$  és  $X_2$  moduló kettővel vett összege és  $X_1, X_2$  egymástól (és a magyarázandó attribútumtól is) teljesen függetlenek. Az (5.) csökkentő módszerek a teljes attribútumhalmazból indulna ki és minden lépésben egy attribútumot dobna ki.

A legközelebbi szomszéd módszer érzékeny a mértékegységre is. Ez logikus, hiszen a legközelebbi szomszéd módszer érzékeny a távolság definíciójára, az pedig nagyban függ az egyes attribútumok mértékegységétől. A problémát a 6.3 ábra szemlélteti.

Az egyik attribútum jelölje egy ember hajhosszát az átlagtól, a másik attribútum a bevételt jelöli dollárban. Az első ábrán a hosszt méterben mérjük a másodikban pedig lábban. Az osztályozandó (zöld) ponthoz egy piros van a legközelebb az első esetben, míg a második esetben kék pont a legközelebbi. A példából következik, hogy a legközelebbi szomszéd módszer nem skálainvariáns.

Az említett problémák nem feltétlenül az osztályozó hibái. A legközelebbi szomszéd módszerben a távolságfüggvény központi szerepet játszik. A helyes távolságfüggvény meghatározásához válasszuk ki a fontos attribútumokat, normalizáljuk, ha szükséges, illetve fontosságuk alapján súlyozzuk őket.

A  $k$ -legközelebbi szomszéd egy módosítását *súlyozott legközelebbi szomszéd módszernek* hívják. Ebben a  $k$  szomszéd minden tagjának akkora a súlya, amekkora az osztályozandó ponttól



6.3. ábra. Mértékegység hatása a legközelebbi szomszéd osztályozóra

vett távolságának inverze (vagy valamilyen antimonoton függvénye). Az osztályozandó ponthoz közel fekvő tanítópontoknak tehát nagyobb szavuk van a végső osztály meghatározásában, mint a távolabb eső pontoknak.




---

**Weka 3.5.7** *A legközelebbi szomszéd módszerét (tehát amikor csak egy szomszédot veszünk figyelembe) a `weka.classifiers.lazy.IB1` osztály implementálja. Két pont távolságának meghatározásánál az euklideszi normát használja. Ha több legközelebbi pontja van egy osztályozandó pontnak, akkor az elsőként megtalált alapján fog osztályozni.*

---

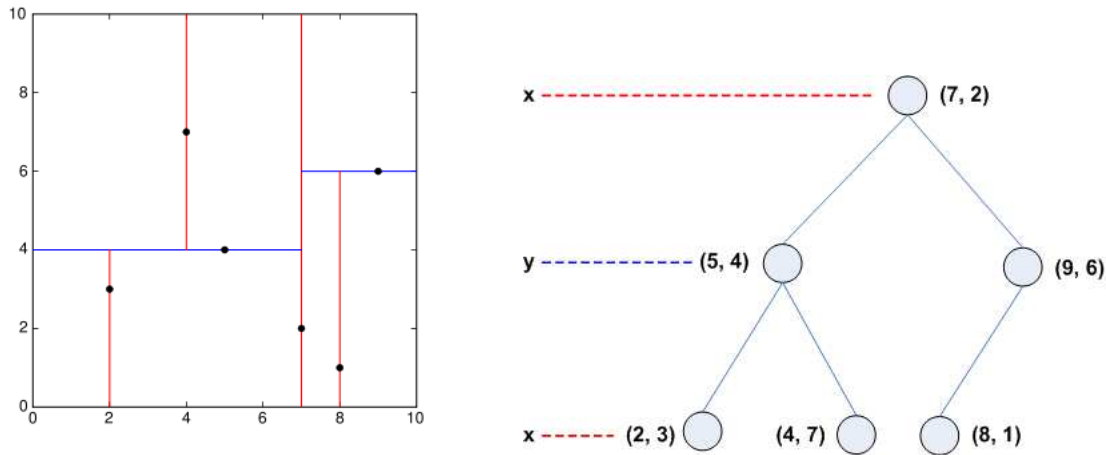


### 6.3.3. Az osztályozás felgyorsítása

Egy új elem osztályozásánál meg kell határoznunk a  $k$ -legközelebbi szomszédot. Ez a teljes adatbázis egyszeri, lineáris végigolvasását jelenti. A mai számítógépes architektúráknak kedvez ez a feladat. A tanítópontok elérnek a memóriában és a prefetch ill. pipeline technikák nagyban gyorsítják a keresést. Talán ez az oka, hogy teszteredmények szerint a kifinomultabb, bonyolultabb módszerek is legfeljebb egy nagyságrendet vernek a lineáris módszerre.

A 70-es évektől egyre több írás született, amelynek témája a lineáris módszernél gyorsabb algoritmus kidolgozása. Az új módszerek általában az ún. *branch-and-bound* technikát alkalmazzák, melynek során a tanító pontok terét felosztják és csak bizonyos részeket vizsgálnak a keresés során. Az előfeldolgozási lépésben gyúrnak valamit az adatbázison, általában egy speciális adatstruktúrába tárolják a tanítópontokat. Amennyiben a tanítás során egyetlen numerikus attribútumot veszünk figyelembe, akkor ha, mint előfeldolgozás, sorba rendezzük az adatokat ( $O(m \log m)$  idő alatt), akkor a legközelebbi szomszédok meghatározásához  $O(\log m)$  lépés elég. Futási idő szempontjából azonos asszimptotikával rendelkező algoritmus adtak két magyarázó változó esetében [10].

Több magyarázó attribútum esetén (nagyobb dimenzióknál) a legismertebb módszer KD-fákat használ [13]. Az algoritmus az előfeldolgozás során a teret hipertéglatestekre osztja, egyes



6.4. ábra. Tartományok meghatározása KD-fa építésénél (bal oldali ábra) és a KD-fa (jobb oldali ábra)

téglatesteket pedig további téglatestekre. A hipertéglatestek oldalai párhuzamosak egymással és a tengelyekkel és egy téglatest kettéosztása mindig egy az oldalfallal párhuzamos sík mentén való kettéosztást jelent. Egy téglatestet nem oszt tovább, ha a téglatestben található pontok száma adott korlát alatt van.

A KD-fa bináris és minden csomópontjának megfelel egy hipertéglatest. A levelekhez hozzá vannak rendelve azok a tanítópontok, amelyek a levél által meghatározott téglatesbe esnek. Tetszőleges csomópont gyermekeihez tartozó hipertéglatest a csomóponthoz tartozó hipertéglatest kettéosztásából jött létre. A 6.4 ábrán néhány tanítópont felosztása látható és a felosztáshoz tartozó KD-fa. Figyelem, a téglatestek által kijelölt terek nem osztályozási tartományoknak felelnek meg. A KD-fát használó algoritmus garantálja, hogy tényleg a legközelebbi szomszédokat fogja megtalálni.

Osztályozásnál nem csak azokat a tanítópontokat veszi figyelembe, amelyek abban a téglatestben vannak, amelyet az osztályozandó pont kijelöl. Az osztályozás menete a következő: A fa csúcsából kiindulva jussunk el addig a levélig, amely téglatestje tartalmazza az osztályozandó pontot. Határozzuk meg a legközelebbi pontot. Amennyiben a legközelebbi pont közelebb van, mint bármelyik oldalfal – másképp fogalmazva, az osztályozandó pontból a legközelebbi ponttól vett távolsággal rajzolt hipergömb nem metsz oldalfalat –, akkor leállunk. Ellenkező esetben meg kell vizsgálni azt a hipertéglatestet, amely fala közelebb van, mint a legközelebbi pont. Ez a téglatest ugyanis tartalmazhat olyan pontot, amely közelebb van az eddig talált legközelebbi ponthoz. A vizsgálandó téglatest nem biztos, hogy az osztályozandó pont által kijelölt téglatest szomszédja a KD-fában. Lehet unokatestvér, vagy sokadrangú unokatestvér. Vegyük észre, hogy az egyszerű konstrukció következtében (érts oldalfalak párhuzamosak a tengelyekkel) nagyon gyorsan el tudjuk dönteni, hogy egy adott téglatestnek lehet-e olyan pontja, amely egy adott ponttól, adott távolságnál közelebb van.

A KD-fa építésénél két cél lebeg a szemünk előtt:

- I. A fa legyen kiegyensúlyozott, abban a tekintetben, hogy minden téglatest ugyanannyi tanítópontot tartalmaz. Ez azért jó, mert ha ki tudunk zárni egy téglatestet a vizsgálatból,



akkor ezzel sok pontot szeretnénk kizárni

- II. A hipertéglalapok legyenek kockák. Ekkor ugyanis nem fordulhat elő, hogy az osztályozandó pont által kijelöl téglatesttel nem érintkező téglatest tartalmazza a legközelebbi szomszédot. Az elnyújtott téglatestek nem kedveznek az algoritmusnak.

Habár a második elvárásnak nem biztos, hogy eleget tesz a következő egyszerű módszer, mégis a gyakorlatban jó eredményt ad. Ki kell jelölni a fal tengelyét, majd meg kell határozni a helyét. A tengely kijelöléséhez nézzük meg mekkora a szórás az egyes tengelyekre nézve. Legyen a fal a legnagyobb szórást eredményező tengelyre merőleges. A fal helyét pedig a medián határozza meg, így garantált hogy a pontok egyik fele az egyik téglatestbe a másik fele a másik téglatestbe fog kerülni.

A KD-fánál vannak újabb adatstruktúrák, ezek közül a legismertebbek a Metric tree (Ball tree)[99, 137] és a Cover Tree [15]. A [73] cikkben a szerzők azonban azt állítják, hogy ezek az új módszerek nem mutatnak számottevő javulást a KD-fához képest. Állításaikat a szerzők számos teszt eredményére alapozzák, melyhez felhasználtak valódi és generált adatbázisokat is.




---

**Weka 3.5.7** *A  $k$ -legközelebbi szomszéd futtatásához  $k > 1$  esetén használjuk a `weka.classifiers.lazy.IBk` osztályt. A  $KNN$  paraméter felel meg a  $k$  értéknek, amelyet nem kell feltétlenül megadnunk. A `weka` a `leave-one-out` módszerrel (lásd a 155 oldal) megpróbálja a megfelelő  $k$  értéket meghatározni, amennyiben a `crossValidate` értéke igaz.*



*Használhatjuk a súlyozott legközelebbi szomszéd módszert is (lásd a 119). Ekkor választanunk kell a `distanceWeighting` paraméterrel, hogy a súly a távolság reciproka, vagy 1-től vett különbsége adja.*

*A `nearestNeighbourSearchAlgorithm` kiválasztóval megadhatjuk, hogy a legközelebbi szomszédok meghatározásához milyen módszert/adatstruktúrát használjon a `weka`. Az alapértelmezett az egyszerű lineáris keresés, de választhatunk `KD-fa`, `Ball tree` és `Cover tree` alapú megoldások közül.*

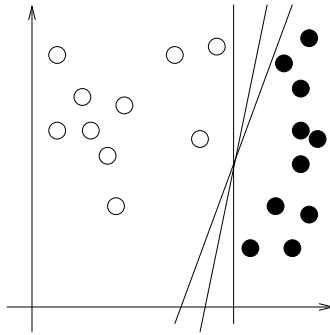
---

## 6.4. Lineárisan szeparálható osztályok

Két osztály lineárisan szeparálható, ha egy hipersík segítségével el tudjuk különíteni a két osztály pontjait. Amennyiben minden pont  $n$  dimenzióban van megadva, akkor  $n - 1$  dimenziós hipersíkot kell meghatározni. Ennek a hipersíknak a képlete:

$$w_1x_1 + w_2x_2 + \dots + w_nx_n = 0.$$

Az osztályozás feladata a  $w$  súlyok meghatározása. Ha ez megvan, akkor jöhet az új elemek osztályozása. Határozzuk meg az új elem attribútumainak  $w$  értékekkel történő súlyozott összegét. Ha az összeg nagyobb nulla, akkor az első osztályba tartozik, ellenkező esetben a másodikba. Kategória típusú magyarázó attribútum esetén az értékeket 0,1, ... számokkal szokás helyettesíteni. Lineárisan szeparálható osztályokra láthatunk példát a 6.5 ábrán.

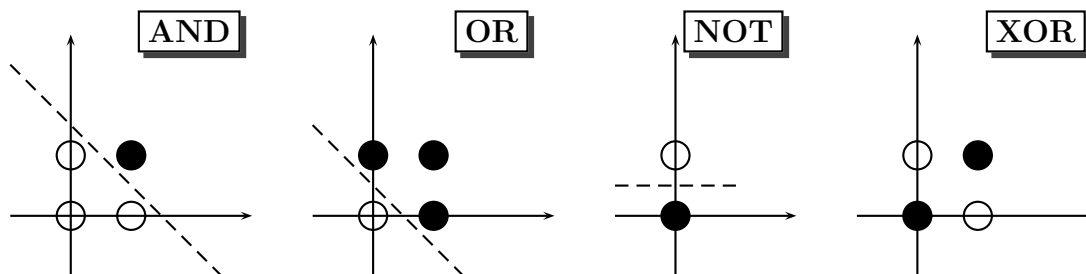


6.5. ábra. Példa lineárisan szeparálható osztályokra

A példából is látszik, hogy adott tanítóhalmazhoz több hipersík is létezik, amellyel kettéválaszthatjuk az osztályokat. A logisztikus regressziónál és az SVM osztályozónál fog felmerülni az a kérdés, hogy melyik hipersík választja el a legjobban a két osztályt, azaz melyik az a sík amelyik jól szeparál és amelytől legtávolabb vannak a pontok.

Mennyire erős megkötés, hogy az osztályok lineárisan szeparálhatók legyenek? Új attribútumok bevezetésével, amelyek az eredeti attribútumok nemlineáris transzformáltjai olyan térbe kerülhetnek, amelyben már lehet lineáris szeparálást végezni.

Amennyiben minden attribútum bináris és 0-1 értékeket vehet fel, akkor a 6.6 ábra jól mutatja, hogy az AND, OR, NOT függvények lineárisan szeparálható osztályokat hoznak létre.



6.6. ábra. AND, OR, NOT logikai függvények tanulása, XOR függvény

Sajnos ugyanez nem mondható el az XOR függvényre. Tehát már egy ilyen egyszerű logikai függvényt, mint az XOR sem tud megtanulni egy lineáris osztályozó. A neurális hálózatoknál vissza fogunk térni az XOR kérdéséhez. Látni fogjuk, hogy a neurális hálózatok, már tetszőleges logikai függvényt képesek megtanulni.

A perceptron és a Winnow módszereket fogjuk először szemügyre venni. Ezek kiindulnak a kezdetben konstans értékeket (perceptronnál nulla, Winnownál egy) tartalmazó súlyvektorból és a tanítópontok hatására a súlyvektort addig módosítják, amíg minden pontot jól szeparál a súlyvektor. A módszerek előnye, hogy jól használható online környezetben is, ahol néha új tanítópont érkezik, amely hatására módosítanunk kell a súlyvektort.

Ismertetjük még a Rocchio-eljárást, amely habár nem állít elő szeparáló hipersíkot mégis

lineáris szeparálást hajt végre. Végül elmélyedünk a logisztikus regresszió rejtelseiben.

### 6.4.1. Perceptron tanulási szabály

Mind az  $n$  attribútumnak valónak kell lennie. A hipersík dimenziója  $n$  lesz, ugyanis fel kell vennünk egy extra attribútumot (az angol irodalomban ezt bias-nak hívják), amelynek értéke minden tanító pontnál egy lesz. A módszer leírása alább olvasható.

---

#### Algorithm 7 Perceptron tanulási szabály

---

**Require:**  $\mathcal{T}$ : tanítópontok halmaza

```

 $\vec{w} = (0, 0, \dots, 0)$ 
while van rosszul osztályozott  $t \in \mathcal{T}$  do
  for all minden  $\vec{t} \in \mathcal{T}$  do
    if  $\vec{t}$  rosszul van osztályozva then
      if  $\vec{t}$  az első osztályba tartozik then
         $\vec{w} = \vec{w} + \vec{t}$ 
      else
         $\vec{w} = \vec{w} - \vec{t}$ 
      end if
    end if
  end for
end while

```

---

Amennyiben az algoritmus során rosszul osztályozott ponttal találkozunk, akkor úgy módosítjuk a hipersíkot, hogy a rosszul osztályozott tanító pont közelebb kerül hozzá, sőt akár át is kerülhet a sík másik oldalára. Ha egy rosszul osztályozott tanító pont az első osztályba tartozik, akkor az attribútum értékeinek súlyozott összege a módosítás során  $\sum w_i t_i$ -ről  $\sum (w_i + t_i) t_i$ -re változik. A különbség, négyzetösszeg lévén, biztosan pozitív. A hipersík a módosítás során helyes irányba mozgott.

A hipersík módosításai egymásnak ellentétesek lehetnek (olyan, mintha a tanítópontoktól jobbról és balról kapná a pofonokat), de szerencsére biztosak lehetünk benne, hogy a sok módosításnak előbb-utóbb vége lesz:

**6.1. lemma.** *Perceptron tanulási algoritmus véges lépésen belül leáll, amennyiben az osztályok lineárisan szeparálhatók.*

Hátrány, hogy ha a tanító pontok nem szeparálhatóak lineárisan, akkor az algoritmus nem áll le. A gyakorlatban ezért egy maximális iterációs számot adnak meg, amelynek elérésekor sikertelen üzenettel leáll az algoritmus.

### 6.4.2. Winnow módszer

Winnow módszerét akkor alkalmazhatjuk, ha minden attribútum bináris. Az eltérés a perceptron tanulástól annyi csak, hogy a rossz osztályozás esetén a súlyvektorhoz nem hozzáadjuk a tanítópont vektorát, hanem a súlyvektor bizonyos elemeit megszorozzuk vagy eloszjuk  $\alpha > 1$

konstanssal, attól függően, hogy melyik csoportba tartozik. Akkor sorol az osztályozó egy  $\vec{x}$  pontot az első osztályba, ha

$$w_1x_1 + w_2x_2 + \dots + w_nx_n > \Theta,$$

ahol  $\Theta$  előre megadott konstans. A szorzást vagy osztást azon elemekre végezzük, amelyre a tanítópont vektora egyest tartalmaz.

Mivel  $\alpha$  pozitív és a kezdeti súlyvektor minden eleme egy, ezért a súlyvektor minden eleme mindig pozitív marad. Vannak alkalmazások, ahol negatív súlyokat is meg kell engedni. Ekkor a *kiegyensúlyozott Winnow* (balanced Winnow) módszert alkalmazhatjuk. Két súlyvektort tartunk karban ( $\vec{w}^+$ ,  $\vec{w}^-$ ). Az osztályozáshoz a  $\vec{w}^+ - \vec{w}^-$  vektort használjuk. A rossz osztályozás esetén a  $\vec{w}^+$ -t ugyanúgy módosítjuk, mint a Winnow alapverziójánál, a  $\vec{w}^-$  elemeit pedig pont ellenkezőképpen, amikor  $w^+_i$ -t szorozzuk  $\alpha$ -val, akkor a  $w^-_i$ -t osztjuk vele.




---

**Weka 3.5.7** A Winnow, illetve a kiegyensúlyozott Winnow módszert a wekában a `weka.classifiers.functions.Winnow` osztály implementálja. A *balanced* paraméter igazra állításával adhatjuk meg, ha kiegyensúlyozott Winnow módszert szeretnénk alkalmazni. A súlyok kezdeti értékét a `defaultWeight` paraméterrel, az iterációk számát a `numIterations` paraméterrel szabályozhatjuk. A  $\Theta$  paraméter a wekában a `threshold` paraméternek felel meg.

---



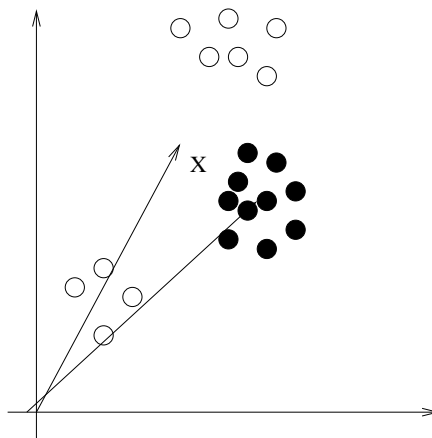
### 6.4.3. Rocchio-eljárás

A *Rocchio-eljárás* klasszikus módszernek számít az információ-visszakeresés területén. Osztályozási feladatra először a [61] munkában adaptálták, és azóta is sok kutatás foglalkozott vele (ld. [122]). Az eljárás feltételezi, hogy minden attribútum valós típusú. Minden  $c$  kategóriához megalkotunk egy *prototípusvektort*, amit a  $\mathcal{D}_c$  tanuló példák átlagaként számítjuk ki (centroid), és ehhez hasonlítjuk az ismeretlen dokumentum vektorát. Az osztályozandó objektum és egy kategória prototípusvektorának távolságát koszinusz- vagy más távolságmértékkel számolhatjuk.

A módszernek kicsiny a számításgénye, ezért a tanulás nagyon gyors. Hátránya viszont, hogy rossz eredményt ad, ha az egy osztályba tartozó pontok nem jellemezhetők egy vektorral (pl. amőba alapú csoportok, vagy két, egymástól jól elkülönülő, csoporthalmaz, amelynek elemei ugyanabba a kategóriába tartoznak). Ezt szemlélteti a 6.7 ábra. Az üres körök az első, a feketével töltött körök a második osztályba tartoznak. Az első osztály prototípusvektora távol esik az üres köröktől. Az  $x$ -szel jelölt osztályozandó pontot a Rocchio az első osztályba sorolná a második helyett.

A módszer hatékonysága lényegesen javítható, ha a prototípusvektorok megalkotásánál a negatív tanulóadatokat is figyelembe vesszük. Ekkor a

$$\vec{c} = w \cdot \sum_{t_j \in C} \vec{t}_j - \gamma \cdot \sum_{t_j \notin C} \vec{t}_j \quad (6.3)$$



6.7. ábra. Példa a Rocchio rossz osztályozására

képlettel számítható a  $c$  prototípusvektora<sup>2</sup>. Ha a második tagban nem az összes negatív tanulópéldát, hanem csak a majdnem pozitív tanulópéldák átlagát vesszük — ezek ugyanis azok, amelyekről a legnehezebb megkülönböztetni a pozitív tanulóadatokat, akkor további lényeges hatékonysági javulás érhető el [119, 141].

#### 6.4.4. Lineáris regresszió

Utazzunk vissza az időben, jussunk el az alapokig, azon belül a lineáris regresszióig. Tesszük ezt azért, mert a kifinomultabb, fejlettebb módszerek a lineáris regresszióból indulnak, illetve azért, mert ha lineárisan szeparálható osztályokkal van dolgunk, akkor a lineáris regressziót felhasználhatjuk a feladat megoldásához. A lineáris regresszió csak abban az esetben használható, ha minden attribútum valós típusú.

Feltételezzük, hogy az  $X$  magyarázó változó  $n$  dimenziós, és a magyarázandó  $Y$  változóval lineáris kapcsolatban áll:

$$\hat{Y} = \hat{w}_0 + \sum_{j=1}^n X_j \hat{w}_j$$

A  $\hat{w}_0$  értéket biasnak hívják. Amennyiben felvesszünk egy extra dimenziót és minden pont ezen dimenziója 1, akkor a vektoros felírást használva tovább egyszerűsíthetjük a képletet:

$$\hat{Y} = X^T \hat{w},$$

ahol a  $T$  felsőindex a mátrix transzponálásnak felel meg. A  $\hat{w}$  oszlopvektort kell meghatároznunk úgy, hogy adott tanítópontok  $(x_i, y_i)$  párok mellett négyzetes hibaösszeg minimális legyen.

$$\sum_{i=1}^{|\mathcal{T}|} (y_i - x_i^T w)^2.$$

<sup>2</sup>A pontok centroidjaként számolt prototípusvektort a  $w = 1$ ,  $\gamma = 0$  paraméterek mellett kapjuk meg.

A fenti függvény a  $w$ -ban négyzetes, így minimuma mindig létezik és egyértelmű. Amennyiben a tanítópontokat egy  $|\mathcal{T}| \times n$ -es  $X$  mátrixszal ábrázoljuk (a tanítópontokhoz tartozó  $y$  értékeket pedig az  $\mathbf{y}$  oszlopvektorral), akkor a fenti függvényt átírhatjuk más formába:

$$(\mathbf{y} - \mathbf{X}w)^T(\mathbf{y} - \mathbf{X}w)$$

Ennek  $w$  szerinti deriváltja:

$$-2\mathbf{X}^T\mathbf{y} + 2\mathbf{X}^T\mathbf{X}w$$

Ha a deriváltat egyenlővé tesszük nullával, akkor egyszerűsítés után a következőhöz jutunk:

$$\mathbf{X}^T(\mathbf{y} - \mathbf{X}w) = 0$$

amelyből nonsingularitást feltételezve kapjuk, hogy

$$\hat{w} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y}.$$




---

### **Weka 3.5.7**

*A `weka.classifiers.functions.`-*

*`SimpleLinearRegression` osztály egyetlen attribútum szerinti lineáris regressziót hajt végre. Azt az attribútumot választja, amely a legkisebb négyzetes hibát adja. Csak szám típusú attribútumokkal tud dolgozni és hiányzó értékeket nem enged meg.*

*A `weka.classifiers.functions.LinearRegression` osztály szintén lineáris regressziót hajt végre, de ez már több attribútumot is figyelembe tud venni. Lehetőség van a regresszióba felhasználandó attribútumok automatikus kiválasztására is az `attributeSelectionMethod` paraméterrel.*

*A négyzetes hibák átlaga (amely ugyanazt az eredményt adja, mintha az összeget minimalizálnánk) helyett a mediánt próbálja minimalizálni a `weka.classifiers.functions.LeastMedSq` osztály.*

---



Lineáris regresszióra visszavezethető számos nemlineáris kapcsolat is. Például az  $y = ax^b$  függvény logaritmusa  $\log x$  és  $\log y$ -ra nézve lineáris kapcsolatot ad. Hasonlóan egyszerű transzformációval kezelhető az  $y = \frac{1}{a+bx}$  függvény is.

### **6.4.5. Logisztikus regresszió**

Ha a lineáris regressziót osztályozásra akarjuk használni (de a magyarázó változók továbbra is valós számok), akkor az egyes osztályoknak egy valós számot kell megfeleltetnünk. Bináris osztályozásnál a nullát és az egyet szokás használni. Ezzel azonban nem oldottuk meg a problémát. A lineáris regresszió egy tanítópont osztályozásnál egy számot fog előállítani és a hibát a tanítópont ettől a számtól vett különbségével definiálja. Tehát egyes típusú tanítópont esetén ugyanakkora lesz a hiba 0 és 2 kimenetek esetén. Ez nem túl jó.

Egy  $x$  oszlopvektorral leírt, pont osztályának jóslásánál meg kell határoznunk az  $x^T w$  értéket. Amennyiben ez nagyobb, mint 0.5, akkor az 1-eshez tartozó osztály a jóslás, ellenkező esetben pedig a nullához tartozó osztály. Az egyszerűség kedvéért jelöljük az  $x^T w - 0.5$

értéket  $\hat{y}$ -nal. Sőt, a könnyebb jelölés érdekében a  $x^T w$  szorzatot jelöljük  $\hat{y}$ -nal és emlékezzünk rá, hogy a 0.5-öt belegyúrtuk a  $w_0$  torzítás értékbe.

Az a függvény, amely nullánál kisebb értékekre 0-át ad, nagyobbakra pedig 1-et eléggé hasonlít az előjel (szignum) függvényre. Ha megengedjük, hogy értelmetlen eredményt kapjunk  $\hat{y} = 0$  esetben – amelyet értelmezhetünk úgy, hogy az osztályozó nem képes dönteni –, akkor a jóslást megkaphatjuk az

$$\frac{1 + \operatorname{sgn}(\hat{y})}{2} \quad (6.4)$$

kiszámításával.

Ha így definiáljuk a kimenetet, akkor a hiba definíciója is megváltozott és a lineáris regresszió nem használható a  $w$  vektor meghatározásához. Egyenlőre most maradjunk annál, hogy egy mágus adott nekünk egy jól szeparáló hipersíkot. Tudunk-e árnyaltabb kimenetet adni, mint pusztán egy osztály (nulla vagy egy)?

Minél közelebb vagyunk a szeparáló hipersíkhhoz, annál bizonytalanabbak vagyunk a döntést illetően. A hipersíkon lévő pontokra már egyáltalán nem tudjuk mit tegyünk. Természetes gondolat hát, hogy az aktuális osztály jóslása helyett az osztály előfordulásának valószínűségét jósoljuk adott bemenet esetén. Ehhez csak annyit kell tennünk, hogy a 6.4 függvényt „lesimítjuk”, azaz egy olyan  $f(\hat{y})$  függvénnyel helyettesítjük, amely

- I. értéke 1-hez közelít, ha  $\hat{y}$  tart végtelenhez,
- II. értéke 0-hoz közelít, ha  $\hat{y}$  tart mínusz végtelenhez,
- III.  $f(0) = 0.5$ ,
- IV. szimmetrikus nullára nézve, tehát  $f(\hat{y}) + f(-\hat{y}) = 1 = 2f(0)$ ,
- V. „sima”, azaz  $f(\hat{y})$  differenciálható minden pontban és
- VI. monoton növekvő (vagy ne legyen lokális szélső érték).

Az ilyen függvényeket nevezzük *szigmoid* függvényeknek.

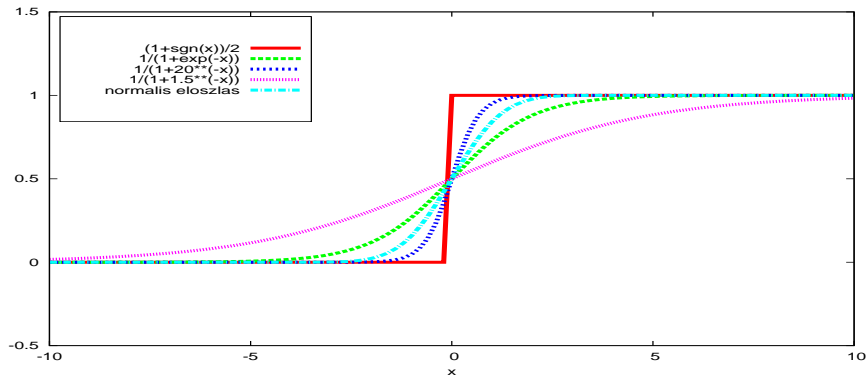
Sok függvény megfelel a fenti elvárásoknak. Könnyű belátni, hogy az  $1/(1+a^{-\hat{y}})$  függvények  $a > 1$  esetében az elvárások tekintetében rendben vannak. Amennyiben  $a = e$ , akkor az ún. *logisztikus függvényt* kapjuk

$$\mathbb{P}(Y = 1|X) = \frac{1}{1 + e^{-\hat{y}}}, \quad (6.5)$$

A logisztikus függvény inverzét ( $\ln(\frac{x}{1-x})$ ) *logit függvénynek* hívják. A logisztikus függvény szépsége, hogy a deriváltja  $f(\hat{y})(1 - f(\hat{y}))$ , amely a mi esetünkben  $\mathbb{P}(Y = 1|X)\mathbb{P}(Y = 0|X)$ -el egyezik meg.

Más függvények is eszünkbe juthatnak. Valószínűségi változók eloszlásfüggvénye is nullából indul mínusz végtelenben és egyhez tart a végtelenben. A harmadik és negyedik feltétel ( $f(0) = 0.5$ ,  $f(\hat{y}) + f(-\hat{y})$ ) megkívánja, hogy a sűrűségfüggvény szimmetrikus legyen, azaz az  $f'(x) = f'(-x)$  teljesüljön minden  $x$  valós számra. A nulla várható értékű normális eloszlás eloszlásfüggvénye megfelel a feltételeknek.

Az előjelfüggvény eltolt változatát, a  $1/(1+a^{-\hat{y}})$  típusú függvényeket különböző  $a$ -kra és a normális eloszlásfüggvényt a 6.8 ábra mutatja.



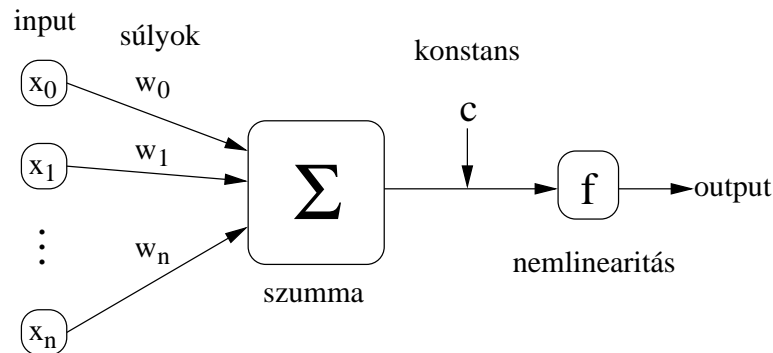
6.8. ábra. az eltolts elöjelfüggvény és néhány „simítása”

Ezzel el is jutottunk a logisztikus regresszió feladatához. Szemben a lineáris regresszióval, lineáris kapcsolat nem  $X$  és  $Y$  között van, hanem  $\ln\left(\frac{\mathbb{P}(Y=1|X)}{1-\mathbb{P}(Y=1|X)}\right)$  és  $x^T w$  között, tehát

$$\mathbb{P}(Y = 1|X) = \frac{1}{1 + e^{-X^T w}}, \quad (6.6)$$

$$\mathbb{P}(Y = 0|X) = 1 - \mathbb{P}(Y = 1|X) = \frac{e^{-X^T w}}{1 + e^{-X^T w}}. \quad (6.7)$$

Meg kell határozni azt a  $\hat{w}$  értéket, amelyik a legkisebb hibát adja.



6.9. ábra. Logisztikus regresszió

Sajnos a  $\hat{w}$  érték meghatározására nincs olyan szép zárt képlet, mint ahogy a lineáris regresszió esetében volt. Iteratív, közelítő módszert használhatunk, amely gradiensképzésen alapul. A hiba minimalizálása helyett a feltételes valószínűségeket maximalizáljuk:

$$\hat{w} \leftarrow \operatorname{argmax}_w \sum_{i=1}^{|\mathcal{I}|} \ln \mathbb{P}(y^i | x^i, w).$$

A fenti képletben a regressziós függvény a szokásos  $\mathbb{P}(Y^i | X^i)$  helyett  $\mathbb{P}(Y^i | X^i, w)$ , hiszen  $w$  most nem mint konstans játszik, hanem mint változó.



Felhasználva, hogy az  $y$  csak nulla vagy egy értéket vehet fel, a maximálandó függvényt átírhatjuk:

$$l(w) = \sum_i y^i \ln \mathbb{P}(y^i = 1|x^i, w) + (1 - y^i) \ln \mathbb{P}(y^i = 0|x^i, w).$$

Kiindulunk valamilyen szabadon megválasztott  $w^{(0)}$  vektorból, majd a  $k$ -adik lépésben a  $w^{(\ell)}$  vektorhoz hozzáadjuk a  $\frac{\delta l(w)}{\delta w}$  vektor  $\lambda$ -szorosát, így megkapjuk a  $w^{(\ell+1)}$  vektort. A  $\lambda$  egy előre megadott konstans, amelynek értékét 0.01-re szokták állítani. A  $\frac{\delta l(w)}{\delta w}$  vektor a  $\frac{\delta l(w)}{\delta w_j}$  parciális deriváltakból áll:

$$\frac{\delta l(w)}{\delta w_j} = \sum_{i=1}^{|T|} x_j^i (y^i - \widehat{\mathbb{P}}(y^i = 1|x^i, w)),$$

ahol  $\widehat{\mathbb{P}}(y^i = 1|x^i, w)$  a logisztikus regresszió által adott jóslás. Az  $y^i - \widehat{\mathbb{P}}(y^i = 1|x^i, w)$  tag a hibát ragadja meg, amely meg van szorozva egy nagyság jellegű tényezővel. Az  $x_j^i$  érték adja meg a jóslásban a  $w_j$  szerepének nagyságát (gondoljunk itt a  $w_j x_j^i$  súlyozott összegre).

Az  $l(w)$  konkáv ezért a gradiens módszer a globális maximumhoz fog konvergálni. A gyakorlat azt mutatja, hogy a konvergencia igen gyors, a  $w_j$  értékek néhány iteráció után már alig-alig változnak.




---

<b>Weka</b>	<b>3.5.7</b>	<i>Logisztikus</i>	<i>regressziót</i>	<i>a</i>
		<i>weka.classifiers.functions.Logistic</i>	<i>és</i>	<i>a</i>
		<i>weka.classifiers.functions.SimpleLogistic</i>	<i>függvények</i>	<i>implementálják.</i>

---



Felmerülhet a kérdés, hogy minek vacakolunk itt mi a logisztikus regresszióval, amikor lineárisan szeparálható osztályokra van már egyszerű megoldásunk (értsd perceptron tanítási szabály). Sőt a logisztikus regresszióban használt gradiens módszer lassú eljárás sok tanítópont esetén. A helyzet az, hogy a logisztikus regresszió zajos környezetben is jól működik, tehát amikor egy-egy pont átkerül a rossz oldalra. Ilyenkor a Perceptron tanulási szabály és a Winnow algoritmus nem áll le (csak ha nyomunk rajta egy STOP gombot). További előnye a logisztikus regresszióknak, hogy ha az osztályok lineárisan szeparálhatóak, akkor nem csak egy „elég jó” hipersíkot talál, hanem megtalálja a legjobbat (a legkisebb négyzetes hibát eredményezőt). Gondoljuk meg, ha csak két tanítópont adott a síkon, akkor a két pontot összekötő szakasz felező merőlegese lesz a szeparáló egyenes.

### Logisztikus regresszió általános osztályozásnál

Az eddigiekben feltettük, hogy bináris osztályozással van dolgunk. Mit tudunk tenni akkor, ha az magyarázandó attribútum  $k > 2$  értéket vehet fel?

A többválaszú logisztikus regresszió (multiresponse/multinomial logisztic regression) a fent említett logisztikus regressziót alkalmazza  $k$ -szor. Veszti az első osztályt és a többi egy kalap alá vonva végrehajt egy logisztikus regressziót. Ez ad egy valószínűséget. Ezután a második osztályt emeli ki és az összes többi osztályt vonja egy kalap alá. Így az összes osztályhoz meg tud határozni egy valószínűséget. Mintha csak egy tagsági függvényt próbálna meghatározni.

Új elem osztályozásánál az osztályozó arra az osztályra teszi le a voksát, amelyik a legnagyobb valószínűséget kapta. Ha csak a jóslott osztály érdekel minket és a kapott valószínűségekre nem vagyunk kíváncsiak, akkor nem szükséges a logisztikus függvény alkalmazása. Mivel a logisztikus függvény monoton növekvő, ezért megőrzi a sorrendet.

A kapott valószínűségeknek az összege nem feltétlenül ad egyet. Ezért a fenti módszer helyett csak  $k-1$  darab  $w$  vektort állítsunk elő úgy, hogy minden  $\ell = 1, \dots, k-1$ -re

$$\mathbb{P}(Y = \ell | X) = \frac{e^{w_0^\ell + x^T w^\ell}}{1 + \sum_{\ell'=1}^{k-1} e^{x^T w^{\ell'}}},$$

és

$$\mathbb{P}(Y = k | X) = \frac{1}{1 + \sum_{\ell'=1}^{k-1} e^{x^T w^{\ell'}}}.$$

A gradiens módszernél alkalmazott vektor, amely  $\lambda$ -szorosát hozzá kell adni az aktuális  $w$  vektorhoz a következő

$$\sum_{i=1}^{|\mathcal{T}|} x_j^i (\delta(y^i = \ell) - \hat{\mathbb{P}}(y^i = \ell | x^i, w)),$$

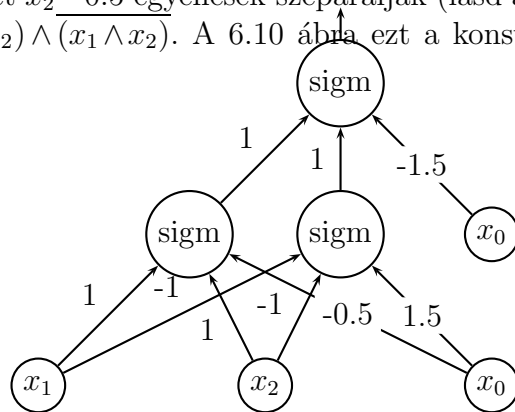
ahol  $\delta(y^i = \ell) = 1$ , ha az  $i$ -edik tanítópont osztálya  $\ell$ , különben 0.

A logisztikus regresszió és a Bayes osztályozó kapcsolatáról a 6.8 részben szólunk.

## 6.5. Mesterséges neurális hálózatok

A logisztikus regresszió modelljét egyrétegű mesterséges neurális hálózatnak is nevezik, sejt-hetjük, hogy ez az alapja a „komolyabb” mesterséges neurális hálózatoknak.

Egy darab logisztikus regresszió elemmel nem sok mindent tudunk kezdeni. Mivel lineáris osztályozó, ezért meg tudja tanulni az *és*, a *vagy* továbbá a *nem* logikai függvényeket, de nem tudja megtanulni az *xor* függvényt. Három logisztikus regresszió felhasználásával azonban az *xor*-t is ki tudjuk fejezni. Idézzük fel, hogy az *és* függvényt a  $x_1 + x_2 - 1.5$ , a *vagy*-ot a  $x_1 + x_2 - 0.5$  egyenes, a *nem*-et  $x_2 - 0.5$  egyenesek szeparálják (lásd a 6.6 ábra). Az *xor* függvény pedig felírható, mint  $(x_1 \vee x_2) \wedge (x_1 \wedge x_2)$ . A 6.10 ábra ezt a konstrukciót mutatja. A felső szignum

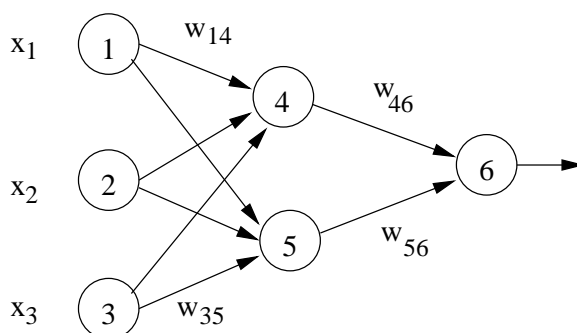


6.10. ábra. xor függvény logisztikus regressziók összekapcsolásával

függvényhez tartozó logisztikus regresszió az *és*-t, a bal alsó a *vagy*-ot a jobb alsó pedig a *nem és*-t adja vissza.

Az építőelemeket ismerve tetszőleges logikai formulát kifejezhetünk logisztikus regressziók összekapcsolásával, ezért a logisztikus regressziók kapcsolata univerzális függvényapproximátornak tekinthető. Ebből a tényből származik a neurális hálózatok elmélete.

A mesterséges neuronhálózatok – némileg az agyműködést utánzó biológiai analógiára is támaszkodva – a logisztikus regressziók kapcsolata. A legnépszerűbb modell a többrétegű előrecsatolt neuronhálózat (lásd 6.11. ábra). Az első réteg csomópontjaiban (neuronok) az input (magyarázó változók, 1-3. neuronok) helyezkedik el, az outputot (magyarázott változókat) a legutolsó réteg kimenete (6. neuroné) adja. A közbenső rétegeket rejtett (hidden) rétegeknek (4-5. neuronok) nevezzük. Minden réteg minden neuronjának kimenete a következő réteg összes neuronjának bemenetével kapcsolatban áll. A kapcsolat szorosságát  $w_{ij}$  súlyok jellemzik. (A 6.11. ábrában 4-6. neuronok helyébe egy 6.9. ábra szerinti logisztikus regressziót kell képzeljünk.)



6.11. ábra. Többrétegű előrecsatolt neurális hálózat

Mind a logisztikus regresszió, mind a neurális hálózatok paramétereikben nem lineáris függvény-approximátornak tekinthetők. A tapasztalatok és az elméleti eredmények (lásd.: [46]) szerint is ugyanannyi paramétert (súlyt) használva nemlineárisan paraméterezett függvényekkel gyakran jobb közelítést érhetünk el, mint lineárisan paraméterezett társaikkal.

Az alkalmas súlyokat nemlineáris optimalizációs technikával, gradiens módszerrel kereshetjük meg szinte ugyanúgy, mint a logisztikus regressziónál tettük. A gradiens eljárások alapelve, hogy egy függvény maximum / minimum helyét úgy keresik meg, hogy egy kezdőpontból kiindulva a gradiens (derivált) irányában / a gradienssel ellentétes irányban mozdulunk el, majd az eljárást ismétlik.

Az előrecsatolt topológiának köszönhetően az *egész* neuronháló hibafüggvényének  $w$  súlyok szerinti gradiensét könnyen kiszámíthatjuk. A súlyok megtalálása a tanító példák alapján az ún. backpropagation (hiba visszaterjedés) eljárás szerint zajlik:

- I. Az inputokból előrehaladva kiszámítjuk az outputok eredményét.
- II. Az utolsó output rétegből rétegről rétegre visszafelé haladva a megfelelő gradiens szabály szerint módosítjuk  $w_{ij}$  értékeket.

Mivel a neuronháló által reprezentált függvénynek lehetnek lokális maximumai ezért a módszer nem biztos, hogy a globális optimumot adja. A backpropagation eljárást ezért többször szokás futtatni különböző kezdeti súlyokkal.

A neuronhálók hátrányaként említhető, hogy a súlyok rendszere közvetlenül nem értelmezhető emberek számára. Nem tudjuk egyszerűen megindokolni, hogy mi alapján hozta meg a neuronháló a döntést. Egy hálózat tulajdonképpen egy fekete doboznak tekinthető a felhasználó szemszögéből. Sok területen nem elfogadható, ha egy módszer nem ad magyarázatot, ezért a neuronhálók alkalmazási köre erősen korlátozott. Ugyanakkor léteznek olyan eljárások, amelyek a neuronhálók súlyaiból emberek számára érthető, a döntéseket indokló szabályokat nyernek ki [55].

Egy városi legenda szerint a 80-as években az amerikai hadsereg szolgálatba akarta állítani a mesterséges intelligenciát és a számítástudományt. Céljuk volt minden tankra egy kamerát tenni, a kamera képét egy számítógépnek továbbítani, amely automatikusan felismeri, ha ellenséges tank búj meg a közeli erdőben. A kutatók neurális hálózat alapú megközelítés mellett döntöttek. A tanításhoz előállítottak 100 darab olyan képet amelyen a fák mögött tank bújt meg és 100 olyat, amelyen tank nem volt látható.

Néhány iteráció után a hálózat tökéletesen osztályozta a képeket. A kutatók és a Pentagon munkatársai nagyon meg voltak elégedve az eredményekkel, ugyanakkor még maguk sem voltak biztosak abban, hogy a neurális hálózat valóban a tank koncepciót tanulta-e meg. Független szakértőktől kért verifikáció során azonban a háló rosszul szerepelt. A pontossága nem haladta meg egy teljesen véletlenszerűen tippelő osztályozó pontosságát.

Valaki aztán rájött a rossz szereplés okára. A tanító képeken az összes tankos képen borult volt az idő, a tank nélküli képeken pedig sütött a nap. Ezt tanulta meg a háló.

Nem lehet tudni, hogy ebből a városi legendából mennyi igaz, az azonban tény, hogy a neurális háló nem ad magyarázatot az osztályozás okára. Ez komoly hátrány például a pénzügyi világban. A befektetők vonakodnak fekete doboz rendszerekre bízni a pénzüket, akkor is, ha ezek nagyon jó eredményeket adnak a tesztek során.




---

**Weka 3.5.7** A *backpropagation* tanító módszert használó neurális hálózatot a `weka.classifiers.functions.MultilayerPerceptron` osztály implementálja. A hálózatot felépíthetjük kézzel vagy automatikusan. A neuronokban használt nemlinearitás a szigmoid függvény.

Az osztálynak számos paramétere van. A *GUI* paraméterrel bekapcsolhatunk egy grafikus interfészt, melyen keresztül láthatjuk, illetve módosíthatjuk a neurális hálózatot. Az *autoBuild* paraméter engedélyezésével a hálózat automatikusan bővíti további rejtett rétegekkel. A *hiddenLayers* paraméter adja meg a neurális hálózat rejtett rétegeinek a számát. Az attribútumok előfeldolgozására ad lehetőséget a *nominalToBinaryFilter* paraméter. A kategória típusú attribútumokat bináris típusúvá alakítja (lásd 3.1 rész). Az attribútumok normalizálását a *normalizeAttributes* paraméterrel tudjuk engedélyezni. A *normalizeNumericClass* paraméter az osztályattribútumot normalizálhatjuk, amennyiben az szám típusú. A *validationSetSize* paraméter a teszhalmaz százalékos méretét adja meg. A tesztelés leállítását szabályozza a *validationThreshold*. Ez az érték adja meg, hogy egymás után hányszor romolhat a tesztelési hiba, mielőtt leállna a tanítás.

---



## 6.6. Döntési szabályok

**6.2. definíció.** Az  $\mathcal{A}$  attribútumhalmaz felett értelmezett döntési szabály alatt olyan  $R: \phi(\mathcal{A}) \rightarrow Y = y$  logika implikációt értünk, amelyek feltételrészében attribútumokra vonatkozó feltételek logikai kapcsolatai állnak, a következményrészben pedig az osztályattribútumra vonatkozó ítélet.

Például a  $H\ddot{O}M\ddot{E}R\ddot{S}\acute{E}K\ddot{L}E\ddot{T} = \text{magas AND SZ}\acute{E}L = \text{nincs} \rightarrow \text{ID}\ddot{O}\_J\ddot{A}T\acute{E}K\ddot{R}A$  alkalmas egy döntési szabály, amely azt fejezi ki, hogy ha magas a hőmérséklet és nincs szél, akkor az idő alkalmas kültéri játékra.

A valószínűségi döntési szabályokban a következményrészben az osztályattribútumra vonatkozó valószínűségi eloszlás szerepel. Ilyen szabályra példa az autóbiztosítás területéről, hogy  $\text{nem} = \text{férfi AND gyerekek száma} = 0 \text{ AND autó teljesítmény} > 150\text{LE} \rightarrow \text{kockázatos} = (80\%, 20\%)$ .

A feltételrészben az és, vagy és a negáció tetszőleges kombinációját felhasználhatjuk ...elvileg. A gyakorlatban ugyanis csak olyan szabályokkal foglalkoznak, amelyben egy alapfeltétel negációja és a feltételek és kapcsolatai szerepelnek. Ez azért nem olyan nagy megszorítás. Ha az azonos következményrésszel rendelkező szabályokból egy szabályt készítünk úgy, hogy a feltételek vagy kapcsolatát képezzük, akkor elmondhatjuk, hogy a szabályok feltételrészében diszjunktív normál formulák állnak. Minden ítéletlogikában megadott formula átírható diszjunktív normál formulává a dupla negáció eliminálásával, a de Morgan és a disztributivitási szabály alkalmazásával.

**6.3. definíció.** Az  $R: \phi(\mathcal{A}) \rightarrow Y = y$  szabályra illeszkedik az  $t$  objektum, ha a feltételrész attribútumváltozóiba az  $t$  megfelelő értékeit helyettesítjük, akkor igaz értéket kapunk.

Amennyiben a szabály következményrésze is igazra értékelődik az objektumon, akkor a szabály fennáll vagy igaz az objektumon.

**6.4. definíció.** Az  $R: \phi(\mathcal{A}) \rightarrow Y = y$  szabály lefedi az  $T$  objektumhalmazt, ha minden objektum illeszkedik a szabályra. Adott  $\mathcal{T}$  tanító halmaz esetén az  $R$  által fedett tanítópontok halmazát  $cover_{\mathcal{T}}(R)$ -rel jelöljük.

Helyesen fedi az  $T$  halmazt az  $R: \phi(\mathcal{A}) \rightarrow Y = y$  szabály, ha  $R$  fedi  $T$ -t és a halmaz összes objektuma az  $y$  osztályba tartozik. Ellenkező esetben helytelen fedésről vagy egyszerűbben rosszul osztályozásról beszélünk. A  $cover_{\mathcal{T}}$ -ben az  $R$  által helyesen fedett pontok halmazát  $cover^+_{\mathcal{T}}(R)$ -rel jelöljük (a helytelenül fedetteket pedig  $cover^-_{\mathcal{T}}(R)$ -rel).

**6.5. definíció.** Az  $R$  szabály relatív fedési hibája megegyezik a rosszul osztályozott pontok számának a fedett tanítópontokhoz vett arányával, tehát

$$Er_{\mathcal{T}}(R) = \frac{cover^-_{\mathcal{T}}(R)}{cover_{\mathcal{T}}(R)}.$$

### Döntési szabályok kifejezőereje

Kifejező erejük szempontjából a döntési szabályok következő típusairól beszélünk:

**ítétekalkulus-alapú döntési szabály** A feltételrészben predikátumok logikai kapcsolata áll (ítétekalkulus egy formulája, amelyben nem szerepelnek a  $\rightarrow$  és  $\longleftrightarrow$  műveleti jelek). Minden predikátum egy attribútumra vonatkozik. Amennyiben az attribútum kategória típusú, akkor  $A = a$  vagy  $A \in \mathcal{A}$  alakú a feltétel, ahol  $a$  egy konstans,  $\mathcal{A}$  pedig az  $A$  értékészletének egy részhalmaza. Sorrend vagy intervallum típusú attribútum esetében emellett  $A \leq a$  és  $a' \leq A \leq a''$  szabályokat is megengedünk.

Az algoritmusok többsége csak olyan egyszerű formulákat tud előállítani, amelyekben a predikátumok és kapcsolataik állnak, például  $\text{MAGASSÁG} \leq 170 \text{ AND HAJSZÍN} = \text{barna AND SZEMSZÍN} \in \{\text{kék, zöld}\}$ .

A csak ítétekalkulus-alapú szabályokat tartalmazó döntési szabályokat/fákat *univariate* (egyváltozós) döntési szabályoknak/fáknak hívjuk.

**reláció-alapú döntési szabály** Ha halmazelméleti szemmel nézzük a predikátumokat, akkor az attribútumokra vonatkozó predikátumot nevezhetünk bináris relációnak, amelynek egyik tagja egy változó, másik tagja egy konstans. A reláció-alapú döntési szabályokban a második tag attribútumváltozó is lehet. Itt például a  $\text{hajszín} = \text{szemszín}$  vagy a  $\text{szélesség} < \text{magasság}$  megengedett feltételek. A reláció-alapú szabályokat tartalmazó döntési szabályokat/fákat multivariate (többváltozós) döntési szabályoknak/fáknak hívjuk. A reláció alapú döntési szabályoknak nem nagyobb a kifejező erejük, amennyiben az attribútumok értékészlete véges. Ekkor ugyanis egy relációs szabály helyettesíthető sok egyváltozós szabálypárral. A fenti példa megfelelője a  $\text{hajszín} = \text{barna AND szemszín} = \text{barna}$ ,  $\text{hajszín} = \text{kék AND szemszín} = \text{kék}$ ,  $\text{hajszín} = \text{mályva AND szemszín} = \text{mályva}$  szabályokkal.

**induktív logikai programozás** Példaként tegyük fel, hogy építőelemek egy kupacát toronynak hívjuk, amelynek legfelső elemére a csúcs, a maradék elemekre pedig a maradék attribútummal hivatkozunk. A  $\text{szélesség} < \text{magasság} \rightarrow \text{ALAK} = \text{álló}$  szabályt úgy is írhatjuk, hogy  $\text{szélesség}(\text{építőelem}) < \text{magasság}(\text{építőelem}) \rightarrow \text{álló}(\text{építőelem})$ . Sőt a szabályt tovább is bonyolíthatjuk:  $\text{szélesség}(\text{torony.csúcs}) < \text{magasság}(\text{torony.csúcs}) \text{ AND } \text{álló}(\text{torony.maradék}) \rightarrow \text{álló}(\text{torony})$ . Ez egy rekurzív kifejezés, amely szerint egy torony akkor álló, ha a legfelső elem magassága nagyobb a szélességénél és a maradék elem álló. A rekurziót le kell zárni:  $\text{torony} = \text{üres} \rightarrow \text{álló}(\text{torony})$ . A rekurzív szabályoknak nagyobb a kifejezőerejük, mint a reláció-alapú döntési szabályhalmaznak, hiszen kifejezve tetszőleges számú predikátumot tartalmazhatnak. A rekurzív szabályokat is tartalmazó szabályhalmazt *logikai programnak* nevezzük, ezekkel továbbiakban nem foglalkozunk.

### 6.6.1. Szabályhalmazok és szabálysorozatok

Beszélünk *szabályhalmazról* és *szabályok sorozatáról*. Halmazok esetén a szabályok függetlenek egymástól. A szabályhalmaz *egyértelmű*, ha tetszőleges objektum csak egy szabályra illeszkedik.

Sorozat esetében egy új objektum osztályattribútumának jóslásánál egyesével sorra vesszük a szabályokat egészen addig, amíg olyat találunk, amelyre illeszkedik az objektum. Ennek a szabálynak a következményrésze adja meg az osztályattribútum értékét.

Egy szabályrendszer (sorozat vagy halmaz) *teljes*, ha tetszőleges objektum illeszthető egy szabályra. Ez köznyelven azt jelenti, hogy az osztályozó minden esetben (tetszőleges osztályozandó elemre) döntést hoz. Sorozatok esetében a teljességet általában az utolsó, ún. *alapértelmezett* szabály biztosítja, amelynek feltételrésze üres, tehát minden objektum illeszkedik rá.

Szabálysorozat esetében nem kell beszélnünk egyértelműségről, hiszen több szabályra való illeszkedés esetén egyértelmű a legelső illeszkedő. A szabályok közötti sorrend (vagy másképp prioritás) biztosításával kerüljük el azt a problémát, hogy milyen döntést hozzunk, ha egy objektumra több, különböző következményrésszel rendelkező szabály illeszkedik.

Sajnos a sorrend definiálásának ára van. Szabályhalmaz esetén ugyanis minden szabály a tudásunk egy töredékét rögzíti. Sorozatok esetében azonban egy szabályt nem emelhetünk ki a környezetéből; egy  $R$  szabály csak akkor süthető el, ha az  $R$ -et megelőző szabályok feltételrészei nem teljesülnek.

A szabályok sorozata átírható szabályok halmazába úgy, hogy egyesével vesszük a szabályokat az elsőtől és a feltételrészhez hozzáfűzzük az előtte álló szabályok feltételrész negáltjainak és kapcsolatát. Az így kapott szabályhalmaz azonban túl olvashatatlan, bonyolult lesz. Sorozattal az összefüggés esetleg egy tömörebb, könnyebben értelmezhetőbb formáját kapjuk.

„Kaliforniai kutatók szerint a marihuána egyik összetevője blokkolni képes a mellrák szétterjedését a szervezetben.” Forrás: <http://velvet.hu/blogok/gumicukor/tags/kutat%C3%A1s>

## 6.6.2. Döntési táblázatok

A döntési táblázat minden oszlopa egy attribútumnak felel meg, az utolsó oszlop az osztályattribútumnak. Az  $A$  attribútumhoz tartozó oszlopban az  $A$  értékére vonatkozó feltétel szerepelhet, leggyakrabban  $A=a$  alakban (ítéltkalkulus-alapú döntési szabály). A táblázat egy sora egy döntési szabályt rögzít. Ha az attribútumok a sorban szereplő feltételeket kielégítik, akkor az osztályattribútum értéke megegyezik a sor utolsó elemének értékével. Elég az elméletből, lássunk egy példát:

időjárás	hőmérséklet	páratartalom	szél	játékidő?
napos	meleg	magas	nincs	nem
napos	meleg	magas	van	nem
borús	meleg	magas	nincs	nem
esős	enyhe	magas	nincs	igen
esős	hideg	magas	nincs	igen
esős	hideg	magas	nincs	igen
esős	hideg	magas	nincs	igen

Egy döntési táblázat tulajdonképpen egy speciális döntési szabályhalmaz, amelyre igaz, hogy a feltételrészben pontosan ugyanazok az attribútumok szerepelnek.

Döntési táblák előállításánál a következő kérdéseket kell tisztázni:

- I. Az attribútumok melyik részhalmazát érdemes kiválasztani? Ideális az lenne, ha minden részhalmazt ki tudnánk értékelni és kiválasztani azt, amelyik a legkisebb hibát (rosszul

osztályozott tanítópontok száma) adja. A gyakorlatban azonban az attribútumok száma nagy ezért az összes részhalmaz kipróbálása sok időbe telik.

- II. Hogyan kezeljük a folytonos attribútumokat? A fenti példában a hőmérsékletet diszkrétizáltuk. Meleg az idő, ha 25 foknál több van, alatta enyhe 5 fokig. Ha a hőmérséklet 5 fok alá megy, akkor hideg van. Ideális az lenne, ha a folytonos attribútumokat az algoritmus automatikusan tudná diszkrétizálni.

### 6.6.3. Az 1R algoritmus

Talán a legegyszerűbb osztályzó algoritmus az 1R. Kiválaszt egy attribútumot és az osztályozásban kizárólag ezt használja. Annyi szabályt állít elő, ahány értéket felvesz a kiválasztott attribútum a tanítóhalmazban. Az  $A = a \rightarrow Y = c$  szabály következményrészében szereplő  $c$  osztály a legtöbbször előforduló osztály az  $A$  attribútumában  $a$  értéket felvevő tanítóminták közül.

Nyilvánvaló, hogy az 1R egyértelmű szabályhalmazt állít elő.

Minden attribútumértékhez meg tudjuk határozni a rosszul osztályozott tanítópontok számát. Ha összeadjuk az  $A$  attribútum értékeihez tartozó rosszul osztályozott tanítópontok számát, akkor megkapjuk, hogy mennyi tanítópontot osztályoznánk rosszul, ha az  $A$  attribútum lenne a kiválasztott. A legkevesebb rosszul osztályozott tanítópontot adó attribútumot választjuk osztályzó attribútumnak. Hiányzó attribútumértékeket úgy kezeljük, mintha az attribútumnak lenne egy különleges, a többitől eltérő értéke.

Sorrend és intervallum típusú attribútumnál  $A \leq a$ ,  $a' \leq A < a''$  és  $a''' \leq A$  típusú szabályokat célszerű előállítani. Ehhez csoportosítsuk az egymást követő értékeket úgy, hogy a hozzájuk tartozó osztályérték szempontjából homogén csoportokat hozzanak létre. Erre diszkrétizálásként is hivatkozunk és az 1R során használt módszert az Előfeldolgozás fejezetben ismertettük (lásd 3.3.5 rész).

Habár a sorrend és intervallum típusú attribútum csoportosításán sokat lehet elmélkedni az 1R módszer nem túl bonyolult. Egyszerűsége ellenére elég jól muzsikál a gyakorlatban. Egy meglepő cikkben [59] a szerző arról írt, hogy az 1R sokkal jobb osztályzó algoritmus, mint azt hinnénk. A szerzők azon a 16 adatbázison értékelték ki a különböző osztályzó módszereket – köztük az 1R-t –, amelyeket a kutatók gyakran használnak cikkeikben. A diszkrétizálásnál 3 helyett 6-ot használt, a módszereket kereszt-validációs eljárással hasonlította össze. Az 1R zavarba ejtően jó helyen végzett, a pontosság tekintetében alig maradt el az újabb és jóval bonyolultabb eljárásoktól.

Az 1R nevében szereplő szám az osztályozás során felhasznált attribútum számára utal. Létezik 0R osztályzó is, amely nem használ fel egyetlen attribútumot sem. Az osztályzó ekkor egy feltétel nélküli szabály, amely ítéletrészében a leggyakoribb osztály áll.

---

#### **Weka 3.5.7**

*A wekában a 0R és 1R módszereket a `weka.classifiers.rules.ZeroR` és a `weka.classifiers.rules.OneR` osztályok implementálják. Az 1R módszer egyetlen paramétere a diszkrétizálás során használt elemszám küszöb.*

---





### 6.6.4. A Prism módszer

A Prism módszer [26] feltételezi, hogy a tanító adatbázisban nincs két olyan elem, amelynek a fontos magyarázó attribútumai megegyeznek, de más osztályba tartoznak. Ha mégis akadnak ilyen objektumok, akkor csak egyet tartunk meg méghozzá olyat, amelyik a leggyakrabban előforduló osztályba tartozik. A leggyakoribb osztályt az azonos attribútumértékkel rendelkező pontok körében kell nézni. A Prism módszer a *fedő módszerek* közé tartozik.

A fedő algoritmus egyesével veszi az osztályattribútum értékeit és megpróbál olyan szabályokat előállítani, amelyek helyesen fedik azon tanítópontokat, amelyek a vizsgált osztályba tartoznak. A szabályok előállításánál a feltételrészhez adunk hozzá egy-egy újabb részfeltételt törekedve arra, hogy olyan részfeltételt vegyünk, amely legnagyobb mértékben növeli a pontosságot. A módszer hasonlít a döntési fák előállítására (lásd következő fejezet) ott is a meglévő szabályhalmazhoz egy új részfeltételt adunk. Döntési szabályoknál más a cél; pontosság növelése helyett az osztályok közötti szeparációt szeretnénk maximalizálni.

A Prism menete a következő. Egyesével sorra vesszük az osztályattribútum értékeit. Minden értéknél kiindulunk egy olyan döntési szabályból, amelynek feltételrészre üres, következményrészében pedig az aktuális osztályérték szerepel. Minden lehetséges  $A$  attribútum,  $a$  érték párra kiszámítjuk, hogy mennyi lenne a helytelenül osztályozott tanítópontok száma, ha az  $A = a$  részfeltételt adnánk a feltételrészhez. Azt a részfeltételt választjuk, amely a legkisebb relatív fedési hibát adó szabályt eredményezi. A részfeltételek hozzáadását addig folytatjuk, amíg olyan szabályt kapunk, amelynek nem nulla a fedése, de nulla a relatív fedési hibája.

Ezután töröljük a tanítópontok közül azokat, amelyeket az újonnan előállított szabály lefed. Ha nincs több olyan tanítópont, amelynek osztályattribútuma az aktuális osztályértéket veszi fel, akkor a következő attribútumértéket vesszük a következményrészbe. Az algoritmus pszeudokódja a 8 ábrán olvasható.

A Prism algoritmus alkotta szabályokat szabálysorozatként célszerű értelmezni. A módszer mindig olyan szabályokat hoz létre, amely lefed néhány tanítópontot. A következő szabály a maradék tanítópontokra szól ezért új objektum osztályozásakor akkor süssük el, ha az előző szabályt nem tudtuk illeszteni. A Prism algoritmusra, mint *separate and conquer* (*leválaszt majd lefed*) módszerre szoktak hivatkozni. A Prism először leválasztja a tanítópontok egy csoportját, majd megpróbálja lefedni azokat szabályokkal.

A Prism csak 100%-os pontosságú szabályokat állít elő. Az ilyen egzakt szabályok mindig a túltanulás veszélyét hordozzák magukban. Az ilyen szabályok sok feltételt tartalmaznak és általában kevés tanítópontot fednek. Hasznosabb lenne kisebb pontosságú, de több pontot fedő szabályokat előállítani. A tökéletességre való törekvés a Prism egy vitathatatlan hibája. Ha például egy feltétel két meghosszabbítása olyan, hogy az első lefed 1000 pontot, de egyet negatívan, a másik pedig csak egy pontot fed le (nyilván helyesen), akkor a Prism a második meghosszabbítást fogja választani. Egy Prism változat a  $\phi$  növelésénél a jelölt AND  $A = a$  taggal a relatív fedési hiba helyett egy információ nyereség jellegű értékkel számol. Jelöljük a  $\phi$  AND  $A = a \rightarrow Y = y$  szabályt  $R$ -rel.

$$\text{hiba}^* = \text{cover}^+(R) \cdot [\log(\text{Er}(R)) - \log(\text{Er}(\phi \rightarrow Y = y))].$$

Az információnyereség-alapú Prism is addig bővíti a feltételrészét, amíg nem sikerül 100%-os pontosságú szabályt előállítani.

Összehasonlítva az információnyereség és a relatív fedési hiba alapján előállított szabályokat a következőket mondhatjuk. A relatív fedési hiba esetén eleinte kis fedésű szabályokat nyes le,

**Algorithm 8 Prism****Require:**  $\mathcal{T}$ : tanítópontok halmaza, $Y$ : osztályattribútum változó,

```

for all  $y \in$  osztályattribútum értékre do
   $E \leftarrow$  az  $y$  osztályba tartozó tanítópontok
   $\phi \leftarrow \emptyset$ 
  while  $E \neq \emptyset$  do
     $R \leftarrow \phi \rightarrow Y = y$ 
    while  $Er_{\mathcal{T}}(R) \neq 0$  do
      hiba  $\leftarrow 1$ 
      for all  $(A, a)$  attribútum-érték párra do
        if  $Er(\phi \text{ AND } A = a \rightarrow Y = y) < \text{hiba}$  then
          hiba  $\leftarrow Er(\phi \text{ AND } A = a \rightarrow Y = y)$ 
           $A^* \leftarrow A$ 
           $a^* \leftarrow a$ 
        end if
      end for
       $\phi \leftarrow \phi \text{ AND } A^* = a^*$ 
    end while
     $\mathcal{T} \leftarrow \mathcal{T} \setminus \text{cover}(R)$ 
  end while
end for

```

hogy a kivételeket jelentő tanító pontokat lefedje. A komoly szabályokat a futás végére hagyja. Az információnyereség-alapú módszer fordítva működik, a speciális eseteket a végére hagyja.



**Weka 3.5.7** *A wekában a Prism módszert a weka.-  
classifiers.rules.Prism osztály implementálja.*

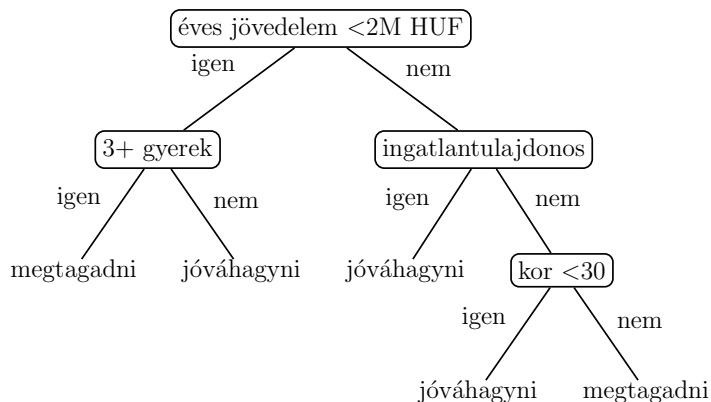


## 6.7. Döntési fák

A döntési fák alapötlete, hogy bonyolult összefüggéseket egyszerű döntések sorozatára vezet vissza. Egy ismeretlen minta klasszifikálásakor a fa gyökeréből kiindulva a csomópontokban feltett kérdésekre adott válaszoknak megfelelően addig lépkedünk lefelé a fában, amíg egy levélbe nem érünk. A döntést a levél címkéje határozza meg. Egy hipotetikus, leegyszerűsített, hitelbírálatra alkalmazható döntési fát mutat be a 6.12. ábra.<sup>3</sup>

A döntési fák nagy előnye, hogy automatikusan felismerik a lényegtelen változókat. Ha egy változóból nem nyerhető információ a magyarázott változóról, akkor azt nem is tesztelik. Ez a tulajdonság azért előnyös, mert így a fák teljesítménye zaj jelenlétében sem romlik, valamint a problémamegértésünket is nagyban segíti, ha megtudjuk, hogy mely változók fontosak, és me-

<sup>3</sup>Az ábrázolt döntési fa sem értékítéletet, sem valós hitelbírálati szabályokat nem tükröz, pusztán illusztráció.



6.12. ábra. Döntési fa hitelbírálatra

lyek nem. Általában elmondható, hogy a legfontosabb változókat a fa a gyökér közelében teszteli. További előny, hogy a döntési fák nagyméretű adathalmazokra is hatékonyan felépíthetők.

A döntési fák egyik fontos tulajdonsága, hogy egy csomópontnak mennyi gyermeke lehet. Nyilvánvaló, hogy egy olyan fa, amely pontjainak kettőnél több gyermeke is lehet mindig átrajzolható bináris fává. A legtöbb algoritmus ezért csak bináris fát tud előállítani.



**Weka 3.5.7** A döntési fákkal kapcsolatos módszereket a `weka.classifiers.trees` csomagban találjuk. A `Classifier output` ablakban a döntési fát szövegesen megjelenítve láthatjuk, amennyiben nem kapcsoljuk ki a `Classifier evaluation options` panelen az `Output model` kapcsolót. A döntési fa grafikus megjelenítéséhez jobb gombbal klikkeljünk a `Result list` ablakban a megfelelő elemre és válasszuk a `Visualize tree` lehetőséget.



### 6.7.1. Döntési fák és döntési szabályok

A döntési fák előnyös tulajdonsága, hogy a gyökérből egy levélbe vezető út mentén a feltételeket összeolvasva könnyen értelmezhető szabályokat kapunk a döntés meghozatalára, illetve hasonlóan egy laikus számára is érthető módon azt is meg tudjuk magyarázni, hogy a fa miért pont az adott döntést hozta.

**6.6. észrevétel.** A döntési fákból nyert döntési szabályhalmazok egyértelműek.

Ez nyilvánvaló, hiszen tetszőleges objektumot a fa egyértelműen besorol valamelyik levelébe. E levélhez tartozó szabályra az objektum illeszkedik, a többire nem.

Vannak olyan döntési feladatok, amikor a döntési fák túl bonyolult szabályokat állítanak elő. Ezt egy példával illusztráljuk.

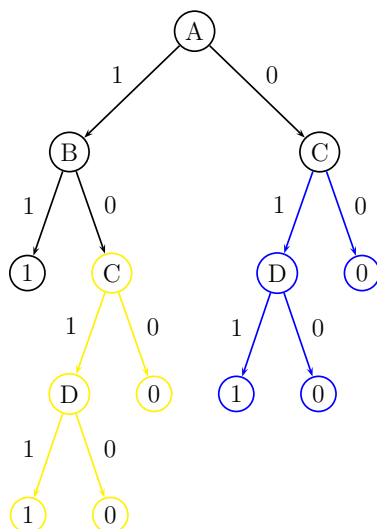
**6.7. példa.** Jelöljük a négy bináris magyarázó attribútumot  $A, B, C, D$ -vel. Legyen az osztályattribútum is bináris és jelöljük  $Y$ -nal. Álljon a döntési szabálysorozat három szabályból:

I.  $A=1$  AND  $B=1 \rightarrow Y=1$

II.  $C=1$  AND  $D=1 \rightarrow Y=1$

III.  $\rightarrow Y=0$

A szabálysorozat teljes, hiszen az utolsó, feltétel nélküli szabályra minden objektum illeszkedik. A fenti osztályozást a 6.13 ábrán látható döntési fa adja.

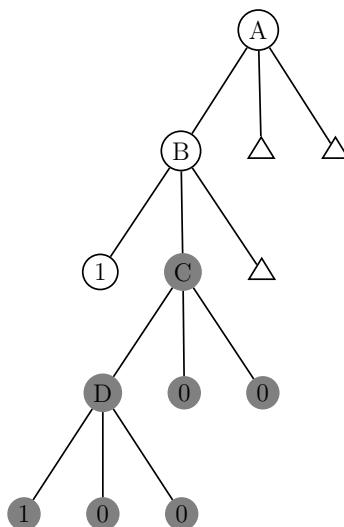


6.13. ábra. Példa adott döntési sorozattal ekvivalens döntési fa

A fenti példában a döntési fa az osztályozás bonyolultabb leírását adja, mint a szabálysorozat. A sárga és kék részfák izomorfak. A részfa által adott osztályozást egyszerűen tudjuk kezelni a döntési szabálysorozatokkal, de a részfák ismételt felrajzolása nem elkerülhető döntési fák esetében. Ezt a problémát az irodalom *ismétlődő részfa problémaként* (*replicated subtree problem*) emlegeti és a döntési fák egy alapproblémájának tekinti. A döntési fák a megoldást nagymértékben elbonyolíthatják. Az előző példában, ha a magyarázó attribútumok nem binárisak, hanem három értéket vehetnek fel, akkor a megadott döntési sorozattal ekvivalens döntési fa a 6.14 ábrán látható. Az a részfa, amelynek pontjait szürkével jelöltük megismétlődik háromszor. Az ismétlődő részfát egy háromszöggel helyettesítettük az áttekinthetőség érdekében. Természetesen a fa jóval egyszerűbb lenne, ha az attribútumot nem csak egy értékkel hasonlíthatnánk össze, hanem olyan tesztet is készíthetnénk, hogy az adott attribútum benne van-e egy adott értékhalmban. Például a gyökérben csak kétfelé célszerű ágazni, attól függően, hogy  $A = 1$  vagy  $A \neq 1$  (másképp  $A \in \{2,3\}$ ). Ha ilyen feltételeket megengednénk, akkor a 6.13 ábrán látható fával izomorf fát kapnánk (ha a címkéket nem vesszük figyelembe).

### 6.7.2. A döntési fa előállítása

A fát a tanító adatbázisból rekurzívan állítjuk elő. Kiindulunk a teljes tanító adatbázisból és egy olyan kérdést keresünk, aminek segítségével a teljes tanulóhalmaz jól szétvágható. Egy szétvágást akkor tekintünk jónak, ha a magyarázandó változó eloszlása a keletkezett részekben



6.14. ábra. Az ismétlődő részfaprobléma szemléltetése

kevésbé szórt, kevésbé bizonytalan, mint a szétvágás előtt. Egyes algoritmusok arra is törekednek, hogy a keletkező részek kb. egyforma nagyok legyenek. A részekre rekurzívan alkalmazzuk a fenti eljárást. Egy csomópont leszármazottjaiban nem vizsgáljuk többé azt az attribútumot, ami alapján szétosztjuk a mintát.

A rekurziót akkor szakítjuk meg valamelyik ágban, ha a következő feltételek közül teljesül valamelyik:

- Nincs több attribútum, ami alapján az elemeket tovább oszthatnánk. A csomópontához tartozó osztály ekkor az lesz, amelyikhez a legtöbb tanítópont tartozik.
- Az adott mélység elért egy előre megadott korlátot.
- Nincs olyan vágás, amely javítani tudna az aktuális osztályzáson. A vágás jóságáról később szólnunk.

Minden levélhez hozzá kell rendelnünk a magyarázandó változó egy értékét, a döntést. Ez általában az ún. többségi szavazás elve alapján történik: az lesz a döntés, amely kategóriába a legtöbb tanítóminta tartozik. Hasonló módon belső csomópontokhoz is rendelhetünk döntést.




---

**Weka 3.5.7** A döntési fa interaktív előállítását teszi lehetővé a `weka.classifiers.trees.UserClassifier` osztály. A módszer elindítása után egy ablak jelenik meg amelynek két füle van. A *Tree Visualizer* fülön az aktuális fát láthatjuk, a *Data Visualizer* pedig a kijelölt fa csomópontjának tanítópontjai jeleníti meg. Itt állíthatjuk elő a vágási függvényt, amelyhez vizuális segítséget kapunk. Az osztály eloszlását láthatjuk két tetszőlegesen kiválasztható attribútum értékeinek függvényében. Ez alapján kijelölhetünk egy téglalapot, poligont vagy összekötött szakaszokat, amely kettéválasztja a pontokat. Akkor jó a kettéválasztás, ha az osztályattribútum szerint homogén csoportok jönnek létre.

---



A döntési fák előállítására a következő három fő algoritmus család ismert:

- I. Iterative Dichotomizer 3 (ID3) család, jelenlegi változat C5.0<sup>4</sup>
- II. Classification and Regression Trees (CART)<sup>5</sup>
- III. Chi-squared Automatic Interaction Detection (CHAID)<sup>6</sup>

### 6.7.3. Az ID3 algoritmus

Az ID3 az egyik legősibb és legismertebb osztályzó algoritmus. A tesztattribútum kiválasztásához az entrópia csökkenését alkalmazza. Ha  $Y$  egy  $\ell$  lehetséges értéket  $p_i$  ( $i = 1, \dots, \ell$ ) valószínűséggel felvevő valószínűségi változó, akkor  $Y$  Shannon-féle entrópiáján a

$$H(Y) = H(p_1, \dots, p_k) = - \sum_{j=1}^l p_j \log_2 p_j$$

számot értjük<sup>7</sup>. Az entrópia az információ-elmélet (lásd [32]) központi fogalma, és  $Y$  változó értékével kapcsolatos bizonytalanságunkat fejezi ki. Ha egy  $X$  változót megfigyelünk és azt tapasztaljuk, hogy értéke  $x_i$ , akkor  $Y$ -nal kapcsolatos bizonytalanságunk

$$H(Y|X = x_i) = - \sum_{j=1}^k \mathbb{P}(Y = y_j | X = x_i) \log_2 \mathbb{P}(Y = y_j | X = x_i)$$

nagyságú. Így ha lehetőségünk van  $X$ -et megfigyelni, akkor a várható bizonytalanságunk

$$H(Y|X) = \sum_{i=1} \mathbb{P}(X = x_i) H(Y|X = x_i)$$

Eszerint  $X$  megfigyelésének lehetősége a bizonytalanság

$$I(Y, X) = H(Y) - H(Y|X)$$

csökkenését eredményezi, azaz  $X$  ennyi információt hordoz  $Y$ -ról. Az ID3 az  $Y$  attribútum szerinti klasszifikálásakor olyan  $X$  attribútum értékei szerint ágazik szét, amelyre  $I(Y, X)$  maximális, azaz  $H(Y|X)$  minimális.




---

**Weka 3.5.7** A wekában az *Id3* algoritmust a *weka.classifiers.treea.Id3* osztály implementálja.

---



<sup>4</sup>Magyarul: Interaktív tagoló / felosztó

<sup>5</sup>Klasszifikáló és regressziós fák

<sup>6</sup>Khi-négyzet alapú automatikus interakció felismerés

<sup>7</sup>Az entrópia képletében  $0 \cdot \infty$  megállapodás szerint 0-val egyenlő.

A kölcsönös entrópia azokat az attribútumokat „kedveli”, amelyek sok értéket vesznek fel és így sokfelé ágazik a fa [111]. Ez terebélyes fákat eredményez. Gondoljuk meg, ha a kiértékelésbe bevesszük az azonosító kódot, akkor az 0 kölcsönös entrópiát fog produkálni, így az algoritmus azt választaná. Hasonló a probléma az 1R módszer diszkretizálásánál említettel (lásd ?? oldal). Egy lehetséges megoldás a nyereségarány mutató (gain ratio) használata [113], amelyre mint normált kölcsönös információ tekintünk. Ez a mutató figyelembe veszi a gyerek csomópontokba kerülő tanítópontok számát és „bünteti” azokat az attribútumokat, amelyek túl sok gyereket hoznak létre. A nyereségarányt úgy kapjuk meg, hogy a kölcsönös információt elosztjuk, az adott attribútum entrópiájával:

$$\text{gain\_ratio}(X) = \frac{I(Y, X)}{H(X)}.$$

Sajnos a nyereségarány sok esetben „túlkompenzál” és olyan attribútumokat részesít előnyben, amelynek az entrópiája kicsi. Egy általános gyakorlat, hogy azt az attribútumot választják, amelyik a legnagyobb nyereségarányt adja, azon attribútumok közül, amelyekhez tartozó kölcsönös információ legalább akkora mint az összes vizsgált attribútumhoz tartozó kölcsönös információk átlaga.

#### 6.7.4. Feltételek a csomópontokban

Az ID3 algoritmus kiválasztja a minimális feltételes entrópiával rendelkező attribútumot és annyi gyerekcsomópont jön létre, amennyi értéket felvesz az attribútum. Leállási feltételként szerepel, hogy egy ágat nem vágunk tovább, ha nincs több vizsgálható attribútum, azaz a fa maximális mélysége megegyezik az attribútumok számával. Az ID3 algoritmus nem feltétlenül bináris fát állít elő.

Ha bináris fa előállítás a cél (továbbá az intervallum típusú attribútum szofisztikáltabb kezelése), akkor a magyarázó  $X$  attribútum típusától függően kétféle feltételt szokás létrehozni. Sorrend típus esetében  $X \geq c$ , ahol  $c$  egy olyan érték, amelyet az  $X$  felvesz valamelyik tanítópont esetén. Intervallum típusú attribútumoknál a  $c$  két szomszédos tanítóérték átlaga. Kategória típus esetében  $X \subseteq K$ , ahol  $K$  az  $X$  értékészletének egy részhalmaza. Az első esetben  $X$  felvett értékeivel lineárisan arányos feltételes entrópiát kell számítani, a másodikban pedig a felvett értékek számával exponenciális számút (ugyanis egy  $n$  elemű halmaznak  $2^n$  darab részhalmaza van).

Sok esetben akkor kapunk jó bináris döntési fát, ha egy gyökérből levélig vezető úton egy attribútumot többször is vizsgálunk (különböző konstansokkal). A fa mélysége ekkor az attribútumok számánál jóval nagyobb is lehet.

#### 6.7.5. Vágási függvények

Miért pont a kölcsönös információt használja az ID3 algoritmus? Milyen jó tulajdonsággal rendelkezik a kölcsönös információ? Van egyéb vágási függvény, amely rendelkezik ezekkel a jó tulajdonságokkal? A válaszok kulcsa a *Taylor-Silverman elvárások* (impurity-based criteria) és a *vágások jósága*.

**6.8. definíció.** *Legyen  $X$  egy olyan diszkrét valószínűségi változó, amely  $k$ -értéket vehet fel. Az eloszlásfüggvény értékei legyenek  $P = (p_1, p_2, \dots, p_k)$ . A  $\Phi : [0,1]^k \mapsto R$  vágási függvényrel szemben támasztott Taylor-Silverman elvárások a következők:*

I.  $\Phi(P) \geq 0$

II.  $\Phi(P)$  akkor veszi fel a minimumát, ha  $\exists j : p_j = 1$

III.  $\Phi(P)$  akkor veszi fel a maximumát, ha  $\forall j : p_j = 1/k$

IV.  $\Phi(P)$  a  $P$  komponenseire nézve szimmetrikus, tehát a  $p_1, p_2, \dots, p_k$  értékek tetszőleges permutációjára ugyanazt az értéket adja.

V.  $\Phi(P)$  differenciálható az értelmezési tartományában mindenhol

Adott  $\mathcal{T}$  tanítóminta esetén a vágási függvény számításakor a  $p_j$  valószínűséget nem ismerjük, így a relatív gyakorisággal közelítjük azaz, ha a  $j$ -edik osztályba tartozó tanítópontok halmazát  $\mathcal{T}^j$ -vel jelöljük, akkor  $p_j = \frac{|\mathcal{T}^j|}{|\mathcal{T}|}$ . A valószínűségvektor empirikus megfelelőjét  $P(\mathcal{T})$ -vel jelöljük ( $P(\mathcal{T}) = (\frac{|\mathcal{T}^1|}{|\mathcal{T}|}, \frac{|\mathcal{T}^2|}{|\mathcal{T}|}, \dots, \frac{|\mathcal{T}^k|}{|\mathcal{T}|}$ ).

**6.9. definíció.** Az olyan  $V$  vágás jósága, amely során a  $\mathcal{T}$  tanítópontokat  $\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_\ell$  diszjunkt tanítóhalmazba osztjuk szét, megegyezik a

$$\Delta\Phi(V, \mathcal{T}) = \Phi(P(\mathcal{T})) - \sum_{i=1}^{\ell} \frac{|\mathcal{T}_i|}{|\mathcal{T}|} \cdot \Phi(P(\mathcal{T}_i))$$

értékkel.

Minél nagyobb egy vágási függvény, annál jobb a vágás. Adott vágási függvény és tanítóponthalmaz esetén célunk megtalálni azt a vágást, amely a maximális vágást eredményezi. Mivel a  $\Phi(P(\mathcal{T}))$  adott tanítóhalmaz esetén rögzített, ezért elég a  $\sum_{i=1}^{\ell} \frac{|\mathcal{T}_i|}{|\mathcal{T}|} \cdot \Phi(P(\mathcal{T}_i))$  értéket minimumát megtalálni.

Amennyiben a vágási függvény csak az osztályok relatív gyakoriságát veszi figyelembe, akkor a vágás jósága 0 lesz abban az esetben, ha az osztályok eloszlása a gyerekekben megegyezik a szülőben található osztályeloszlással. Ez megfelel elvárásainknak, nem nyerünk semmit az olyan vágással, amely során az egyes osztályba tartozó pontok relatív száma egymáshoz viszonyítva mit sem változik.

Most már látható Taylor és Silverman miért fogalmazta meg az elvárásait. A lényeg a második és a harmadik elvárás. Azt szeretnénk, hogy a gyermekekben található tanítóminták minél homogénebbek legyenek. Ideális esetben olyan gyerekek jönnek létre, amelyekhez tartozó tanítópontok egy osztályba tartoznak. Ehhez az osztályhoz tartozó relatív gyakoriság 1, a többi osztályé 0 és a vágási függvény a minimumát veszi fel. A legrosszabb esetben az összes osztály relatív gyakorisága megegyezik, azaz a vágás során olyan gyerek jött létre, amelyben az osztályattribútum teljesen megjósolhatatlan. A harmadik elvárás szerint ezt az esetet büntetni kell, pontosabban a vágási függvény vegye fel a minimumát. Értelemszerűen a minimum és a maximum között a vágási függvény „normális és kezelhető” legyen, azaz legyen deriválható legalábbis minden pontban.

Nem meglepő, hogy az entrópia teljesíti az öt feltételt.

**6.10. lemma.** Az entrópia, mint vágási függvény, megfelel a Taylor-Silverman elvárásoknak [112].



Különböző kutatók különböző vágási függvényeket vezettek be. Például a CART algoritmusban a Gini indexet [21, 47] használták:

$$Gini(\mathcal{P}) = 1 - \sum_{j=1}^k p_j^2.$$

A DKM vágási függvényt [34][72] bináris osztályozás esetén ajánlják:

$$DKM(\mathcal{P}) = 2 \cdot \sqrt{p_1 p_2}$$

**6.11. lemma.** *A Gini és a DKM vágási függvények megfelelnek a Taylor-Silverman elvárásoknak.*

Elméletileg bizonyították [72], hogy a DKM vágási függvény ugyanakkora hiba mellett kisebb döntési fákat állít elő, mintha entrópia vagy Gini index alapján választanánk ki a vágást.

Itt szeretnénk visszautalni az ID3 algoritmus ismertetése végén leírtakra. Az entrópia alapú vágási függvények azokat a vágásokat részesítik előnyben, amelyek sokfelé vágnak, azaz sok gyereket hoznak létre. Általában is igaz, hogy ha a vágás jóságát a fenti módon definiáljuk és a vágási függvény kielégíti a Taylor-Silverman elvárásokat, akkor olyan vágások jönnek létre, amelyekhez sok gyerek tartozik. Természetesen ez a probléma nem jelentkezik bináris döntési fák esetében. Ott minden belső csúcson pontosan két gyereke van.

A megoldást a vágás jóságának normalizálása jelenti. Például az információnyereség helyett a nyeresóarányt (gain ratio) célszerű használni, amelyet megkapunk, ha az információnyereséget elosztjuk az entrópiával. Általános esetben is hasonlóat teszünk. A [90] cikk szerint a vágás jóságának normáltját a következőképpen célszerű képezni:

$$\|\Delta\Phi(V, \mathcal{T})\| = \frac{\Delta\Phi(V, \mathcal{T})}{-\sum_{i=1}^{\ell} \sum_{j=1}^k p_{ij} \log p_{ij}},$$

ahol  $p_{ij} = |\mathcal{T}_i^j|/|\mathcal{T}|$ . Az  $\mathcal{T}_i^j$  az  $i$ -edik gyermek  $j$  osztályba tartozó tanítópontjainak halmazát jelöli.

### 6.7.6. Továbbfejlesztések

Míg az ID3 családba tartozó fák csak klasszifikációra, addig a CHAID és a CART klasszifikációra és előrejelzésre is alkalmazható. A C4.5 (amelynek kereskedelmi, javított változata a C5.0) és a CHAID fák kizárólag egyetlen attribútumra vonatkozó egyenlő, kisebb, nagyobb tesztek használják a csomópontokban a döntésekhez (egyváltozós fák), azaz a jellemzők terét téglatestekre vágják fel. A CART fák ferdén is tudnak vágni, attribútumok lineáris kombinációját is tesztelik (relációs fák). Míg a CART eljárás mindig bináris döntéseket használ a csomópontokban, addig egy nominális attribútumra egy C4.5 fa annyi felé ágazik, ahány lehetséges értéket az attribútum felvehet.

Talán a leglényegesebb különbség a különböző fák között, hogy mit tekintenek jó döntésnek, vágásnak. Nominális magyarázott változó esetén a CHAID eljárás – nevének megfelelően – a  $\chi^2$ -tesztet használja. A CART metodológia a Gini-indexet minimalizálja. A Gini-index alapján mindig olyan attribútumot keresünk, amely alapján a legnagyobb homogén osztályt tudjuk leválasztani.

Ha a magyarázandó  $Y$  változó intervallum skálán mért, akkor a CART eljárás egyszerűen a varianciájának csökkentésére törekszik, a CHAID pedig  $F$ -tesztet használ.

A CHAID konzervatív eljárás, csak addig növeli a fát, amíg a csúcsban alkalmazható legjobb szétvágás  $\chi^2$ -, vagy  $F$ -teszt szerinti szignifikanciája meghalad egy előre adott küszöböt. A CART és C4.5 eljárások nagyméretű fát építenek, akár olyat is, amelyik tökéletesen működik a tanuló adatbázison vagy olyan heurisztikus leállási szabályokat alkalmaznak, hogy a fa nem lehet egy előre adott korlátnál mélyebb, vagy hogy egy csúcsot nem szabad már szétvágni, ha egy korlátnál kevesebb eset tartozik bele. Mindenesetre a kialakuló fa nagy és terebélyes lesz, túl speciális, amely nem csak az alappopuláció jellemzőit, hanem a mintában előforduló véletlen sajátosságokat is modellezi. Ezért a fát felépítése után egy ellenőrző adatbázist használva megszozták metszeni (pruning) és elhagyják a felesleges döntéseket.

Tanácsos megvizsgálni, hogy nem fordul-e elő, hogy a generált C5.0 vagy CHAID fa egymás után ismételten kevés (2-3) attribútum értékét teszteli. Ez arra utalhat, hogy az attribútumok valamely függvénye (pl.: hányadosa - egy főre eső jövedelem) bír magyarázó erővel és a fa ezt a kapcsolatot próbálja ismételt vagdosással közelíteni.




---

**Weka 3.5.7** A C4.5 egy továbbfejlesztett változatának java implementációja a `weka.classifiers.trees.J48` osztály. Talán ez a legelismertebb döntési fa előállító módszer a wekában.

---



Láttuk, hogy többféle mutatószám létezik a vágási kritérium kiválasztására. Ezek között nem létezik a legjobb. Bármelyikhez lehet készíteni olyan adatbázist, amelyet rosszul osztályoz a vágási kritériumot használó algoritmus. A következőkben néhány ismert vágási függvény egységes leírását mutatjuk be.

### 6.7.7. Súlyozott divergenciafüggvények alapján definiált vágási függvények

Bináris vágási függvények esetén a szülő csomópont  $N$  tanító pontját osztjuk szét úgy, hogy a bal oldali gyerekebe  $N_b$  tanító pont jut, a jobboldaliba pedig  $N_j$ . Az  $N_i, i \in \{B, J\}$  pontból  $N_{ji}$  tartozik a  $j$ -edik osztályba. Legyen  $\pi_i = N_i/N$  és  $p_{ji} = N_{ji}/N$ . A  $j$ -edik osztály gyakoriságát a szülőben  $p_j$ -vel jelöljük.

A fenti jelölésekkel a  $\chi^2$  statisztika átírható az alábbi formába:

$$\chi^2/N = \pi_B \sum_{j=1}^k p_{jB}(p_{jB}/p_j - 1) + \pi_J \sum_{j=1}^k p_{jJ}(p_{jJ}/p_j - 1)$$

Legyen  $\mathbf{u} = (u_1, u_2, \dots, u_k)$  és  $\mathbf{v} = (v_1, v_2, \dots, v_k)$  két diszkrét eloszlásfüggvény. Amennyiben a divergencia-függvényüket az alábbi módon definiáljuk

$$d(\mathbf{u} : \mathbf{v}) = \sum_{j=1}^k u_j(u_j/v_j - 1),$$

akkor a  $\chi^2$  statisztika átírható a következőképpen (alkalmazzuk a  $u_j(u_j/v_j - 1) = 0$  konvenciót  $u_j = v_j = 0$  esetén):

$$\chi^2 = N(\pi_B d(p_B : \mathbf{p}) + \pi_J d(p_J : \mathbf{p})).$$

Ha a divergenciafüggvénynek a következőt használjuk

$$d(\mathbf{u} : \mathbf{v}) = 2 \sum_{j=1}^k u_j \log(u_j/v_j),$$

akkor az entrópiához jutunk. Továbbá  $d(\mathbf{u} : \mathbf{v}) = 2 \sum_{j=1}^k (u_j^2 - v_j^2)$  esetén a Gini index  $N$ -edrészét kapjuk.

A közös magot az *erő divergencia függvény* adja [126]:

$$d_\lambda(\mathbf{u} : \mathbf{v}) = \frac{1}{\lambda(\lambda+1)} \sum_{j=1}^k u_j ((u_j/v_j)^\lambda - 1),$$

ahol  $-1 < \lambda \leq \infty$ . A  $d_\lambda$  függvény értékét a  $\lambda = 0$  helyen, az ugyanitt vett határértéke adja  $d_\lambda$ -nak. Az erő divergencia függvény alapján definiáljuk a vágási függvények egy családját:

$$C(\lambda) = \pi_B d_\lambda(p_B : \mathbf{p}) + \pi_J d_\lambda(p_J : \mathbf{p})$$

Láttuk, hogy  $\lambda = 1$  estén a  $\chi^2$  statisztikát kapjuk,  $\lambda = 0$ -nál pedig az entrópiát. További ismert vágási függvényeket is megkaphatunk az erő divergencia függvényből. Freeman-Tuckey statisztika adódik  $\lambda = -1/2$ -nél és a Cressie-Read  $\lambda = -2/3$ -nál [114].

**6.12. tétel.** *A  $C(\lambda)$  vágási függvényosztályba tartozó vágási függvények teljesítik a Taylor-Silverman elvárásokat.*

Ismert vágási függvény az MPI index, amelyet az alábbi módon definiálnak:

$$M = \pi_B \pi_J \left( 1 - \sum_{j=1}^k p_{jB} \cdot p_{jJ} / p_j \right)$$

Egy kis kézimunkával az MPI index átalakítható az alábbi formára:

$$M = \pi_B 2\pi_J^2 d_1(p_J : \mathbf{p}) + \pi_J 2\pi_B^2 d_1(p_B : \mathbf{p}),$$

amely a  $D(\lambda) = \pi_B 2\pi_J^2 d_\lambda(p_J : \mathbf{p}) + \pi_J 2\pi_B^2 d_\lambda(p_B : \mathbf{p}) = \pi_B \pi_J C(\lambda)$  vágási függvényosztály tagja. Szerencsére ez a függvényosztály is rendben van az elvárásaink tekintetében:

**6.13. tétel.** *A  $D(\lambda)$  vágási függvényosztályba tartozó vágási függvények teljesítik a Taylor-Silverman elvárásokat.*

### 6.7.8. Döntési fák nyesése

A döntési fák nyesése során a felépített fát kicsit egyszerűsítjük, azaz belső csomópontokat távolítunk el. Feltételezzük, hogy a felépítés során a fa túltanult, tehát megtanult olyan esetiségeket is, amelyek csak a tanítóhalmazra jellemző. A nyesést ezért egy külön teszthalmazon szokás elvégezni, vagy ha felhasználtunk minden elérhető tanítópontot a fa felépítése során, akkor használhatunk keresztvalidációt is (lásd a ?? rész).

Beszélhetünk *előnyesésről* (prepruning) és *utónyesésről*. Az előnyesés nem más mint egy intelligens – általában statisztikai megfontolásokon alapuló – STOP feltétel. Habár lenne olyan vágási függvény, amely megfelel minden feltételnek (és nem értük el az előre megadott maximális mélységet) mégsem osztjuk tovább a tanítópontokat egy új feltétel bevezetésével.

Az utónyesés során nagy fát növesztünk, majd elkezdjük azt zsugorítani. A kutatók többsége az utónyesést tartja jobb megoldásnak. Gondoljuk meg, ha a bináris magyarázandó attribútum két teljesen véletlen bináris magyarázó attribútum moduló kettes összege, akkor egy attribútum alapján végzett osztályozás statisztikailag semmivel sem jobb, mint nulla attribútum alapján végzet osztályozás, így az előnyesés nem engedélyez semmilyen feltételt a gyökérbe. Az utónyesés során viszont tökéletes osztályozást lehet létrehozni, hiszen az először nagy méretű fát növeszt, amelyben mindkét attribútum előfordulhat.

A két legismertebb utónyesési eljárás a *részfa helyettesítés* (subtree replacement) és a *részfa felhúzás* (subtree raising).

A részfa helyettesítés során egy belső pontból induló, minden útjában levélig érő fát egyetlen levéllel helyettesítünk. Kérdés, hogy a nyesett fa jobban osztályoz-e, mint az eredeti fa. Osztályozók összehasonlításáról a 6.11 részben írunk. Az összehasonlításnál elég azokat a tesztpontokat figyelembe venni, amelyeket az döntési fa a részfa gyökerébe sorol. Ezen tanítópontokon kell összehasonlítani a 0R osztályozót (csak egy gyökérpontot tartalmazó döntési fa tulajdonképpen egy 0R osztályozó) és a részfát, mint külön döntési fát.

FOLYT. KÖV.

### 6.7.9. Döntési fák ábrázolása

A döntési fa előállítás után két fontos kérdés szokott felmerülni. Egyrészt tudni szeretnénk, hogy melyik levélbe esik sok tanító pont, azaz melyek azok a szabályok, amelyek sok tanító pontra érvényesek. Másrészt látni szeretnénk, hogy a levelek mennyire jól osztályoznak; a tesztpontok közül (ha vannak tesztpontok) milyen arányban osztályozott rosszul az adott levél. Az első kérdés tehát azt vizsgálja, hogy mennyire jelentős az adott levél, a második pedig azt, hogy mennyire jó, mennyire igaz a levélhez tartozó szabály. Ezeket az értékeket azonnal látni szeretnénk, ha ránézünk egy döntési fára.

Elterjedt módszer (ezt használják például a SAS rendszerében is), hogy minden levelet egy körkikkely reprezentál. A körkikkely nagysága arányos a levélhez tartozó tanító pontokkal, a színe pedig a levélhez tartozó szabály jóságát adja meg. Például minél sötétebb a szín, annál rosszabb az osztályozás aránya. Egy ilyen ábrázolásra láthatunk példát a következő ábrán.

FOLYT. KÖV.

„ A grapefruit mindennapos fogyasztása harmadával növelheti a mellrák veszélyét – figyelmeztetnek amerikai kutatók.”  
 Forrás: <http://www.macosz.hu/grapefruit-daganat.html>

### 6.7.10. Hanyag döntési fák

A hanyag döntési fák olyan döntési fák, amelyben az azonos szinten elhelyezkedő pontokban ugyanazt az attribútumot vizsgáljuk.

FOLYT. KÖV.

## 6.8. Bayesi hálózatok

A Bayes hálózatok két fontos elvre építenek. A maximum likelihood szerint egy elem osztályozásánál azt az osztályt fogjuk választani, amelynek a legnagyobb a valószínűsége a megfigyelések és az elem további attribútumai alapján. A bayes-tétel szerint pedig meghatározhatjuk a feltételes valószínűséget, ha ismerünk néhány másik valószínűséget.



**Weka 3.5.7** Számos bayes-háló alapú módszer található a `weka.classifiers.bayes` csomagban.



A Bayes-tétel segítségével meghatározható az optimális (lásd 6.2. szakasz) klasszifikációs szabály. Jelöljük  $Y_i$ -vel azt, amikor a klasszifikálandó eset az  $i$ -edik osztályba tartozik ( $Y = y_i$ ). Az elemek megfigyelhető tulajdonságait az  $X$  vektor írja le. Az egyszerűség kedvéért a tévedés költsége legyen minden esetben azonos. Ekkor egy ismeretlen,  $X$  tulajdonságú példányt abba az osztályba ( $i$ ) érdemes (optimális) sorolni, amelyekre  $\mathbb{P}(Y_i|X)$  maximális. A Bayes-szabály alapján

$$\mathbb{P}(Y_i|X) = \frac{\mathbb{P}(X, Y_i)}{\mathbb{P}(X)} = \frac{\mathbb{P}(X|Y_i) \mathbb{P}(Y_i)}{\mathbb{P}(X)}.$$

Mivel  $\mathbb{P}(X)$  minden  $i$ -re konstans, ezért elegendő  $\mathbb{P}(X|Y_i) \mathbb{P}(Y_i)$ -t maximalizálni.  $\mathbb{P}(Y_i)$  vagy a priori adott, vagy pedig a mintából a relatív gyakoriságokkal egyszerűen becsülhető. Így már „csak”  $\mathbb{P}(X|Y_i)$ -t kéne meghatározni.

Amennyiben  $k$  darab bináris magyarázó attribútumunk van, az  $Y$  pedig  $\ell$  értéket vehet fel, akkor  $\ell(2^k - 1)$  darab  $\mathbb{P}(X|Y_i)$  értéket kellene megbecsülnünk. A 3.3.7 részben láttuk, hogy egy valószínűség megbecsléséhez relatív gyakorisággal mennyi tanítópontot kell vennünk. A gyakorlati esetek többségében ennyi tanítópont nem áll rendelkezésünkre, ezért valamilyen feltétellel kell élnünk a modell kapcsán. A naív bayes-hálók feltételezik, hogy az egyes attribútumok feltételesen függetlenek egymástól.

### 6.8.1. Naív Bayes-hálók

A naív bayes-hálók olyan feltételezéssel élnek, amelynek segítségével a  $\ell(2^k - 1)$  darab megbecsülendő paraméter száma  $\ell \cdot k$ -ra csökken.

**6.14. definíció.** Legyen  $X, Y$  és  $Z$  három valószínűségi változó. Az  $X$  feltételesen független  $Y$ -től adott  $Z$  esetén, ha

$$\mathbb{P}(X = x_i | Y = y_j, Z = z_k) = \mathbb{P}(X = x_i | Z = z_k)$$

minden lehetséges  $x_i, y_j, z_k$  hármasra.

Ha például az  $es\emptyset$ , vihar, villámlás diszkrét valószínűségi változót tekintjük, akkor a vihar feltételesen független az  $es\emptyset t\emptyset l$ , ha a villámlást ismerjük. A villámlás ugyanis vihart okoz (a villámlás hiánya pedig azt jelenti nincs vihar), ezért az  $es\emptyset$  ténye semmilyen további információval nem szolgál a viharra vonatkozóan. Természetesen van összefüggés a vihar és az  $es\emptyset$  között, de nincs köztük feltételes összefüggés, ha a villámlás értékét ismerjük.

A naív Bayes klasszifikáló feltételezése szerint egy osztályon belül az attribútumok feltételesen függetlenek egymástól. Ekkor a  $\mathbb{P}(X|Y)$  valószínűség kifejezhető a  $\mathbb{P}(X_j|Y)$  valószínűségek szorzataként, hiszen

$$\mathbb{P}(X_1, X_2|Y_i) = \mathbb{P}(X_1|X_2, Y_i) \mathbb{P}(X_2|Y_i) = \mathbb{P}(X_1|Y_i) \mathbb{P}(X_2|Y_i)$$

Az első egyenlőségnél a valószínűségek általános tulajdonságát használtuk fel, a másodiknál pedig a feltételes függetlenséget. Könnyű belátni, hogy  $k$  magyarázó változó esetén a következőt kapjuk

$$\mathbb{P}((X_1, X_2, \dots, X_k) = (x_1, x_2, \dots, x_k) | Y_i) = \prod_{j=1}^k \mathbb{P}(X_j = x_j | Y_i).$$

A  $\mathbb{P}(X_j = x_j | Y_i)$  valószínűségek a mintából becsülhetők.

### kategória típusú attribútum

Amennyiben az  $X_j$  kategória típusú, akkor  $\mathbb{P}(X_j = x_j | Y_i)$  valószínűséget a relatív gyakorisággal közelítjük, tehát meghatározzuk a relatív arányát az  $X_j$  attribútumában  $x_j$  értéket felvevő elemeknek a  $Y_i$  osztályú elemek között.

Problémát jelent, ha valamelyik relatív gyakoriság nulla, mert ekkor a szorzat is nulla lesz a többi tagtól függetlenül. Legegyszerűbb megoldás, hogy az adott attribútum minden értékének előfordulásához hozzáadunk egyet. Ha volt elég mintánk, akkor a valószínűségek alig torzulnak, viszont sikerül kiküszöbölnünk, hogy a nulla tag miatt a többi relatív gyakoriságot nem vesszük figyelembe. Ha egy adott osztályba tartozó elemek egy attribútuma három értéket vehet fel és az előfordulások: 0, 150, 250. Akkor 0, 150/400, 250/400 helyett 1/403, 151/403, 251/403 értékeket használunk. Erre a technikára az irodalomban, mint *Laplace estimation* hivatkoznak. Egy kifinomultabb módszer, ha egy helyett  $p_k$ -t adunk a relatív gyakorisághoz, ahol  $p_k$ -val jelöljük a  $k$ -adik attribútumérték relatív gyakoriságát a teljes tanítóhalmazban (tehát nem csak a  $Y_i$  kategóriájú tanítóhalmazban).

### szám típusú attribútum

Amennyiben  $X_j$  szám típusú és tudjuk a  $\mathbb{P}(X_j|Y_i)$  eloszlásának típusát, akkor a keresett valószínűséghez szükséges eloszlásparámétereket statisztikai módszerrel becsüljük. Ha például normális eloszlással van dolgunk, akkor elég meghatározni a várható értéket és a szórást, ezekből tetszőleges értékhez tartozó valószínűség a sűrűségfüggvényből közvetlen adódik. A várható értéket a mintaátlaggal (empirikus közép:  $\bar{X}_{ij} = \sum_{k=1}^{|Y_i|} x_{ij}^k / |Y_i|$ ), a szórásnégyzetet a korrigált empirikus szórásnégyzettel ( $s_{ij}^{*2} = \sum_{k=1}^{|Y_i|} (x_{ij}^k - \bar{X}_{ij})^2 / (|Y_i| - 1)$ ) becsüljük. A keresett valószínűséget a

$$\mathbb{P}(X_j = x_j | Y_i) = \frac{1}{s_{ij}^* \sqrt{2\pi}} e^{-(x_j - \bar{X}_{ij})^2 / 2s_{ij}^{*2}}$$

képlet adja.



**Weka 3.5.7** `weka.classifiers!bayes.NaiveBayesSimple` A naív Bayes osztályozót, amely a szám típusú attribútumoknál normális eloszlást feltételez a `weka.classifiers.bayes.NaiveBayesSimple` osztály implementálja.



A `weka.classifiers.bayes.NaiveBayes` a normalitásra tett feltételt enyhíti. Ez az osztályozó ún. kernel becslőt használ a keresett valószínűségek meghatározásához. Ha pedig a `useSupervisedDiscretization` paramétert igazra állítjuk, akkor a szám típusú attribútumokat kategória típusúvá alakítja egy felügyelt diszkrétizáló módszerrel (`weka.filters.supervised.attribute.Discretize` szűrő lásd a 48 oldal).

A naiv Bayes osztályozó hátrányra, hogy az feltételes függetlenséget és egyenlőséget feltételez. Sokat javíthatunk a naív Bayes osztályozók pontosságán, ha előfeldolgozás során meghatározzuk a fontos attribútumokat, tehát azokat, amelyekről úgy gondoljuk, hogy nem függetlenek az osztályattribútumtól. Több kutató arról számol be, hogy a megfelelő attribútumkiválasztással párosított naív Bayes osztályozó felveszi a versenyt a bonyolultabb, újabb módszerekkel.

## 6.8.2. Naív Bayes-hálók és a logisztikus regresszió kapcsolata

Ebben a részben belátjuk, hogy amennyiben minden magyarázó attribútum valós típusú, akkor a normális eloszlást feltételező naív bayes osztályozó (GNB – Gaussian Naive Bayes) GNB egy lineáris osztályozó, amely nagyon hasonlít a logisztikus regresszióra.

Foglaljuk össze milyen feltételezésekkel él a GNB:

- Az  $Y$  bináris valószínűségi változó, melynek eloszlása  $p_Y$  paraméterű binomiális eloszlás.
- Minden  $X_j$  magyarázó változó valós típusú.

- $X_j|Y = y_i$  feltételes valószínűségi változó  $\mu_{ji}, \sigma_j$  paraméterű normális eloszlással, tehát

$$\mathbb{P}(X_j = x_j | Y = y_i) = \frac{1}{\sqrt{2\pi\sigma_j^2}} e^{-\frac{(x_j - \mu_{ji})^2}{2\sigma_j^2}}$$

- a magyarázó változók adott  $Y$  esetén feltételesen függetlenek egymástól.

Vegyük észre, hogy az  $X_j|Y = y_i$  feltételes valószínűségi változó szórása attribútumról attribútumra más lehet és nem függ  $Y$ -tól.

Célunk belátni, hogy ezek a feltevések hasonló alakú  $\mathbb{P}(Y|X)$ -t adnak, mint azt a logisztikus regresszió teszi (emlékeztetőként:  $\mathbb{P}(Y = 1|X) = \frac{1}{1 + e^{-x^T w}}$ ). Induljunk ki a bayes szabályból

$$\begin{aligned} \mathbb{P}(Y = 1|X) &= \frac{P(Y = 1)P(X|Y = 1)}{P(Y = 1)P(X|Y = 1) + P(Y = 0)P(X|Y = 0)} \\ &= \frac{1}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} = \frac{1}{1 + \exp \ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}} \end{aligned}$$

most használjuk ki a feltételes függetlenséget:

$$\begin{aligned}\mathbb{P}(Y = 1|X) &= \frac{1}{1 + \exp \ln \frac{P(Y=0)}{P(Y=1)} + \sum_j \ln \frac{P(X_j|Y=0)}{P(X_j|Y=1)}} \\ &= \frac{1}{1 + \exp \ln \frac{1-p_Y}{p_Y} + \sum_j \ln \frac{P(X_j|Y=0)}{P(X_j|Y=1)}}\end{aligned}\quad (6.8)$$

Vizsgáljuk meg közelebbről a szummában szereplő tagot:

$$\begin{aligned}\ln \frac{P(X_j|Y=0)}{P(X_j|Y=1)} &= \ln \frac{\frac{1}{\sqrt{2\pi\sigma_j^2}} \exp -\frac{(X_j-\mu_{j0})^2}{2\sigma_j^2}}{\frac{1}{\sqrt{2\pi\sigma_j^2}} \exp -\frac{(X_j-\mu_{j1})^2}{2\sigma_j^2}} = \ln \exp \frac{(X_j-\mu_{j1})^2 - (X_j-\mu_{j0})^2}{2\sigma_j^2} \\ &= \frac{(2X_j(\mu_{j0}-\mu_{j1}) + \mu_{j1}^2 - \mu_{j0}^2)}{2\sigma_j^2} = \frac{\mu_{j0}-\mu_{j1}}{\sigma_j^2} X_j + \frac{\mu_{j1}^2 - \mu_{j0}^2}{2\sigma_j^2}\end{aligned}$$

Ha ezt visszahelyettesítjük a 6.8 egyenletbe, akkor látható, hogy  $\mathbb{P}(Y = 1|X)$  tényleg az  $X_j$  attribútumok súlyozott összegének nemlineáris függvényeként adódik:

$$\mathbb{P}(Y = 1|X) = \frac{1}{1 + e^{w_0 + x^T w}},$$

ahol a súlyok

$$w_j = \frac{\mu_{j0} - \mu_{j1}}{\sigma_j^2},$$

a torzítás pedig:

$$w_0 = \ln \frac{1-p_Y}{p_Y} + \sum_j \frac{\mu_{j1}^2 - \mu_{j0}^2}{2\sigma_j^2}$$

Összegezzük a hasonlóságokat és a különbségeket a GNB és a logisztikus regresszió között. Legfőbb hasonlóság, hogy mind a két módszer lineáris szeparálást végez, azaz az osztályozáshoz a magyarázó attribútumok súlyozott összegét veszi alapul. Különbség van azonban a súlyok meghatározásában. A logisztikus regresszió közvetlenül becsli a súlyokat, míg a GNB normális eloszlást feltételezve megbecsli a várható értéket és a szórást, majd ez alapján számít egy súlyt. A logisztikus regresszió tehát közvetlenül becsli  $\mathbb{P}(Y|X)$ -et, míg a Bayes osztályozó ezt közvetve teszi,  $\mathbb{P}(Y)$  és  $\mathbb{P}(X|Y)$  becslésével. Be lehet látni, hogy amennyiben fennáll a normalitásra tett feltétele a GNB-nek, akkor a GNB és a logisztikus regresszió ugyanazt az osztályozót (azaz ugyanazokat a súlyokat) eredményezik.

A logisztikus regresszió – mivel nem él semmilyen feltételezéssel az adataira vonatkozóan – egy általánosabb módszernek tekinthető, mint a GNB. Ha nem teljesül a normalitásra tett feltétel, akkor a GNB torz eredményt ad, míg a logisztikus regresszió „adaptálódik a helyzethez”.

A legközelebbi szomszéd módszernél már megtanultuk, hogy az általánosabb módszer nem tekinthető jobb módszernek (ha ez nem így lenne, akkor mindenre a legközelebbi szomszéd módszert használnánk, hiszen ez a legáltalánosabb módszer). Az általános módszerek hajlamosak a túltanulásra. Számos írás született, amely pont a logisztikus regresszió túltanulásának hajlamát próbálja visszafogni különféle büntetőfüggvények bevezetésével. Az általános módszerek



további hátránya, hogy jóval több tanítópontonra van szükségük, mint azoknak, amelyek valamilyen feltételezéssel élnek a háttérben megbújó modellel kapcsolatban.

Nem meglepő ezért, hogy különbség van a tanulás konvergenciájának sebességében is. A logisztikus regresszióknak  $O(n)$  a Bayes hálónak csak  $O(\log n)$  tanítópontonra van szüksége ugyanakkora pontosság eléréséhez (amennyiben a normalitásra tett feltétel teljesül).

### 6.8.3. Bayes hihetőségi hálók

A Bayes hihetőségi hálók (Bayesian belief networks) a függetlenségre tett feltételt enyhítik. Lehetővé teszik az adatbányásznak, hogy egy irányított, körmentes gráf segítségével a változók közötti függőségi struktúrát előre megadja. A gráf csomópontjai megfigyelhető és nem megfigyelhető, de feltételezett (rejtett) változók lehetnek. Úgy gondoljuk, hogy a gráf a függőségeket jól leírja, azaz

$$\mathbb{P}((Z_1, Z_2, \dots, Z_s) = (z_1, z_2, \dots, z_s)) = \prod_{j=1}^s \mathbb{P}(Z_j = z_j | \text{par}(Z_j))$$

teljesül, ahol  $\text{par}(Z_j)$  a  $Z_j$  csúcs szüleit (a gráfban közvetlenül belemutató csúcsok halmazát jelöli). Minthogy a háló struktúrája a teljes eloszlást leírja, ezért tetszőleges  $Z_j$  csúcsokat kiemelhetünk outputnak / előrejelzendőnek. Ha nincsenek rejtett változók, akkor a szükséges

$$\mathbb{P}(Z_j = z_j | \text{par}(Z_j))$$

valószínűségek közvetlen becsülhetők a mintából. Ha a háló rejtett változókat is tartalmaz, akkor a gradiens módszer egy változata alkalmazható. Végül olyan eljárások is ismertek, amelyek segítségével a hálózat topológiája a tanuló példákból kialakítható, nem feltétlenül szükséges azt előre megadni.

## 6.9. Osztályozók kombinálása

### 6.9.1. Bagging

### 6.9.2. Randomizálás

### 6.9.3. Boosting

## 6.10. Osztályozók kiértékelése

Az adatbányászati modellekre – sajnos – ritkán állnak rendelkezésre olyan módszerek, amelyek segítségével az illeszkedés jószágáról statisztikai teszttel dönthetünk (a kivételeket lásd például [18] cikkben). Egy lehetséges általános megközelítést Rissanen adott meg [1]. A „legrövidebb leírás”<sup>8</sup> elve szerint egy adathalmazt magyarázó elméletek közül az a leginkább elfogadhatóbb, amelynél összesen a legkevesebb bit szükséges a modell és az adatoknak a modell segítségével való leírásához.

---

<sup>8</sup>Minimum Description Length, MDL.

Már említettük, hogy a túltanulás miatt nem célszerű a tanítóhalmazt használni az osztályozó pontosságának megállapításához. A tanítóhalmazon számított hibát *resubstitution error*, azaz *visszahelyettesítési hibának* nevezzük.

A leggyakrabban alkalmazott módszer a következő. Adatainkat három részre osztjuk (általában 70%-20%-10% arányban). Ugyanazon tanító adatokon több konkurens modellt építünk, majd az ellenőrző adathalmaz segítségével kiválasztjuk a legjobbat, amelyet alkalmazni fogunk. A végső modell teljesítményét, pedig egy – az előző kettőtől diszjunkt – teszt adatbázison mérjük. Ismételt mintavételezési technikákkal csökkenthetjük a fenti eljárás adatigényét, illetve több klasszifikáló eredményeinek kombinálásával is javítható az előrejelzés pontossága [55].

A három részre osztós technikánál a tanító halmazba csak az adatok 70%-a kerül. Minél kisebb a tanító adathalmaz, annál kevésbé lehetünk biztosak, hogy az osztályozónál nem lépett fel túltanulás. Továbbá minél kisebb az adat, annál kevésbé reprezentatív a rejtett információ, így annál nehezebb megtanulni. A reprezentativitásnál nem elég figyelembe venni az adatok méretét. Bonyolultabb, sok szabályt tartalmazó modelleknek nagyobb tanítóhalmazra van szükségük. Érezzük, hogy kevesebb tanítópontra van szükségünk bináris osztályozás esetén, mint akkor, amikor 20 különböző osztályt hozhatunk létre.

Honnan tudjuk eldönteni, hogy az adathalmazunk egy része reprezentatív-e? Általánosan sehog. Van azonban egy egyszerű vizsgálat, amelyet érdemes elvégezni. A tanító és a tesztelő adathalmazban az egyes osztályok eloszlása nagyjából meg kell egyezzenek. Nem várhatunk jó osztályozást, ha a tanítóhalmazba nem került valamely osztályból egyetlen elem sem. Az eredeti adathalmaz olyan particionálását (tanító és teszthalmazra), amelyre teljesül, hogy az osztályok relatív előfordulása a tanítóhalmazban és a teszthalmazban megegyeznek, *rétegzett (stratified) particionálásnak/mintavételezésnek* hívjuk.

Nem mindig áll rendelkezésünkre annyi adat, hogy a három részre osztás után is azt tudjuk mondani, hogy a tanító adathalmaz elég reprezentatív. Kisebb adathalmazok esetén ismételt mintavételezéssel szoktak segíteni a helyzeten. A következőkben ezeket a technikákat tekintjük át.

Az osztályozók legfontosabb mérőszáma a *hibaarány*, amely a tévesen osztályozott objektumok számát adja meg. A hibaarány inverze a pontosság. Érdekes lehet tudni, hogy mennyi a hibaarány a tanítóhalmazon, de a túltanulás veszélye miatt a hibaarányt máshogy szokás mérni.

### Ismételt mintavételezés

Az eredeti adathalmaz nagyobb részét (általában kétharmadát) válasszuk tanítóhalmaznak, a maradékon határozzuk meg a hibaarányt. Ismételjük többször az eljárást különböző, véletlenszerűen választott tanítóhalmazokon. Az osztályozás végső hibaarányát az egyes felosztásokból származó hibaarányok átlagaként adjuk meg.

### Kereszt-validáció és a leave-one-out

Osszuk fel a tanítóhalmazt  $N$  részre. Az adott osztályozó módszerrel  $N$  különböző tanítást fogunk végezni. Minden tanításnál egy rész lesz a tesztelőhalmaz a többi uniója pedig a tanítóhalmaz. Minden tanításnál más tesztelőhalmazt választunk. A végső hibaarányt megint

az egyes hibarányok átlaga adja. Igen elterjedt (habár elméletileg nem megalapozott), hogy  $N$  értékének 10-et adnak meg.

A kereszt-validáció egy speciális esete, amikor a  $N$  értéke megegyezik a tanítópontok számával, azaz csak egy elemet tesztelünk. Ezt a módszert *leave-one-out*-nak (egy kimarad) hívják. Ennek a módszernek két előnye és két hátránya van. Előny, hogy a módszer teljesen determinisztikus, továbbá a tanításhoz a lehető legtöbb információt használja. Hátrány ugyanakkor, hogy a tanítást sokszor kell elvégezni, ami nagyon költséges lehet, továbbá a teszteléshez használt adathalmaz biztos, hogy nem rétegzett.

Egyes kutatók úgy vélik, hogy a kereszt-validáció jelentősége túl van értékelve, hiszen elméletileg nem lehet bizonyítani, hogy megbízhatóbb eredményt szolgál, mint az egyszerű oszd ketté (taníts, majd tesztelj) módszer.

## Bootstrap

Az eddigi megoldásokban egy tanítópontot egyszer használtunk fel a résztanítások során. A bootstrap visszatevéses mintavételezésen alapul, amely eredményeképpen a tanítóhalmazban (a halmaz szó használata itt most helytelen) ugyanaz az elem többször is előfordulhat. A bemeneti adathalmaz méretét jelöljük  $n$ -nel. A módszer egyszerű. Válasszunk visszatevéses mintavételezéssel  $n$  elemet. Lesznek olyan elemek, amelyeket többször választottunk és lesznek olyanok is, amelyeket egyszer sem. Azok az elemek adják a teszhalmazt, amelyeket egyszer sem választottunk. Annak a valószínűsége, hogy egy elem a tanítóhalmazban lesz, közel 63% nagy  $n$  esetén, hiszen annak valószínűsége, hogy egy elemet nem választunk:

$$\left(1 - \frac{1}{n}\right)^n \rightarrow e^{-1} \approx 0.368.$$

A bootstrap esetében a hibarányt a tanító és a teszhalmazon számított hibarányok súlyozott összege adja, ahol a súlyok  $1 - e^{-1}$  és  $e^{-1}$ :

$$e = (1 - e^{-1})e_{\text{teszt}} + e^{-1}e_{\text{tanító}}$$

A bootstrap nem feltétlenül rétegzett tanító mintát hoz létre. Sőt valószínű, hogy a tanító minta torz lesz, hiszen lesznek olyan tanító pontok, amelyek nem kerülnek bele a tanító halmazba és olyanok is lesznek, amelyek többször is szerepelni fognak. A bootstrapot is többször ismételtjük, különböző mintákkal és a végső hibát az egyes hibák átlagaként számítjuk. Jogos az a kétely a bootstrappal kapcsolatban, hogy torz adatokon torz osztályozók fognak létrejönni. Ezek hibáinak súlyozott átlaga pedig nem feltétlenül közelíti a helyes osztályozás hibáját.




---

**Weka 3.5.7** A wekában az osztályozás kiértékelésének módját a *Test options* panelen adhatjuk meg. *Use training set* esetén a hibát és egyéb paramétereket a tanítóhalmazon mérjük. *Supplied test set* esetén külön teszhalmazt adhatunk meg, *Cross-validation* választásakor kereszt-validációt használunk. A *Folds* paraméterrel adhatjuk meg, hogy hány részre ossza a weka a tanítóhalmazt. Ha a hagyományos tanítóhalmaz, teszhalmaz kettéosztást kívánjuk használni, akkor válasszuk a *Percentage split* opciót. Ilyenkor megadhatjuk tanítóhalmazba kerülő elemek százalékos arányát.

---



### 6.10.1. Értekezés

A fenti módszerek az egyszerű „oszd ketté, taníts majd számíts hibát” megközelítés azon gyenge pontját próbálják orvosolni, hogy a tanítóhalmaz vagy a teszhalmaz (vagy mindkettő) torz lehet, valamely szabály szempontjából. Ebből adódóan hamis hibaarányt fog szolgáltatni. Sajnos a fenti módszerek ugyanúgy adhatnak rossz eredményt. Vegyünk egy egyszerű példát, amelyben bináris osztályozót készítünk, de az adatok teljesen véletlenszerűek nincs semmilyen összefüggés az attribútumok és az osztály között. Döntési fa ebben az esetben egyetlen gyökércsomópontot tartalmazna és minden objektumot abba az osztályba sorolna, amelyikbe több tanítópont tartozik (többségi szavazás). Ha a tanítópontok száma páros és a tanítópontok fele-fele tartozik az egyik ill. a másik osztályba, akkor a leave-one-out értékelő 100%-os hibát állapítana meg az elvárt 50% helyett. Ha az osztályozónk olyan, hogy teljesen megtanulja (megjegyzí) az elem, osztály hozzárendelést, azaz a tanítóhalmazon 100%-os pontosságot produkál, akkor a bootstrap szerint a hiba  $(1 - e^{-1}) \cdot 0.5 + e^{-1} \cdot 0$ , ami az 50%-nál  $e^{-1}/2$ -vel kisebb.

Tegyük fel, hogy a bináris osztályozónk  $p$  valószínűséggel ad helyes eredményt, tehát a pontossága  $p$ . Adott  $N$  tesztpont mellett a helyesen osztályozott pontok számát jelöljük  $f$ -fel. A helyesen osztályozott pontok száma egy  $N, p$  paraméterű, binomiális eloszlású valószínűségi változó. Tetszőleges  $\alpha$  értékhez ( $1 - \alpha$  a próba szintje) meg tudjuk határozni az elfogadási tartományt a helyesen osztályozott pontok számára vonatkozóan (ezt hívják bináris tesztnek, lásd a 2.6.2 rész). Határozzuk meg azt az  $f/N$ -nél kisebb, legkisebb  $p$ -t (jelöljük  $p_l$ -el), amelyre  $f$  az elfogadási tartományba esik. Határozzuk meg ezenkívül azt az  $f/N$ -nél nagyobb, legnagyobb  $p$ -t (jelöljük  $p_u$ -val), amelyre  $f$  az elfogadási tartományba esik. A tesztelő pontok alapján csak azt tudjuk elmondani, hogy az igazi  $p$  a  $[p_l, p_u]$  intervallumba esik.

A fenti módszer meglehetősen számításigényes. A  $p_l, p_u$  értékek meghatározásának nehézsége abból adódik, hogy a  $p$  a valós számok halmazából kerül ki, a valószínűségi események (és így  $f$  és az elfogadási tartományok korlátai is) azonban egész számok. A pontosság rovására a  $[p_l, p_u]$  intervallumok meghatározása közvetlen számítható, amennyiben a binomiális eloszlást  $Np, Np(1 - p)$  paraméterű normális eloszlással közelítjük.

### 6.10.2. Hiba mérése regresszió esetében

Amikor a magyarázandó attribútum szám típusú, akkor a leggyakrabban használt hiba a négyzetes hibaátlag (vagy annak gyöke). Az elterjedt használat oka, hogy a négyzetes hibaösszeg könnyen kezelhető matematikailag – gondoljuk csak a lineáris regresszióra, amely sok regressziós módszer kiindulópontjaként szolgál. Ha csökkenteni szeretnék a külön pontok által okozott hiba mértékét, akkor használhatunk átlagos hibakülönbséget is.

Többször láttuk, hogy nem az abszolút hiba érdekel minket, hanem a relatív hiba. Azt gondoljuk, hogy ugyanakkora hibát vétünk, ha 200 helyett 220-at jósolunk, mint amikor 1 helyet 1.1-et. A fenti hibamértékek relatív változatainak pontos képlete a következő táblázatban látható. Valamely teszhalmazban (amely származhat kereszt-validációból vagy bootstrapből) az  $i$ -edik pont osztályértékét  $y_i$ -vel, a modell által jósolt osztályértéket  $\hat{y}_i$ -vel jelöljük.

hibamérték	képlet
átlagos négyzetes hiba	$\frac{(y_1 - \hat{y}_1)^2 + \dots + (y_n - \hat{y}_n)^2}{n}$
átlagos négyzetes hibagyök	$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + \dots + (y_n - \hat{y}_n)^2}{n}}$
abszolút hibaátlag	$\frac{ y_1 - \hat{y}_1  + \dots +  y_n - \hat{y}_n }{n}$
relatív négyzetes hiba	$\frac{(y_1 - \hat{y}_1)^2 + \dots + (y_n - \hat{y}_n)^2}{(y_1 - \bar{y})^2 + \dots + (y_n - \bar{y})^2}$
relatív négyzetes hibagyök	$\sqrt{\frac{(y_1 - \hat{y}_1)^2 + \dots + (y_n - \hat{y}_n)^2}{(y_1 - \bar{y})^2 + \dots + (y_n - \bar{y})^2}}$
relatív abszolút hiba	$\frac{ y_1 - \hat{y}_1  + \dots +  y_n - \hat{y}_n }{ y_1 - \bar{y}  + \dots +  y_n - \bar{y} }$
korrelációs együttható	$\frac{(y_1 - \bar{y})(\hat{y}_1 - \bar{\hat{y}}) + \dots + (y_n - \bar{y})(\hat{y}_n - \bar{\hat{y}})}{\sqrt{((y_1 - \bar{y})^2 + \dots + (y_n - \bar{y})^2)((\hat{y}_1 - \bar{\hat{y}})^2 + \dots + (\hat{y}_n - \bar{\hat{y}})^2)}}$

A korrelációs együttható (amely mínusz egy és plusz egy közé esik) kilóg a sorból két dolog miatt. Egyrészt ez a mérték skála invariáns, azaz, ha minden jóslt értéket megszorozunk egy adott konstanssal, akkor a korrelációs együttható nem változik. Legtöbb alkalmazásban nem ezt szeretnénk. Másrészt minél jobb az osztályozó módszer, annál közelebb lesz az együttható egyhez. A többi mérték értéke 0 lesz a tökéletes osztályozó estében.

Az alkalmazási terület tippet adhat arra, hogy melyik hibamértéket használjuk, de a gyakorlat azt mutatja, hogy osztályozók rangsorálásánál ugyanazt a sorrendet szokták adni az egyes mutatók. Talán a fő kérdés, hogy az teljes hiba mérésénél az egyes hibák abszolút értékét vagy négyzetét használjuk.

### 6.10.3. Hiba mérése valószínűségi döntési rendszerek esetén

Valószínűségi döntési rendszerek esetén a kimenet egy valószínűségi eloszlás, nem pedig egy konkrét osztály. Nem azt mondjuk, hogy egy adott feltétellel rendelkező ügyfél kockázatos, hanem azt, hogy 80%-ot adunk annak valószínűségére, hogy kockázatos és 20%-at arra, hogy nem. Ha az osztályok száma  $k$ , akkor az osztályozás eredménye egy  $k$  dimenziós valószínűségi (elemeinek összege 1) vektor. Hogyan határozzuk meg a hibát valószínűségi osztályozás esetében?

#### Négyzetes veszteségfüggvény

Tetszőleges elem konkrét osztályát is leírhatjuk egy valószínűségi vektorral. Ha az elem a  $j$ -edik osztályba tartozik, akkor a valószínűségi vektor  $j$ -edik eleme legyen 1, a többi pedig nulla. Az osztályozás hibája, ekkor az elem osztályához tartozó vektor és az osztályozás eredményeként kapott vektor különbségének normája lesz. Általában az euklideszi normát használjuk és a négyzetgyök számításától eltekintünk:

$$Er(\mathbf{p}, \mathbf{a}) = \sum_{i=1}^k (p_i - a_i)^2.$$

Az  $a_i$ -k közül egyetlen érték 1, a többi nulla, ezért a négyzetes veszteségfüggvény átírható  $1 - 2p_j \sum_{i=1}^k p_i^2$ , ahol  $j$ -vel az osztály sorszámát jelöltük.

Ha az osztályattribútum teljesen független a többi attribútumtól, akkor a négyzetes veszteségfüggvény azokat az osztályozásokat fogja jutalmazni, amelyek a bemenettől függetlenül olyan valószínűségi vektorokat állítanak elő, amely megfelel az osztályattribútum eloszlásfüggvényének, azaz a kimeneti vektor  $i$ -edik eleme adja meg az  $i$ -edik osztály előfordulásának valószínűségét. Nem nehéz ezt az állítást belátni. Jelöljük az  $i$ -edik osztály előfordulásának valószínűségét  $p_i^*$ -vel. A várható értéke a négyzetes veszteségfüggvénynek egy adott teszteleme esetén:

$$\mathbb{E}\left[\sum_{i=1}^k (p_i - a_i)^2\right] = \sum_{i=1}^k (\mathbb{E}[p_i^2] - 2\mathbb{E}[p_i a_i] + \mathbb{E}[a_i^2]) = \sum_{i=1}^k (p_i^2 - 2p_i p_i^* + p_i^*) = \sum_{i=1}^k ((p_i - p_i^*)^2 + p_i^*(1 - p_i^*)).$$

Felhasználtuk, hogy az  $a_i$  várható értéke  $p_i^*$ , továbbá, hogy  $a_i^2 = a_i$  hiszen  $a_i$  értéke csak egy vagy nulla lehet. A végső képletből látszik, hogy a várható érték akkor lesz minimális, ha  $p_i = p_i^*$  minden  $i$ -re.

#### 6.10.4. Osztályozók hatékonyságának mutatószámai

A legfontosabb mutatószám az osztályozó pontossága, amely a jól osztályozott pontok számának arányát adja meg az összes pont számához viszonyítva. Többet mond az ún. *keveredési mátrix* (*confusion matrix*), amely annyi sorból és oszlopból áll, amennyi az osztályok száma. Az  $i$ -edik sor  $j$ -edik eleme adja meg azoknak a pontoknak a számát, amelyeket az osztályozó a  $j$ -edik osztályba sorol, holott azok az  $i$ -edik osztályba tartoznak. A diagonálisban található elemek adják meg a helyesen osztályozott pontok számát. Egy keverési mátrixot láthatunk a következő ábrán.

		Jósolt osztály			
		a	b	c	$\Sigma$
	a	88	10	2	100
tényleges	b	14	40	6	60
osztály	c	18	10	12	40
	$\Sigma$	120	60	20	



**Weka 3.5.7** *Osztályozás esetén a weka alapértelmezés szerint kirajzolja a keveredési mátrixot a kimeneti panelen (Classifier output). Ha erre nem vagyunk kíváncsiak, akkor a Test options panelen klikkeljünk a More options... feliratú gombra. Ez felhossa a Classifier evaluation options panelt, itt töröljük az output confusion matrix kijelölését. Itt állíthatjuk többek között azt is, hogy megjelenjen-e az osztályozó által előállított modell ( döntési szabályok, fák, feltételes valószínűségi táblák – Bayes osztályozók esetében, stb.).*



A pontosság megtévesztő lehet. A magas pontosság nem biztos, hogy a szofisztikált módszerünk eredménye. Ha például bináris osztályozás esetében az egyik osztály előfordulásának

valószínűsége 90%, akkor egy 88% pontosságú osztályozó rossz osztályozó, hiszen pontossága rosszabb, mint annak a butuska osztályozónak, amely mindig a gyakori osztályra tippel. Másik butuska osztályozó a véletlen osztályozó, amely a  $C$  osztály  $p_c$  valószínűséggel választja, ahol  $p_c$  a  $C$  osztály előfordulásának valószínűsége. A valószínűséget relatív gyakorisággal közelítik. A véletlen osztályozó várható pontossága  $0.9 * 0.9 + 0.1 * 0.1 = 82\%$ .

Egy osztályozó kappa statisztikája az osztályozó pontosságát a véletlen osztályozóhoz hasonlítja. Tegyük fel, hogy a tanítóhalmazon az egyes osztályok relatív gyakoriságai  $p_1, p_2, \dots, p_k$  és a tanítóhalmazon az osztályok előfordulása  $n_1, n_2, \dots, n_k$ . Legyen  $N = \sum_{i=1}^k n_i$  és  $M = \sum_{i=1}^k n_i p_i$ . A kappa statisztikát ekkor a

$$\frac{T - M}{N - M}$$

adja, ahol  $T$ -vel a helyesen osztályozott pontokat jelöljük. A kappa statisztika nulla és egy közé esik. A véletlen osztályozó kappa statisztikája nulla, a tökéletes osztályozóé pedig egy.



**Weka 3.5.7** *A kimeneti panelen (Classifier output) mindig megjelenik a jól/rosszul osztályozott és a nem osztályozható elemek száma (Correctly/Incorrectly Classified Instances, UnClassified Instances) és ezen értékek összes tanítóponthoz vett aránya, a kappa statisztika, az abszolút hibák átlaga (Mean absolute error – lásd 157 oldal), a négyzetes hibaátlag (Root mean squared error), a relatív abszolút hibaátlag (Relative absolute error), a relatív négyzetes hibaátlag (Root relative squared error).*



Bináris osztályozás esetén, amikor az osztályozó kimenete nulla vagy egy (igaz/hamis, vagy pozitív/ negatív) további elnevezéseket kell megismernünk. A jól osztályozott pontok számát TP-vel (True Positiv) és TN-nel (True Negative) jelöljük attól függően, hogy melyik osztályba tartoznak. A rosszul osztályozott pontok jelölése FP, FN (False Positive, False Negative). A következő keveredési mátrix összefoglalja a jelöléseket:

	tényleges osztály	
	1	0
jósolt	1 TP	FP
osztály	0 FN	TN

A *felidézést* vagy *megbízhatóságot* (angolul recall vagy true positive rate), amelyet bináris osztályozásnál érzékenységgnek (sensitivity) is hívnak az  $R = \frac{TP}{TP+FN}$  hányados adja. A selejtet (fallout) a  $\frac{FP}{FP+TN}$ , a *precisiont* (a precision és az accuracy angol szavak - amelyek az adatbányászatban mást jelentenek - magyar megfelelője a pontosság) a  $P = \frac{TP}{TP+FP}$  adja. E két érték parametrikus harmonikus közepét  $F$ -mértéknek nevezzük:

$$F = \frac{1}{\alpha \frac{1}{P} + (1 - \alpha) \frac{1}{R}}$$

Az *accuracy* definíciója  $\frac{TP+TN}{N}$ .




---

**Weka 3.5.7** *Ha a Classifier evaluation options panelen (elérhetjük a Test options panel More options... feliratú gombján keresztül) bejelöljük az Output per-class stats opciót, akkor minden osztályhoz megkapjuk a TP és FP arányt, a precíziót, a felidézést, az F-mértéket ( $\alpha = 0.5$  mellett) és a ROC görbe alatti területet.*

---



## 6.11. Osztályozók összehasonlítása

Számos módszert ismertünk meg ebben a fejezetben. A módszerek célja ugyanaz volt, minél pontosabb modellt készíteni. A bőség zavarában persze tanácstalanok lehetünk, hogy most melyik osztályozót használjuk egy adott feladat megoldásánál. Szerencsére néhány módszert azonnal kizárhatunk, mert például az adatok nem lineárisan szeparálhatók, vagy vannak valós típusú attribútumok, stb. Néhány módszer kizárása után valószínűleg nem egy fog maradni; melyik ezek közül a legjobb. Az osztályozók összehasonlításának helyes módja a kutatói világban is gyakran felmerül, amikor valaki új módszert javasol, és meg kell mutatnia, hogy az ő megoldása miért jobb a többiekénél.

Első gondolatunk az lehetne, hogy minek időzünk ilyesmivel; vegyünk egy teszhalmazt és nézzük meg, hogy melyik osztályozónak lesz nagyobb a pontossága. Csakhogy egy adathalmaz a valódi pontosság egy becslése és lehet, hogy a különbség a becslés hibájának eredménye.

A következő részekben legyen adott  $B$  és  $L$  osztályozók (gondolhatjuk, hogy a  $B$  egy Bayes osztályozóra, a  $L$  pedig egy logisztikus regresszió alapuló módszerre utal) és jelöljük az egyes pontosságokat  $p_B$ ,  $p_L$ -lel. Null hipotézisünk, hogy

$$H_0 : p_B = p_L.$$

### 6.11.1. Binomiális próbán alapuló összehasonlítás

Amennyiben a két osztályozó és egy teszhalmaz a rendelkezésünkre áll, tehát minden tesztpontra el tudjuk dönteni, hogy az egyes osztályozók helyes eredményt adnak-e, akkor az összehasonlításhoz binomiális próbát célszerű használni.

Jelöljük  $d_i$ -vel azt a valószínűségi változót, amely 1 értéket vesz fel, amennyiben az  $i$ -edik olyan tesztpontnál, ahol a két osztályozó különböző eredményt adott, a  $B$  osztályoz helyesen. Ha az  $L$  ad jó eredményt, akkor  $d_i = 0$ . Tehát azokat a tesztpontokat nem vesszük figyelembe a teszt során, amelyekre  $B$  és  $L$  ugyanazt az eredményt adja.

Amennyiben a null hipotézis fennáll, akkor ugyanannyi esélye van annak, hogy  $d_i$  egyet vesz fel, mint annak, hogy nullát, azaz  $\mathbb{P}(d_i = 1) = 1/2$ . Ha  $n$  olyan tesztpont áll rendelkezésünkre, amelyre a két osztályozó eltérő választ ad és  $D = \sum_{i=1}^n d_i$ , akkor  $D$  egy  $(n, 1/2)$  paraméterű binomiális eloszlású valószínűségi változó. Jelöljük  $D^*$ -gal azon tanítópontok számát, amelyekre  $B$  jó eredményt adott, de  $L$  rosszat, és legyen  $p_{D^*} = \binom{n}{D^*} 1/2^n$ .

A binomiális proba során szolgáltatott  $p$ -érték megegyezik azon  $[0, n]$  intervallumhoz tartozó pontok valószínűségeinek összegével, amely valószínűségek kisebb vagy egyenlőek  $p_{D^*}$ -nál.



### 6.11.2. Student próbán alapuló összehasonlítás

Amennyiben nem tudjuk kiértékelni minden tanítóponton az osztályozónkat, csak azt tudjuk, hogy mennyi pontot osztályozott jól az egyik és mennyit a másik osztályozó, akkor a Student próbát alkalmazhatjuk. Ez egy pontatlanabb próba, mint a binomiális próba, mert közelítően alapul.

Legyen adott  $N$  darab tesztalmaz. A két osztályozó  $i$ -edik tesztalmazon mért pontosságát jelöljük  $b_i$  és  $l_i$ -vel, a pontosságok átlagát pedig  $\bar{b}$  és  $\bar{l}$ -vel és legyen  $d_i = b_i - l_i$ .

Tegyük fel, hogy egy  $\nu$  pontosságú osztályozó pontossága egy adott tesztalmazon megegyezik egy  $\nu$  várható értékű és ismeretlen szórású normális eloszlású valószínűségi változó egy megfigyelésével. Ezek a megfigyelt pontosságok rendelkezésünkre állnak, azt kéne eldönteni, hogy az eredeti pontosságok statisztikailag eltérnek-e egymástól. Nullhipotézisünk tehát, hogy  $\nu_B = \nu_L$  vagy  $\nu_D = 0$ , ahol  $D = B - L$ .

A student t-próba (lásd a 32) egy ismeretlen várható értékű és szórású normális eloszlású valószínűségi változó várható értékére tett feltételt próbál eldönteni. Számítsuk hát ki a

$$\frac{\bar{d} - 0}{\sqrt{\sigma_d^{*2}/N}}$$

értéket és vessük össze azzal az értékkel, ahol a student eloszlás megegyezik  $1 - \alpha$ -val. A szokásos módon a  $\bar{d}$ -gal a mintaátlagot, a  $\sigma_d^*$ -gal az empirikus korrigált szórást, a  $1 - \alpha$ -val az próba szintjét jelöltük. Ha a teszt elutasítja a nullhipotézist, akkor a nagyobb átlaghoz tartozó osztályozó statisztikailag is jobb a másiknál.

### 6.11.3. Osztályozó módszerek összehasonlítása

Az előbbi részekben két osztályozót hasonlítottunk össze. Osztályozót egy osztályozó módszer alkalmazásával kapunk egy tanító halmazon. Például az ID3 algoritmus egy osztályozó módszer, míg egy döntési fa, amelyet az ID3 alkalmazásával kapunk egy adott tanítóhalmazon egy osztályozó.

Tudományos munkákban a kutatók új módszereket fejlesztenek ki és a módszerek hatékonyságát nyilvánosan elérhető adatbázisokon végzett eredményekkel kívánják bemutatni. Ehhez általában keresztvalidációt használnak. Vesznek  $n$  darab adatbázist, minden adatbázist  $r$  diszjunkt részre osztanak véletlenszerűen és a véletlen felosztást  $k$ -szor megismétlik. Összesen  $nrk$  tanítást majd tesztelést végeznek és az ezekből származó pontosságot használják a student t-próba során.



**Weka 3.6.0** A wekában az osztályozók összehasonlítása során a fent vázolt módszert használhatjuk. Ehhez az **Experimenter** ablakot kell kiválasztani. A **Datasets** panelen megadhatjuk az adatbázisokat, az **Algorithms** panelen pedig az összehasonlítandó osztályozó módszereket. Ha a keresztvalidáció helyett egyszerű kettéosztást (tanító- és tesztelőhalmazra) szeretnénk alkalmazni, akkor az **Experiment type** panelen válasszuk a **Train/Test Percentage Split** lehetőséget.

A módszerek alkalmazását a **Run** fül **Start** gombjának lenyomásával érhetjük el. Az eredményt elmenthetjük arff formátumba, így az be-menetként szolgálhat tetszőleges adatbányászati módszernek. Ha csv



*formátumba mentünk, akkor az eredményt kényelmesen böngészhetjük excel vagy openoffice táblázatkezelőkkel.*

*Az **Analyse** fülön végezhetjük el a Student próbát. A student próbához nem csak az osztályozó pontosságát használhatjuk (ez az alapértelmezett érték), hanem bármilyen kiszámított értéket (pl. kap-pa statisztika, F-mérték, ROC görbe alatti terület, stb.) használhatunk. A **Test output** panelen az Student teszt eredményét láthatjuk. Meg van adva minden osztályozóhoz az összesített pontosság. Amennyiben az érték mellett egy \* szerepel, akkor a Student próba alapján az osztályozó módszer szignifikánsan rosszabb, mint a legelső osztályozó. A szám melletti v betű esetén pedig a Student próba szignifikánsan jobb eredményt adott.*

---

## 7. fejezet

# Klaszterezés

Klaszterezésen elemek csoportosítását értjük. Úgy szeretnénk a csoportosítást elvégezni, hogy a hasonló elemek ugyanazon, míg az egymástól eltérő elemek külön csoportba kerüljenek. Sajnos a „jó” csoportok kialakítása nem egyértelmű feladat, hiszen az emberek gyakran más-más szempontokat vesznek figyelembe a csoportosításnál. Ugyanazt azt adathalmazt, alkalmazástól és szokásoktól függően, eltérően klasztereznék az emberek. Például az 52 darab francia kártyát sokan 4 csoportra osztanák (szín szerint), sokan 13-ra (figura szerint). A Black Jack játékosok 10 csoportot hoznának létre (ott a 10-es, bubi, dáma, király között nincs különbség), míg a Pikk Dáma játékot kedvelők hármat (pikk dáma, a kőrök és a többi lap). Klaszterezéskor tehát az adathalmaz mellett meg kell adnunk, hogy miként definiáljuk az elemek hasonlóságát, továbbá, hogy mi alapján csoportosítsunk (összefüggő alakzatokat keressünk, vagy a négyzetes hibát minimalizáljuk stb.).

A jóság egzakt definíciójának hiánya mellett nagy problémát jelent az óriási keresési tér. Ha  $n$  pontot akarunk  $k$  csoportba sorolni, akkor a lehetséges csoportosítások számát a Stirling számok adják meg:

$$\mathcal{S}_n^{(k)} = \frac{1}{k!} \sum_{i=0}^k (-1)^{k-i} \binom{k}{i} i^n.$$

Még egy egészen kicsi adathalmaz mellett is megdöbbentően sokféleképpen csoportosíthatunk. Például 25 elemet 5 csoportba  $\mathcal{S}_{25}^{(5)} = 2,436,684,974,110,751$  különböző módon partícionálhatunk. Ráadásul, ha a csoportok számát sem tudjuk, akkor a keresési tér még nagyobb ( $\sum_{k=1}^{25} \mathcal{S}_{25}^{(k)} > 4 \cdot 10^{18}$ ).

Szükség van azonban az elemek automatikus csoportosítására, így a problémákon túl kell lépni. Objektív definíciót kell adnunk az elemek hasonlóságának mértékére és a klaszterezés minőségére. Amennyiben megfelelő matematikai modellbe ágyaztuk a feladatot, lehetőség nyílik olyan algoritmusok megkeresésére, amelyek jól és gyorsan oldják meg a feladatot. Ezekről az algoritmusokról és a hasonlóság megállapításának módjáról szól ez a fejezet.

Klaszterezés során csoportokba, osztályokba soroljuk az elemeket, tehát osztályozást végzünk. Az eredeti osztályozási feladattól (lásd előző fejezet) az különbözteti meg a klaszterezést, hogy nincs megadva, hogy melyik elem melyik osztályba tartozik (tehát nincs egy tanító, aki helyes példákkal segíti a tanulásunkat), ezt nekünk kell meghatároznunk. Ezért hívják a klaszterezést *felügyelet nélküli tanulásnak* (unsupervised learning) is.

A klaszterezés az adatbányászat legrégebbi és leggyakrabban alkalmazott része. Számos helyzetben alkalmazzák, így csoportosítanak weboldalakat, géneket, betegségeket stb. Az egyik

legdinamikusabban fejlődő terület azonban a személyre szabott szolgáltatásoké, ahol az ügyfeleket, ill. vásárlókat kategorizálják, és az egyes kategóriákat eltérően kezelik. A klaszterezésre azért van szükség, mert az ügyfelek számossága miatt a kézi kategorizálás túl nagy költséget jelentene.

Gyakran nem az a fontos, hogy az egyes elemeket melyik csoportba soroljuk, hanem az, hogy mi jellemző a különböző csoportokra. Például egy banki stratégia kialakításánál nem érdekel bennünket, hogy Kis Pista melyik csoportba tartozik, hanem csak az, hogy milyen ügyfélcsoportokat célszerű kialakítani és ezekre a csoportokra mi jellemző. A klaszterezés segítségével egy veszteséges tömörítést végeztünk. A teljes ügyfeleket tartalmazó adatbázist egy kisebb, átláthatóbb, emészthetőbb ügyfélcsoport adatbázissá alakítottuk.




---

**Weka 3.5.7** A klaszterező módszereket az *Experimenter* alkalmazás *Cluster* fülén keresztül érhetjük el.

---



A fejezet további részében először egy meghökkentő kutatási eredményről számolunk be, majd a hasonlóság meghatározásáról beszélünk végül rátérünk a legismertebb klaszterező algoritmusokra.

## 7.1. Egy lehetetlenség-elmélet

A klaszterezés az egyik legnehezebben átlátható adatbányászati terület. Napról napra újabb és újabb cikkek jelennek meg különböző „csodaalgoritmusokról”, amelyek szupergyorsan és helyesen csoportosítják az elemeket. Elméleti elemzésekről általában kevés szó esik – azok is gyakran elnagyoltak, sőt hibásak –, viszont az algoritmust igazoló teszteredményekből nincs hiány. Mintha minden algoritmusnak illetve szerzőnek létezne a maga adatbázisa, amivel az eljárás remek eredményeket hoz.

Ebben a káoszban kincset érnek a helyes irányvonalak megvilágításai és a megalapozott elméleti eredmények. Egy ilyen gyöngyszem Jon Kleinberg munkája, amit az „An Impossibility Theorem for Clustering (A Klaszterezés Lehetetlenség-elmélete)” című cikkében publikált 2002-ben [74]. A cím már sejteti az elszomorító eredményt, miszerint *nem létezik jó, távolság alapú<sup>1</sup> klaszterező eljárás!* Ezt a meglepő állítást úgy bizonyítja, hogy három tulajdonságot mond ki, amellyel egy klaszterező eljárásnak rendelkeznie kell, majd belátja, hogy nem létezhet klaszterező eljárás, amelyre ez igaz. A tulajdonságok az alábbiak:

**Skála-invariancia:** Ha minden elempár távolsága helyett annak az  $\alpha$ -szorosát vesszük alapul (ahol  $\alpha > 0$ ), akkor a klaszterező eljárás eredménye ne változzon!

**Gazdagság (richness):** Tetszőleges előre megadott csoportosításhoz tudjunk megadni távolságokat úgy, hogy a klaszterező eljárás az adott módon csoportosítson.

**Konzisztencia:** Tegyük fel, hogy a klaszterező eljárás valahogy csoportosítja az elemeket. Ha ezután tetszőleges, azonos csoportban lévő elempárok között a távolságot csökkentem,

---

<sup>1</sup>A különbözőség megállapításához használt távolságfüggvénynek szemi-metrikának kell lennie, tehát a háromszög egyenlőtlenségnek nem kell teljesülnie

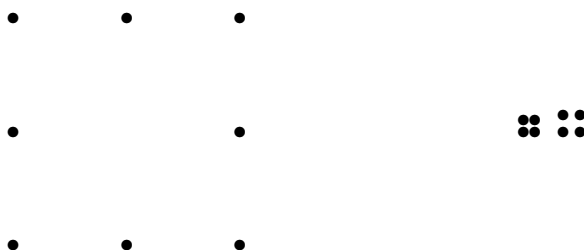
illetve külön csoportban lévő elempárok távolságát növelem, akkor az újonnan kapott távolságok alapján működő eljárás az eredetivel megegyező csoportosítást adja.

A fenti tulajdonságok teljesen természetesek, azt gondolnánk, hogy minden algoritmus ilyen. Ezért nem túl biztató a következő tétel:

**7.1. tétel.** *Amennyiben az elemek száma nagyobb 1-nél, akkor nem létezik olyan klaszterező eljárás, ami rendelkezik a Skála-invariancia, a Gazdagság és a Konzisztencia tulajdonságokkal.*

Kleinberg azt is bebizonyítja, hogy bármely két tulajdonsághoz létezik klaszterező eljárás, amely rendelkezik a választott tulajdonságokkal. Például a single-linkage eljárás (lásd 7.7.1. rész) skála-invariáns és konzisztens. Ezen kívül az is igaz, hogy a partíciónáló algoritmusok (pl.: k-means, k-medoid), ahol a cél a középpontoktól vett távolság függvényének (például négyzetes hiba összege) minimalizálása, nem konzisztensek.

Vitatkozhatunk azon, hogy a konzisztencia jogos elvárás-e egy klaszterező algoritmussal szemben. Nézzük a következő ábrát. Bal oldalon láthatjuk az eredetileg megadott pontokat, jobb oldalon pedig az átmozgatás során kapottakat.



Legtöbbször a bal oldali pontokat egy csoportba vennék (nagy négyzetet reprezentáló pontok), a jobb oldalon láthatókat viszont két külön csoportba sorolnák (két kis négyzethez tartozó pontok). A klaszteren belüli távolságokat tehát csökkentettük, a klaszterezés mégis megváltozott, azaz klaszterezési eljárásunk nem rendelkezik a konzisztencia tulajdonsággal.

Sajnos Kleinberg erre az észrevételre is tud elszomorítóan reagálni. A konzisztencia fogalmát lazíthatjuk. Amennyiben a klasztereken belüli távolságokat csökkentjük, a klaszterek közötti távolságokat növeljük, és ezáltal bizonyos klaszterek kisebb klaszterekké bomlanak, akkor a klaszterező eljárás *finomítás-konzisztens*. Belátható, hogy nem létezik olyan klaszterező eljárás, ami skála-invariáns, gazdag és finomítás-konzisztens.

Ha viszont a gazdagságból is engedünk egy kicsit, nevezetesen, hogy a klaszterező algoritmus sose tudjon minden pontot külön klaszterbe sorolni – de tetszőleges más módon tudjon partíciónálni –, akkor létezik klaszterező eljárás, amely kielégíti a három tulajdonságot.

Mielőtt továbblépnénk gondolkodjunk el azon, hogy jogos-e a hasonlóságot és különbözőséget pusztán egy távolság alapján definiálni. A klaszterezés eredeti célja az, hogy a hasonló elemek egy csoportba, míg a különböző elemek eltérő csoportba kerüljenek. Ebből következik, hogy egy tetszőleges elem különbsége (távolsága) a saját csoportbeli elemeitől kisebb lesz, mint a különbsége más csoportban található elemektől. Biztos, hogy jó ez? Biztos, hogy az ember is így csoportosít, tehát ez a természetes klaszterezés? Sajnos nem lehet a kérdésre egyértelmű választ adni. Van amikor az ember így csoportosít, van, amikor máshogy. Tekintsük a következő ábrán elhelyezkedő pontokat.

.....  
 .....

Valószínűleg kivétel nélkül minden ember két csoportot hozna létre, az alsó szakaszhoz tartozó pontokét és a felső szakaszhoz tartozó pontokét. Mégis, ha megnézzük, akkor az alsó szakasz bal oldali pontja sokkal közelebb van a felső szakasz bal oldali pontjaihoz, mint azokhoz a pontokhoz, amelyek az alsó szakasz jobb oldalán helyezkednek el. Mégis ragaszkodunk ahhoz, hogy a bal- és jobboldali pontok egy csoportba kerüljenek. Úgy érezzük, egymáshoz tartoznak, mert mindannyian az alsó szakasz elemei.

Következésképpen a klaszterezés célja – az eredetivel szemben – gyakran az, hogy úgy csoportosítsunk, hogy egy csoportba kerüljenek az elemek akkor, ha ugyanahhoz az absztrakt objektumhoz tartoznak, és különbözőbe, ha más absztrakt objektum részei. A klaszterezés nehézsége pont abban rejlik, hogy automatikusan kell felfedezni objektumokat az elemek alapján, ami ráadásul nem egyértelmű feladat (például Rubin vázájának esete).

Ha a klaszterezés során az absztrakt objektumokat összefüggő alakzatok formájában keressük (pl. vonal, gömb, amőba, pácikaember stb.) akkor van esély jól megoldani a feladatot. Összességében tehát tökéletes klaszterezés nem létezik, ugyanakkor a lehetetlenség elmélet nem zárja ki az összefüggő alakzatokat felfedező eljárás létezését.

## 7.2. Hasonlóság mértéke, adatábrázolás

Adott  $n$  elem (vagy más néven objektum, egyed, megfigyelés stb.). Tetszőleges két elem  $(x,y)$  között értelmezzük a *hasonlóságukat*. Mi a hasonlóság helyett annak inverzével, a *különbözőséggel* dolgozunk  $(d(x,y))$ . A  $d(x,y)$ -től elvárjuk (amellett, hogy  $d(x,y) \geq 0$ ) azt, hogy

I. reflexív:  $d(x,x) = 0$ ,

II. szimmetrikus legyen:  $d(x,y) = d(y,x)$ ,

III. és teljesüljön a háromszög-egyenlőtlenség:  $d(x,z) \leq d(x,y) + d(y,z)$ ,

tehát a különbözőség metrika (távolság) legyen<sup>2</sup>. A továbbiakban elemek különbözősége helyett gyakran mondunk elemek *távolságát*.

A klaszterezés legáltalánosabb esetében minden egyes elempár távolsága előre meg van adva. Az adatokat ekkor egy ún. távolság mátrixszal reprezentáljuk:

$$\begin{bmatrix} 0 & d(1,2) & d(1,3) & \cdots & d(1,n) \\ & 0 & d(2,3) & \cdots & d(2,n) \\ & & 0 & \cdots & d(3,n) \\ & & & \ddots & \vdots \\ & & & & 0 \end{bmatrix},$$

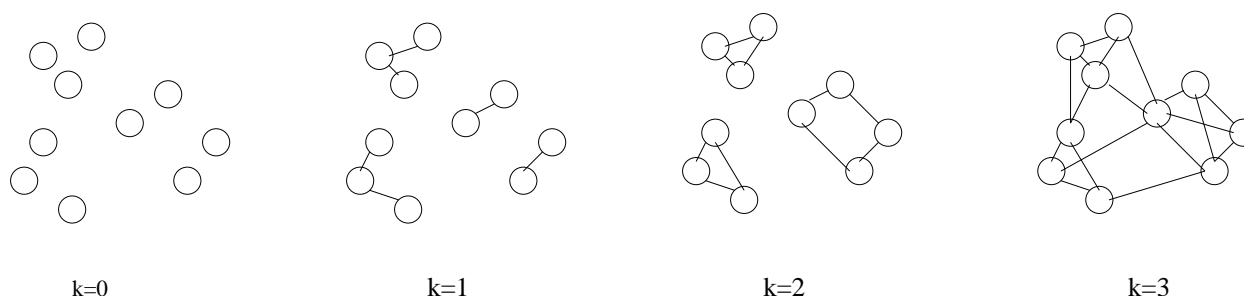
<sup>2</sup>Megjegyzés: Ha a 3. tulajdonság nem teljesül, akkor szemi-metrikáról beszélünk, ha az erősebb  $d(x,y) \leq \max\{d(x,z), d(y,z)\}$  tulajdonság áll fenn, akkor pedig ultrametrikusról (más néven nem-archimédeszi).

ahol  $d(i, j)$  adja meg az  $i$ -edik és a  $j$ -edik elem különbözőségét.

A gyakorlatban az  $n$  elem (vagy objektum) attribútumokkal van leírva, és a különbözőséget az attribútumok alapján definiálhatjuk valamilyen *távolságfüggvénnyel*. Ha megadjuk a távolságfüggvényt, akkor elvben felírhatjuk a fenti mátrixot. Sok esetben azonban az elemek száma olyan nagy, hogy a mátrix rengeteg helyet foglalna. Modellünkben ezért rendelkezésünkre állnak az attribútumokkal megadott elemek halmaza és a távolságfüggvény. Az  $n$  értéke nagy lehet, így nem tehetjük fel, hogy az adatok elférnek a memóriában.

Sokszor fogjuk a klaszterezést gráfparticionálási feladatként vizsgálni. Az elemekre tekintetünk úgy, mint egy  $G = (V, E)$  súlyozott, irányítatlan, teljes gráf pontjaira, ahol az éleken található súlyok a távolságot, vagy éppen a hasonlóságot adják meg. Az  $(u, v) \in E$  él súlyát  $w(u, v)$ -vel jelöljük.

Vannak algoritmusok, amelyek nem az eredeti gráfon dolgoznak, hanem az úgynevezett *k-legközelebbi szomszéd gráfon*, amit  $G_k$ -val jelölünk.  $G_k$ -ban is a pontoknak az elemek, az éleken található súlyok pedig a hasonlóságoknak felelnek meg, de itt csak azokat az éleket tároljuk, amelyek egyik pontja a másik pont  $k$  legközelebbi pontjai között szerepel. Az alábbi ábrán ilyen gráfokat láthatunk:



7.1. ábra. Példa  $k$ -legközelebbi szomszéd gráfokra  $k=0,1,2,3$  esetén

Ha az adathalmazt a  $k$ -legközelebbi szomszéd gráffal ábrázoljuk, akkor ugyan veszünk némi információt, de a lényeg megmarad, és jóval kevesebb helyre van szükségünk. Az egymástól nagyon távoli elemek nem lesznek összekötve  $G_k$ -ban. További előny, hogy amennyiben egy klaszter sűrűségét a benne található él összsúlyával mérjük, akkor a sűrű klasztereknél ez az érték nagy lesz, ritkákánál pedig kicsi.

### 7.3. A klaszterek jellemzői

A  $C$  klaszter elemeinek számát  $|C|$ -vel jelöljük. A klaszter „nagyságát” próbálja megragadni a klaszter *átmérője* ( $D(C)$ ). A két legelterjedtebb definíció az elemek közötti átlagos, illetve a maximális távolság:

$$D_{avg}(C) = \frac{\sum_{p \in C} \sum_{q \in C} d(p, q)}{|C|^2},$$

$$D_{max}(C) = \max_{p, q \in C} d(p, q).$$

Ízlés kérdése, hogy a klaszter átmérőjének számításakor figyelembe vesszük-e a pontok önmaguktól vett távolságát (ami 0). Nyugodtan használhatjuk az átmérő  $D'_{avg}(C) = \frac{\sum_{p,q \in C, p \neq q} d(p,q)}{\binom{|C|}{2}} = 2 \frac{N}{N-1} D_{avg}(C)$  definícióját is. A *klaszterek közötti távolságot* ( $d(C_i, C_j)$ ) is többféleképpen értelmezhetjük.

**Minimális távolság:**  $d_{min}(C_i, C_j) = \min_{p \in C_i, q \in C_j} d(p, q)$ .

**Maximális távolság:**  $d_{max}(C_i, C_j) = \max_{p \in C_i, q \in C_j} d(p, q)$ .

**Átlagos távolság:**  $d_{avg}(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{p \in C_i} \sum_{q \in C_j} d(p, q)$ , ami a külön klaszterben lévő pontpárok átlagos távolságát adja meg.

**Egyesített klaszter átmérője:**  $d_D(C_i, C_j) = D(C_i \cup C_j)$

A vektortérben megadott elemeknél gyakran használt fogalmak a klaszter *középpontja* ( $\vec{m}_C$ ) és a *sugara* ( $R_C$ ).

$$\vec{m}_C = \frac{1}{|C|} \sum_{p \in C} \vec{p},$$

$$R_C = \frac{\sum_{p \in C} |\vec{p} - \vec{m}_C|}{|C|}.$$

A klaszterek közötti távolság mérésére pedig gyakran alkalmazzák a középpontok közötti távolság értékét:

$$d_{mean}(C_i, C_j) = |\vec{m}_i - \vec{m}_j|.$$

Az átlagok kiszámításánál – például átmérő, sugár esetében – számtani közepet használtunk. Bizonyos cikkekben négyzetes közepet alkalmaznak helyette. Tulajdonképpen tetszőleges közép használható, egyik sem rendelkezik elméleti előnnyel a többivel szemben. Gondoljuk meg azonban, hogy a hatvány alapú közepeknél jóval nagyobb számokkal dolgozunk, így ezek számítása esetleg nagyobb átmeneti tárat kíván.

A négyzetes középnek előnye a számtani középpel szemben, hogy könnyű kiszámítani, amennyiben vektortérben dolgozunk. Ezt a BIRCH algoritmusnál (7.7.3. rész) is kihasználják, ahol nem tárolják a klaszterekben található elemeket, hanem csak 3 adatot:  $|C|$ ,  $\vec{L}S_C = \sum_{p \in C} \vec{p}$ ,  $SS_C = \sum_{p \in C} \vec{p}\vec{p}^T$ . Könnyű belátni, hogy a fenti három adatból két klaszter ( $C_i, C_j$ ) közötti átlagos távolság (és hasonlóan az egyesített klaszter átmérője) közvetlenül adódik:  $d_{avg}(C_i, C_j) = \frac{SS_{C_i} + SS_{C_j} - 2\vec{L}S_{C_i}\vec{L}S_{C_i}^T}{|C_i||C_j|}$ .

## 7.4. A klaszterezés „jósága”

Mint már említettük, a klaszterezés jóságára nem lehet minden szempontot kielégítő, objektív mértéket adni. Ennek ellenére néhány függvény minimalizálása igen elterjedt a klaszterező algoritmusok között.

A továbbiakban  $n$  darab elemet kell  $k$  rögzített számú csoportba sorolni úgy, hogy a csoportok diszjunktak legyenek, és minden csoportba kerüljön legalább egy elem.



### 7.4.1. Klasszikus mértékek

Az alábbi problémákat különböztetjük meg a minimalizálandó célfüggvény alapján:

**Minimális átmérő probléma:** Célunk itt a legnagyobb klaszterátmérő minimalizálása. Átmérőnek ez esetben  $D_{max}$ -ot szokás használni.

**$k$ -medián probléma:** Válasszuk ki az  $n$  elem közül  $k$  ún. reprezentáns elemet, amelyek a minimális hibaösszeget adják. Egy elem hibája a hozzá legközelebbi reprezentáns elem távolsága. A feladat NP-nehéz, még akkor is, ha olyan síkba rajzolható gráfokra szorítkozunk, amelyeknek a maximális fokszáma 3 (ha a gráf fa, akkor már lehet polinomrendű algoritmust adni,  $p = 2$  esetében a feladat lineáris időben megoldható)[70]. A feladat NP-nehéz marad, ha a gráf Euklideszi térbe képezhető, sőt, konstans szorzó erejéig közelítő megoldást adni, még ilyenkor is, nehéz feladat [92]!

**$k$ -center probléma:** Ez a feladat a  $k$ -medián módosítása, csak itt a legnagyobb hibát kell minimalizálni.

**$k$ -klaszter probléma:** Célunk itt a klaszteren belüli távolságösszegek  $(\sum_{i=1}^k \sum_{p,q \in C_i} d(p,q) = \sum_{i=1}^k |C_i|^2 D_{avg}(C_i))$  minimalizálása. A feladat (és konstans szorzó erejéig annak közelítése) NP-nehéz  $k \geq 2$  ( $k \geq 3$ ) esetén [116].

**Legkisebb (négyzetes) hibaösszeg:** Csoportosítsuk úgy a pontokat, hogy a középpontoktól való távolság összege  $(E = \sum_{i=1}^k \sum_{p \in C_i} (|\vec{p} - \vec{m}_{C_i}|))$  minimális legyen. Nyilvánvaló, hogy ez a megközelítés csak olyan esetekben használható, amikor értelmezni tudjuk a klaszterek középpontját ( $\vec{m}_{C_i}$ -t).

Sok esetben a középpontoktól való távolságösszeg helyett a távolság négyzeteinek összegét minimalizálják.

Legkisebb (négyzetes) hibaösszeg probléma eléggé hasonlít a  $k$ -klaszter problémához.

**7.2. észrevétel.**  $\sum_{i=1}^k \sum_{p,q \in C_i} d(p,q) = \sum_{i=1}^k \sum_{p \in C_i} \|p - \vec{m}_{C_i}\|^2$ , ahol  $\vec{m}_C = \frac{1}{|C|} \sum_{q \in C} \vec{q}$ .

*Bizonyítás:*

$$\begin{aligned} \sum_{i=1}^k \sum_{p \in C_i} \|p - \vec{m}_{C_i}\|^2 &= \sum_{i=1}^k \sum_{p \in C_i} \left\| p - \frac{1}{|C_i|} \sum_{q \in C_i} \vec{q} \right\|^2 = \sum_{i=1}^k \sum_{p \in C_i} \sum_{q \in C_i} \frac{1}{|C_i|} \|p - q\|^2 \\ &= \sum_{i=1}^k \frac{1}{|C_i|} \sum_{(p,q) \in C_i} \|p - q\|^2 = \sum_{i=1}^k |C_i| D_{avg}(C_i) \end{aligned}$$

■

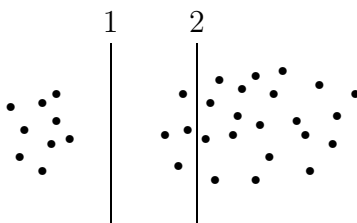
Azok az algoritmusok, amelyek a fenti célfüggvényeket minimalizálják, az elemeket kis kompakt felhőkbe csoportosítják. Ez valamennyire elfogadhatónak tűnik, azonban ezeknek a megközelítéseknek számos súlyos hátránya van.

- I. Legfontosabb, hogy csak elliptikus klasztereket generál, tehát tetszőleges amőba alakú, de kompakt klasztert felvág kisebb kör alakú klaszterekre.

II. Rosszul csoportosít, ha a klaszterek között nagyok a méretkülönbségek. Ennek oka az, hogy a nagy klaszterben lévő pontok távol esnek a középponttól, ami nagy hibát eredményez. Tehát hiába kompakt egy nagy klaszter, a hibát minimalizáló algoritmusok kis részekre fogják felosztani.

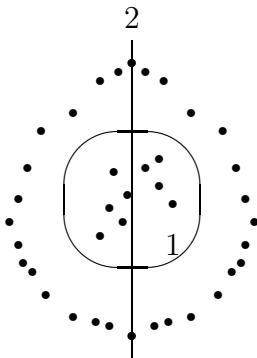
A négyzetes hibaösszeget minimalizáló eljárások további hibája, hogy érzékeny a távol eső (outlier) pontokra, hiszen egy távoli pont a klaszter középpontját nagyon „elhúzhatja”.

Elrettentő példaként nézzük a következő két ábrán látható pontokat. Ha a maximális átmérőt minimalizáljuk, akkor a 2. egyenes alapján osztjuk ketté a pontokat. Ennek ellenére minden „rendszerető” ember a két csoportot inkább az 1-es egyenes mentén szeparálná. A gyenge klaszterezés oka, hogy a klasztereken belüli maximális eltérést annak árán minimalizáljuk, hogy sok különböző pont egy klaszterbe kerül (Megjegyzés: ugyanezt a rossz eredményt kapnánk, ha a közepektől való távolságot, esetleg a távolságösszeget akarnánk minimalizálni.)



7.2. ábra. Hibás klaszterezés: eltérő méretű klaszterek esetén

A 7.3. ábrán látható pontokat a 2-medián problémát megoldó algoritmusok a 2-es egyenes szerint csoportosítanák.



7.3. ábra. Hibás klaszterezés: egymást tartalmazó klaszterek esetén

### 7.4.2. Konduktancia alapú mérték

A klasszikus mértékek igen közkedveltek a matematikusok körében, köszönhetően az egyszerűségüknek és annak, hogy remekül elemezhetőek. Rengeteg írás született, amelyek matematikus szemmel kitűnőek, a gyakorlatban azonban – ahogy azt az előző két példa is illusztrálta – haszontalanok. 30-40 évig mégis ezek a problémák álltak a középpontban. A kutatók szép eredményeket értek el, a hasznosság igénye azonban sokáig kimondatlan maradt.

Az adatbányászat népszerűsödésével egyre fontosabb szerepük lett a "Mi haszna van?", "Mi a jó klaszterezés?" kérdéseknek. Hamar kiderült, hogy a klasszikus mértékek a gyakorlati esetek többségében egyszerűen rossz eredményt adnak. Új mértékek, megközelítések születtek. Ezek közül talán a legígéretesebb a konduktancia alapú mérték [69].

Tekintsünk az adathalmazunkra, mint egy  $G = (V, E)$  gráfra, de most az éleken található súly a hasonlósággal legyen arányos, ne pedig a távolsággal (különbözőséggel). Jelöljük egy  $T \subseteq E$  élhalmazban található élek súlyainak összegét  $w(T)$ -vel ( $w(T) = \sum_{e \in T} w(e)$ ),  $C \subseteq V$  klaszterben található elemek számát  $|C|$ -vel,  $E(S) = E(V \setminus S)$ -el (edge(S)) pedig az  $(S, V \setminus S)$  vágást keresztező élek halmazát:  $E(S) = \{(p, q) | p \in S, q \in V \setminus S\}$ .

Vizsgáljuk azt az egyszerű esetet, amikor  $k = 2$ , tehát a gráf pontjait két részre akarjuk osztani. Klaszterezésnél az egyik célunk, hogy az elemeket úgy csoportosítsuk, hogy a különböző elemek külön klaszterbe kerüljenek. Ez alapján mondhatnánk, hogy egy minimális összsúlyú vágás jól osztaná ketté a pontokat. Sajnos ez a módszer legtöbb esetben kiegyensúlyozatlan (nagyon eltérő méretű) csoportokat hozna létre. Gondoljuk meg, hogy ha az egyik klaszterben csak 1 elem található, akkor  $n - 1$  súlyt kell összegezni, míg egyenletes kettéosztásnál ugyanez az érték  $(\frac{n}{2})^2$ .

A vágás helyett célszerű olyan mértéket bevezetni, amely figyelembe veszi valahogy a gráf kiegyensúlyozottságát is, és kisebb jelentőséget tulajdonít az olyan vágásnak, amely kis elemszámú részhez tartozik. Egy gráf *kiterjedése* (expansion) megadja az összes vágás közül azt, amelyiknél a legkisebb az arány a vágás súlya és a vágást alkotó két ponthalmaz közül a kisebbik elemszáma között. Formálisan az  $(S, V \setminus S)$  vágás kiterjedése:

$$\varphi(S) = \frac{w(E(S))}{\min(|S|, n - |S|)}.$$

Látható, hogy a számláló a kis vágásértéket, míg a nevező a kiegyensúlyozottságot preferálja. Egy gráf kiterjedése pedig a vágások minimális kiterjedése, egy klaszter kiterjedését pedig a hozzá tartozó részgráf kiterjedésével definiálhatjuk. A klaszterezés jóságát, ez alapján, a klaszterek minimális kiterjedésével adhatjuk meg.

Sajnos a kiterjedés képletében a nevező nem veszi figyelembe az élek súlyait. Azt szeretnénk, hogy azok a pontok, amelyek nagyon különbözőek az összes többi ponttól, kisebb összsúllyal szerepeljenek a „jóság” definíciójában, mint azok a pontok, amelyeknek jóval több ponthoz hasonlítanak. A kiterjedés általánosítása a *konduktancia* (conductance).

**7.3. definíció.** Legyen  $G = (V, E)$  gráf egy vágása  $(S, V \setminus S)$ . A vágás konduktanciáját a következőképpen definiáljuk:

$$\phi(S) = \frac{w(E(S))}{\min(a(S), a(V \setminus S))},$$

ahol  $a(S) = \sum_{p \in S, q \in V} w(p, q)$ .

A gráf konduktanciája pedig a vágások minimális konduktanciája:  $\phi(G) = \min_{S \subseteq V} \phi(S)$ .

A konduktancia könnyen általánosítható  $k$  klaszter esetére. Egy  $C \subseteq V$  klaszter konduktanciája megegyezik a vágásai legkisebb konduktanciájával, ahol az  $(S, C \setminus S)$  vágás konduktanciája:  $\phi(S) = \frac{\sum_{p \in S, q \in C \setminus S} w(p, q)}{\min(a(S), a(C \setminus S))}$ . Egy klaszterezés konduktanciája a klaszterek minimális konduktanciájával egyezik meg. A klaszterezés célja tehát az, hogy keressük meg azt a klaszterezést, ami a legnagyobb konduktanciát adja. A 7.2 és a 7.3 ábrákon látható pontokat a konduktancia alapú klaszterező eljárások helyesen csoportosítják.

Sajnos a konduktancia alapú mérték még nem tökéletes. Ha például egy jó minőségű klaszter mellett van néhány pont, amelyek mindentől távol esnek, akkor a klaszterezés minősége igen gyenge lesz (hiszen a minőség a leggyengébb klaszter minősége). A probléma egy lehetséges kiküszöbölése, ha a klaszterezés minősítésére két paramétert használunk. A konduktancia mellett bevezethetjük azt a mértéket, amely megadja, hogy az összes él súlyának hányad részét nem fedik le a klaszterek.

**7.4. definíció.** A  $\{C_1, C_2, \dots, C_k\}$  a  $(V, E)$  gráf egy  $(\alpha, \epsilon)$ -partíciója, ha:

I. minden  $C_i$  klaszter konduktanciája legalább  $\alpha$ ,

II. a klaszterek közötti élek súlya legfeljebb  $\epsilon$  hányada az összes él súlyának.

A klaszterezés célja ekkor az lehet, hogy adott  $\alpha$  mellett találjunk olyan  $(\alpha, \epsilon)$ -partíciót, amely minimalizálja  $\epsilon$ -t, vagy fordítva (adott  $\epsilon$  mellé találjunk olyan  $(\alpha, \epsilon)$ -partíciót, amely maximalizálja  $\alpha$ -t). A feladat NP-nehéz.

## 7.5. Klaszterező algoritmusok típusai

A szakirodalomban jónéhány klaszterező algoritmus található. Nem létezik ideális klaszterező algoritmus, mivel az eredmények összehasonlítására nincs objektív mérték. Az egyes alkalmazások jellegétől függ, hogy melyik algoritmust célszerű választani.

A klaszterező algoritmusokat 5 kategóriába soroljuk.

**Partíciós módszer:** A partíciós módszerek a pontokat  $k$  diszjunkt csoportra osztják úgy, hogy minden csoportba legalább egy elem kerüljön. A csoportok a klasztereknek felelnek meg. Egy kezdeti particionálás után egy újraparticionálási folyamat kezdődik, mely során egyes pontokat más csoportba helyezünk át. A folyamat akkor ér véget, ha már nem „mozognak” az elemek.

**Hierarchikus módszer:** A hierarchikus módszerek a klaszterekből egy hierarchikus adatszerkezetet (általában fát, amit a szakirodalomban *dendogram*nak neveznek) építenek fel.

**Spektrál módszerek:** Spektrál módszerek közé soroljuk az olyan algoritmusokat, amelyek a csoportok meghatározásához az adathalmazt reprezentáló mátrix sajátértékeit, illetve sajátvektorait használja fel.

**Sűrűség-alapú módszerek:** A legtöbb klaszterező algoritmus csak elliptikus alakú klasztereket tud kialakítani. A sűrűség-alapú módszerek ennek a hibának a kiküszöbölésére születtek meg. Az alapvető ötlet az, hogy egy klasztert addig növesztenek, amíg a sűrűség a „szomszédságban” meghalad egy bizonyos korlátot. Pontosabban egy klaszteren belüli elemekre mindig igaz, hogy adott sugarú körön belül mindig megtalálható bizonyos számú elem. A sűrűség-alapú módszereket a klaszterezés mellett kivételek, kívülálló elemek felderítésére (outlier analysis) is alkalmazzák.

**Grid-alapú módszerek:** A grid-alapú módszerek az elemeket rácspontokba képezik le, és a későbbiekben már csak ezekkel a rácspontokkal dolgoznak. Ezeknek az algoritmusoknak a gyorsaság a fő előnyük.

Klaszterező algoritmusokkal Dunát lehetne rekeszteni. Szinte bármilyen „butuska” klaszterező algoritmushoz tudunk generálni olyan adathalmazt, amit az fog a legjobban csoportosítani. Sajnos ezt a tényt a cikkek szerzői is gyakran kihasználják. A végeredményen kívül akadnak még szempontok, amelyeket meg lehet vizsgálni az egyes klaszterező algoritmusoknál. A legfőbb elvárásaink az alábbiak lehetnek:

**Skálázhatóság:** Sok algoritmus csak akkor hatékony, ha az elemek elférnek a memóriában. Sajnos a gyakorlatban gyakran olyan nagy adatbázisokat kell feldolgozni, hogy ez a feltétel nem tartható.

**Adattípus:** Vannak algoritmusok, amelyek csak intervallum típusú attribútumokkal megadott elemeken működnek. Nyilvánvaló, hogy ez a feltétel szűkíti az alkalmazások körét.

**Tetszőleges alakú, méretű és sűrűségű klaszterek:** A legtöbb klaszterező algoritmus csak elliptikus klasztereket képes felfedezni. A gyakorlati életben azonban ritkán elliptikusak a klaszterek. Jogos elvárás, hogy az algoritmus akár amőba alakú, sőt egymásba ágyazódó klasztereket is meg tudjon határozni. Emellett jól tudjon csoportosítani eltérő méretű és sűrűségű elemhalmazokat.

**Előzetes ismeretek:** Elvárjuk, hogy az algoritmusok automatikusan meghatározzák a szükséges klaszterek számát. Sajnos vannak algoritmusok, amelyeknek előre meg kell adni ezt a paramétert.

**Zajos adatok, távol eső elemek kezelése:** A legtöbb adatbázis tartalmaz valamekkora zajt, kivételes, a többségtől távol eső elemeket. Rossz tulajdonsága egy algoritmusnak, ha ezeknek az elemeknek nagy hatása van a klaszterek kialakítására.

**Adatok sorrendjére való érzékenység:** Miért fogadnánk el az algoritmus eredményét, ha az teljesen megváltozik, mihelyt más sorrendben vesszük az elemeket? Az eredményként kapott klaszterek nem függhetnek az adatok feldolgozásának sorrendjétől.

**Dimenzió:** Bizonyos algoritmusok csak alacsony dimenzió esetén hatékonyak. Vannak azonban olyan alkalmazások, ahol az elemek nagyon sok paraméterrel vannak leírva, azaz az elemeket egy magas dimenziójú tér elemeiként kell felfognunk.

**Értelmezhetőség:** A felhasználók azt várják el a klaszterező algoritmusoktól, hogy olyan klasztereket találjanak, amelyek jól meghatározott jegyekkel bírnak, és viszonylag könnyű magyarázatot adni arra, hogy milyen tulajdonságú elemek tartoznak az egyes klaszterbe.

## 7.6. Particionáló eljárások

A particionáló algoritmusoknál a csoportok száma előre adott ( $k$ ). Azért nevezzük ezeket az eljárásokat particionáló eljárásoknak, mert a legelső lépésben particionáljuk az elemeket, és a továbbiakban csak áthelyezgetünk bizonyos elemeket az egyik részből a másikba. Akkor kerül egy elem egy másik részbe, ha ezáltal javul a klaszterezés minősége. A klaszterezés minőségére az egyes partíciós algoritmusok eltérő célfüggvényt használnak. Egy lépés során a célfüggvény javítására általában több lehetőség is kínálkozik. Általában az algoritmusok a legjobbat választják ezek közül, tehát a „legmeredekebb lejtő” irányába lépnek a célfüggvény völgyekkel teli görbéjén.

### 7.6.1. Forgó $k$ -közép algoritmus

A  $k$ -közép algoritmus ( $k$ -means algorithm) [42] az egyik legrégebbi (1965-ből származik) és legegyszerűbb klaszterező algoritmus vektortérben megadott elemek csoportosítására. A klaszterezés minőségének jellemzésére a négyzetes hibafüggvényt használja.

Az algoritmus menete a következő: kezdetben választunk  $k$  darab véletlen elemet. Ezek reprezentálják eleinte a  $k$  klasztert. Ezután besorolunk minden pontot ahhoz a klaszterhez, amely reprezentáns eleméhez az a leginkább hasonló. A besorolás után új reprezentáns pontot választunk, éspedig a klaszter középpontját. A besorolás, új középpont választás iterációs lépéseket addig ismételjük, amíg történik változás.

Jancey 1966-ban Forgó-tól teljesen függetlenül ugyanezt az algoritmust javasolta egy apró módosítással [66]. Az új reprezentáns pont ne az új középpont legyen, hanem a régi és az új középpontot összekötő szakaszon, például a középponton. Ez egy visszafogottabb, kisebb lépésekben haladó algoritmus. A kisebb lépések előnye, hogy esetleg nem lesznek túllövéses, túl nagy oszcillációk. Elméletileg azonban egyikről sem lehet elmondani, hogy jobb lenne a másiknál, bizonyos helyzetekben az egyik, máskor a másik ad jobb eredményt.

Az algoritmus szép eredményeket hoz, amennyiben a klaszterek egymástól jól elszigetelődő kompakt „felhők”. Előnye, hogy egyszerű és jól skálázható, futási ideje  $O(nkt)$ , ahol  $t$  az iterációk számát jelöli. A  $k$ -közép algoritmusnak azonban számos hátránya van.

- I. Lehet, hogy az algoritmus lokális optimumban áll meg, tehát az iterációk során nem változik semmi, mégis létezik olyan csoportosítás, ahol a négyzetes hiba kisebb.
- II. Csak olyan elemek csoportosítására használható, amelyek vektortérben vannak megadva, hiszen értelmezni kell az elemek középpontját. Ezek szerint a  $k$ -közép nem használható olyan alkalmazásokban, ahol az elemek attribútumai között például kategória típusú is szerepel.
- III. Rendelkezik a négyzetes hibát minimalizáló algoritmusok minden hibájával (lásd a 7.4.1-es részt).

A lokális optimumba kerülés esélyének csökkentése érdekében érdemes az algoritmust többször futtatni különböző kezdeti pontokkal. Azt a csoportosítást fogadjuk el, amelynek legkisebb a négyzetes hibája. Vegyük észre, hogy ez a megoldás erős heurisztika! Minél nagyobb  $n$ , elvben annál több lokális optimum lehet, annál nagyobb az esélye, hogy lokális optimumban kötünk ki. Ismereteink szerint nincs olyan képlet, ami megmondja, hogy adott elemszám esetén hányszor kell futtatni az algoritmust, hogy biztosan (vagy adott valószínűséggel) megtaláljuk a globális optimumot.

---

**Weka 3.5.7** A  $k$ -közép algoritmust a `weka.clusterers.SimpleKMeans` osztály implementálja.

---



### 7.6.2. A $k$ -medoid algoritmusok

Ezek az algoritmusok a  $k$ -közép két hibáját próbálják kiküszöbölni. Egyrészt az eredmény kevésbé érzékeny a kívülálló pontokra, másrészt csak a hasonlósági értékeket használja. Tehát

nem feltétel, hogy az elemek vektortérben legyenek megadva.

A  $k$ -medoid algoritmusokban egy klasztert nem a középpont reprezentál, hanem a leginkább középen elhelyezkedő elem, a medoid. Továbbra is egy négyzetes hiba jellegű függvényt próbálunk minimalizálni, de a négyzetes hiba itt a medoidoktól való távolságok összegét jelenti ( $k$ -medián probléma, lásd 7.4.1 rész).

## A PAM algoritmus

A PAM (Partitioning Around Medoids) algoritmus [71] menete teljesen megegyezik a  $k$ -közép menetével. Kezdetben választunk  $k$  véletlen elemet, ezek lesznek először a medoidok. Ezután elkezdődik az elemhözrendelés medoidokhoz, új medoid választása iteratív folyamat. Egy elemet a legközelebbi medoidhoz rendelünk.

Abban az esetben választunk új medoidot egy klaszterben, ha ezzel csökken a négyzetes hiba. Határozzuk meg az összes nem medoid, medoid párra  $(x, x_m)$ , hogy mennyivel változna a négyzetes hiba, ha  $x_m$ -nek átvinné a szerepét  $x$ . Nyilvánvaló, hogy nincsenek hatással a négyzetes hiba változására azok az elemek, amelyekhez van  $x$  és  $x_m$ -nél közelebbi medoid. A négyzetes hiba változásánál nem csak  $x_m$  klaszterébe tartozó elemeket kell megvizsgálnunk, hiszen lehet, hogy a medoidváltás hatására egyes elemek új klaszterbe kerülnek. Ha vannak olyan párok, amelyeknél a négyzetes hiba változása negatív, akkor cseréljen szerepet annak a párosnak a két eleme, amelyhez tartozó négyzetes hiba változás abszolút értékben a legnagyobb.

Sajnos az algoritmus lassú, hiszen egy iteratív lépés időigénye  $O(k(n-k)^2)$ . Összehasonlítva a  $k$ -közép algoritmussal elmondhatjuk, hogy lassabb, viszont kevésbé érzékeny a távol eső pontokra.

## A CLARA és CLARANS algoritmusok

A PAM algoritmus nem alkalmas nagy adathalmazok klaszterezésére, mert lassú. A CLARA és CLARANS algoritmusok a PAM algoritmus kiterjesztései. Nem vizsgálnak meg minden medoid, nem medoid párt, hanem csak néhányat. Így az iterációs lépésben elvégzendő ellenőrzések száma kisebb, ami által az algoritmusok nagyobb adathalmazokon is használhatók.

A CLARA algoritmus [71] az eredeti adatbázis helyett egy véletlenszerűen választott mintán dolgozik. A medoidok csak ebből a véletlen mintából kerülhetnek ki, de a négyzetes hibát a teljes adatbázisra nézve számítja. Az iterációs lépés időigénye így  $O(k(n'-k)(n-k))$ , ahol  $n'$  a minta nagysága.

Ha a legkisebb négyzetes hibát eredményező  $k$  elem nincs a mintában, akkor a CLARA nem fogja megtalálni az optimális megoldást. Célszerű ezért az algoritmust több véletlenszerűen kiválasztott mintán lefuttatni, és a legkisebb négyzetes hibát szolgáltatót elfogadni.

A CLARANS algoritmus [97] az eredeti adathalmazon dolgozik. Nem az összes lehetséges csere által eredményezett négyzetes hiba változását számítja ki, hanem csak egy, véletlenszerűen választott párét. Ha a változás negatív (sikeres választás), akkor a pár elemei szerepet cserélnek, ellenkező esetben új párt sorsolunk. A felhasználó egy paraméterrel szabályozhatja a sikertelen választások maximális számát, amely elérésével az algoritmus véget ér.

A CLARANS sem biztos, hogy megtalálja a legkisebb négyzetes hibát adó  $k$  medoidot mielőtt a sikertelen próbálkozások száma elérné a küszöböt. Ezért többször futtassuk az algoritmust mindig más kezdeti medoidokkal.

Vegyünk észre egy fontos tulajdonságot: eredményezhetik ugyanazt a klaszterezést különböző medoidok. Valószínű, hogy az optimális közeli megoldás ugyanazt a csoportosítást adja, mint a legkisebb négyzetes hibát szolgáltató medoidok. Ezért javasolt a fenti heurisztikák alkalmazása nagy adathalmazok esetén.

A CLARANS nagy adathalmazokon való alkalmazhatóságával foglalkoznak a [37] cikkben.  $R^*$ -fák használatával feloldják azt a kényszert, miszerint a pontok férjenek el a memóriában. A PAM és rokonainak hibája, hogy a  $k$ -medián problémát próbálja megoldani, így csak egyszerű, eliptikus klasztereket képes felfedezni.

## 7.7. Hierarchikus eljárások

A hierarchikus eljárások onnan kapták a nevüket, hogy az elemeket, klasztereket, alklasztereket hierarchikus adatszerkezetbe rendezik el. Két fajta hierarchikus eljárást különböztetünk meg: a *lentről építkezőt*, más néven *egyesítőt* és a *fentről építkezőt*, avagy az *osztót*. A lentről építkező eljárásoknál kezdetben minden elem külön klaszter, majd a nagyon közeli klasztereket egyesíti, amennyiben azok teljesítenek bizonyos feltételt. A fentről építkezők fordítva működnek: kezdetben egyetlen klaszter létezik, amit kisebb alklaszterekre osztunk, majd ezeket finomítjuk tovább.

A hierarchikus algoritmusok kényes pontja az egyesítendő vagy osztandó klaszterek kiválasztása. Miután egy egyesítés (osztás) megtörténik, az összes további műveletet az új klaszteren végezzük. Ezek a műveletek tehát nem megfordítható folyamatok, így rossz választás rossz minőségű klaszterezéshez vezet.

### 7.7.1. Single-, Complete-, Average Linkage Eljárások

A legegyszerűbb egyesítő, hierarchikus eljárás az alábbi.

- I. Kezdetben minden pont külön klaszterhez tartozik.
- II. Keressük meg, majd egyesítsük a legközelebbi klasztereket.
- III. Ha a klaszterek száma lecsökkent  $k$ -ra, akkor álljunk meg, ellenkező esetben ugorjunk a 2. lépésre.

Ez az egyszerű eljárás a távolság mátrixszal dolgozik, feltételezi, hogy az elfér a memóriában. A távolság mátrix egy felső háromszög-mátrix, amelynek  $i$ -edik sorának  $j$ -edik eleme tárolja az  $i$  és  $j$  elemek távolságát. Célszerű kiegészíteni minden sort két extra információval: a legközelebbi klaszter sorszámával és annak távolságával.

Az eljárás tár- és időigénye (az összehasonlítások száma)  $O(n^2)$ . A mai tárkapacitások mellett (1-2 Gbyte memóriával rendelkező gép teljesen hétköznapi számítás) ez azt jelenti, hogy az elemek száma 30-40 ezernél nem lehet több.

Amennyiben két klaszter távolságát a legközelebbi pontjaik távolságával definiáljuk, akkor az eljárást *single linkage eljárásnak* nevezzük. A single linkage rendkívül alkalmas jól elkülönülő, tetszőleges alakú klaszterek felfedezésére. Mivel a minimális távolságon alapszik, ezért ha a klaszterek túl közel kerülnek egymáshoz, akkor hajlamos a single linkage összekötni őket, és esetleg a klaszteren belül egy vágást képezni.



További hibája, hogy érzékeny az outlierekre. Általában az outlierek távol esnek a többi ponttól, így ezeket a pontokat nem fogja egyesíteni. Például két jól elszeparálódó, sok pontot tartalmazó klasztert és egy nagyon távoli pontot úgy oszt két részre, hogy az egyik részben lesz a távoleső pont, a másikban pedig az összes többi. Ha tudjuk, hogy olyan adathalmazt kell klasztereznünk, ahol a (tetszőleges alakú) csoportok jól elkülönülnek egymástól, továbbá nincsenek outlierok, akkor a single eljárás jó munkát végez.

Ha az eljárást gráfelméleti szemszögből nézzük (teljes gráfban a pontoknak az elemek, az éleken lévő súlyoknak pedig a távolságok felelnek meg), akkor a single-linkage eljárás egy minimális feszítőfát fog találni, amennyiben a klaszterek számának egyet adunk meg értékül. Ha  $k$  darab csoportba akarjuk sorolni a pontokat, akkor ezt a minimális feszítőfa  $k - 1$  legnagyobb élének elhagyásával nyerhetjük. Azon elemek lesznek egy klaszterben, amelyek egy komponensbe kerültek. Rájöhetünk arra is, hogy a single-linkage eljárás nem más, mint Kruskal algoritmus más köntösben.

Ha két klaszter távolságának megállapításához a minimális távolság helyett a maximális távolságot használjuk, akkor *complet linkage* eljárásról beszélünk, ha pedig az átlagos hasonlóságot, vagy az egyesített klaszter átmérőjét, akkor *average linkage* eljárásról.

### 7.7.2. Ward módszere

Ward módszere a particionáló eljárásokhoz hasonlóan a legkisebb négyzetes hibát próbálja minimalizálni (tehát csak vektortérben megadott pontoknál alkalmazható). Az egyszerű hierarchikus eljárástól csak az egyesítendő klaszterek kiválasztásának módjában különbözik. Azt a két klasztert egyesíti, amelyek a legkisebb négyzetes hibanövekedést okozzák (nyilvánvalóan kezdetben – amikor minden pont külön klaszter – a négyzetes hibaösszeg 0).

A módszer rendelkezik a négyzetes hibát minimalizáló eljárások minden hibájával. Emellett nem skálázható, hiszen a távolságmátrixszal dolgozik, és végül nem garantált, hogy megtalálja a globális minimumot.

### 7.7.3. A BIRCH algoritmus

Ezt az algoritmust nagy adathalmazok klaszterezésére találták ki. Csak vektortérben adott elemeket tud klaszterezni. Egy klasztert a  $CF = (N, \vec{L}S, SS)$  hármas (Cluster Feature) jellemez, ahol  $N$  a klaszterben található elemek száma,  $\vec{L}S = \sum_{i=1}^N \vec{x}_i$  és  $SS = \sum_{i=1}^N |\vec{x}_i|$ . Egy klaszter CF-je a klaszter statisztikai jellemzőit tárolja: a nulladik, első és második momentumokat. Az algoritmus során a klaszterek CF-értékeit tároljuk, a benne lévő elemeket nem. Egy klaszter CF-jéből ki tudjuk számolni a klaszter átmérőjét vagy akár két klaszter átlagos távolságát.

A CF-ekből az algoritmus egy ún. CF-fát épít fel. A CF-fa egy gyökeres, (lefelé) irányított fa. A fa leveleiben CF-értékek vannak, egy belső pont pedig a pontból induló alfához tartozó klaszterek egyesítéséből kapott CF-értéket tárolja. A fának két paramétere van. Az első határozza meg a belső pontból induló ágak maximális számát (ágszám korlát). A második paraméter adja meg, hogy mekkora lehet maximálisan a levélhez tartozó klaszterek átmérője. Ennek a paraméternek nagy hatása van a fa méretére: minél kisebb a maximális átmérő, annál több kis klasztert kell létrehozni, tehát annál nagyobb lesz a fa.

A BIRCH algoritmus két fázisra oszlik. Az elsőben az elemek egyszeri végigolvasása során felépítjük a CF-fát. Ez már eleve egyfajta klaszterezést eredményez. A második fázisban minden

elemet besorolunk a hozzá legközelebbi klaszterbe, majd esetleg ebből kiindulva lefuttatunk egy particionáló algoritmust (például a  $k$ -közepet).

Az első fázis során kapott CF-fa – az ágszám korlát miatt – nem valószínű, hogy meg fog felelni a természetes klaszterezésnek. Lehet, hogy bizonyos pontok, amelyeknek egy klaszterben kellene lenniük, szét lesznek választva, és a fordítottja is előfordulhat. Sőt, az is megtörténhet, hogy ugyanazok az elemek a fa építésének különböző fázisaiban különböző levelekbe fognak kerülni!

A cikk szerzői javaslatot adnak az outlierok kiszűrésére: ha a CF-fában valamely levélben található pontok száma „jóval kevesebb”, mint a levelekben található pontok átlagos száma, akkor töröljük a levelet, mert valószínűleg outlierokat tartalmaz. A felhasználónak kell megadni az elemszámra vonatkozó küszöbszámot, ami alatt töröljük a leveleket. Vegyük észre, hogy ez a megoldás erős heurisztika. Számtalan példát lehetne mutatni, amikor fontos pontokat is töröl, miközben valódi outlierokat a fában hagy.

A BIRCH algoritmus tehát tényleg meg tud bírkozni igazán nagy méretű adathalmazokkal, azonban rendelkezik szinte az összes rossz klaszterezési tulajdonsággal. Mivel a második szakaszban egyfajta  $k$ -közép algoritmust futtatunk, ezért a BIRCH-re is igazak a  $k$ -középről mondottak. De ezen kívül érzékeny az adatok sorrendjére, és nem skála-invariáns, hiszen a CF-fában lévő klaszterek maximális átmérőjét a felhasználónak kell megadnia.

#### 7.7.4. A CURE algoritmus

A CURE (Clustering Using REpresentatives) algoritmus [51] átmenet a BIRCH és a single linkage eljárás között a reprezentáns elemek számát illetően (a BIRCH-ben a középpont a reprezentáns pont, a single linkage-ben a klaszter összes elemét számon tartjuk.). Egy klasztert  $c$  darab (ahol  $c$  előre megadott konstans) elem jellemez. Ezek az elemek próbálják megragadni a klaszter alakját. Ennek következménye, hogy a CURE nem ragaszkodik az eliptikus klaszterekhez.

Hierarchikus eljárás lévén, kezdetben minden elem külön klaszter, majd elkezdődik a klaszterek egyesítése. Az egyesítésnek három lépése van.

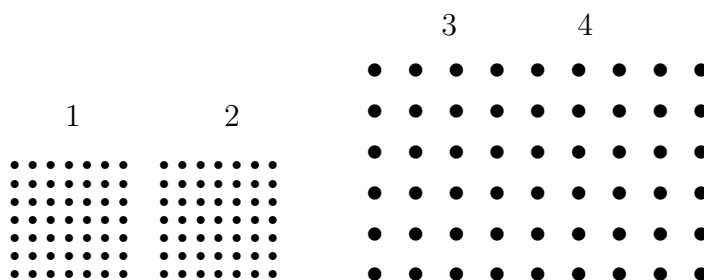
- I. A reprezentáns pontok alapján kiválasztja a két legközelebbi klasztert. Két klaszter távolságát reprezentáns pontjaik távolságának minimuma adja.
- II. Egyesíti a klasztereket, majd a  $2c$  reprezentáns pont közül kiválaszt  $c$  darabot, amelyek várhatóan jól fogják reprezentálni az egyesített klasztert. Ennek módja a következő. Legyen az első reprezentáns pont a középponttól legtávolabbi elem. A következő  $c - 1$  lépésben mindig az előzőleg kiválasztott ponthoz képest a legtávolabbit válasszuk reprezentáns elemnek.
- III. A reprezentáns pontokat „összehúzza”, azaz az általuk kijelölt középpont felé mozdulnak úgy, hogy az új távolság a középponttól az  $\alpha$ -szorosa legyen az eredeti távolságnak. Ennek a lépésnek a célja az outlierok hatásának csökkentése.

Az egyesítés akkor ér véget, amikor a klaszterek száma eléri  $k$ -t. Az eljárás végeztével rendelkezésünkre áll  $c$  reprezentáns ponttal jellemzett  $k$  darab klaszter. Ezután a teljes adatbázis egyszeri végigolvasásával az elemeket besoroljuk a hozzá legközelebbi klaszterbe a legközelebbi reprezentáns pont alapján.

A CURE algoritmust felkészítették, hogy kezelni tudjon nagy adathalmazokat is. Véletlen mintát vesz, azt felosztja részekre, majd az egyes részeket klaszterezi (ez a rész tehát párhuzamosítható). A kapott klaszterekből végül kialakítja a végső klasztereket. A részletek és az algoritmus során felhasznált adatszerkezetek ( $k$ -d fa és kupac) iránt érdeklődőknek ajánljuk a [51]-t.

A CURE algoritmus számos hibával rendelkezik:

- I. Az elemeknek vektortérben adottnak kell lenniük.
- II. Minden klasztert rögzített számú reprezentáns pont jellemez. De miért jellemezzon ugyanannyi pont egy kis kör alakú klasztert és egy nagy amőba alakút? Természetesebb lenne, ha a reprezentáns pontok száma függene a klaszter méretétől és alakjától.
- III. A reprezentáns pontok kiválasztása sem túl kifinomult. A módszer nem a klaszter alakját fogja meghatározni, hanem inkább egy konvex burkot. Gondoljuk meg, ha egy kör alakú klaszterben van egy bemélyedés, akkor a bemélyedés környékén található pontokat sosem fogja az eljárás reprezentánsnak választani, hiszen ők közel vannak a többi ponthoz. Amőba alakú klaszternél tehát a reprezentáns pontok alapján nem lehet megállapítani, hogy hol vannak a bemélyedések, így azt sem dönthetjük el, hogy egy nagy ellipszissel van dolgunk, vagy egy érdekes alakú amőbával.
- IV. Klaszter egyesítése után a reprezentáns pontokat összehúzzuk a középpont felé. Nagy klaszter esetében sok egyesítés, így sok összehúzás van. Az összehúzás által távolabb kerülnek a reprezentáns pontjai más reprezentáns pontoktól, így egyre ritkábban lesz kiválasztva egyesítésre. Mintha az algoritmus „büntetné” a nagy klasztereket.
- V. Rosszul kezeli az eltérő sűrűségű pontokat. Ezt leginkább az alábbi ábra illusztrálja. A



7.4. ábra. Hibás klaszterezés: eltérő sűrűségű klaszterek esetén

CURE az 1-es és 2-es klasztereket fogja egyesíteni (azok reprezentáns pontjai vannak egymáshoz legközelebb) a 3-as és 4-es klaszterek egyesítése helyett.

Megjegyezzük, hogy az algoritmust bemutató cikk hosszú bevezetőjében éppen arra hívta fel a figyelmet, hogy a mások által publikált algoritmusok mennyire rosszul kezelik a különböző méretű és amőba alakú klasztereket. Ennek ellenére a tesztekben bemutatott adathalmazban nagyságrendileg azonos méretűek voltak a klaszterek és alakjuk elliptikus volt.

### 7.7.5. A Chameleon algoritmus

A Chameleon két nagy fázisra oszlik. Kiindulásként előállítja a  $k$ -legközelebbi gráfot, majd ezt részekre osztja. A második fázisban bizonyos részeket egyesít, előállítva ezzel a végleges csoportokat.

A Chameleon az első olyan hierarchikus algoritmus, amely a klaszterek egyesítésénél nem csak a klaszterek távolságát ( $d(C_i, C_j)$ ) veszi figyelembe, hanem az egyes klasztereken belüli távolságokat is, pontosabban a *relatív belső kapcsolódásukat* ( $RI(C_i, C_j)$ ) is (relative interconnectivity). Abban az esetben egyesít két klasztert, amennyiben  $d(C_i, C_j)$  és  $RI(C_i, C_j)$  is nagy értéket vesz fel. Ennek az ötletnek köszönhető, hogy a Chameleon – szemben az eddigi algoritmusokkal – jól tud klaszterezni eltérő sűrűségű adathalmazokat is. Nézzük meg, hogyan definiálja az algoritmus két klaszter relatív belső kapcsolódását és relatív távolságát.

#### Relatív távolság

#### Relatív belső kapcsolódás

## 7.8. Sűrűség-alapú módszerek

A sűrűség alapú módszerek (density based methods) szerint egy klaszteren belül jóval nagyobb az elemek sűrűsége, mint a klaszterek között. Ez alapján lehet elválasztani a klasztereket és azonosítani az outliereket.

### 7.8.1. A DBSCAN algoritmus

A DBSCAN a legelső sűrűség-alapú eljárás [38]. A sűrűség meghatározásához két paramétert használ, egy sugár jellegű mértéket ( $eps$ ) és egy elemszám küszöböt ( $minpts$ ). A  $p$  elem *szomszédai* ( $N_{eps}(p)$ ) azok a elemek, amelyek  $p$ -től legfeljebb  $eps$  távolságra vannak. A  $q$  elem a  $p$ -ből *sűrűség alapon közvetlen elérhető*, ha  $q \in N_{eps}(p)$  és  $|N_{eps}(p)| \geq minpts$ . Naivan azt gondolhatnánk, hogy egy klaszterben található elemek sűrűség alapon közvetlen elérhetőek egymásból. Ez nem így van, hiszen a klaszter határán lévő elemek  $eps$  távolságán belül nincs mindig  $minpts$  darab elem.

A  $q$  elem *sűrűség alapon elérhető*  $p$ -ből, ha léteznek  $p_1 = p, p_2, \dots, p_n = q$  elemek úgy, hogy  $p_{i+1}$  sűrűség alapon közvetlen elérhető  $p_i$ -ből. A  $p$  és  $q$  elemek *sűrűség alapon összekötöttek*, ha létezik olyan  $o$  elem, amelyből  $p$  és  $q$  sűrűség alapon elérhető. A klaszter definíciója ezek alapján:

**7.5. definíció.** *Az elemek egy  $C$  részhalmaza klaszter, amennyiben*

*I. Ha  $p \in C$  és  $q$  sűrűség-alapon elérhető  $p$ -ből, akkor  $q \in C$  (maximalitás).*

*II. Ha  $p, q \in C$ , akkor  $p$  és  $q$  sűrűség alapon összekötöttek.*

Egy elemet *zajnak* (noise) hívunk, ha nem tartozik egyetlen klaszterbe sem.

Legyen a  $C$  klaszter egy  $p$  eleme olyan, hogy  $|N_{eps}(p)| \geq minpts$ . Ekkor könnyű belátni, hogy  $C$  megegyezik azoknak a elemeknek a halmazával, amelyek  $p$ -ből sűrűség alapján elérhetőek. E tulajdonságot használja az algoritmus. Válasszunk egy tetszőleges elemet ( $p$ ) és határozzuk meg a sűrűség alapján elérhető elemeket. Amennyiben  $|N_{eps}(p)| \geq minpts$  feltétel teljesül, akkor

meghatároztunk egy klasztert. A feltétel nemteljesülés nem jelenti azt, hogy  $p$  zaj, lehet, hogy egy klaszter határán helyezkedik el.  $|N_{eps}(p)| < minpts$  esetén egyszerűen válasszunk egy új elemet. Ha már nem tudunk új elemet választani, akkor az algoritmus véget ér. Azokat az elemeket tekintjük zajnak (outliernek), amelyeket nem soroltunk semelyik klaszterbe.

A DBSCAN algoritmus előnye, hogy tetszőleges alakú klasztert képes felfedezni, és ehhez csak az elemek távolságát használja. Hátránya, hogy rendkívül érzékeny a két paraméterre ( $eps$ ,  $minpts$ ). Sőt amennyiben a klaszterekben található elemek sűrűsége eltérő, akkor nem biztos, hogy lehet olyan paramétereket adni amivel a DBSCAN jó eredményt ad.

---

***Weka 3.5.7*** A *DBScan* algoritmust a *weka.clusterers.DBScan* osztály implementálja.

---



## 8. fejezet

### Idősorok elemzése

## 9. fejezet

# Webes adatbányászat

Az Internetről történő automatikus információkinyerő alkalmazások gombamód szaporodnak napjainkban. A terület mélyebb áttekintése túlmutat ezen írás keretein, ezért csak a talán legfontosabb két témát járjuk körül: a weboldalak rangsorolását és az intelligens Internetes keresést.

Az oldalak közötti megfelelő rangsor felállítása napjaink kritikus feladata. A keresőrendszerek mindennapos eszközökké váltak. Naponta milliók használják, így a helyes működésük mindenki érdeke.

Minden honlapkészítő álma, hogy az oldala elsőként jelenjen meg a keresők által visszaadott listában. Ez cégeknek sok látogatót és így sok potenciális ügyfelet jelent, másfelől a gyakran látogatott oldalakon elhelyezett reklámok is jó bevételt jelentenek. Központi szerepük miatt fokozott támadásoknak vannak kitéve. Egy rangsoroló algoritmus elkészítésekor ezért fontos megvizsgálni, hogy az milyen trükkökkel lehet azt becsapni, és ezek ellen hogyan kell védekezni.

### 9.1. Oldalak rangsorolása

Képzeld el azt a ritkának nem mondható helyzetet, amikor egy keresőrendszer a feltett kérdésünkre rengeteg oldalt talál, olyan sokat, hogy kivitelezhetetlen feladat egyesével átnézni azokat és kiválasztani a fontosakat. Mégis tudjuk azt, hogy a talált oldalaknak valamilyen köze van a kérdésünkhöz: egyeseknek több, másoknak kevesebb.

Szükség van tehát az oldalak automatikus rangsorolására, aminek alapkövetelménye, hogy formálisan is tudjuk definiálni egy weboldal „fontosságát”. Felmerül a kérdés, hogy objektíve a fontosság definiálása. A válasz egyszerű: nem. A kérdésre kiadott dokumentumok között ugyanis különböző emberek nem ugyanazt a sorrendet állítanák fel.

A feladatot meg kell oldani, akkor is ha tökéletes megoldást még elméletileg sem tudunk adni. Megelégszünk ezért a fontosság valamilyen heurisztikán alapuló közelítésével. Algoritmikusan szemlélve két algoritmuscsalád létezik, az egyik családba tartozó algoritmusok a *gráf összes pontját súlyozzák*, majd a súlyok rendezésével állapítják meg a sorrendet. Ezek a *globális algoritmusok*, míg a másik családot *kérdésfüggő rangsorolásoknak* nevezhetjük, ami azt jelenti, hogy a rangsoroló algoritmus minden kérdésnél lefut, és ekkor csak egy *részgráf* csúcsait pontozza. Mindkét családnak megvan a maga előnye, az elsőnek csak egyszer kell lefutni, és utána csak memóriából való olvasás a keresőszoftver feladata, míg a második figyelembe tudja venni azt a tényt, hogy egy weboldal különböző témájú kereséseknél különböző módon szignifikáns.

### 9.1.1. Az egyszerű Page Rank

A Google keresőrendszerben 1998-ban implementált Page Rank (Brin-Page) algoritmus a gyakorlati alkalmazások során nagyon jó eredményt hozott [102]. A továbbiakban az ő módszerüket mutatjuk be.

Rendelkezésünkre áll  $N$  darab weboldal a hagyományos weboldalak minden tulajdonságával. Feladat lenne ezek között egyfajta fontossági sorrendet felállítani. Egy oldal fontosságát, hasznosságát jól tükrözi az oldalt meglátogató emberek száma. A legtöbb oldal készítői azonban a letöltések számát illetően semmilyen auditálást nem végez, így rangsoroló algoritmust nem alapozhatunk ezen információkra.

A linkeknek nagy szerepük van a fontosságban. Ha valaki saját oldalán egy másik oldalra mutató linket helyez el, akkor azt azért teszi, mert szerinte, a másik oldal hasznos információt tartalmaz, kapcsolódik az oldal témájához, valamilyen szempontból fontos. A Page Rank algoritmus (és minden kifinomult kereső rendszer) az oldalak közötti linkstruktúra alapján definiálja a fontosságot.

Egy oldal fontosságát az oldal rangja adja meg. Elképzeléseinknek megfelel az az állítás, hogy ha valahova sok link mutat, akkor az fontos oldal, továbbá, ha egy oldal fontos, akkor az általa mutatott lapok is azok. Informálisan egy rekurzív definíciót adnánk a fontosságnak: „egy oldal fontos, ha fontos oldalak mutatnak rá”.

A rang meghatározásához szükségünk van az oldalak közötti linkstruktúra ismeretére. Defináljuk az  $N$  weblaphoz  $A$   $N \times N$ -es sor-sztochasztikus mátrixot az alábbiak szerint: amennyiben az  $i$ -edik lapon  $n$  link található, akkor

$$A_{ij} = \begin{cases} \frac{1}{n} & \text{ha } j\text{-re mutat link } i\text{-ről} \\ 0 & \text{egyébként} \end{cases}$$

Az  $A$  mátrix (sor-)sztochasztikus, azaz  $\forall i$ -re  $\sum_{j=1}^N A_{ij} = 1$ ,  $A_{ij} \geq 0$ . A sor-sztochasztikus mátrixokra igaz a következő tétel:

**9.1. tétel.** *Legyen  $A$  sor-sztochasztikus mátrix ( $N \times N$ -es),  $j = (\frac{1}{N}, \dots, \frac{1}{N})$ . Ekkor*

$$p = \lim_{m \rightarrow \infty} jA^m$$

*létezik és  $pA = p$ .*

**9.2. definíció.** *A  $p = (p_1, \dots, p_N) \in \mathbb{R}_+^N$  vektor a lapok rang-vektora (tehát az  $i$ -edik lap rangja  $p_i$ ).*

Az algoritmus menete a következő:

- I. Készítsük el az  $A$  mátrixot az adott weblapok topológiájából.
- II. Kezdetben minden oldal rangja  $\frac{1}{N}$ , tehát  $p = (\frac{1}{N}, \dots, \frac{1}{N})$ ,
- III. végezzük el  $p_{i+1} \leftarrow p_i A$  iterációt,
- IV. ha teljesülnek a leállási feltételek akkor STOP, ellenkező esetben ugrás az előző utasításra.



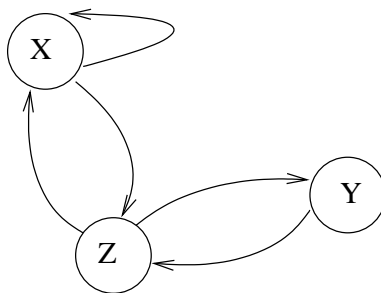
Leállási feltétel lehetne az, hogy a  $p$  rang-vektor egy adott küszöbnél kisebbet változik. Az eredeti célunk azonban az egyes oldalak rangsorolása, nem pedig a pontos rangértékek meghatározása. Ezért sokkal ésszerűbb, ha akkor állítjuk le az iterációt, ha a rang-vektor alapján felállított sorrend nem változik egy adott számú iteráció után. A fent kimondott tétel a garancia arra, hogy az iteráció során  $p$  rang-vektor egy vektorhoz konvergál, amiből következik, hogy az algoritmus minden esetben le fog állni.

Képzeljük azt, hogy kezdetben minden oldal fontossága  $\frac{1}{N}$ , és minden lap a következő lépést hajtja végre: a saját fontosságát egyenlő mértékben szétosztja az általa mutatott oldalak között. Könnyű végiggondolni, ha a fenti lépést hosszú időn keresztül folytatják, akkor minden lap fontossága meg fog egyezni a fent definiált rang-vektor laphoz tartozó rangértékével.

A fenti algoritmus elfogadásához egy másik intuitív magyarázat lehetne az alábbi: Tegyük fel, hogy a „sztochasztikus szörfölő” egy olyan, az Interneten barangoló lény, aki a kiindulási lapot, egyenletes eloszlás szerint, véletlenszerűen választja ki, valamint minden következő oldalt az aktuálisról elérhetőek közül választja ki hasonlóan véletlenszerűen. Belátható, hogy annak a valószínűsége, hogy végtelen sok lépés után a szeszélyes szörfölő az  $i$ -edik lapra kerül,  $p_i$ .

Az algoritmus vitathatatlan előnye, hogy gyors ( $N \cdot j \cdot A^M$  jól számítható) és könnyen programozható.

Nézzünk egy nagyon egyszerű példát az algoritmusra. 3 oldalt kell rangsorolnunk, amelyek linkstruktúrája a következő ábrán látható.



9.1. ábra. Példa az egyszerű Page Rank algoritmusra

A topológia alapján az  $A$  mátrix:

$$A = \begin{pmatrix} \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 1 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$

Az első három iteráció után a rang vektor  $N$ -szerese:

$$N \cdot p_1 = (1, 1, 1)$$

$$N \cdot p_2 = \left(1, \frac{1}{2}, \frac{3}{2}\right)$$

$$N \cdot p_3 = \left(\frac{9}{8}, \frac{1}{2}, \frac{11}{8}\right)$$

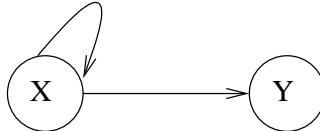
$$N \cdot p_4 = \left(\frac{5}{4}, \frac{11}{16}, \frac{17}{16}\right)$$

Megmutatható, hogy  $p = (\frac{6}{5}, \frac{3}{5}, \frac{6}{5})$

Ennek az egyszerű algoritmusnak két nagy hibája van, melyeket zsákutca, illetve pókháló problémának hívunk.

### Zsákutca probléma

Zsákutcának nevezzük azt az oldalt, amiről nem mutat link semmilyen más lapra, de más lapról mutat rá. Amennyiben az oldalak között zsákutca van, akkor az  $A$  mátrix ehhez az oldalhoz tartozó sora csupa 0 elemet fog tartalmazni. Ekkor az  $A$  mátrix nem lesz sor-sztocasztikus, és oldalak fontossága „kiszivárog” a rendszerből. A probléma szemléltetésére nézzük a következő ábrán látható lapstruktúrát.



9.2. ábra. Példa zsákutcára

A hozzá tartozó mátrix:  $A = \begin{pmatrix} \frac{1}{2} & \frac{1}{2} \\ 0 & 0 \end{pmatrix}$ . Könnyen ellenőrizhető, hogy  $A^2 = \begin{pmatrix} \frac{1}{4} & \frac{1}{4} \\ 0 & 0 \end{pmatrix} = \frac{1}{2}A$ , továbbá  $A^m = \frac{1}{2^{m-1}}A$ , amiből adódik, hogy a rangvektor a 0 vektorhoz fog tartani.

### Pókháló probléma

Lapok olyan rendszerét, amelyben minden link csak e rendszerbeli lapra mutat, pókhálónak nevezzük. Jellemző rájuk, hogy az iteráció során magukba gyűjtik (esetleg az összes) a fontosságot. Ez komoly visszaélésekhez adhat alapot és SPAM-eléshez vezethet, hiszen linkek eltávolításával bárki alakíthat ki pókhálót, amennyiben van arra az oldalra mutató link.

Példaként térjünk vissza a 9.1 laptopológiához, csak most tegyük fel, hogy  $Y$  a  $Z$ -re mutató linkjét átállítja úgy, hogy ezentúl saját magára mutasson. Ekkor  $A$  mátrix a következőképpen módosul:

$$A = \begin{pmatrix} \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 1 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}$$

a rang vektor  $N$ -szerese az első négy iteráció során:

$$N \cdot p_1 = (1, 1, 1)$$

$$N \cdot p_2 = (1, \frac{3}{2}, \frac{1}{2})$$

$$N \cdot p_3 = (\frac{3}{4}, \frac{7}{4}, \frac{1}{2})$$

$$N \cdot p_4 = (\frac{5}{8}, 2, \frac{3}{8})$$

$$N \cdot p_5 = (\frac{1}{2}, \frac{35}{16}, \frac{5}{16})$$

belátható, hogy a rang vektor a  $p = (0, 1, 0)$  vektorhoz fog tartani.

### 9.1.2. Az igazi Page Rank

A fenti két probléma kiküszöbölésére az oldalak megadóztatását javasolták. Ennek ötlete az, hogy szedjük be mindenkitől fontosságának bizonyos százalékát, majd a beszedett adót osszuk el egyenlően. Amennyiben  $\epsilon$ -nal jelöljük a befizetendő adót, akkor a fentiek alapján  $A$  mátrix

helyett a  $B = \epsilon \cdot U + (1 - \epsilon) \cdot A$  mátrixot használjuk, ahol  $U = \begin{pmatrix} \frac{1}{N} & \frac{1}{N} \\ \dots & \dots \\ \frac{1}{N} & \frac{1}{N} \end{pmatrix}$ .

Könnyen ellenőrizhető, hogy a  $B$  mátrix sor-sztocasztikus, így alkalmazhatjuk rá a 9.1-es tételt, ami ismét garantálja, hogy az algoritmus le fog állni.

Az igazi Page Rank algoritmusban az egyes lapok nem csak szomszédjaiknak osztják szét fontosságukat, hanem először befizetik az adót a királyi kincstárba, és csak a maradékot osztják szomszédjaiknak. Fontosságot pedig kapnak a rá mutató oldalak mellett a kincstárban található beszedett adóból is, egyenlő mértékben.

Amennyiben  $A$  mátrix helyett  $B$  mátrixot alkalmazzuk, a sztochasztikus szörfölőre nem lesz igaz az, hogy  $p_i$  annak valószínűsége, hogy  $i$ -edik oldalra lép. Igaz lesz viszont a „szeszélyes sztochasztikus szörfölőre”, akire  $\epsilon$  valószínűséggel rájön a szeszély, és ilyenkor a következő állomását, egyenletes eloszlást követve, véletlenszerűen választja a lapok közül.

Az igazi Page Rank algoritmust a kezdeti Google(<http://www.google.com>) keresőrendszer használta a talált oldalak rangsorolásához. A keresőrendszerről részletesebb leírás található a [23] cikkben.

## 9.2. Webes keresés

Internetes keresés során egy keresőrendszertől két típusú kérdésre kérhetünk választ:

**tág kérdés** A választ tartalmazó, vagy a kérdéshez kapcsolódó oldalak száma nagy. Ilyen kérdés lehet, hogy információt szeretnénk a java nyelvről, vagy a gépkocsigyártókról.

**szűk kérdés** Ezen olyan specifikus kérdést értünk, amelyre a választ kevés oldal tartalmazza. Ilyen kérdés lehet, hogy „A 2001. *Űrodüsszeia* hányadik percében hangzik el az első emberi szó?”

Szűk kérdésre a válaszadás automatikus módja jóval nehezebb feladat, mint tág kérdésre. Szűk kérdésnél annak veszélye fenyeget, hogy egyáltalán nem találunk választ pusztán hasonló szavakon alapuló kereséssel. Tág kérdéseknél ezzel szemben a probléma éppen a válaszhoz kapcsolódó lapok túl nagy száma lehet. Ebben a részben arra keresünk választ, hogy miként tudjuk kiválasztani a tág kérdésre kapott nagy mennyiségű oldalból a kérdéshez leginkább kapcsolódó oldalakat.

### 9.2.1. Gyűjtőlapok és Tekintélyek – a HITS algoritmus

Az 1999-ben Jon Kleinberg által publikált Gyűjtőlapok és Tekintélyek (Hubs and Authorities) módszere [75] a lapok linkstruktúráját használja fel. A linkstruktúra mellett számos információ állhat rendelkezésünkre, amelyek segítségünkre lehetnek az oldalak fontosságának meghatározásában. A látogatások számát már említettük. Probléma vele, hogy az oldalak elenyésző részét figyelik auditáló szoftverek.

Az oldalon elhelyezett metaadatok, kulcsszavak, az oldal leírása, de ezenkívül a szövegben kiemelt szavak (dőlt betű, vastag betű, villogó betű ...) szintén segíthetnek a kérdéshez kapcsolódás mértékének eldöntésében. A tanulmányban ezek szerepét nem vesszük figyelembe. Jelöljük  $\sigma$ -val a kérdést, amire a választ keressük. Az algoritmus fázisai a következők:

- I.  $M_\sigma$  (mag)laphalmaz kiválasztása hagyományos keresővel.
- II.  $M_\sigma$  bővítésével bázis lap-részgráf konstruálása. Jelöljük ezt a bázist  $B_\sigma$ -val.
- III. A  $\sigma$ -hoz tartozó gyűjtőlapok és tekintélyek (szimultán) kiszűrése  $B_\sigma$ -ból.

A gyűjtőlapoknak és tekintélylapoknak nem adunk pontos matematikai definíciót. Minden oldalhoz egy gyűjtőlap- és egy tekintélyértéket fogunk rendelni. Minél nagyobbak ezek az értékek, annál inkább tekintünk egy oldalt az adott kérdéshez tartozó gyűjtő-, illetve tekintélylapnak. Intuitív definíciója a két fogalomnak a következő lehetne: gyűjtőlap az olyan lap, ami sok tekintélylapra mutat, tekintélylapok pedig azok, amire sok gyűjtőlap mutat. Ezek szerint a gyűjtőlapok a  $\sigma$  szempontjából értékes linkek gyűjteménye, a tekintélylapok pedig a  $\sigma$  kérdéshez kapcsolódó értékes információkat tartalmazó lapok. Például az AMS honlapja egy matematikai gyűjtőlap, Jeffrey D. Ullman adatbányászatról szóló jegyzetvázlata pedig tekintélylap, amennyiben  $\sigma =$  "adatbányászati algoritmusok". Amikor egy kérdést felteszünk, akkor elsősorban a válasz érdekel bennünket, nem pedig az olyan oldalak, amik sok hasznos oldalra mutatnak. Az eredmény szempontjából a tekintélyoldalak a fontosak. Ezek megtalálásához gyakran a gyűjtőoldalakon keresztül vezet az út, így érdemes őket együtt keresni. Most pedig nézzük részletesen az algoritmus egyes lépéseinek működését.

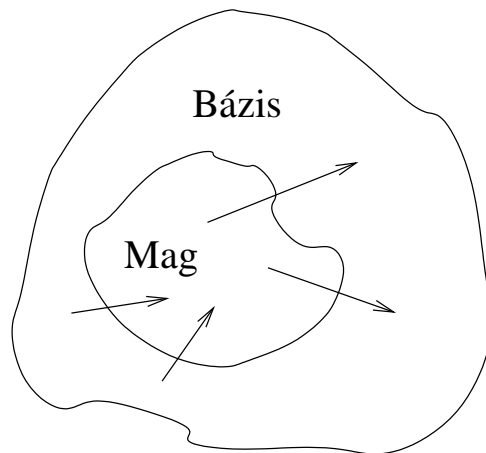
### $M_\sigma$ mag meghatározása

Az algoritmus kiindulását képező weboldaloknak egy hagyományos kereső által  $\sigma$  kérdésre kiadott első  $t$  darab lapját vesszük. Ez a kezdőkészlet azonban nem mentes a hagyományos keresőrendszerek által adott hibáktól. Egyrésztől lehet, hogy fontos oldalak nincsenek benne a találati listában. A "gépkocsi gyártók" kérdésre például nem fogják kiadni a Honda honlapját, mert a lapon ilyen szóösszetétel nincsen. Másrésztől sok olyan oldalt is generálni fog, amelyek nem kapcsolódnak a témához. Ennek több oka is lehet, például az, hogy a kérdésnek több értelme is van (gondoljunk itt a Java nevű szigetre), vagy az egyes oldalak „hazudnak”, azaz olyan tartalmat állítanak magukról, amelyek nem igazak(pl.:mp3, free holiday ...). A fenti hátrányok ellenére elmondhatjuk, hogy ennek a magnak a „környezete” már hasznos információkban gazdag lesz.

### $B_\sigma$ bázis létrehozása

A gyűjtőlapokat és a tekintélyoldalakat a bázisból fogjuk kinyerni, így ezzel szemben az alábbi elvárásaink vannak:

- I. Ne legyen túl nagy!
- II. Legyen fontos lapokban gazdag!
- III. Tartalmazza a  $\sigma$ -hoz releváns lapokat (vagy azok legtöbbszörét)!



9.3. ábra. Bázis generálása a magból

A tesztelés során kapott eredmények azt mutatták, hogy az alábbi egyszerű algoritmus a gyakorlatban jól működik. Induljunk ki az előző pontban definiált magból (azaz legyen  $B_\sigma = M_\sigma$ ), majd adjuk hozzá az összes olyan oldalt, amelyre mutat link valamely  $B_\sigma$ -beli oldalról. Ezen kívül vegyük  $B_\sigma$ -hoz azokat az oldalakat, amelyekről mutat link valamely  $B_\sigma$ -beli lapra. Elképzelhető, hogy népszerű oldal is van  $B_\sigma$ -ban, amelyre rengeteg oldal mutathat, ezért egy oldal maximum egy előre meghatározott konstans ( $d$ ) számú új lap felvételét „okozhatja”. Ezért ha egy lapra  $d$ -nél több lap mutat, akkor válasszunk ezek közül véletlenszerűen  $d$  darabot. Töröljük a bázisból a navigációt szolgáló éleket (pl.: vissza az előző oldalra) úgy, hogy csak a különböző hosztok közötti élek maradjanak. Itt azt a feltételezést tettük, hogy a hosztokat meg lehet különböztetni URL-jük alapján (Ez nyilván nem tökéletes megoldás, gondoljunk csak a unix alapú rendszerekre, ahol az egyes felhasználók honlapjának domainnevei megegyeznek. Nem könnyű kérdés az, hogy egy adott domain mikor tekintsünk csak egy oldalnak, illetve mikor osszuk fel többre.

Kleinberg tapasztalata szerint a  $t = 200$ ,  $d = 50$  mellett a bázis mérete 1000 és 5000 között lesz.

### Tekintélyek kinyerése

A tesztek alapján a bázis tartalmazni fogja a tekintélyek nagy részét. Hogyan leljük meg ezeket a több ezer oldal közül? Első ötlet lehetne, hogy a nagy be-fokú csúcsok reprezentálják a kereséshez kapcsolódó fontos oldalakat. Ez a megoldás azonban felemás eredményt ad: a jó oldalak mellett lesznek úgynevezett „univerzálisan népszerű” oldalak is. Ezekre jellemző, hogy  $\sigma$ -tól függetlenül a legtöbb kérdéshez tartozó bázisban megtalálhatóak. Például, ha  $\sigma = \text{„java”}$ , akkor a  $B_\sigma$ -ban a legnagyobb be-fokú csúcsokhoz tartozó oldalak a

- I. **www.gamelan.com**
- II. **java.sun.com**
- III. **amazon.com**
- IV. karibi vakációkat hirdető oldal

Az utolsó két oldalt valamilyen automatikus módon ki kellene szűrni.

Kleinbergnek a következő szűrő ötlete támadt. A  $\sigma$  kérdéshez tartozó tekintélyeknek nagy befokon kívül jellemzője, hogy nagy az átfedés azokban a laphalmazokban, amik rájuk mutatnak. Ezekben benne lesznek a téma gyűjtőlapjai. A következő ábra szemlélteti a tekintélyek és az univerzálisan népszerű lapok közötti különbséget. A téma gyűjtőlapjai és tekintélyei általában



9.4. ábra. Topológiai különbség a tekintélyek és az univerzálisan népszerű lapok között

egy sűrű páros gráfot alkotnak, míg az univerzálisan népszerű lapokra szabálytalanul, összevissza mutatnak a linkek.

A sűrű páros gráf megtalálása a következőképpen történik. Legyen  $C$  a  $B_\sigma$  weblaphalmazhoz tartozó szomszédossági mátrix, tehát  $c_{ij} = 1$  ha  $i \rightarrow j$ , 0 különben. Ez hasonlít a Page Rank algoritmusnál ismertetett  $A$  mátrixra, azzal a különbséggel, hogy nincs sztochasztikusan skálázva. Rendeljünk minden laphoz egy gyűjtőlap, illetve egy tekintélylap értéket, tehát vezessük be a

$$g = (\dots, g_i, \dots), g_i \geq 0$$

$$t = (\dots, t_i, \dots), t_i \geq 0$$

gyűjtő-, illetve tekintély vektorokat, amelyek legyenek normált vektorok, tehát  $\|g\| = \|t\| = 1$ . A két vektorra a tekintély és gyűjtőlap intuitív definíciója miatt legyen érvényes a következő két szabály:

$$g = \lambda C t$$

$$t = \mu C^T g$$

azaz egy lap gyűjtőértéke az általa mutatott tekintélyértékeinek összege-  $\lambda$ -val skálázva, és egy lap tekintélyértéke azon lapok gyűjtőértékeinek összege, amelyek rá mutatnak- $\mu$ -vel skálázva.

A két egyenletet egymásba írva:

$$g = \lambda \mu C C^T g$$

$$t = \lambda \mu C^T C t$$

Hasonlóan, mint az oldalak rangját a Page Rank algoritmusnál, a  $g$  és  $t$  vektorokat is iteratíván határozzuk meg. A lépések:

$$\text{I. } t^{(0)} = g^{(0)} = \begin{pmatrix} \frac{1}{|B_\sigma|} \\ \vdots \\ \frac{1}{|B_\sigma|} \end{pmatrix}$$

$$\text{II. } \hat{t}^{(i+1)} \leftarrow C^T C t^{(i)} \text{ és } \hat{g}^{(i+1)} \leftarrow C C^T g^{(i)}$$

$$\text{III. } t^{(i+1)} \leftarrow \frac{\hat{t}^{(i+1)}}{\|\hat{t}^{(i+1)}\|} \text{ és } g^{(i+1)} \leftarrow \frac{\hat{g}^{(i+1)}}{\|\hat{g}^{(i+1)}\|}$$

IV. ha teljesül a leállási feltétel, akkor STOP, ha nem GOTO 2

A leállási feltételről hasonló mondható el, mint a Page Rank algoritmusnál: nem  $g$  és  $t$  pontos értéke érdekel bennünket, hanem az első néhány, legnagyobb tekintélyértékkel rendelkező oldal. A tapasztalati eredmények azt mutatták, hogy 20 iteráció után a legnagyobb 5-10 tekintélyértékkel rendelkező oldal már stabilizálódik.

A kísérleti eredmények mellett mindig hasznos, ha matematikai tételek is igazolják azt, hogy az algoritmus véget fog érni, azaz  $t^{(i)}$  és  $g^{(i)}$  konvergálnak valahova. A következő tétel ezt a matematikai megalapozást nyújtja. A tétel bizonyítása a B függelékben található.

**9.3. tétel.** *A fent definiált  $t^{(i)}$  és  $g^{(i)}$  sorozatok konvergálnak nemnegatív értékű vektorokhoz.*

Kleinberg módszere igen jó eredményt ért el lényeges oldalak kiszűrésénél nagy találati halmazokból. Például a  $\sigma$ ="Gates"-re, a legfontosabb oldalnak a <http://www.roadahead.com>-ot találta, majd ezek után jöttek a Microsofthoz kapcsolódó oldalak. A győztes oldal Bill Gates könyvének hivatalos weblapja, amit az AltaVista csak a 123. helyre rangsorolt.

### 9.2.2. A SALSA módszer

Az algoritmus ([84], Stochastic Approach for the Link-Structure Analysis) a már megismert Mag és Bázis halmazokon dolgozik, és egy véletlen sétát valósít meg az alább definiált gráfokon, amely az eredeti gráf pontjainak Gyűjtőlap és Tekintély tulajdonságait emeli ki.

A  $G_t$  ill.  $G_g$  gráfok csúcsai legyenek az eredeti gráf csúcsai (a weboldalak), az  $i$  és  $j$  pont között pedig annyi él van, ahány olyan csúcs (Gyűjtőlap) van, amiből  $i$ -be és  $j$ -be is mutat link, ill. hány olyan csúcs (Tekintély) van, amibe  $i$ -ből és  $j$ -ből is van él.

Megjegyzésként elmondható, hogy a HITS algoritmusban egy (dupla) lépés alatt ezen gráfok összes élén továbbadtuk az induló csúcs pontszámát, míg a SALSA algoritmusnál nem az egészet, hanem figyelembe vesszük azt, hogy minden csúcs ugyanannyit továbbbítson, így egy-egy Markov láncot definiálunk a gráfokon.

Az  $M_t$  és  $M_g$  Markov láncok formális definíciójához a  $B(i) = \{k : k \rightarrow i\}$  mellett szükségünk lesz a  $F(i) = \{k : i \rightarrow k\}$  jelölésre. Az előző bekezdés szerint a megfelelő átmenetvalószínűségek a következők:

$$P_t(i, j) = \sum_{k: k \in B(i) \cap B(j)} \frac{1}{|B(i)|} \frac{1}{|F(k)|} \text{ ill.}$$

$$P_g(i, j) = \sum_{k: k \in F(i) \cap F(j)} \frac{1}{|F(i)|} \frac{1}{|B(k)|}.$$

Az egyensúlyi súlyokat kiszámító iteráció indítása

$$[\underline{t}]_0 := [\underline{g}]_0 := \frac{1}{N}(1, \dots, 1)^T,$$

majd az iteráció lépése:

$$[\underline{t}(i)]_k := \sum_j P_t(j, i) [\underline{t}(j)]_{k-1}, \text{ ill.}$$

$$[\underline{g}(i)]_k := \sum_j P_g(j, i) [\underline{g}(j)]_{k-1}.$$

Feltéve egy pillanatra, hogy a Markov láncaink irreducibilisek, azaz a két fent definiált gráf összefüggő, az állítható, hogy az egyensúlyi eloszlásokban két pont tekintély ill. gyűjtőlap súlyának aránya megegyezik az eredeti gráfban vett be- ill. ki-fokszámainak arányával. Az állítás abból következik, hogy irreducibilis Markov láncnak egyértelmű a stacionárius eloszlása, és a fenti súlyarányokat feltéve az állítás ellenőrizhető a következőképpen:

Az irreducibilitás miatt egyértelmű stacionárius eloszlás ki kell elégítse, hogy

$$\forall i \ \underline{t}(i) = \sum_j P_t(j, i) \underline{t}(j).$$

Most  $B$ -vel az élek halmazát jelölve és feltéve az előzőek szerint, hogy

$$\forall i \ \underline{t}(i) = \frac{|B(i)|}{|B|},$$

így számolhatunk:

$$\begin{aligned} \underline{t}(i) &= \sum_j \underline{t}(j) P_t(j, i) = \sum_j \underline{t}(j) \sum_{k \in B(j) \cap B(i)} \frac{1}{|B(j)|} \frac{1}{|F(k)|} = \\ &= \sum_j \frac{|B(j)|}{|B|} \sum_{k \in B(j) \cap B(i)} \frac{1}{|B(j)|} \frac{1}{|F(k)|} = \\ &= \sum_j \frac{|B(j)|}{|B|} \frac{1}{|B(j)|} \sum_{k \in B(j) \cap B(i)} \frac{1}{|F(k)|} = \\ &= \frac{1}{|B|} \sum_j \sum_{k \in B(j) \cap B(i)} \frac{1}{|F(k)|} = \\ &= \frac{1}{|B|} \sum_{k \in B(i)} \sum_{j \in F(k)} \frac{1}{|F(k)|} = \\ &= \frac{1}{|B|} \sum_{k \in B(i)} 1 = \frac{|B(i)|}{|B|} \end{aligned}$$

Ennek megfelelően a leírt iteratív algoritmus lefuttatására tulajdonképpen nincs szükség, hiszen a stacionárius eloszlás az előbbi elméleti eredmény felhasználásával *közvetlenül számítható*. Ezzel együtt természetesen az is igaz, hogy az algoritmus könnyen becsapható, hiszen – az előző fejezetben leírtak szerint – az oldalunkra mutató linkek száma tetszőlegesen növelhető.

Itt jegyezzük meg, hogy a SALSA az egyenletes eloszlással indított HITS algoritmus első lépésének felel meg, ezután az első lépés után a SALSA stacionárius eloszlásának súlyai jelennek meg.



Több komponensből álló gráf esetén az algoritmus csak abban a komponensben dolgozik, ahonnan indult a séta. Mivel az indulás egyenletesen lett választva, ezért egy adott komponensből való indulás valószínűsége a komponens méretével arányos, azaz az alapgráfot  $G$ -vel, komponenseit  $G_k$ -val, az  $i$  csúcs komponensének indexét  $j$ -vel jelölve

$$a_i = \frac{|G_j|}{|G|} \frac{|B(i)|}{\sum_{\alpha \in G_j} |B(\alpha)|},$$

ahol a nevezőben levő összeg a komponens összes éleinek száma.

### 9.2.3. Gyűjtőlapok, Tekintélyek és véletlen séták

Mint láthattuk, a SALSA algoritmus ötlete az eredeti gráfból egyszerűen származtatható másik gráf(ok)on megvalósított véletlen séta volt. Láttuk továbbá, hogy az eredeti Gyűjtőlap és Tekintély algoritmusunk első lépése ekvivalens a SALSA-val. Jogosan kérdezhetjük tehát, hogy az eredeti algoritmusnak létezik-e véletlen séta analogonja, illetve átfogalmazva a kérdést, hogy az eredeti algoritmus is kapcsolatba hozható-e Markov láncok stacionáris eloszlásaival? A válasz persze igenlő, hiszen minden sztochasztikus vektorhoz létezik olyan Markov lánc, amelynek stacionáris eloszlása éppen az a vektor. A kérdés már csak az, hogy létezik-e olyan ezzel a tulajdonsággal bíró Markov lánc is, amelynek átmenetvalószínűségei az eredeti gráfból származtathatók?

A válasz – kissé meglepő módon – az, hogy az algoritmus összes közbülső eredménye előáll az eredeti gráfból származtatható Markov lánc stacionáris eloszlásaként, bár az, hogy az első fél lépés után már igaz ez (ott a SALSA súlyok jelennek meg), már előrevetíti az eredményt. Vezessük be a következő jelöléseket: egy  $B$  illetve egy  $F$  lépésnek egy a webgráfban levő link követését nevezzük hátra illetve előre irányban. Ezek kombinációit is definiáljuk, például  $BFBF = (BF)^2$  egy négylépéses sétát jelent a webgráfban. Az  $i$  pontból a  $j$  pontba vezető  $(BF)^n$  séták halmazát jelölje  $(BF)^n(i, j)$ , az  $i$  pontból induló  $(BF)^n$  séták halmazát jelölje  $(BF)^n(i)$ , továbbá az összes  $(BF)^n$  séták halmazára használjuk magát a  $(BF)^n$  jelölést! Az  $(FB)^n$  séták halmazai hasonlóan értendők.

Most definiáljuk a következő két Markov láncot: az állapotok halmaza az összes csúcs, amely magában az alapgráfban is benne volt, míg két csúcs között pontosan akkor van él, ha az alapgráfban van köztük legalább egy  $(BF)^n$  illetve  $(FB)^n$  séta. Az átmenetvalószínűségek pedig legyenek:

$$P_t(i, j) := \frac{|(BF)^n(i, j)|}{|(BF)^n(i)|}, \quad \text{illetve}$$

$$P_g(i, j) := \frac{|(FB)^n(i, j)|}{|(FB)^n(i)|}.$$

A definíciókból az látható, hogy

$$|(BF)^n(i, j)| = (C^T C)^n(i, j) \quad \text{és}$$

$$|(FB)^n(i, j)| = (C C^T)^n(i, j), \quad \text{és ezekből}$$

$$|(BF)^n(i)| = \sum_j (C^T C)^n(i, j) \quad \text{és}$$

$$|(FB)^n(i)| = \sum_j (C C^T)^n(i, j).$$

Az eredeti HITS algoritmus  $n$ . iterációja után a pontszám vektorok normálás nélkül  $(C^T C)^n \mathbf{1}$  illetve  $(C C^T)^n \mathbf{1}$ , azaz összeg normában ez ugyanaz, mint a megfelelő Markov láncok stacionáris eloszlása:

$$\underline{t}(i) = \frac{|(BF)^n(i)|}{|(BF)^n|}$$

$$\underline{g}(i) = \frac{|(FB)^n(i)|}{|(FB)^n|}$$

Elmondható tehát, hogy az algoritmus végső pontszámarányai a csúcsokból induló hosszú  $BF$  illetve  $FB$  séták számainak arányától függ, aminek az a következménye, hogy nagyon erősen kötött alakzatok (teljes páros részgráfok) környékét az algoritmus kiemeli.

### 9.2.4. Automatikus forrás előállító - Gyűjtőlapok és Tekintélyek módosításai

Gyűjtőlapok és Tekintélyek alapú keresést sikerrel alkalmazták automatikus forrás előállítás során (automatic resource compilation, röviden ARC) [27]. A továbbiakban erről szólnak pár szót.

Általánosabb fogalmak keresésénél gyakran használunk előre szerkesztett hierarchikus fogalomtárakat. A legismertebb fogalomtárak a Yahoo! vagy az Infoseek oldalán találhatóak. Ha például információra van szükségünk a tangóról, akkor a Yahoo! főoldaláról a Recreation & Sports (kikapcsolódás és sport) linket választva eljuthatunk egy újabb oldalra. Itt már választhatjuk a dance (tánc) linket, majd a Ballroom-ot (társas) és végül a tangót. Innen már nem léphetünk tovább újabb alkategória kiválasztásával, hanem a tangóval foglalkozó legfontosabb weboldalak listáját láthatjuk. Mind a fogalomhierarchia felépítése, mind az egyes fogalmakhoz tartozó legfontosabb weboldalak megkeresése manuális úton történik, tehát emberek járják a világhálót és keresik az olyan oldalakat, amelyek tényleg hasznos információval szolgálnak a fogalomról.

Az ARC-nál a második lépést próbálták automatizálni: adott egy tág fogalom, keressük meg a hasznos információkat tartalmazó weboldalakat. Ehhez a Gyűjtőlapok és Tekintélyek keresést használták, két módosítással.

Egyrészt kétszer, nem pedig egyszer alkalmazták azt a lépést, amely során a Magból( $M_\sigma$ ) a Bázis( $B_\sigma$ ) előállították. Emiatt a Bázis mérete nőtt, viszont nem veszünk el olyan oldalt, amely a hagyományos kereső által kiadott oldalaktól 2 link távolságra van.

Másrészt módosították az iteráció során használt  $C$  mátrixot is. Tudjuk, hogy a weboldalak HTML kódjában a  $\langle A \text{ HREF}="" \rangle_i / A_i$  tag jelent egy linket. Ha például egy oldalon a  $\langle A \text{ HREF}="" \rangle$  tag található, akkor ha a szörfölő az *ingyen sms* szóra kattint, akkor a **www.mtnsms.com** oldalra kerül.

Megfigyelték, hogy nagyon gyakran a HREF tag környezetében az oldalt jellemző szavak találhatóak. Ez nem meglepő, hiszen az oldalak készítői minél jobban próbálják segíteni az oldalt látogatóinak navigációját. A tag környezete tehát fontos, mert ha megtalálható ott a kérdéses fogalom, akkor várható, hogy a link egy hasznos oldalra mutat.

Szomszédossági mátrix helyett ezért olyan mátrixot javasoltak, amely elemei a következőképp számíthatók:

$$c_{ij} = \begin{cases} 1 + n(f) & \text{ha } j\text{-re mutat link } i\text{-ről} \\ 0 & \text{egyébként} \end{cases}$$

ahol  $n(f)$  a fogalom előfordulásának száma egy adott szélességen belül a HREF tagtól.

A szélességet kísérleti úton próbálták meghatározni: azt vizsgálták, hogy pár ismert oldalra mutató több ezer oldalon hol található meg az ismert oldalakat jellemző szó. A tesztek eredményeként megállapították, hogy ha az oldalon megtalálható a jellemző szó, akkor 97%-ban az a HREF 50 byte-os környezetében is megtalálható.

Az algoritmust implementálták, és széleskörű felmérést készítettek, amelyben a megkérdezetteknek arra kellett válaszolniuk, hogy szerintük adott fogalmakra a három kereső(ARC, Infoseek, Yahoo!) közül melyik találta meg a legjobb oldalakat. A felmérésből kiderült, hogy a teljesen automatikus, emberi munkát nem igénylő ARC ugyanolyan jól teljesített, mint a másik két rendszer [27].

### 9.2.5. Gyűjtőlapok és Tekintélyek módszerének hátrányai

Vizsgálatok kimutatták, hogy a Gyűjtőlapok és Tekintélyek módszerének három hátránya van [16].

- I. Előfordulhat, hogy egy hoszton található dokumentumhalmaz minden eleme egy másik hoszton található dokumentumra mutató linket tartalmaz. Ez növelni fogja a dokumentumhalmaz elemeinek gyűjtőlap értékét és a másik hoszton található dokumentum tekintélyértékét. Ennek ellenkezője is könnyen előfordulhat, nevezetesen: egy hoszton található dokumentum több olyan dokumentumra mutat, amelyek egy másik hoszton találhatóak. Látható, hogy a hosztpárok létrehozásával a gyűjtőlap- és tekintélyértékek növelhetők, ami visszaélésre ad lehetőséget. Egy igazságos algoritmustól elvárjuk, hogy egyik hoszt se növelhesse túlzott mértékben mások fontosságát.
- II. A weboldalakat gyakran automatikusan állítják elő valamilyen segédeszköz segítségével. Ezek az eszközök sokszor linkeket helyeznek el a generált oldalakon. Például a Hypernews rendszer USENET cikkeket konvertál weblapokká úgy, hogy a Hypernews honlapjára mutató linket szúr az oldal végére. Ezekre a linkekre nem igaz a fejezet elején elhangzott állítás, miszerint az oldal szerzője azért helyezi el a linket oldalán, mert a másik oldal a saját oldal témájára nézve hasznos információkat tartalmaz.
- III. Bázis laphalmaz létrehozása során a Mag laphalmazhoz új oldalakat veszünk fel a linkstruktúra alapján. Az új oldalak között sok olyan lehet, amelyek nem kapcsolódnak a kérdéses témához. Amennyiben ezeket az oldalakat szoros linkstruktúra köti össze, akkor a „témásodródás” problémája mutatkozik: a legnagyobb tekintélyértékkel rendelkező oldalak csak tágabb értelemben fognak a témához kapcsolódni. Egy egyszerű teszt megmutatta, hogy a „jaguar *and* car” kérdésre a legjobb tekintélyoldalak (amelyek különböző autógyártó cégek honlapjai lettek) az általánosabb fogalomhoz (car) kapcsolódtak.

Az első esetben a problémát az okozza, hogy egy hoszton elhelyezett több dokumentum összbefolyása túl nagy lehet: minél több dokumentum található egy hoszton, annál inkább képes növelni más hoszton található dokumentum tekintély- vagy gyűjtőlapértékét.

Ideális esetben azt várnánk, hogy egy hoszton található dokumentumhalmaznak összesen akkora befolyása legyen, mintha ezen a hoszton csak egyetlen dokumentum lenne található. Ehhez módosítanunk kell az iteráció során használt mátrixot: amennyiben egy hosztrol  $k$  darab

dokumentum tartalmaz linket egy másik hoszton található dokumentumra, akkor a  $C$  mátrix ezen dokumentumaihoz tartozó értéke 1 helyett  $\frac{1}{k}$  legyen.

Az [16] cikkben a másik két problémára is javasoltak megoldást. Szövegelemzés felhasználásával a Bázisban található oldalakhoz relevancia értéket társítanak, ami megadja, hogy az adott oldal mennyire kapcsolódik a témához. A relevancia értéknek több szerepe van. Egyrészt a témához kis mértékben kapcsolódó (kis relevancia értékű) lapokat töröljük a Bázisból, másrészt a tekintély- illetve gyűjtőlapérték meghatározásához a lap relevanciaértékét is figyelembe vesszük: a relevanciaértékkel arányosan nő egy lap tekintély- illetve gyűjtőlapértéke. A szövegelemzéssel bővített Gyűjtőlapok és Tekintélyek módszerét a továbbiakban nem tárgyaljuk, a részletek megtalálhatók a [16] cikkben.

A fejezetben bemutatott két fő algoritmusról pár összehasonlító tesztet találhatunk a [8] cikkben.

## 10. fejezet

# Gyakori minták kinyerése

A fejlett társadalmakra jellemző, hogy számos, a mindennapi életünk során gyakran használt terméket és szolgáltatást nélkülözhetetlennek tartunk. Minél sokszínűbb a felhasználói csoport, annál nehezebb egy olyan üzenetet eljuttatni részükre, ami mindenki számára egyértelmű, ám ha valakinek ez sikerül, az nagy haszonnal járhat, hiszen pár százalékpontos növekedés is szignifikáns a nagy volumenben értékesített termékeknél. A piaci stratégiák kialakításánál is elsősorban a sokaságra, illetve a sokaság jellemzőire vagyunk kíváncsiak. Egyedi, külön elemek akkor érdekesek, ha például csalásokat akarunk felderíteni. Fenti eseteken kívül vizsgálhatjuk a gyakori balesetet okozó helyzeteket, a számítógépes hálózatban gyakran előforduló, riasztással végződő eseménysorozatokat, vagy pl. azt, hogy az egyes nyomtatott médiumoknak milyen az olvasói összetétele, és amennyiben több magazinnak, újságnak hasonló a célcsoportja, érdemes üzenetünket több helyen is elhelyezni, hogy hatékonyabban ösztönözzük meglévő és potenciális vásárlóinkat.

Oldalakon keresztül lehetne sorolni azon példákat, amikor a gyakran előforduló „dolgok” értékes információt rejtenek magukban. A szakirodalomban a dolgokat mintáknak nevezzük, és *gyakori minták kinyeréséről* beszélünk.

A minta típusa többféle lehet. Vásárlói szokások felderítésénél gyakori *elemhalmazokat* keresünk, ahol az elemek a termékeknek felel meg. Utazásokkal kapcsolatos szokásoknál a gyakran igénybe vett, költséges szolgáltatások sorrendje is fontos, így gyakori *sorozatokat* keresünk. Telekommunikációs hálózatokban olyan feltételek (predikátumok) gyakori fennállását keressük, amelyek gyakran eredményeznek riasztást. Ezeket a gyakori *bool formulákat* megvizsgálva kaphatjuk meg például a gyakori téves riasztások okait. A böngészési szokások alapján fejleszthetjük oldalaink struktúráját, linkjeit, így a látogatók még gyorsabban és hatékonyabban találják meg a keresett információkat. A böngészés folyamatát *címkézett gyökeres fákkal* jellemezhetjük. Gyakori mintákat kinyerő algoritmusokat a rákkutatásban is alkalmaztak. Azt vizsgálták, hogy a rákkeltő anyagokban vannak-e gyakran előforduló molekula-struktúrák. Ezeket a struktúrákat címkézett gráfokkal írjuk le.

A példákból következik, hogy a minta típusa sokféle lehet. Sejtethetjük, hogy más technikákat kell majd alkalmazni pl. címkézett gráfok keresésénél, mintha csak egyszerű elemhalmazokat keresünk. Ebben a részben egy általános leírást adunk, egy egységes matematikai keretbe helyezzük a gyakori minta kinyerésének feladatát. Emellett ismertetjük a legfontosabb módszerek általános – a minta típusától független – leírását.

## 10.1. A gyakori minta definíciója

E rész megértéséhez feltételezzük, hogy az olvasó tisztában van a 2.1 részben definiált fogalmakkal (rendezések, korlát, valódi korlát, maximális korlát, predikátum,).

**10.1. definíció.** *A  $H$  halmaz a  $\preceq$  rendezésre nézve lokálisan véges, ha minden  $x, y \in H$  elemhez, ahol  $x \preceq y$ , véges számú olyan  $z$  elem létezik, amelyre  $x \preceq z \preceq y$ .*

**10.2. definíció.** *Az  $\mathcal{MK} = (\mathcal{M}, \preceq)$  párost, ahol  $\mathcal{M}$  egy alaphalmaz,  $\preceq$  az  $\mathcal{M}$ -en értelmezett részben rendezés, mintakörnyezetnek nevezzük, amennyiben  $\mathcal{M}$ -nek pontosan egy minimális eleme van,  $\mathcal{M}$  halmaz a  $\preceq$  rendezésre nézve lokálisan véges és rangszámozott (graded), azaz létezik a  $|| : \mathcal{M} \rightarrow \mathbb{Z}$  ún. méretfüggvény, amire  $|m| = |m'| + 1$ , ha  $m$ -nek maximális valódi alsó korlátja  $m'$ . Az  $\mathcal{M}$  elemeket mintáknak (pattern) nevezzük és  $\mathcal{M}$ -re, mint mintahalmaz vagy mintatér hivatkozunk.*

Az  $m' \preceq m$  esetén azt mondjuk, hogy  $m'$  az  $m$  részmintája, ha  $m' \prec m$ , akkor valódi részmintáról beszélünk. A  $\preceq$ -t tartalmazási relációnak is hívjuk. Az általánosság megsértése nélkül feltehetjük, hogy a minimális méretű minta mérete 0. Ezt a mintát üres mintának hívjuk.

Íme az egyik legegyszerűbb példa mintakörnyezetre, amelyet vásárlói szokások feltárása során alkalmaztak először. Legyen  $\mathcal{J}$  véges halmaz. Gyakori elemhalmazok keresésénél a  $(2^{\mathcal{J}}, \subseteq)$  lesz a mintakörnyezet, ahol  $\subseteq$  a halmazok tartalmazási relációját jelöli. A méretfüggvény egy halmazhoz az elemszámát rendeli. Az elemhalmazokon túl kereshetünk gyakori sorozatokat, epizódokat (véges halmazon értelmezett részben rendezéseket), bool formulákat, címkézett gyökeres fákat vagy általános gráfokat. Ezen mintakörnyezetek pontos definícióját a következő fejezetekben találjuk.

**10.3. definíció.** *Legyen  $(H_1, \preceq_1)$   $(H_2, \preceq_2)$  két részben rendezett halmaz. Az  $f: H_1 \rightarrow H_2$  függvény rendezés váltó vagy más szóval anti-monoton, amennyiben tetszőleges  $x, y \in H_1$ ,  $x \preceq_1 y$  elemekre  $f(y) \preceq_2 f(x)$ .*

**10.4. definíció.** *A gyakori minta kinyerésnek feladatában adott egy  $\mathcal{B}$  bemeneti (vagy feldolgozandó) adathalmaz,  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet, egy  $\text{supp}_{\mathcal{B}}: \mathcal{M} \rightarrow \mathbb{N}$  anti-monoton függvény és egy  $\text{min\_supp} \in \mathbb{N}$  küszöbszám. Feladat, hogy megkeressük azon mintákat, amelyekre a  $\text{supp}$  függvény  $\text{min\_supp}$ -nál nagyobb vagy egyenlő értéket ad:*

$$GY = \{gy : gy \in \mathcal{M}, \text{supp}_{\mathcal{B}}(gy) \geq \text{min\_supp}\}.$$

A  $\text{supp}_{\mathcal{B}}$  függvényt támogatottsági függvénynek (support function),  $\text{min\_supp}$ -ot támogatottsági küszöbnek, a  $GY$  elemeit pedig gyakori mintáknak hívjuk. A nem gyakori mintákat ritkáknak nevezzük. Az érthetőség kedvéért a  $\mathcal{B}$  tagot gyakran elhagyjuk, továbbá a  $\text{supp}(m)$ -re mint a minta támogatottsága hivatkozunk. A támogatottsági függvény értéke adja meg, hogy egy minta mennyire gyakori a bemenetben.

Az elemhalmazok példájánál maradva a bemenet lehet például elemhalmazok sorozata. Ekkor egy  $H$  halmaz támogatottságát úgy értelmezhetjük, mint a sorozat azon elemeinek száma, amelyek tartalmazzák  $H$ -t. Például az  $\langle \{A, D\}, \{A, C\}, \{A, B, C, D\}, \{B\}, \{A, D\}, \{A, B, D\}, \{D\} \rangle$  bemenet esetén  $\text{supp}(\{A, D\}) = 4$ . Ha  $\text{min\_supp}$ -nak 4-et adunk meg, akkor  $GY = \{\{A\}, \{D\}, \{A, D\}\}$ .

A támogatottság anti-monotonitásából következik az alábbi egyszerű tulajdonság.

**10.5. tulajdonság.** *Gyakori minta minden részmintája gyakori.*

A mintákat elemhalmazok, sorozatok, gráfok, stb. formájában fogjuk keresni, azaz a minták mindig valamilyen alaphalmazon definiált struktúrák lesznek. Ha az alaphalmazon definiálunk egy teljes rendezést, akkor az alapján – könnyebben vagy nehezebben – a mintákon is tudunk teljes rendezést adni. Ezt például elemhalmazok esetében a lexikografikus rendezés, gráfok esetében a kanonikus címkézés segítségével fogjuk megtenni. A mintákon értelmezett teljes rendezés egyes algoritmusoknál (pl.: APRIORI) a hatékonyság növelésére használható, másoknak pedig alapfeltétele (pl.: Zaki). Sokszor fog felbukkanni a *prefix* fogalma is, amihez szintén egy teljes rendezésre lesz szükség.

„Amerikai kutatás során megállapították, hogy 1 óra tévénézés hatására 200 dollárral többet költünk a reklámok miatt.” Forrás: Sláger rádió, 2007. október 25., 17 óra 48 perc

**10.6. definíció.** Legyen  $\preceq$  a  $H$  halmazon értelmezett részben rendezés. A  $\prec'$  teljes rendezést a  $\prec$  lineáris kiterjesztésének hívjuk, ha minden  $x \prec y$  párra  $x \prec' y$  teljesül.

A lineáris kiterjesztéseknek azon csoportja érdekes számunkra, amelyek *mérettartóak*. Ez azt jelenti, hogy  $|x| < |y|$  esetén a  $x \prec' y$  feltételnek is fenn kell állnia. Amikor tehát a  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet  $\preceq$  tagjának egy mérettartó lineáris kiterjesztését akarjuk megadni, akkor az azonos méretű elemek között definiálunk egy sorrendet. A továbbiakban a mérettartó jelzőt elhagyjuk, és minden lineáris kiterjesztés alatt mérettartó lineáris kiterjesztést értünk.

**10.7. definíció.** Legyen  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet és  $\preceq'$  a  $\preceq$  egy lineáris kiterjesztése. Az  $m$  minta  $\ell$ -elemű részmintái közül az  $\preceq'$  szerinti legelsőt hívjuk az  $m$  minta  $\ell$ -elemű prefixének.

Például, ha  $\mathcal{J} = \{A, B, C, D, E\}$ , és az azonos méretű mintákon az abc rendezés szerinti lexikografikus rendezést vesszük a teljes rendezésnek, akkor például az  $\{A, C, D, E\}$  minta 2-elemű prefixe az  $\{A, C\}$  halmaz.

### 10.1.1. Hatékonysági kérdések

A bemeneti adat és a minták halmaza általában nagy. Például bemeneti sorozatok esetében nem ritkák a  $10^9$  nagyságrendű sorozatok, a mintatér pedig általában  $10^5$  nagyságrendű halmazok hatványhalmaza. Ilyen méretek mellett a naív algoritmusok (például határozzuk meg a mintahalmaz minden elemének támogatottságát, majd válogassuk ki a gyakoriakat) túl sok ideig futnának, vagy túl nagy lenne a memóriaigényük. Hatékony, kifinomult algoritmusokra van szükség, amelyek speciális adatstruktúrákat használnak.

Egy algoritmus hatékonyságát a futási idővel (ami arányos az elemi lépések számával) és a felhasznált memóriával jellemezzük. Például megmondhatjuk, hogy adott méretű bemenet esetén átlagosan, vagy legrosszabb esetben mennyi elemi lépést (összehasonlítás, értékadás), illetve memóriát használ. Sajnos a gyakori mintát kinyerő algoritmusok mindegyike legrosszabb esetben a teljes mintatér megvizsgálja, ugyanis a támogatottsági küszöb függvényében a mintatér minden eleme gyakori lehet.

A gyakori minta-kinyerés korszakának első 10-15 évében az algoritmusok hatékonyságát – elméleti elemzések híján – minden esetben teszteredményekkel igazolták. Szinte minden algoritmusozóhoz lehet találni olyan bemeneti adatot, amit az algoritmus nagyon hatékonyan képes feldolgozni. Ennek eredményeként például, csak a gyakori elemhalmazokat kinyerő algoritmusok

száma meghaladja a 150-et, és a mai napig nem tudunk olyan algoritmusról, amelyik az összes többit legyőzné futási idő vagy memóriefogyasztás tekintetében.

A jövő feladata ennek a káosznak a tisztázása. Ehhez a legfontosabb lépés a bemeneti adat karakterisztikájának formális leírása lenne. Sejtjük, hogy legjobb gyakori mintakinyerő algoritmus nem létezik, de talán van esélyünk értelmes megállapításokra, ha a bemenetre vonatkozóan különböző feltételezésekkel élünk (szokásos feltétel például az, hogy a bemenet olyan sorozat, melynek elemei kis méretű halmazok vagy az, hogy csak nagyon kevés magas támogatottságú minta van) és ezekhez próbáljuk megtalálni az ideális algoritmust.

## 10.2. További feladatok

A gyakori mintakinyerés egyik nagy kritikája, hogy sokszor túl nagy a kinyert minták száma. Vannak olyan feladatok, ahol nem az összes gyakori mintát kívánjuk kinyerni, hanem csak egy részüket. Erre példa az ún. *top-k* mintakinyerés, melynek során a  $k$  legnagyobb támogatottságú mintát keressük. Emellett az alábbi feladatok léteznek.

### 10.2.1. Nem bővíthető és zárt minták

**10.8. definíció.** *Az  $m$  gyakori minta  $\mathcal{B}$ -re nézve nem bővíthető (maximal), ha nem létezik olyan  $m'$  gyakori minta  $\mathcal{B}$ -ben, amelynek  $m$  valódi részmintája.*

**10.9. definíció.** *Az  $m$  minta  $\mathcal{B}$ -re nézve zárt, amennyiben nem létezik olyan  $m'$  minta  $\mathcal{B}$ -ben, amelynek  $m$  valódi részmintája, és  $m'$  támogatottsága megegyezik  $m$  támogatottságával ( $\text{supp}(m') = \text{supp}(m)$ ).*

Az ember azonnal láthatja, hogy mi értelme van annak, hogy csak a nem bővíthető mintákat keressük meg: egyértelműen meghatározzák a gyakori mintákat és számuk kevesebb. Sajnos a nem bővíthető minták alapján csak azt tudjuk megmondani, hogy egy minta gyakori-e, a támogatottságot nem tudjuk megadni (legfeljebb egy alsó korlátot).

Nem ilyen triviális, hogy mi értelme van a gyakori zárt mintáknak. Azt látjuk, hogy a zárt gyakori minták a gyakori minták részhalmazai, és a zárt minták részhalmaza a nem bővíthető minták, hiszen

**10.10. tulajdonság.** *Minden nem bővíthető minta zárt.*

Mégis mi célt szolgálnak a gyakori zárt minták? Ennek tisztázásához két új fogalmat kell bevezetnünk.

**10.11. definíció.** *Az  $m'$  minta az  $m$  minta lezártja, ha  $m \preceq m'$ ,  $\text{supp}(m) = \text{supp}(m')$  és nincs  $m'' : m' \prec m''$ , melyre  $\text{supp}(m') = \text{supp}(m'')$ .*

Nyilvánvaló, ha  $m$  zárt, akkor lezártja megegyezik önmagával.

**10.12. definíció.** *Az  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet a zártságra nézve egyértelmű, amennyiben minden  $m \in \mathcal{M}$  minta lezártja egyértelmű.*



Látni fogjuk, hogy sorozat típusú bemenet esetén például az elemhalmazokat tartalmazó mintakörnyezet zártságra nézve egyértelmű, míg a sorozatokat tartalmazó nem az. A zártságra nézve egyértelmű mintakörnyezetekben a zárt minták jelentősége abban áll, hogy ezek ismeretében tetszőleges mintáról el tudjuk dönteni, hogy gyakori-e, és ha igen, meg tudjuk pontosan mondani támogatottságát. Szükségtelen tárolni az összes gyakori mintát, hiszen a zárt mintákból ezek egyértelműen meghatározhatók. Az  $m$  minta gyakori, ha része valamely gyakori zárt mintának, és  $m$  támogatottsága megegyezik a legkisebb olyan zárt minta támogatottságával, amelynek része  $m$  (ez ugyanis az  $m$  lezártja).

### 10.2.2. Kényszerek kezelése

Nem mindig érdekes az összes gyakori minta. Előfordulhat, hogy például a nagy méretű, vagy bizonyos mintákat tartalmazó, vagy nem tartalmazó, stb. gyakori minták nem fontosak. Általánosíthatjuk a feladatot úgy, hogy a felhasználó kényszereket, predikátumokat ad meg, és azokat a mintákat kell meghatároznunk, amelyek kielégítik az összes kényszert.

A feladat egyszerű megoldása lenne, hogy – mint utófeldolgozás – a gyakori mintákat egyesével megvizsgálva törölnénk azokat, amelyek nem elégítenek minden kényszert. Ez a megoldás nem túl hatékony. Jobb lenne, ha a kényszereket minél „mélyebbre” tudnánk helyezni a gyakori mintákat kinyerő algoritmusokban. Ez bizonyos kényszereknél megtehető, másoknál nem. Nézzük, milyen osztályokba sorolhatjuk a kényszereket.

Tulajdonképpen az is egy kényszer, hogy gyakori mintákat keresünk. A gyakoriságra vonatkozó predikátum igaz, ha a minta gyakori, ellenkező esetben hamis. Ez a predikátum anti-monoton:

**10.13. definíció.** Legyen  $(H, \preceq)$  egy részben rendezett halmaz. A  $p : H \rightarrow \{\text{igaz}, \text{hamis}\}$  predikátum anti-monoton, amennyiben tetszőleges  $x \in H$  elem esetén, ha  $p(x) = \text{igaz}$ , akkor  $p(y)$  is igazat ad minden  $y \preceq x$  elemre.

Ha a fenti definícióba  $y \preceq x$  helyett  $x \preceq y$  írunk, akkor a *monoton* predikátumok definícióját kapjuk. Egy predikátum akkor és csak akkor monoton és anti-monoton egyben, ha a mintatér minden eleméhez igaz (vagy hamis) értéket rendel. Az ilyen predikátumot *triviális predikátumnak* hívjuk.

**10.14. definíció.** Legyen  $(H, \preceq)$  egy részben rendezett halmaz. A  $p : H \rightarrow \{\text{igaz}, \text{hamis}\}$  predikátum prefix anti-monoton, amennyiben megadható a  $\prec$ -nek egy olyan  $\preceq'$  lineáris kiterjesztése amire, ha  $p(m) = \text{igaz}$ , akkor  $p$  az  $m$  minden prefixén is igaz.

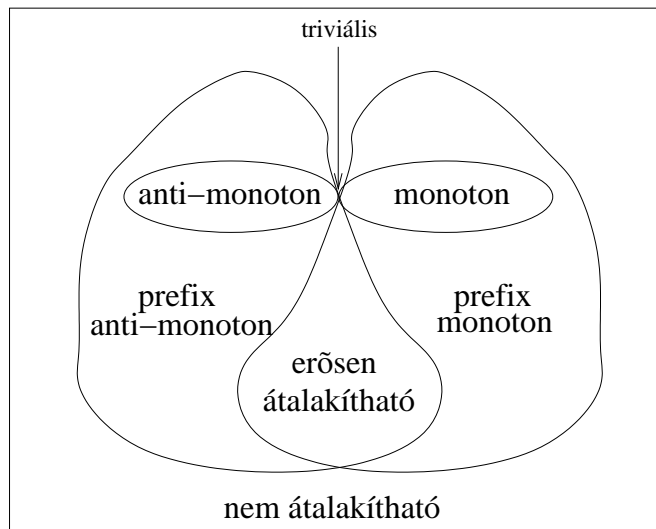
**10.15. definíció.** Legyen  $(H, \preceq)$  egy részben rendezett halmaz. A  $p : H \rightarrow \{\text{igaz}, \text{hamis}\}$  predikátum prefix monoton, amennyiben megadható a  $\prec$ -nek egy olyan  $\preceq'$  lineáris kiterjesztése amely, ha  $p(m) = \text{igaz}$ , és az  $m'$  mintának  $m$  prefixe. akkor  $p(m')$  is igaz.

Minden anti-monoton (monoton) predikátum egyben prefix anti-monoton (prefix monoton) is.

**10.16. definíció.** A  $p$  predikátum erősen átalakítható, amennyiben egyszerre prefix anti-monoton és prefix monoton.

A 10.1 ábrán látható a kényszerek kapcsolata [107].

Sejthetjük, hogy az anti-monoton predikátumok lesznek a legegyszerűbben kezelhetők. Ilyen anti-monoton predikátumok például a következők:



10.1. ábra. A kényszerek (predikátumok) osztályozása

- A minta mérete ne legyen nagyobb egy adott küszöbnél.
- A mintának legyen része egy rögzített minta.

Vásárlói szokások vizsgálatánál – amikor a vásárlói kosarakban gyakran előforduló termékhalmozokat keressük – monoton kényszer például az, hogy a termékhalmozban lévő elemek profitjának összértéke (vagy minimuma, maximuma) legyen nagyobb egy adott konstansnál.

Prefix monoton predikátum például, hogy a termékhalmozban található termékek árának átlaga nagyobb-e egy rögzített konstansnál. Rendezzük a termékeket áruk szerint növekvő sorrendbe. Ezen rendezés szerinti lexikografikus rendezés legyen a teljes rendezés. Nyilvánvaló, hogy ekkor a prefixben található termékek árai nagyobbak, mint a prefixben nem szereplő termékek árai. Ez a kényszer prefix monoton, hiszen a prefix a legolcsóbb termékeket tartalmazza, így átlaga nem lehet kisebb. Érdeemes átgondolni, hogy ez a predikátum ráadásul erősen átalakítható.

### 10.2.3. Többszörös támogatottsági küszöb

Vannak olyan alkalmazások, amelyekben a gyakoriság egyetlen, univerzális támogatottsági küszöb alapján történő definiálása nem megfelelő. Ha például vásárlási szokások elemzésére gondolunk, akkor a nagy értékű termékekkel kapcsolatos tudás legalább annyira fontos, mint a nagy mennyiségben értékesített, de kis haszonnal járó termékekkel kapcsolatos információ. Kézenfekvő megoldás, hogy annyira lecsökkentjük a támogatottsági küszöböt, hogy ezek a ritka elemek is gyakoriak legyenek, ami azzal a veszéllyel jár, hogy (ezen fontos elemek mellett) a mintatér nagy része gyakorivá válik. *Többszörös támogatottsági küszöbnél* a mintatér minden eleméhez egyedileg megadhatunk egy támogatottsági küszöböt, azaz létezik egy  $min\_supp : \mathcal{M} \rightarrow \mathbb{N}$  függvény, és az  $m$  akkor gyakori, ha  $supp(m) \geq min\_supp(m)$ .

Többszörös támogatottsági küszöb esetén nem igaz a 10.5 tulajdonság. Hiába nagyobb ugyanis egy rész minta támogatottsága, a rész mintához tartozó támogatottsági küszöb még nagyobb lehet, és így a rész minta nem feltétlenül gyakori.

### 10.2.4. Dinamikus gyakori mintakinyerés

Egyre népszerűbb adatbányászati feladat a gyakori minták ún. dinamikus kinyerése. Adott egy kiindulási  $\mathcal{B}$  bemenet a hozzá tartozó gyakori mintákkal és támogatottságokkal és egy másik  $\mathcal{B}'$  bemenet. Általában a  $\mathcal{B}'$ -t valami apró módosítással kapjuk  $\mathcal{B}$ -ből. Feladat, hogy minél hatékonyabban találjuk meg a  $\mathcal{B}'$ -ben gyakori mintákat, azaz minél jobban használjuk fel a meglévő tudást (a  $\mathcal{B}$ -ben gyakori mintákat). Gondolhatunk itt egy on-line áruházra, ahol kezdetben rendelkezésünkre állnak az elmúlt havi vásárlásokhoz tartozó gyakori termékhalmozok, miközben folyamatosan érkeznek az új vásárlások adatai. Hasznos, ha az újonnan felbukkanó gyakori mintákat minél hamarabb felfedezzük, anélkül, hogy a bővített adatbázisban off-line módon lefuttatnánk egy gyakori mintákat kinyerő algoritmust.

## 10.3. Az algoritmusok jellemzői

*Helyes* vagy *helyesen működő* jelzővel illetjük azokat az algoritmusokat, amelyek nem hibáznak, tehát csak gyakori mintákat nyernek ki és azok támogatottságát jól határozzák meg. *Teljes* egy algoritmus, ha be lehet bizonyítani, hogy az összes gyakori mintát és támogatottságait meghatározza. Helyesen működő és teljes algoritmusokról fogunk beszélni, de szó lesz olyan algoritmusokról is, amelyekről csak azt tudjuk, hogy (bizonyos feltételezésekkel élve) kicsi annak a valószínűsége, hogy nem talál meg minden gyakori mintát.

*Szélességi bejárást* valósítanak meg azok az algoritmusok<sup>1</sup>, amelyek a legkisebb mintákból kiindulva egyre nagyobb méretű gyakori mintákat nyernek ki. Egy ilyen algoritmusra igaz, hogy az  $\ell$ -elemű gyakori mintákat hamarabb találja meg, mint az  $\ell$ -nél nagyobb elemű mintákat. *Mélységi bejárást* megvalósító algoritmusokra ez nem igaz; ezek minél gyorsabban próbálnak eljutni a nem bővíthető mintához. Ha ez sikerül, akkor egy újabb, nem bővíthető mintát vesznek célba.

A következőkben ismertetjük a három legfontosabb gyakori mintákat kinyerő módszert az APRIORI-t, Zaki módszerét és a mintanövelő módszert. Ennek a három algoritmusnak a szerepe abban áll, hogy szinte az összes többi algoritmus ezeknek a továbbfejlesztése, vagy ezen algoritmusok keveréke. Jelentőségüket tovább növeli az a tény, hogy ezek a módszerek alkalmazhatóak akármilyen típusú mintákat keresünk, legyenek azok elemhalmozok, sorozatok vagy gráfok. Nem pontos algoritmusokat adunk, hanem csak egy általános módszerleírást. Egyes lépéseket csak a minta típusának ismeretében lehet pontosan megadni.

„Fogszuvasodás ellen zöld tea  
Távol-keleti felmérések szerint  
azoknak az iskolásgyerekeknek,  
akik napi egy csésze cukor nélküli  
zöld teát isznak, feleannyi od-  
vas foguk van, mint az átlagnak.”  
Forrás: <http://www.terebess.hu/szorolapok/zoldtea.html>

## 10.4. Az APRIORI módszer

Az eredeti APRIORI algoritmust gyakori elemhalmozok kinyerésére használták, és mint az AIS algoritmus [3] továbbfejlesztett változata adták közre. Rakesh Agrawal és Ramakrishnan Srikant [5] publikálták 1993-ban, de tőlük függetlenül, szinte ugyanezt az algoritmust javasolta

<sup>1</sup>A szélességi bejárást megvalósító algoritmusokat szintenként haladó (levelwise) algoritmusoknak is hívják.

Heikki Mannila, Hannu Toivonen és A. Inkeri Verkamo [89]-ben. Az 5 szerző végül egyesítette a két írást [4]. Kis módosítással az algoritmust gyakori sorozatok kinyerésére is (APRIORIALL, GSP algoritmusok), sőt, alapelvét bármely típusú gyakori minta (epizód, fa stb.) keresésénél is alkalmazhatjuk.

Az algoritmus rendkívül egyszerű, mégis gyors és kicsi a memóriaigénye. Talán emiatt a mai napig ez az algoritmus a legelterjedtebb és legismertebb gyakori mintakinyerő algoritmus.

Az APRIORI szélességi bejárást valósít meg. Ez azt jelenti, hogy a legkisebb mintából kiindulva szintenként halad előre a nagyobb méretű gyakori minták meghatározásához. A következő szinten (iterációban) az eggyel nagyobb méretű mintákkal foglalkozik.

Az algoritmusban központi szerepet töltenek be az ún. *jelöltek*. Jelöltnek hívjuk egy adott iterációban azt a mintát, amelynek támogatottságát meghatározzuk, azaz, aminek figyelmünket szenteljük. *Hamis jelölteknek* hívjuk azokat a jelölteket, amelyekről ki fog derülni, hogy ritka minták, *elhanyagolt minták* pedig azok a gyakori minták, amelyeket nem választunk jelöltnek (nem foglalkozunk velük ☺). Nyilvánvaló, hogy csak azokról a mintákról tudjuk eldönteni, hogy gyakoriak-e, amelyeknek meghatározzuk a támogatottságát, tehát amelyek jelöltek valamikor. Ezért elvárjuk az algoritmustól, hogy minden gyakori mintát felvegyen jelöltnek. A teljesség feltétele, hogy ne legyen elhanyagolt minta, a hatékonyság pedig annál jobb, minél kevesebb a hamis jelölt.

A jelöltek definiálásánál a 10.5 tulajdonságot használjuk fel, ami így szól: „Gyakori minta minden részmintája gyakori.”. Az állítást indirekten nézve elmondhatjuk, hogy egy minta biztosan nem gyakori, ha van ritka részmintája! Ennek alapján ne legyen jelölt az a minta, amelynek van ritka részmintája. Az APRIORI algoritmus ezért építkezik letről. Egy adott iterációban pontosan tudjuk, hogy a részminták gyakoriak vagy sem! Az algoritmus onnan kapta a nevét, hogy az  $\ell$ -elemű jelölteket a bemeneti adat  $\ell$ -edik átolvasásának megkezdése előtt (a priori) állítja elő.

Az algoritmus pszeudokódja a következő ábrán látható. Kezdeti értékek beállítása után belépünk egy ciklusba. A ciklus akkor ér véget, ha az  $\ell$ -elemű jelöltek halmaza üres. A cikluson belül először a `támogatottság_meghatározás` eljárást hívjuk meg, amely a jelöltek támogatottságát határozza meg. Ha ismerjük a jelöltek támogatottságát, akkor ki tudjuk választani belőlük a gyakoriakat. A `jelölt_előállítás` függvény az  $\ell$ -elemű gyakori mintákból  $(\ell + 1)$ -elemű jelölteket állít elő.

Az APRIORI elvet adaptáló algoritmusok mind a fenti lépéseket követik. Természetesen a különböző típusú mintáknál különböző módon kell elvégezni a támogatottság-meghatározás, gyakoriak kiválogatása, jelöltek előállítás lépéseket.

Az algoritmus hatékonyságának egyik alapfeltétele, hogy a jelöltek elférjenek a memóriában. Ellenkező esetben ugyanis rengeteg idő menne el az olyan I/O műveletekkel, amelynek során a jelölteket a háttér és a memória között ide-oda másolgatjuk. A fenti pszeudokód az eredeti APRIORI egyszerűsített változatát írja le. Valójában ugyanis addig állítjuk elő az  $\ell$ -elemű jelölteket, amíg azok elférnek a memóriában. Ha elfogy a memória, akkor  $\ell$  növelése nélkül folytatjuk az algoritmust, majd a következő iterációban ott folytatjuk a jelöltek előállítását, ahol abbahagytuk.

### 10.4.1. Jelöltek előállítása

Az  $\ell$ -elemű jelöltek előállításának egyszerű módja az, hogy vesszük az összes  $\ell$ -elemű mintát, és azokat választjuk jelöltnek, amelyekre teljesül, hogy minden részmintájuk gyako-

**Algorithm 9** Az APRIORI módszer**Require:**  $B$ : bementei adat $min\_supp$ : támogatottsági küszöb $\ell \leftarrow 0$  $J_\ell \leftarrow \{ \text{az üres minta} \} \cup \{ J_\ell : \text{Az } \ell\text{-elemű jelöltek} \}$ **while**  $|J_\ell| \neq 0$  **do**    támogatottság\_meghatározás( $B, J_\ell$ )    **for all**  $j \in J_\ell$  **do**        **if**  $supp(j) \geq min\_supp$  **then**             $GY_\ell \leftarrow GY_\ell \cup \{j\}$         **end if**    **end for**     $J_{\ell+1} \leftarrow \text{jelölt-előállítás}(GY_\ell)$      $GY \leftarrow GY \cup GY_\ell$      $\ell \leftarrow \ell + 1$ **end while****return**  $GY$ : gyakori minták

ri. Szükségtelen az összes részmintát ellenőrizni, ugyanis a támogatottság anti-monotonitásából következik az, hogy ha az összes  $(\ell-1)$ -elemű rész minta gyakori, akkor az összes valódi rész minta is gyakori.

Ez a módszer azonban nem túl hatékony, vagy úgy is megfogalmazhatnánk, hogy túl sok felesleges munkát végez, túl sok olyan mintát vizsgál meg, amelyek biztosan nem gyakoriak. Hívjuk *potenciális jelölteknek* azon mintákat, amelyeket előállítunk, majd ellenőrizzük, hogy rész mintáik gyakoriak-e. Ha egy potenciális minta átesik a teszten, akkor jelölté válik.

Tudjuk, hogy ha egy minta jelölt lesz, akkor minden  $(\ell-1)$ -elemű rész mintája gyakori, tehát célszerű az  $(\ell-1)$ -elemű gyakori mintákból kiindulni. Egy egyszerű megoldás lenne, ha sorra vennénk az  $(\ell-1)$ -elemű gyakori minták minimális valódi felső korlátait, mint potenciális jelöltek. Még jobb megoldás, ha a  $(\ell-1)$ -elemű gyakori mintapárokra vesszük a minimális valódi felső korlátait. Ekkor ugyanis csak olyan potenciális jelöltet állítunk elő, amelynek van két  $(\ell-1)$ -elemű gyakori rész mintája. A minimális valódi felső korlátot egy *illesztési művelettel* fogjuk előállítani. A két gyakori mintát a potenciális jelölt generátorainak hívjuk. Az illesztési műveletet a  $\otimes$ -el fogunk jelölni. Akkor illesztünk két mintát, ha van  $(\ell-2)$ -elemű közös rész mintájuk. Ezt a rész mintát *magnak* (core) fogjuk hívni.

Ha az előállítás módja olyan, hogy nem állíthatjuk elő ugyanazt a potenciális jelöltet két különböző módon, akkor ezt a jelölt-előállítást *ismétlés nélkülinek* nevezzük. Nézzünk egy példát. Legyenek a mintatér elemei elemhalmazok. Akkor állítsuk elő két  $(\ell-1)$ -elemű gyakori elemhalmaznak a minimális valódi korlátját, ha metszetük  $(\ell-2)$ -elemű. A minimális valódi korlátok halmaza csak egy elemet fog tartalmazni, a két halmaz unióját. Ez a jelölt-előállítás nem ismétlés nélküli, ugyanis például az  $(\{A, B\}, \{A, C\})$  párnak ugyanaz a legkisebb felső korlátja, mint az  $(\{A, B\}, \{B, C\})$  párnak.

Az ismétlés nélküli jelölt-előállítást mindig a minta elemein értelmezett teljes rendezés fogja garantálni, ami a  $\preceq$  rendezés egy lineáris kiterjesztése lesz. A teljes rendezésnek megfelelően végigmegyünk az  $(\ell-1)$ -elemű gyakori mintákon és megnézzük, hogy mely sorban utána követ-

kező  $(\ell - 1)$ -elemű gyakori mintával illeszthető, illetve az illesztésként kapott potenciális jelölt minden  $(\ell - 1)$ -elemű részmintája gyakori-e. Sok esetben a ismétlés nélküliségnek elégséges feltétele az lesz, hogy a két gyakori minta  $(\ell - 2)$ -elemű prefixeik megegyezzenek. A minta típusának ismeretében a teljességet (minden minimális valódi felső korlátbeli elemet előállítunk) és az ismétlés nélküliséget könnyű lesz bizonyítani.

---

**Algorithm 10** Jelöltek előállítása
 

---

**Require:**  $GY_{\ell-1}$ :  $(\ell - 1)$ -elemű gyakori minták

```

for all  $gy \in GY_{\ell-1}$  do
  for all  $gy' \in GY_{\ell-1}, gy \preceq gy'$  do
    if  $gy$  és  $gy'$  illeszthető then
       $\hat{J} \leftarrow \text{minimális\_valódi\_felső\_korlát}(gy, gy')$   $\{\hat{J}: \text{potenciális jelöltek halmaza}\}$ 
      for all  $\hat{j} \in \hat{J}$  do
        if minden_részalmaz_gyakori( $\hat{j}, GY_{\ell-1}$ ) then
           $J_\ell \leftarrow \hat{j}$ 
        end if
      end for
    end if
  end for
end for
return  $J_\ell$ :  $\ell$ -elemű jelöltek
  
```

---

### 10.4.2. Zárt minták kinyerése, az APRIORI-CLOSE algoritmus

A zárt minták jelentőségét a 10.2.1 részben már tárgyaltuk. Itt most két feladat megoldásával foglalkozunk. Megnézzük, hogy az összes gyakori mintából hogyan tudjuk előállítani a zártakat, illetve bemutatjuk az APRIORI-CLOSE [104–106] algoritmust, amely már eleve csak a zárt mintákat határozza meg. Mindkét módszerhez az alábbi észrevételt használjuk fel:

**10.17. észrevétel.** *Ha az  $m$  minta nem zárt, akkor van olyan  $m$ -et tartalmazó eggyel nagyobb méretű minta, amelynek támogatottsága megegyezik  $m$  támogatottságával.*

Tegyük fel, hogy a legnagyobb méretű gyakori minta mérete  $k$ . A  $GY_k$  elemei zártak. Egy egyszerű algoritmus menete a következő:

Nézzük sorban  $GY_{k-1}, GY_{k-2}, \dots, GY_0$  elemeit. Ha  $m \in GY_\ell$ -hez találunk olyan  $m' \in GY_{\ell+1}$  elemet, amelynek támogatottsága megegyezik  $m$  támogatottságával, akkor  $m$  nem zárt. Ha nincs ilyen tulajdonságú  $m'$ , akkor  $m$  zárt.

Az APRIORI-CLOSE menete teljes mértékben megegyezik az APRIORI algoritmus menetével. Az egyetlen különbség, hogy az  $\ell$ -elemű gyakori minták meghatározása után törli az  $(\ell - 1)$ -elemű nem zártakat. Miután eldöntötte, hogy az  $\ell$ -elemű  $m$  minta gyakori, megvizsgálja az összes  $(\ell - 1)$ -elemű részmintáját  $m$ -nek. Amennyiben van olyan részalmaz, aminek támogatottsága egyenlő  $m$  támogatottságával, akkor ez a rész minta nem zárt, ellenkező esetben zárt.

## 10.5. Sorozat típusú bemenet

A legáltalánosabb eset leírásánál nem tettünk semmi megkötést a bemenet típusára és a támogatottsági függvényre vonatkozóan. Az esetek többsége azonban egy speciális családba tartozik. Ennek a problémacsaládnak a jellemzője, hogy a bemenet egy véges sorozat, és a támogatottságot azon elemek száma adja, amelyek valamilyen módon *illeszkednek* a mintára<sup>2</sup>. Az illeszkedést egy illeszkedési predikátummal adhatjuk meg, melynek értelmezési tartománya a mintatér.

$$\text{bemenet} : \mathcal{S} = \langle s_1, s_2, \dots, s_n \rangle$$

A támogatottság definíciója megköveteli, hogy ha egy minta illeszkedik egy sorozatelemre, akkor minden részmintája is illeszkedjen. A legtöbb esetben a sorozat elemei megegyeznek a mintatér elemeivel és az  $m$  minta akkor illeszkedik egy sorozatelemre, ha annak  $m$  a részmintája.

A szakirodalomban igen elterjedt a sorozatok helyett a halmazokkal leírt bemenet, ahol minden egyes elem egyedi azonosítóval van ellátva. A jegyzetben a sorozatos leírást fogjuk használni, akinek ez szokatlan, az tekintse azonosítóknak a sorozat elemeinek sorszámát.

Az  $m$  minta *gyakoriságát* (jelölésben:  $\text{freq}_{\mathcal{S}}(m)$ , ami a frequency szóra utal) az  $m$  támogatottsága és az  $\mathcal{S}$  hosszának hányadosával definiáljuk. A gyakorisági küszöböt ( $\frac{\text{min-suppl}}{|\mathcal{S}|}$ ) következetesen  $\text{min\_freq}$ -el jelöljük. Az értelmesen megválasztott gyakorisági küszöb mindig 0 és 1 között van. Az esetek többségében támogatottsági küszöb helyett gyakorisági küszöböt adnak meg.

Sorozat típusú bemenet esetén merül fel azon elvárás az algoritmusokkal szemben, hogy ne legyen érzékeny a bemenet *homogenitására*. Intuitíve akkor homogén egy bemenet, ha nincsenek olyan részei, amelyben valamely minta gyakorisága nagyon eltér a teljes bemenet alapján számított gyakoriságától. Sok alkalmazásban ez a feltétel nem áll fenn, így azokat az algoritmusokat kedveljük, amelyek hatékonysága független a bemenet homogenitásától. Könnyű átgondolni, hogy az APRIORI algoritmus rendelkezik ezzel a tulajdonsággal.

### 10.5.1. Apriori

Amennyiben a támogatottságot illeszkedési predikátum alapján definiáljuk, akkor az APRIORI algoritmus a bemeneti elemeken egyesével végigmegy és növeli azon jelöltek számlálóját, amelyek illeszkednek az éppen aktuális bemeneti elemre. Azonos bemeneti elemeknél ez a művelet ugyanazt fogja csinálni ezért célszerű az azonos bemeneti elemeket összegyűjteni és csak egyszer meghívni az eljárást. A bemenet azonban túl nagy lehet, ezért ezt gyorsítási lépést csak akkor szokás elvégezni, amikor már rendelkezésünkre állnak az egyelemű gyakori minták. Ezek alapján további szűréseket lehet végezni. Például elemhalmaz/elemsorozat típusú bemeneti elemeknél töröljük a halmazból/sorozatból a ritka elemeket. Ez duplán hasznos, hiszen csökkentjük a memóriafogyasztást és mivel az azonos szűrt elemek száma nagyobb lehet, mint az azonos elemek száma a támogatottságok meghatározása még kevesebb időbe fog telni.

Vannak olyan minták, amelyeknél a illeszkedés eldöntése drága művelet. Például gráf típusú mintáknál az illeszkedés meghatározásához egy részgráf izomorfia feladatot kell eldönteni, ami bizonyítottan NP-teljes. Ilyen mintáknál hasznos, ha minden jelöltnél rendelkezésünkre állnak

<sup>2</sup>Ha csak a matematikai definíciókat tekintjük, akkor törekedhetünk volna a legegyszerűbb leírásra és használhattunk volna sorozatok helyett multihalmazokat. A valóságban azonban a bemenet tényleg sorozatok formájában adott, így nem tehetjük fel, hogy az azonos bemeneti elemek össze vannak vonva.

azon bemeneti elemek sorszámait, amelyekre illeszkednek a generátorok (nevezzük ezt a halmazt illeszkedési halmaznak). Az illeszkedési predikátum anti-monoton tulajdonságából következik, hogy a jelölt csak azon bemeneti elemekre illeszkedhet, amelyekre generátoraink is illeszkednek. A támogatottság meghatározása során a jelöltek illeszkedési halmazát is meg kell határoznunk hiszen a jelöltek lesznek a generátorok a következő iterációban. Természetesen a generátorok illeszkedési listáit törölhetjük miután meghatároztuk a jelöltek illeszkedési listáit.

## DIC

A DIC (Dynamic Itemset Counting) algoritmus [22] az APRIORI továbbfejlesztése. Gyakori elemhalmazok kinyerésére javasolták, de minden olyan gyakori mintákat kereső feladatban alkalmazható, amelyben a bemenet sorozat típusú, és a támogatottságot illeszkedési predikátum alapján definiáljuk. Az algoritmus nem tisztán szélességi bejárást valósít meg; a különböző elemszámú minták együtt vannak jelen a jelöltek között. Ha  $k$  a legnagyobb gyakori minta mérete, akkor várhatóan  $(k+)$ -nél kevesebbszer, de legrosszabb esetben  $(k+1)$ -szer kell végigolvasni a bemenetet.

A DIC algoritmusban – szemben az APRIORI-val – nem válik szét az egyes iterációkban a jelöltek előállítás, a támogatottságok meghatározása és a ritka minták törlése. Miközben vesszük a bemeneti elemeket és határozzuk meg a jelöltek támogatottságát, új jelölteket vehetünk fel és törölhetünk (azaz dinamikus elemszámlálást alkalmazunk, ahogyan erre az algoritmus neve is utal). Akkor veszünk fel egy mintát a jelöltek közé, ha minden valódi részmintájáról kiderült, hogy gyakori. Akárhon veszünk fel egy jelöltet, egy iterációval később ugyanott, ahol felvettük, törölnünk kell a jelöltek közül, hiszen a pontos támogatottság meghatározásához a teljes bemenetet át kell néznünk. Ha a törölt jelölt gyakori, akkor természetesen a mintát felvesszük a gyakori minták halmazába. Minden jelöl esetén tárolnunk kell, hogy hányadik bemeneti elemnél lett jelölt. A kiindulási állapotban minden egyelemű minta jelölt és akkor ér véget az algoritmus, amikor nincs egyetlen jelölt sem.

Elemhalmazok példáját nézve, ha az  $A$  és  $B$  elemek olyan sokszor fordulnak elő, hogy támogatottságuk, már a bemenet egyharmadának átolvasása után eléri  $min\_supp$ -ot, akkor az  $\{A, B\}$  két elemű halmaz már ekkor jelölt lesz, és előfordulásait el kell kezdeni összeszámolni. A bemenet végigolvasása után ismét az első bemeneti elemre lépünk és az egyelemű jelöltek törlése után folytatjuk a jelöltek támogatottságának meghatározását. Az  $A, B$  jelöltet a bemenet egyharmadánál töröljük a jelöltek közül. Ha nincs más jelölt, akkor az algoritmus véget ér. Látható, hogy ekkor a DIC algoritmus  $1+1/3$ -szor olvassa végig a bemenetet, amit az APRIORI kétszer tesz meg.

Az algoritmus hátránya, hogy minden bemeneti elemnél meg kell vizsgálni, hogy vannak-e törlendő jelöltek. Ez költséges művelet ezért célszerű „állomásokat” létrehozni. Például minden ezredik bemeneti elem lehet egy állomás. Csak az állomásoknál nézzük meg, hogy egy jelölt támogatottsága elérte-e  $min\_supp$ -ot, így csak állomásnál veszünk fel, illetve törölünk jelölteket.

A DIC algoritmus, szemben az APRIORI-val, érzékeny az adatok homogenitására. Amennyiben egy minta a felszállóhelyétől nagyon távol koncentrálva gyakori, akkor az összes, őt részmintaként tartalmazó minta is csak sokára lesz jelölt. Ekkor a DIC hatékonysága rosszabb az APRIORI algoritmusénál, hiszen ugyanannyiszor járja végig a bemenetet, mint az APRIORI, de eközben olyan munkát is végez, amit az APRIORI nem (minden állomásnál ellenőrzi, hogy mely jelölteket kell törölni). Összességében elmondhatjuk, hogy a DIC csak abban az esetben lesz gyorsabb az APRIORI-nál, ha a bemenet olyan nagy, hogy a futási időben nagy szerepet



játszik a bemenet beolvasása. A mai memóriakapacitások mellett ez ritkán áll fenn.

A következőkben olyan algoritmusokat ismertetünk, amelyek sorozat típusú bemenet és illeszkedés alapú támogatottság esetén tudják meghatározni a gyakori mintákat.

### 10.5.2. Zaki módszere

Zaki módszere [151] szintén jelölteket használ a keresési tér bejárásához, de a bejárás típusa – szemben az APRIORI-val – mélységi. A  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet esetén csak akkor használható, ha tudjuk definiálni a  $\preceq$ -nek egy lineáris kiterjesztését, ugyanis az algoritmus építőelemei a prefixek.

A prefix alapján definiálhatunk egy ekvivalencia relációt. Adott  $\ell$  esetén két minta ekvivalens, ha  $\ell$ -elemű prefixük megegyezik. A  $P$  prefixű minták halmazát  $[P]$ -vel jelöljük. A prefixek segítségével a minták halmazát diszjunkt részekre osztjuk, azaz a feladatot kisebb részfeladatokra vezetjük vissza.

Nézzük például az elemhalmazok esetét. Legyen  $\mathcal{J} = \{A, B, C, D\}$  és  $\mathcal{M} = (2^{\mathcal{J}}, \subseteq)$ , akkor  $I' \prec I''$ , ha  $|I'| < |I''|$  vagy, ha  $|I'| = |I''|$  és  $I'$  lexikografikusan megelőzi  $I''$ -t. Például  $\{D\} \prec \{A, C\}$  és  $\{A, B, D\} \prec \{B, C, D\}$ . Amennyiben  $\ell = 1$ , akkor például a  $\{A, B\}$ ,  $\{A, C\}$ ,  $\{A, D\}$  egy ekvivalencia osztályba tartozik, aminek például a  $\{B, C\}$  nem eleme.

A prefix mellett Zaki módszerének központi fogalma az *illeszkedési lista*. Egy mintához tartozó illeszkedési lista tárolja a minta illeszkedéseit. Az illeszkedési lista két fontos tulajdonsággal bír:

- I. Az illeszkedési listából könnyen megkapható a támogatottság.
- II. Egy jelölt illeszkedési listája megkapható a generátorainak illeszkedési listáiból.

Például elemhalmaz típusú minták esetében (ha az illeszkedést a tartalmazási reláció alapján definiáljuk) egy elemhalmaz illeszkedési listája egy olyan lista lesz, amely a bemeneti sorozat azon elemeinek sorszámát tárolja, amelyeknek része az adott elemhalmaz. Például  $\langle \{A, D\}, \{A, C\}, \{A, B, C, D\}, \{B\}, \{A, D\}, \{A, B, D\}, \{D\} \rangle$  bemenet esetén az  $\{A, C\}$  illeszkedési listája:  $\langle \{1, 2\} \rangle$ .

Zaki algoritmusának pszeudokódja az alábbi.

---

#### Algorithm 11 Zaki módszere

---

**Require:**  $B$ : bementei adat

$min\_supp$ : támogatottsági küszöb

$J \Leftarrow 1$  elemű minták halmaza  $\{J: \text{jelöltek}\}$

$ILL(J) \Leftarrow \text{ILL\_felépítés}(B, J)$   $\{ILL(J): \text{jelöltek illeszkedési listája}\}$

**for all**  $j \in J$  **do**

**if**  $|ILL(j)| \geq min\_supp$  **then**

$GY_1 \Leftarrow GY_1 \cup j$

**end if**

**end for**

$zaki\_segéd(GY, ILL(GY), min\_supp)$   $\{GY = [0]^+\}$

**return**  $GY$ : gyakori minták

---

Először felépítjük az egyelemű minták illeszkedési listáit. Ezek alapján meghatározzuk a gyakori mintákat. A későbbiekben nem használjuk a bemenetet csak az illeszkedési listákat, ezekből ugyanis a támogatottságok egyértelműen meghatározhatók. Az algoritmus lényege a `zaki_segéd` rekurziós eljárás, amelynek pszeudokódja a ?? ábrán látható.

---

**Algorithm 12** `zaki_segéd` eljárás
 

---

**Require:**  $[P]^+$ :  $P$  prefixű,  $P$ -nél eggyel nagyobb gyakori minták

$ILL[P]^+$ :  $[P]^+$ -beli minták illeszkedési listája

$min\_supp$ : támogatottsági küszöb

**for all**  $m \in [P]^+$  **do**

**for all**  $m' \in [P]^+, m \preceq m'$  **do**

$J, ILL(J) \leftarrow \text{minimális\_valódi\_felső\_korlát}(m, m', ILL(m, m'))$

$\{J: \text{jelöltek}, ILL(J): \text{jelöltek illeszkedési listája}\}$

**for all**  $j \in J$  **do**

**if**  $|ILL(j)| \geq min\_supp$  **then**

$GY' \leftarrow GY' \cup \{j\}$

**end if**

**end for**

**end for**

$\text{zaki\_segéd}(GY', ILL(GY'), min\_supp)$

$GY \leftarrow GY \cup GY'$

**end for**

**return**  $GY$ :  $P$  prefixű összes gyakori minta

---

A Zaki féle jelölt előállításnak két feladata van. Természetesen az egyik a jelöltek előállítása, de emellett az illeszkedési listákat is előállítja. A jelölt-előállítás megegyezik az APRIORI jelölt előállításának első lépésével (potenciális jelöltek előállítása). A második lépést nem is tudnánk elvégezni, ugyanis nem áll rendelkezésünkre az összes rész minta, így nem is tudjuk ellenőrizni, hogy az összes rész minta gyakori-e.

Nézzünk erre egy gyors példát. Amennyiben a mintákat elemhalmazok formájában keressük, akkor az APRIORI és Zaki módszere is először meghatározza a gyakori elemeket. Legyenek ezek az  $A, C, D, G, M$  elemek. Az APRIORI ezek után előállítana  $\binom{5}{2}$  darab jelöltet, majd meghatározná támogatottságaikat. Zaki ehelyett csak az  $A$  prefixű kételemű halmazok támogatottságát vizsgálja. Ha ezek közül gyakori például az  $\{A, C\}$ ,  $\{A, G\}$ , akkor a következőkben az  $\{A, C, G\}$ -t nézi, és mivel további jelöltet nem tud előállítani, ugrik a  $C$  prefixű elemhalmazok vizsgálatára, és így tovább.

Látnunk kell, hogy Zaki módszere csak több jelöltet állíthat elő, mint az APRIORI. A mélységi bejárás miatt ugyanis egy jelölt előállításánál nem áll rendelkezésünkre az összes rész minta. Az előző példa esetében például az  $\{A, C, G\}$  támogatottságát hamarabb vizsgálja, mint a  $\{C, G\}$  halmazét, holott ez utóbbi akár ritka is lehet. Ebben a tekintetben tehát Zaki módszere rosszabb az APRIORI-nál, ugyanis több hamis jelöltet állít elő.

Zaki módszerének igazi ereje a jelöltek támogatottságának meghatározásában van. A minták illeszkedési listáinak előállítása egy rendkívül egyszerű és nagyon gyors művelet lesz. Emellett ahogy haladunk egyre mélyebbre a mélységi bejárás során, úgy csökken az illeszkedési listák hossza, és ezzel a támogatottság meghatározásának ideje is.

A bemenet szűrésének ötletét az APRIORI algoritmusnál is elsűthetjük, de nem ilyen mértékben. Ha ismerjük a gyakori egyelemű mintákat, akkor törölhetjük azon sorozatelemeket, amelyek nem illeszkednek egyetlen gyakori egyelemű mintára sem. Sőt ezt a gondolatot általánosíthatjuk is: az  $\ell$ -edik lépésben törölhetjük a bemeneti sorozat azon elemeit, amelyek nem illeszkednek egyetlen  $(\ell - 1)$ -elemű mintára sem. Ez a fajta bemeneti tér szűkítés azonban nem lesz olyan hatékony, mint amilyen a Zaki módszerében. Ott ugyanis egyszerre csak 1 prefixet vizsgálunk, az APRIORI-nál azonban általában sok olyan minta van, aminek csak az üres minta a közös részmintája.

Összességében tehát az APRIORI kevesebb jelöltet generál, mint Zaki módszere, de a jelöltek támogatottságának meghatározása több időt vesz igénybe. Általánosságban nem lehet megmondani, hogy melyik a jobb módszer. Egyes adatbázisok esetén az APRIORI, másoknál a Zaki módszer. Sőt könnyen lehet olyan példát mutatni, amikor az egyik algoritmus nagyságrendileg több időt tölt a feladat megoldásával, mint a másik.

Zaki módszerénél könnyű kezelni a anti-monoton és a prefix anti-monoton kényszereket. A nem gyakori minták törlésekor töröljük azokat a mintákat is, amelyek nem elégítenek ki minden anti-monoton kényszert. A prefix anti-monoton kényszereket a jelöltek előállítását után kell figyelembe vennünk: törölhetjük azokat a generátorokat, amelyekre nem teljesül az anti-monoton kényszer. A `zaki_segéd` eljárásból következik, hogy ilyen  $m$  mintát legfeljebb olyan jelölt előállításánál fogunk felhasználni, aminek  $m$  a prefixe. Természetesen itt is bajban vagyunk, ha több prefix anti-monoton kényszer van adva, hiszen ezek  $\prec$ -nek különböző lineáris kiterjesztéseit használhatják.

### 10.5.3. Mintanövelő algoritmusok

A mintanövelő (pattern growth) algoritmus olyan mintakeresés esetén alkalmazható, amikor a bemenet minták sorozataként van megadva, és az illeszkedést a tartalmazás alapján definiáljuk, értelmezhető a prefix, és a minták *egyértelműen növelhetők*. Például a növelés művelet halmazok esetén az unió, sorozatok esetében a konkatenáció képzésének felel meg (és ebből látszik, hogy a növelés művelete nem feltétlenül kommutatív).

**10.18. definíció.** Az  $\mathcal{MK}=(\mathcal{M}, \preceq)$  mintakörnyezet mintái egyértelműen növelhetők, ha létezik egy olyan  $+$  „növelő” művelet, amellyel az  $\mathcal{M}$  félcsoportot alkot.

A növelés inverze a *csökkentés*, jelölése:  $-$ . Az  $m - m'$  művelet eredménye az az  $m''$  minta, amivel  $m'$ -t növelve  $m$ -et kapjuk.

A mintanövelő módszerek csak egyelemű jelölteket használnak, és emellett a bemeneten végeznek olyan műveleteket, amelyek eredményeként megkapjuk a gyakori mintákat. A két művelet a *szűrés* és a *vetítés*, amelyek az eredeti  $\mathcal{S}$  bemenetből egy „kisebb”  $\mathcal{S}'$  bemenetet állítanak elő. A szűrés a gyakori egyelemű mintákat használja és olyan  $\mathcal{S}'$  bemenetet állít elő amelyben a gyakori minták megegyeznek az  $\mathcal{S}$ -beli gyakori mintákkal. Az  $\mathcal{S}$  bemenet  $m$  mintára vetítése (jelölésben  $\mathcal{S}|m$ ) pedig olyan  $\mathcal{S}'$  bemenetet állít elő, amelyre igaz, hogy ha  $m$ -et az  $\mathcal{S}'$ -beli gyakori mintákkal növeljük, akkor megkapjuk az  $\mathcal{S}$ -beli,  $m$ -et tartalmazó gyakori mintákat. A  $m$ -et tartalmazó gyakori minták meghatározásához csak azokra a bemeneti elemekre van szükség, amelyekre illeszkedik  $m$ , ezért a vetítés első lépése mindig ezen elemek meghatározása lesz.

Ha például a bemenet elemei elemhalmazok és akkor illeszkedik egy elemhalmaz a bemenet egy elemére, ha annak része, akkor szűrés művelet az lesz, hogy a bemeneti

elemekből töröljük a ritka elemeket. Nyilvánvaló, hogy ritka elem nem játszik szerepet a gyakori elemek meghatározásában. A bemenet  $X$  halmazra vetítését megkapjuk, ha töröljük azon bemeneti elemeket, amelyeknek nem része  $X$ , majd a kapott elemekből töröljük  $X$ -et. Legyen  $\mathcal{S} = \langle \{A, C, F\}, \{B, G\}, \{A, C, D\}, \{A, C\}, \{B, C\}, \{C, D, E\}, \{A, B, C\} \rangle$  amelynek szűrése 2-es támogatottsági küszöb esetén az  $\tilde{\mathcal{S}} = \langle \{A, C\}, \{B\}, \{A, C, D\}, \{A, C\}, \{B, C\}, \{C, D\}, \{A, B, C\} \rangle$  sorozat és  $\tilde{\mathcal{S}}|_{\{A, C\}} = \langle \{D\}, \{B\} \rangle$ .

A mintanövelő módszer rendkívül egyszerű, tulajdonképpen a feladatot rekurzívan kisebb részfeladat megoldására vezeti vissza. A rekurziós eljárást a bemenet szűrésével és különböző mintákra vett vetítéseivel hívja meg, miközben a mintateret is csökkenti. Jelöljük  $\mathcal{M} \setminus \bar{m}$ -el azt a mintateret, amit úgy kapunk  $\mathcal{M}$ -ből, hogy töröljük azon mintákat, amelynek  $m$  részmintája ( $\bar{m}$ ). Ha az  $m$  minta támogatottsága  $\mathcal{S}$ -ben  $\text{supp}_{\mathcal{S}}(m)$  és az  $m' \in \mathcal{M} \setminus \bar{m}$  támogatottsága  $\mathcal{S}|_m$ -ben  $\text{supp}_{\mathcal{S}|_m}(m')$ , akkor  $m+m'$  támogatottsága is  $\text{supp}_{\mathcal{S}|_m}(m')$ . A módszer pszeudokódja a ?? ábrán látható.

---

**Algorithm 13** Mintanövelő módszer
 

---

**Require:**  $B$ : bemeneti adat

$\text{min\_supp}$ : támogatottsági küszöb

$\mathcal{M}$ : mintatér

$J_1 \leftarrow$  1-elemű minták  $\{J_1$ : egyelemű jelöltek}

támogatottság\_meghatározás( $B, J_1$ )

$GY_1 \leftarrow$  gyakoriak\_kiválogatása( $J_1, \text{min\_supp}$ )

$\tilde{B} \leftarrow$  szűrés( $B$ )

**for all**  $gy \in GY_1$  **do**

$GY' \leftarrow$  mintanövelő\_módszer( $\tilde{B}|_{gy}, \text{min\_supp}, \mathcal{M} \setminus \tilde{g}y$ )

**for all**  $gy' \in GY'$  **do**

$GY \leftarrow GY \cup \{gy + gy'\}$

**end for**

**end for**

**return**  $GY$ : gyakori minták

---

A módszer előnye abban rejlik, hogy szűrést, vetítést és az egyelemű jelöltek támogatottságát hatékonyan tudjuk megvalósítani. A hatékonyság növelése érdekében a vetített tranzakciók azonos elemeit csak egyszer tároljuk, általában egy fa-szerű struktúrában.

Az anti-monoton kényszerek kezelése a mintanövelő algoritmusok esetében is egyszerű. Ne folytassuk a rekurziót, ha a minta nem elégíti ki minden anti-monoton kényszert.

Az egyes mintatípusok esetében úgy fogjuk megadni a növelés műveletet, hogy tetszőleges minta csökkentése a minta prefixét fogja adni. Ez azt eredményezi, hogy törölhetjük azt a mintát, amelyik nem elégíti ki a prefix anti-monoton kényszert, és leállhatunk a rekurzióval. Hasonlóan az APRIORI és a Zaki módszeréhez itt sincs mód több prefix anti-monoton kényszer hatékony kezelésére. Az algoritmus menetét ugyanis egyértelműen megadja a növelés művelet, amit a prefix anti-monoton kényszerben felhasznált teljes rendezés alapján definiálunk.

### 10.5.4. Kétlépcsős technikák

A szélességi bejárást megvalósító algoritmusok az adatbázist legalább annyiszor olvassák végig, amekkora a legnagyobb gyakori minta mérete. Előfordulhatnak olyan alkalmazások, amelyeknél az adatbázis elérése drága művelet. Ilyenre lehet példa, amikor az adatbázis egy elosztott hálózatban található, vagy lassú elérésű háttértárolón.

A kétlépcsős algoritmusok [118, 136] a teljes adatbázist legfeljebb kétszer olvassák végig. I/O tekintetében tehát legyőzik például az APRIORI algoritmust, azonban olyan futási környezetben, ahol a futási időt nem szinte kizárólag az I/O műveletek határozzák meg (ha a bemenet elfér a memóriában akkor ez a helyzet áll fenn), az APRIORI algoritmus gyorsabban ad eredményt.

#### Naiv mintavételező algoritmus

Olvassuk be a teljes bemenet egy részét a memóriába (a rész nagyságára nézve lásd 87. oldal). Erre a kis részre futtassuk le az APRIORI algoritmust az eredeti *min\_freq* gyakorisági küszöbvel. A kis részben megtalált gyakori minták lesznek a jelöltek a második fázisban, amelynek során a jelöltek támogatottságát a teljes adatbázisban meghatározzuk. Ezáltal ki tudjuk szűrni azokat a mintákat, amelyek ritkák, de a kis részben gyakoriak. Előfordulhat azonban a fordított helyzet, azaz a kis adatbázisban egy minta ritka, viszont globálisan gyakori, tehát nem kerül a jelöltek közé, és így nem is találhatjuk azt gyakorinak. Javíthatunk a helyzeten, ha csökkentjük a kis részben a gyakorisági küszöböt, amivel növeljük a jelöltek számát, de csökkentjük annak veszélyét, hogy egy gyakori mintát ritkának találunk.

Ennek az egyszerű algoritmusnak két hátránya van. Egyrészt nem ad arra garanciát, hogy minden gyakori mintát megtalálunk (azaz nem teljes), másrészt a gyakorisági korlát csökkentése miatt a hamis jelöltek száma túlzottan nagy lehet. A fenti két problémát küszöböli ki a partíciós, illetve a Toivonen-féle algoritmus.

Mivel a kétlépcsős algoritmusok egy kis rész kiválasztásán alapulnak, így nagyon érzékenyek az adatbázis homogenitására. Gondoljunk itt a szezonális elemekre, amelyek lokálisan gyakoriak, de globálisan ritkák. Például a kesztyűk eladása tél elején nagy, de mégis a kesztyű önmagában ritka elem. Amennyiben a kis rész kiválasztása a bemenet egy véletlen pontjáról történő szekvenciális olvasást jelentene, akkor az nagy eséllyel sok hamis és hiányzó jelöltet eredményezne.

#### Partíciós algoritmus

A partíciós algoritmus [118] kétszer olvassa végig a teljes adatbázist. Páronként diszjunkt részekre osztja a bemenetet ( $\mathcal{S} = \langle \mathcal{S}_1, \mathcal{S}_2, \dots, \mathcal{S}_r \rangle$ ), majd az egyes részekre meghívja az APRIORI algoritmust, ami megadja az egyes részekben gyakori mintákat (hívjuk őket lokálisan gyakori mintáknak). A második végigolvasásnál egy minta akkor lesz jelölt, ha valamelyik részben gyakori volt. Könnyen látható, hogy az algoritmus teljes, hiszen egy gyakori mintának legalább egy részben gyakorinak kell lennie, és ezt az APRIORI ki fogja szűrni (mivel az APRIORI is teljes).

Kérdés, hogy hány részre osszuk a teljes adatbázist. Nyilvánvaló, hogy minél nagyobb az egyes részhalmazok mérete, annál jobb képet ad a teljes adatbázisról, tehát annál kevesebb lesz a hamis jelölt. A részek nagy mérete azonban azt eredményezi, hogy azok nem férnek el a memóriában, és így az APRIORI algoritmus sok időt tölt el partíciórészek ideiglenes háttérbe másolásával és visszaolvasásával. Habár globálisan csak kétszer olvassuk végig a teljes

adatbázist, azonban az egyes partíciók I/O igényének összege legalább akkora, mintha a teljes adatbázisra futtatnánk le az APRIORI algoritmust. Végeredményben a második végigolvasás miatt a partíciós algoritmus I/O igénye nagyobb lesz, mint az APRORI algoritmusé.

Ha az egyes részek elférnek a memóriában, akkor nem lép fel a fenti probléma, hisz az APRIORI algoritmus nem fog I/O műveletet igényelni (feltéve, ha a jelöltek a számlálóikkal együtt is elférnek még a memóriában). Túl kis méret választása azonban azt eredményezheti, hogy a partíció nem ad hű képet a teljes adatbázisról, így a lokális gyakori minták mások (is!) lesznek, mint a globális gyakori minták, ami túl sok hamis jelöltet eredményezhet.

A helyes partícióméret tehát a rendelkezésünkre álló memóriától függ. Legyen minél nagyobb, de úgy, hogy a jelöltek számlálóikkal együtt is elférjenek a memóriában. Természetesen a jelöltek száma a gyakori minták méretétől függ, amiről a partícióméret meghatározásakor még nincs pontos képünk.

A partíciós algoritmus szintén érzékeny a bemenet homogenitására. Ezt az érzékenységet csökkenthetjük, ha módosítjuk egy kicsit az algoritmust. Ha egy  $m$  minta gyakori az  $\mathcal{S}_i$  részben, akkor a rákövetkező  $\mathcal{S}_{i+1}, \mathcal{S}_{i+2}, \dots, \mathcal{S}_{i+l}$  részekben is határozzuk meg a támogatottságát egészen addig, amíg  $\text{freq}_{\bigcup_{j=i}^{i+l} \mathcal{S}_j}(m) \geq \text{min\_freq}$ . Ha ezalatt eljutunk az utolsó részig, akkor vegyük fel  $m$ -et a második végigolvasás jelöltjei közé. Ellenkező esetben felejtsük el, hogy  $m$  gyakori volt ezen részekben. Ha egy mintát az összes részben vizsgáltunk, akkor ezt szintén szükségtelen felvenni jelöltnek a második végigolvasásnál, hiszen támogatottsága megegyezik az egyes résztámogatottságok összegével.

A partíciós algoritmus további előnye, hogy remekül párhuzamosítható. Saját memóriával rendelkező feldolgozó egységek végezhetik az egyes részek gyakori mintakeresését, és ezáltal mind az első, mind a második fázis töredék idő alatt elvégezhető.

### Toivonen algoritmusa

Az naív mintavételező algoritmus nagy hátránya, hogy még csökkentett  $\text{min\_freq}$  mellett sem lehetünk biztosak abban, hogy nem vesztettünk el gyakori mintát. Toivonen algoritmusa [136] az adatbázist egyszer olvassa végig, és ha jelenti, hogy minden mintát megtalál, akkor bizonyítható, hogy ez igaz. Az algoritmus nem más, mint a naív mintavételező algoritmus továbbfejlesztett változata. Az egyszerű algoritmusnál azonban több információt ad, ugyanis jelenti, ha biztos abban, hogy minden gyakori mintát előállított, és azt is jelenti, amikor lehetséges, hogy van hiányzó jelölt (olyan gyakori minta, ami nem jelölt, és így nem találhatjuk azt gyakorinak). A lehetséges hiányzó jelöltekről információt is közöl.

Alapötlete az, hogy ne csak a kis részben található gyakori minták előfordulását számoljuk össze a teljes adatbázisban, hanem azok minimális valódi felső korlátait is. Mit jelent az, hogy az  $m$  minta tetszőleges  $M \subseteq \mathcal{M}$  mintahalmaz minimális valódi felső korlátai közé tartozik (jelölésben  $m \in MVFK(M)$ )? Először is a valódi felső korlát formálisan:  $m' \prec m$  minden  $m' \in M$ . A minimalitás pedig azt jelenti, hogy nem létezik olyan  $m''$  minta, amely  $M$ -nek valódi felső korlátja és  $m'' \prec m$ . A gyakori minták minimális valódi felső korlátjai azok a ritka minták, amelyek minden részmintája gyakori.

Például elemhalmaz típusú minta esetén, ha  $\mathcal{M} = 2^{\{A,B,C,D,E,F\}}$  és  $M = \{\{A\}, \{B\}, \{C\}, \{F\}, \{A, B\}, \{A, C\}, \{A, F\}, \{C, F\}, \{A, C, F\}\}$ , akkor  $MVFK(M) = \{\{B, C\}, \{B, F\}, \{D\}, \{E\}\}$ .

Toivonen algoritmusában a teljes adatbázisból egy kis részt veszünk. Ebben meghatározzuk a gyakori minták halmazát és ennek minimális valódi felső korlátját. A teljes adatbázisban ezek

támogatottságát vizsgáljuk, és gyűjtjük ki a globálisan gyakoriakat. A következő egyszerű tétel ad információt arról, hogy ez az algoritmus mikor teljes, azaz mikor lehetünk biztosak abban, hogy minden gyakori mintát meghatároztunk.

**10.19. tétel.** *Legyen  $\mathcal{S}'$  az  $\mathcal{S}$  bemeneti sorozat egy része. Jelöljük  $GY$ -vel az  $\mathcal{S}$ -ben,  $GY'$ -vel az  $\mathcal{S}'$ -ben gyakori mintákat és  $GY^*$ -al azokat az  $\mathcal{S}$ -ben gyakori mintákat, amelyek benne vannak  $GY' \cup MVFK(GY')$ -ben ( $GY^* = GY \cap (GY' \cup MVFK(GY'))$ ). Amennyiben*

$$GY^* \cup MVFK(GY^*) \subseteq GY' \cup MVFK(GY')$$

*teljesül, akkor  $\mathcal{S}$ -ben a gyakori minták halmaza pontosan a  $GY^*$ , tehát  $GY^* \equiv GY$ .*

*Bizonyítás:* Indirekt tegyük fel, hogy létezik  $m \in GY$ , de  $m \notin GY^*$ , és a feltétel teljesül. A  $GY^*$  definíciója miatt ekkor  $m \notin GY' \cup MVFK(GY')$ . Vizsgáljuk azt a legkisebb méretű  $m' \preceq m$ -t, amire  $m' \in GY$  és  $m' \notin GY^*$  (ilyen  $m'$ -nek kell lennie, ha más nem, ez maga az  $m$  minta). Az  $m'$  minimalitásából következik, hogy minden valódi részmintája eleme  $GY' \cup MVFK(GY')$ -nek és gyakori. Ebből következik, hogy  $m'$  minden részmintája eleme  $GY^*$ -nak, amiből kapjuk, hogy  $m' \in MVFK(GY^*)$ . Ez ellentmondást jelent, hiszen a feltételnek teljesülnie kell, azonban van olyan elem ( $m'$ ), amely eleme a bal oldalnak, de nem eleme a jobb oldalnak. ■

Tetszőleges  $GY'$  halmaz esetén az  $MVFK(GY') \cup GY'$ -t könnyű előállítani. Sőt, amennyiben a gyakori mintákat APRIORI algoritmussal határozzuk meg, akkor  $MVFK(GY')$  elemei pontosan a ritka jelöltek lesznek (hiszen a jelölt minden része gyakori).

Nézzünk egy példát Toivonen algoritmusára. Legyen a mintatér a  $\{A,B,C,D\}$  hatványhalmaza. A kis részben az  $\{A\}, \{B\}, \{C\}$  elemhalmazok gyakoriak. Ekkor a minimális valódi felső korlát elemei az  $\{A,B\}, \{A,C\}, \{B,C\}, \{D\}$  halmazok. Tehát ennek a 7 elemhalmaznak fogjuk a támogatottságát meghatározni a teljes adatbázisban. Ha például az  $\{A\}, \{B\}, \{C\}$   $\{A,B\}$  halmazokat találjuk gyakorinak a teljes adatbázisban, akkor a tételbeli tartalmazási reláció fennáll, hiszen az  $\{A\}, \{B\}, \{C\}, \{A,B\}$  halmaz minimális valódi felső korlátai közül mind szerepel a 7 jelölt között. Nem mondható ez, ha  $\{D\}$ -ről derül ki, hogy gyakori. Ekkor Toivonen algoritmus jelent, hogy előfordulhat, hogy nem biztos, hogy minden gyakori elemhalmazt megtalált. Az esetleg kimaradtak csak (!) az  $\{A,D\}, \{B,D\}, \{C,D\}$  halmazok lehetnek.

### 10.5.5. A zárt minták „törékenysége”

Tagadhatatlan, hogy a zárt mintákon alapuló memóriacsökkentés egy szép elméleti eredmény. Ne foglaljunk helyet a memóriában a gyakori, nem zárt mintáknak, hiszen a zárt, gyakori mintákból az összes gyakori minta meghatározható.

Ez a technika ritkán alkalmazható azon esetekben, amikor a bemenet sorozat formájában adott, a támogatottságot pedig egy illeszkedési predikátum alapján definiáljuk. És, mint azt már említettük, a legtöbbször ez áll fenn.

Ennek oka, hogy gyakori mintákat általában nagy, zajokkal terhelt adatbázisokban keresnek. Ilyen adatbázisban szinte az összes elemhalmaz zárt, így a módszerrel nem nyerünk semmit. Gondoljuk meg, hogy ha egy adatbázist úgy terhelünk zajjal, hogy véletlenszerűen beszurunk egy-egy új elemet, akkor folyamatosan növekszik az esélye annak, hogy egy minta zárt lesz. A nemzártág tehát egy „sérülékeny” tulajdonság. Tetszőleges nem zárt  $m$  mintát zárttá tehetünk egyetlen olyan tranzakció hozzáadásával, amely illeszkedik  $m$ -re, de nem illeszkedik egyetlen olyan mintára sem, amelynek  $m$  valódi részmintája.

### 10.5.6. Dinamikus gyakori mintabányászat

Nagy adatbázisok esetén a gyakori minták kinyerése még a leggyorsabb algoritmusokat felhasználva is lassú művelet. Az adatbázisok többségében a tárolt adatok nem állandóak, hanem változnak: új elemeket veszünk fel, egyeseket módosítunk, vagy törölünk. Ha azt szeretnénk, hogy a kinyert gyakori minták konzisztensek legyenek az adatbázisban tárolt adatokkal, akkor bizonyos időközönként a gyakori minták adatbázisát is frissíteni kell.

A konzisztenciát elérhetjük úgy, hogy lefuttatjuk valamelyik ismert (APRIORI, Zaki stb.) algoritmust minden módosítás után. Ennek az a hátránya, hogy lassú, hiszen semmilyen eddig kinyert tudást nem használ fel. Szükség van tehát olyan algoritmusok kifejlesztésére [11, 29, 30, 98, 117, 134], ami felhasználja az adatbázis előző állapotára vonatkozó információkat és így gyorsabban ad eredményt, mint egy nulláról induló, hagyományos algoritmus.

Itt most azt az esetet nézzük, amikor csak bővíthetjük a bemenetet, de a leírt módszerek könnyen általánosíthatók arra az esetre, amikor törölhetünk is a bemenetből. Adott tehát  $\mathcal{S}$  bemeneti sorozat, amelyben ismerjük a gyakori mintákat ( $GY^{\mathcal{S}}$ ) és azok támogatottságát. Ezen kívül adott az új bemeneti elemek sorozata  $\mathcal{S}'$ . A feladat a  $\langle \mathcal{S}, \mathcal{S}' \rangle$ -ben található gyakori minták ( $GY^{\langle \mathcal{S}, \mathcal{S}' \rangle}$ ) és azok támogatottságának meghatározása.

#### FUP algoritmus

A FUP (Fast Update) [29] a legegyszerűbb szabály- karbantartó algoritmus. Tulajdonképpen nem más, mint az APRIORI algoritmus módosítása.

Kétféle jelöltet különböztetünk meg: az első csoportba azok a minták tartoznak, melyek az eredeti adatbázisban gyakoriak voltak, a másodikba azok, amelyek nem. Nyilvánvaló, hogy az új adatbázisban mindkét csoport elemeinek támogatottságát meg kell határozni, a régi adatbázisban azonban elég a második csoport elemeit vizsgálni.

A FUP az alábbi trivialisításokat használja fel.

- I. Ha egy minta  $\mathcal{S}$ -ban gyakori volt és  $\mathcal{S}'$ -ben is az, akkor az  $\langle \mathcal{S}, \mathcal{S}' \rangle$ -ben is biztos gyakori, előfordulása megegyezik  $\mathcal{S}'$ -beni és  $\mathcal{S}$ -beni előfordulások összegével.
- II. Amennyiben egy elemhalmaz  $\mathcal{S}$ -ban ritka, akkor  $\langle \mathcal{S}, \mathcal{S}' \rangle$ -ben csak abban az esetben lehet gyakori, ha  $\mathcal{S}'$ -ben gyakori. Ezek szerint ne legyen jelölt olyan elemhalmaz, amely sem  $\mathcal{S}$ -ban, sem  $\mathcal{S}'$ -ben nem gyakori.

Ezekből következik, hogy csak olyan elemhalmazok lesznek jelöltek  $\mathcal{S}$  végigolvasásánál, amelyek  $GY^{\mathcal{S}}$ -ban nem szerepeltek, de  $\mathcal{S}'$ -ben gyakoriak voltak.

Az algoritmus pszeudokódja a 14-es ábrán látható.

A támogatottság meghatározás, gyakoriak kiválogatása és a jelölt-előállítás lépések teljes egészében megegyeznek a APRIORI ezen lépéseivel.

A FUP algoritmust könnyű módosítani arra az esetre, amikor nem csak hozzáadunk új elemeket az eredeti bemeneti sorozathoz, hanem törölünk is néhányat a régi elemek közül (FUP2 algoritmus [30]).

A  $FUP$  és  $FUP_2$  algoritmusok nem mentesek az APRIORI algoritmus legfontosabb hátrányától, attól, hogy a teljes adatbázist annyiszor kell átolvasni, amekkora a legnagyobb gyakori jelölminta mérete. Ezen a problémán próbáltak segíteni a később publikált algoritmusok.



**Algorithm 14** FUP algoritmus**Require:**  $S$ : régi bemeneti adat $S'$ : új bemeneti adat $GY^S$ : régi gyakori minták $min\_freq$ : gyakorisági küszöb $\ell \leftarrow 0$  $J_\ell^1 \leftarrow GY_\ell^S \setminus \{J_\ell^1: 1\text{-es csoportbeli jelöltek}\}$  $J_\ell^2 \leftarrow \{\text{üres minta}\} \setminus GY_\ell^S \setminus \{J_\ell^2: 2\text{-es csoportbeli jelöltek}\}$ **while**  $|J_\ell^1| + |J_\ell^2| \neq 0$  **do**    támogatottság\_meghatározás( $S', J_\ell^1 \cup J_\ell^2$ )     $J_\ell^* \leftarrow$  gyakoriak\_kiválogatása( $J_\ell^2, min\_freq$ )    **if**  $|J_\ell^*| \neq 0$  **then**        támogatottság\_meghatározás( $S, J_\ell^*$ )    **end if**     $GY_\ell \leftarrow$  gyakoriak\_kiválogatása( $J_\ell^1 \cup J_\ell^*, min\_freq$ )     $J_{\ell+1}^{**} \leftarrow$  jelölt\_előállítás( $GY_\ell$ )     $\ell \leftarrow \ell + 1$      $J_\ell^1 \leftarrow J_\ell^{**} \cap GY_\ell^S$      $J_\ell^2 \leftarrow J_\ell^{**} \setminus GY_\ell^S$     delete( $J_\ell^{**}$ )**end while****return**  $GY^{(S, S')}$ : gyakori minták**Esélyes jelölteken alapuló dinamikus algoritmus**

A [134] cikkben Toivonen algoritmusában használt minimális valódi felső korlátokat használják annak érdekében, hogy csökkentsék a nagy adatbázist átolvasásának számát. Az adatbázis növekedése során először a minimális valódi felső korlátok válnak gyakorivá. Ha nem csak a gyakori minták előfordulását ismerjük a régi adatbázisban, hanem azok minimális valódi felső korlátait is, akkor lehet, hogy szükségtelen a régi adatbázist végigolvasni. Ha ugyanis az új tranzakciók felvételével egyetlen minimális valódi felső korlát sem válik gyakorivá, akkor biztos, hogy nem keletkezett új gyakori minta. A 10.19-as tétel ennél erősebb állítást fogalmaz meg: még ha bizonyos minimális valódi felső korlátok gyakorivá váltak, akkor is biztosak lehetünk abban, hogy nem kell a régi adatbázist átvizsgálunk, mert nem keletkezhetett új gyakori minta. Átültetve a tételt a jelenlegi környezetbe: ha  $GY^{S \cup S'} \cup MVFK(GY^{S \cup S'}) \subseteq GY^S \cup MVFK(GY^S)$ , akkor biztosak lehetünk, hogy nem keletkezett új gyakori minta, és csak a támogatottságokat kell frissíteni.

# 11. fejezet

## Gyakori sorozatok, bool formulák és epizódok

A kutatások középpontjában a gyakori elemhalmazok állnak. Tovább léphetünk, és kereshetünk bonyolultabb típusú mintákat is. Erről szól ez a fejezet.

### 11.1. Gyakori sorozatok kinyerése

Napjainkban az elektronikus kereskedelem egyre nagyobb méretet ölt. A vevők megismerésével és jobb kiszolgálásával célunk a profitnövekedés mellett a vásárlói elégedettség fokozása. Az elektronikus kereskedelem abban különbözik a hagyományos kereskedelemtől, hogy az egyes tranzakciókhoz hozzárendelhetjük a vásárlókat. Eddig a tranzakciók (kosarak) óriási halmaza állt rendelkezésünkre, most ennél több: pontosan tudjuk, hogy ki, mikor, mit vásárol.

Az újabb adatok újabb információkinyeréshez adhatnak alapot. Nem csak általános vásárlási szabályokat állíthatunk elő, hanem ennél többet: személyre szabhatjuk a vásárlási szokásokat, vevők csoportjait alakíthatjuk ki, megkereshetjük a sok, illetve kevés profitot hozó vásárlási csoportokat, stb.

Ebben a fejezetben a vevők között gyakran előforduló vásárlói minták kinyerésével foglalkozunk. Két példa: sok vevő a „Csillagok háborúja” DVD megvétele után a „Birodalom visszavág” című filmet, majd később a „Jedi visszatér” című filmet is megveszi DVD-n, vagy a vevők 30%-a új mobiltelefon és új tok vásárlása után új előlapot is vásárol.

Kereskedelmi cégek a kinyert gyakori mintákat, epizódokat újabb profitnövekedést hozó fogásokra használhatják. Például kiderülhet, hogy a videomagnót vásárlók nagy aránya a vásárlást követő 3-4 hónappal kamerát is vásárolnak. Ekkor ha valaki videomagnót vesz, küldjük ki postán kamerákat reklámozó prospektusokat a vásárlást követően 2-3 hónappal. A szekvenciális mintakinyerés (és egyéb epizódkutató algoritmusok) nem csak az on-line áruházakra jellemző. Felhasználási területük egyre bővül, a további kutatásokat a gyakorlatban előforduló problémák is igénylik. Jellemző terület a direkt marketing, de további felhasználási területre lehet példa az alábbi: páciensek tüneteit és betegségeit tartalmazó adatbázisokból kinyert minták nagy segítségre lehetnek az egyes betegségek kutatásánál, nevezetesen, hogy az egyes betegségeket milyen tünetek, vagy más betegségek előzik meg gyakran.

Mielőtt rátérünk arra, hogy miként lehet kinyerni elemhalmazokat tartalmazó sorozatokból a gyakoriakat, egy egyszerűbb esettel foglalkozunk, ahol a sorozat elemei atomi események.

### 11.1.1. A Gyakori Sorozat Fogalma

A gyakori sorozatok kinyerésének feladata annak a feladatkörnek egy esete, amikor a támogatottságot a tartalmazási predikátum alapján definiáljuk. Feltételezzük, hogy az olvasó tisztában van a 10.5 részben definiált fogalmakkal.

Adott  $(\mathcal{J} = \{i_1, i_2, \dots, i_m\}$  elemek (vagy termékek) halmaza és  $v$  darab  $\mathcal{J}$  felett értelmezett sorozat. Tehát a bemenet sorozatoknak egy sorozata:

$$\text{bemenet} : \mathcal{S} = \langle S_1, S_2, \dots, S_v \rangle,$$

ahol  $S_k = \langle i_1^{(k)}, i_2^{(k)}, \dots, i_{n^{(k)}}^{(k)} \rangle$ , és  $i_j^{(k)} \in \mathcal{J}$ .

Definiáljuk a  $\mathcal{M} = (\mathcal{M}, \preceq)$  mintakörnyezet tagjait sorozatok esetében. Az  $\mathcal{M}$  elemei az  $\mathcal{J}$  felett értelmezett sorozatok:

**11.1. definíció.**  $S = \langle i_1, \dots, i_m \rangle$  sorozat tartalmazza  $S' = \langle i'_1, \dots, i'_n \rangle$  sorozatot (jelöléssel  $S' \preceq S$ ), ha léteznek  $j_1 < j_2 < \dots < j_n$  egész számok úgy, hogy  $i'_1 = i_{j_1}, i'_2 = i_{j_2}, \dots, i'_n = i_{j_n}$ .

Amennyiben  $S' \preceq S$ , akkor  $S'$  az  $S$  részsorozata. Például a  $\langle G, C, I, D, E, H \rangle$  sorozat tartalmazza a  $\langle C, D, H \rangle$  sorozatot. Ebben a mintakörnyezetben  $\|\cdot\|$  függvény a sorozat hosszát adja meg. A fentiek alapján a fedés, TID lista, támogatottság, gyakoriság, gyakori sorozat definíciója egyértelmű.

Egy alap mintakinyerési feladatban adott  $S_i$  sorozatok sorozata, továbbá  $\text{min\_supp}$  támogatottsági küszöb, elő kell állítani a gyakori sorozatokat.

### 11.1.2. APRIORI

A fent definiált feladat a gyakori mintakinyerés egy speciális esete, így alkalmazhatók rá az általános algoritmusok, például az APRIORI. Az általános leírást megadtuk a 10.4 részben, itt most csak azon speciális részleteket vizsgáljuk, amelyek sorozat típusú mintatér esetén érvényesek. Két lépést vizsgálunk közelebbről a jelöltek előállítását és a támogatottság meghatározását.

#### Jelöltek előállítása

Az APRIORI jelöltelőállítás két lépésből áll: potenciális jelöltek előállítása, majd a potenciális jelöltek részmintáinak vizsgálata. Akkor lesz egy  $\ell$ -elemű potenciális jelölből jelölt, ha minden  $\ell - 1$  elemű részsorozata gyakori. Általánosan annyit mondtunk el, hogy egy potenciális jelölt két  $\ell - 1$  elemű gyakori mintáknak (ezeket hívtuk a jelölt generátorainak) a minimális valódi felső korlátja. Sorozat típusú minta esetén akkor lesz két  $\ell - 1$  elemű gyakori mintáknak a minimális valódi felső korlátja  $\ell$  elemű, ha van  $\ell - 2$  elemű közös részsorozatuk.

A hatékonyság szempontjából fontos lenne, ha a jelöltek előállítása ismétlés nélküli lenne. Ehhez szükségünk van a sorozatokon értelmezett teljes rendezésre. Az  $\mathcal{J}$  elemein tudunk egy tetszőleges teljes rendezést definiálni, ami szerinti lexikografikus rendezés megfelel a célnak. A rendezés alapján értelmezhetjük egy sorozat tetszőleges elemű prefixét. Két  $\ell - 1$  elemű gyakori mintákból akkor képezek potenciális jelöltet, ha  $\ell - 2$  elemű prefixük megegyeznek (hasonlóan a halmazok eseténél). A minimális valódi felső korlát a az utolsó elemmel bővített sorozatok lesznek.

A generátorok lehetnek azonos sorozatok is. Például az  $\langle G, C, I \rangle$  sorozat önmagával a  $\langle G, C, I, I \rangle$  jelöltet fogja előállítani. Látnunk kell, hogy ez a jelöltelőállítás ismétlés nélküli, ugyanis tetszőleges jelölteknek egyértelműen meg tudjuk mondani a generátorait.

### Támogatottság meghatározása

A jelölt sorozatok támogatottságának meghatározás szinte megegyezik a jelölt halmazok támogatottságának meghatározásával. Erről részletesen szoltunk a 4.2.2 részben. Itt csak az apró különbségekre térünk ki.

A kételemű jelölteknel nem csak a kétdimenziós tömb egyik felét fogjuk használni, hanem a teljes tömböt. Ez abból következik, hogy számít a sorrend, tehát például az  $\langle A, B \rangle$  sorozat különbözik az  $\langle B, A \rangle$  sorozattól.

Kettőnél nagyobb jelölteket célszerű szófában tárolni. A szófa felépítése, a jelöltek támogatottságának meghatározása 1 apró részlettől eltekintve teljesen megegyezik a halmazoknál leírtakkal. A szófa bejárásakor ügyelni kell arra, hogy a sorozatban lehetnek ismétlődő elemek, illetve az elemek nincsenek sorba rendezve. A rekurziós lépés nem két rendezett lista közös elemeinek meghatározását jelenti, hanem egy rendezett lista (az adott belső pontból kiinduló élek címkéi) azon elemeinek meghatározását, amelyek szerepelnek egy másik listában (az aktuális bemeneti sorozat vizsgálandó része).

„Kínában, ahol sokan fogyasztják rendszeresen, lehetőség volt hosszas kísérletek folytatására, melyek során bebizonyosodott, hogy azok a férfiak és nők, akik hetente legalább egyszer isznak teát, kevesebb eséllyel betegednek meg végbél, hasnyálmirigyés vastagbél-daganatban, illetve a betegség esetleges kialakulása során lelassul a rákos sejtek burjánzása.” Forrás: <http://www.vital.hu/themes/alter/bio9.htm>

### 11.1.3. Elemhalmazokat tartalmazó gyakori sorozatok

Az előző részben definiált feladat általánosítása, amikor a bemeneti sorozat és a mintahalmaz elemei nem elemek sorozata, hanem elemhalmazoké. Azaz megenegedünk  $\langle AB, B, ABC, E \rangle$  típusú sorozatokat is. Vásárlásoknál például nem csak egy terméket vásárolnak az emberek, hanem termékek egy halmazát.

#### Formális leírás

A bemeneti sorozatok és a mintatér elemei a  $2^J$  felett értelmezett sorozatok, azaz a sorozat elemei az  $J$  részhalmazai. A bemeneti sorozat elemeit szokás *vásárlói sorozatoknak* is hívni, utalva arra, hogy először vásárlói sorozatok esetén került elő a feladat.

Hasonlóan az eddigiekhez a támogatottságot a tartalmazási reláció alapján definiáljuk.

**11.2. definíció.**  $S = \langle I_1, \dots, I_m \rangle$  sorozat tartalmazza  $S' = \langle I'_1, \dots, I'_n \rangle$  sorozatot (jelöléssel  $S' \preceq S$ ), ha léteznek  $j_1 < j_2 < \dots < j_n$  egész számok úgy, hogy  $I'_1 \subseteq I_{j_1}, I'_2 \subseteq I_{j_2}, \dots, I'_n \subseteq I_{j_n}$ .

Ezzel a tartalmazási relációval egy sorozat mérete a sorozat elemeinek méretösszege (tehát például a  $\langle AB, B, ABC, E \rangle$  sorozat mérete 7). A támogatottság, gyakoriság, TID lista, gyakori sorozat fogalmi megegyeznek az eddigiekkel. Feladatunk kinyerni az elemhalmazokból felépülő gyakori sorozatokat [6].

## APRIORIAL

Ismét APRIORI! De minek törjük az agyunkat új módszereken, ha van már módszer, ami jól megoldja a feladatot. Csak a jelöltek előállítását kell tisztázni (pontosabban csak annak első lépését), és készen is vagyunk, mehetünk pihenni (☺). Ennél még kényelmesebb megoldást javasoltak az APRIORIAL kitalálói<sup>1</sup>. Visszavezették ezt a feladatot az előző részben bemutatott APRIORI megoldásra.

Bevezethetjük a gyakori elemhalmaz fogalmát. Az  $I$  elemhalmaz támogatottsága megegyezik azon sorozatok számával, amelyek valamelyik eleme tartalmazza  $I$ -t. Az  $I$  gyakori, ha támogatottsága nagyobb  $min\_supp$ -nál. Nyilvánvaló, hogy gyakori sorozat minden eleme gyakori elemhalmaz. Ezeket a gyakori elemeket tekinthetjük atomi elemeknek, és használhatjuk az előző részben bemutatott algorimust. A gyakori elemhalmazok meghatározásához pedig tetszőleges gyakori elemhalmazt kinyerő algoritmust használhatunk. Ügyelnünk kell azonban arra, hogy a támogatottság meghatározásánál egy sorozat csak eggyel növelheti egy jelölt méretét akkor is ha több elemének része a jelölt.

A feladat visszavezetése az előző feladat APRIORI megoldására nem jelenti azt, hogy ez a megoldás megegyezik az absztrakt APRIORI adaptálásával elemhalmazokat tartalmazó sorozatokra. Az APRIORIAL ugyanis az iterációk során eggyel hosszabb jelöltsorozatokhoz létre, amelyek mérete nem feltétlenül eggyel nagyobb generátoraiknál. Az APRIORIAL nagyobb léptékben halad, így kevesebb iterációs lépést hajt végre, de ugyanakkor jóval több hamis jelöltet generálhat. Ez tehát egy kényelmes, de veszélyes megoldás.

### Időkényszerek Bevezetése

A gyakori sorozatok kinyerését – hasonlóan a gyakori minták kinyeréséhez – a marketingesek igénye keltette életre. A kapott eredmények azonban nem elégtették ki őket, újabb feladattal álltak elő [130] [148]!

- I. **Időkényszerek bevezetése.** A felhasználók gyakran specifikálni akarják a sorozatban található szomszédos elemek között eltelt idő maximális és minimális megengedett értékét. Például nem tulajdonítunk túl nagy jelentőséget annak, ha valaki vesz egy tusfürdőt majd három év múlva egy ugyanolyan márkájú szappant.
- II. **Kosarak definíciójának lazítása.** Sok alkalmazásnál nem számít ha a sorozat egy elemét 2 (vagy több) egymás utáni kosár tartalmazza, ha azok vásárlási ideje bizonyos időablakon belül van. Amennyiben egy vevő 5 perc múlva visszatér az áruházba, akkor valószínű, hogy ezt nem az előző vásárlásának hatására tette (még kicsomagolni sem volt ideje az árut), hanem inkább elfelejtett valamit. Logikus, hogy a két vásárlást összehasonlíthatjuk, és lehet, hogy az összevont kosárhalmazban már megtalálható lesz a sorozat egy eleme, míg az eredeti kettőben külön-külön nem. A tranzakciók definíciójának ilyen lazításánál a sorozatok elemeit kosarak uniója tartalmazhatja, ahol az unióban szereplő kosarak vásárlási idejeinek egy előre megadott időablakon belül kell lenniük.

---

<sup>1</sup>Ez nem meglepő, hiszen sem az ismétlés nélküli jelöltelőállítás sem a támogatottság meghatározása nem triviális feladat. Érdemes elgondolkozni azon, hogy miért nem.

## A Gyakori Sorozat Fogalma Időkényszerek Esetén

Ismét vásárlási sorozatok sorozataként adott a bemenet, de most a vásárlási sorozatok elemei nem pusztán elemhalmazok, hanem olyan párok, amelyek első tagja egy elemhalmaz, második tagja pedig egy időbélyeg. Tehát, legyen ismét  $\mathcal{J} = \{i_1, i_2, \dots, i_m\}$  elemek (vagy termékek) halmaza. Egy vásárlói sorozat most  $\mathcal{T} = \langle \hat{t}_1, \hat{t}_2, \dots, \hat{t}_n \rangle$  tranzakciók sorozata, ahol  $\hat{t}_j = (t_j, TIME_j)$ ,  $t_j \subseteq \mathcal{J}$ ,  $TIME_j \in \mathbb{R}$ . A  $\hat{t} = (t, TIME)$  tranzakció tartalmazza  $I \subseteq \mathcal{J}$  elemhalmazt (jelölésben  $I \subseteq \hat{t}$ ), ha  $I \subseteq t$ . A  $\hat{t}$  tranzakció idejére a továbbiakban  $\hat{t}.TIME$ -al hivatkozunk, tranzakciójára  $\hat{t}.t$ -vel.

A mintakörnyezet definíciója megegyezik a hagyományos, sorozatokat tartalmazó mintakörnyezettel. Mivel ebben az esetben a bemenet és a mintatér elemeinek típusa különbözik (párokból álló sorozat, illetve elemhalmazokból álló sorozat) ezért definiálnunk kell a támogatottságot.

**11.3. definíció.** A  $\mathcal{T} = \langle \hat{t}_1, \hat{t}_2, \dots, \hat{t}_n \rangle$  vásárlói sorozat tartalmazza az  $M = \langle I_1, \dots, I_m \rangle$  mintasorozatot, ha léteznek  $1 \leq l_1 \leq u_1 < l_2 \leq u_2 < \dots < l_m \leq u_m \leq n$  egész számok úgy, hogy

$$I. I_j \subseteq \bigcup_{k=l_j}^{u_j} \hat{t}_k.t, 1 \leq j \leq m,$$

$$II. \hat{t}_{u_i}.TIME - \hat{t}_{l_i}.TIME \leq \text{idő\_ablak}, 1 \leq i \leq m,$$

$$III. \hat{t}_{l_i}.TIME - \hat{t}_{u_{i-1}}.TIME > \text{min\_eltelt\_idő}, 2 \leq i \leq m$$

$$IV. \hat{t}_{u_i}.TIME - \hat{t}_{l_{i-1}}.TIME \leq \text{max\_eltelt\_idő}, 2 \leq i \leq m$$

A fentiekből látszik, hogy a 11.1 definícióval ellentétben tetszőleges elemhalmazt tartalmazó tranzakciók elemhalmazainak uniója tartalmazhat, ahol a tranzakcióknak *idő\_ablak*on belül kell lenniük (2. feltétel).

Ez alapján az  $M$  mintasorozat *támogatottsága* legyen az  $M$ -et tartalmazó vásárlói sorozatok száma. Egy mintasorozat gyakori, ha támogatottsága nem kisebb egy előre megadott támogatottsági küszöbnél (*min\_supp*).

Definiáltunk egy gyakori mintákat kinyerő problémát, amit nyilvánvalóan meg tudunk oldani egy APRIORI algoritmussal. A jelöltek előállításának módja egyezzen meg az APRIORIALL jelöltelőállításának módjával (lévén a mintakörnyezet ugyanaz), a támogatottságok meghatározásánál pedig vegyük figyelembe az időkénszereket, annak érdekében, hogy a helyes támogatottságokat kapjuk. Ha lefuttatnánk így az algoritmus, és vizsgálnánk az eredményt, akkor megdöbbenve vennénk észre, hogy az APRIORI algoritmus nem állította elő az összes gyakori sorozatot. Mi az oka ennek? Bizonyítottuk, hogy az APRIORI teljes, de akkor hol bújik el a hiba? A következő részben eláruljuk a megoldást.

## GSP algoritmus

A GSP (Generalized Sequential Patterns) algoritmus alkalmas olyan sorozatok kinyerésre, amelynél időkénszereket alkalmazhatunk és lazíthatjuk a tranzakciók definícióját. A most következő leírás látszólag teljesen eltér a GSP-t publikáló írástól. Ennek oka az, hogy ragaszkodunk az egységes leíráshoz, amit a 10.1 részben adtunk. Ennek a leírásnak nagy előnye az, hogy ha a problémát meg tudjuk fogalmazni ebben a keretben, akkor a megoldás is azonnal adódik.

Térjünk vissza arra a kérdésre, hogy hol a hiba. Tekintsük a következő mintát:  $M = \langle A, B, C \rangle$ , és nézzük a következő vásárlói sorozatot:  $\mathcal{T} = \langle (A, 1.0), (B, 2.0), (C, 3.0) \rangle$ . Ha  $\text{max\_eltelt\_idő} = 1.5$ , akkor  $\mathcal{T}$  tartalmazza  $M$ -et, de nem tartalmazza annak  $M' = \langle A, C \rangle$  részmintáját, ugyanis az  $A$  és  $C$  elem időbélyege között nagyobb a különbség  $\text{max\_eltelt\_idő}$ -nél. Ezek szerint az  $M$  támogatottsága nagyobb, mint  $M'$  részmintájának támogatottsága. Azaz a fent definiált támogatottsági függvény nem teljesít a támogatottsági függvénnyel szembeni elvárásunkat! Hát ez a hiba, ezért nem fog helyes eredményt adni az APRIORI.

Ahelyett, hogy új problémát definiálnánk és új algoritmus keressünk, próbálkozzunk azzal, hogy átírjuk a feladatot úgy, hogy az új feladat megoldásai megegyezzenek az eredeti feladat megoldásaival, és az új feladat beilleszkedjen egységes keretünkbe. A bemenet, a keresett minta típusa és a támogatottsági függvény adott, így csak a  $\mathcal{MK} = (\mathcal{M}, \prec)$  mintakörnyezet második tagját változtathatjuk meg.

**11.4. definíció.** Az  $M = \langle I_1, \dots, I_n \rangle$  sorozatnak  $M'$  részsorozata (vagy az  $M$  tartalmazza  $M'$ -t,  $M' \prec M$ ), amennyiben az alábbi 3 feltétel közül teljesül valamelyik:

I.  $M'$ -t megkaphatjuk  $M$ -ből  $I_1$  vagy  $I_n$  törlésével.

II.  $M'$ -t megkaphatjuk  $M$ -ből egy legalább 2 elemű  $I_i$  valamely elemének törlésével.

III.  $M'$  részsorozata  $M''$ -nek, ahol  $M''$  részsorozata  $M$ -nek.

Ebben a mintakörnyezetben a  $\|$  függvény ismét a sorozat elemei méretének összegét adja meg. Nézzünk példákat részsorozatokra. Legyen  $M = \langle AB, CD, E, F \rangle$ . Ekkor a  $\langle B, CD, E \rangle$ ,  $\langle AB, C, E, F \rangle$  és a  $\langle C, E \rangle$  mind részsorozatai  $M$ -nek, de a  $\langle AB, CD, F \rangle$  és  $\langle A, E, F \rangle$  sorozatok nem azok.

**11.5. észrevétel.** A fenti tartalmazási relációra nézve a támogatottsági függvény rendelkezik a monotonitás tulajdonságával.

Ha visszatérünk ahhoz a példához, amelyen bemutattuk, hogy az eredeti támogatottsági függvény nem igazi támogatottsági függvény, akkor láthatjuk, hogy nem baj, ha  $\langle A, B, C \rangle$  támogatottsága nagyobb, mint az  $\langle A, C \rangle$  támogatottsága, ugyanis  $\langle A, C \rangle$  nem része az  $\langle A, B, C \rangle$  sorozatnak.

Most már alkalmazhatjuk az APRIORI algoritmust. Ezzel kapcsolatban egyetlen kérdést kell tisztáznunk, mégpedig az, hogyan és mikor állítsunk elő két  $\ell - 1$  elemű gyakori sorozatból  $\ell$  elemű jelöltet.

Két  $k$ -méretű sorozatból ( $S_1, S_2$ ) potenciális jelöltet generálunk akkor, ha törölnénk  $S_1$  első elemének legkisebb sorszámú elemét ugyanazt a sorozatot kapnánk, mintha  $S_2$ -ből az utolsó elem legnagyobb sorszámú elemét törölnénk. A jelölt sorozat az  $S_2$  utolsó elemének legnagyobb sorszámú elemével bővített  $S_1$  sorozat lesz. Az új elem külön elemként fog megjelenni a jelöltben, amennyiben  $S_2$ -ben is külön elem volt, ellenkező esetben  $S_1$  utolsó eleméhez csatoljuk. A fentiek alól kivétel az 1-elemes sorozatok illesztése, ahol az új elemet mind a kétféleképpen fel kell venni, tehát mint új elem, és mint bővítés is. Ezek szerint  $\langle (i) \rangle$  és  $\langle (j) \rangle$  illesztésénél  $\langle (i, j) \rangle$ , és  $\langle (j), (i) \rangle$  is bekerül a jelöltek közé (egyértelmű, hogy mindkét jelöltnek mindkét 1-elemes sorozat részsorozata).

A fenti táblázat egy példát mutat a jelöltek előállítására. Az  $\langle (A, B), (C) \rangle$  sorozatot a  $\langle (B), (C, D) \rangle$  és a  $\langle (B), (C), (E) \rangle$  sorozathoz is illeszthetjük. A többi sorozatot egyetlen

3 méretű gyakorik	4 méretű jelöltek	
	potenciális jel.	jelölt
$\langle(A, B), (C)\rangle$	$\langle(A, B), (C, D)\rangle$	$\langle(A, B), (C, D)\rangle$
$\langle(A, B), (D)\rangle$	$\langle(A, B), (C), (E)\rangle$	
$\langle(A), (C, D)\rangle$		
$\langle(A, C), (E)\rangle$		
$\langle(B), (C, D)\rangle$		
$\langle(B), (C), (E)\rangle$		

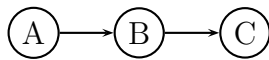
11.1. táblázat. Példa: GSP jelöltgenerálás

másik sorozathoz sem tudjuk illeszteni. Például az  $\langle(A, B), (D)\rangle$  illesztéséhez  $\langle(B), (Dx)\rangle$  vagy  $\langle(B), (D), (x)\rangle$  alakú sorozatnak kéne szerepelnie a gyakorik között, de ilyen nem létezik. A törlési fázisban az  $\langle(A, B), (C), (E)\rangle$  sorozatot töröljük, mert az  $\langle(A), (C), (E)\rangle$  részsorozata nem gyakori.

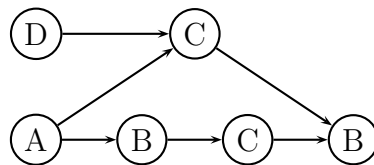
A jelöltek támogatottságának meghatározását nem részletezzük.

#### 11.1.4. Sorozat típusú minta általánosítása

Tetszőleges elemsorozatot ábrázolhatunk egy gráffal. Például a  $\langle A, B, C \rangle$  sorozat megfelelője



a gráf. Az általunk definiált sorozatot, mindig egy nagyon egyszerű gráffal ábrázolnánk, ami egy irányított, körmentes, címkézett út. Mi sem természetesebb, hogy a sorozat általánosítása egy olyan valami, amit teljesen általános irányított, körmentes, címkézett gráffal ábrázolunk. Például lehet egy általános mintához tartozó gráf a 11.1 ábrán látható.



11.1. ábra. Példa: sorozat általánosítása

Érezzük, hogy ezt a mintát tartalmazzák például a  $\langle A, D, C, B, C, B \rangle$  vagy az  $\langle E, D, A, B, B, CC, B \rangle$  sorozatok, de nem tartalmazzák a  $\langle A, D, C, C, B, B \rangle$  illetve a  $\langle A, D, B, C, B, C \rangle$  sorozatok.

Ugyanezt az általános leírást kapnánk, ha egy sorozatra nem mint út tekintünk, hanem mint olyan halmazon értelmezett teljes rendezés, amelynek elemei azonosító, elem párok. A teljes rendezés általánosítása ugyanis a részben rendezés, amit körmentes, irányított gráffal szokás ábrázolni.

Nézzük formálisan. Legyen  $\mathcal{J}$ , illetve  $TID$  elemek és azonosítók halmaza. A mintatér elemei ekkor  $(tid, i)$  párokon értelmezett részben rendezés, ahol  $tid \in TID, i \in \mathcal{J}$ . A  $tid$  címkéjén az  $i$



elemet értjük.

**11.6. definíció.** Az  $m = (\{(tid_1, i_1), \dots, (tid_m, i_m)\}, \leq)$  minta tartalmazza az  $m' = (\{(tid'_1, i'_1), \dots, (tid'_n, i'_n)\}, \leq')$  mintát (jelöléssel  $m' \preceq m$ ), ha létezik  $f : \{tid'_1, \dots, tid'_m\} \rightarrow \{tid_1, \dots, tid_n\}$  injektív függvény úgy, hogy  $tid'_j$  címkéje megegyezik  $f(tid'_j)$  címkéjével ( $1 \leq j \leq m$ ), és  $(tid'_k, i'_k) \leq' (tid'_l, i'_l)$  esetén  $(f(tid'_k), i'_k) \leq (f(tid'_l), i'_l)$  is teljesül minden  $(1 \leq k, l \leq m)$  indexre.

Az általános minta keresésénél a bemenet  $J$  felett értelmezett elemsorozatok sorozataként adott. Egy bemeneti sorozat tulajdonképpen felfogható általános mintának, ahol a rendezés teljes rendezés. Egy minta *támogatottsága* megegyezik azon sorozatok számával, amelyek tartalmazzák a mintát.

## 11.2. Gyakori bool formulák

Legyenek a bemenet  $n$ -esek halmaza. A felhasználó megad predikátumokat, amelyek a bemenet elemein vannak értelmezve, és akár többváltozósak is lehetnek. A mintatér elemei ezen predikátumokon értelmezett bool formula. A formulában megengedjük az *és*, *vagy* illetve *negáció* operátorokat [86], de hatékonysági okok miatt célszerű csak a diszjunktív normál formulákra szorítkozni.

Nézzünk példákat. Tegyük fel, hogy egy telekommunikációs hálózatban egy eseménynek 4 attribútuma van: típus, modul, szint, időbélyeg. Az első megadja egy riasztás típusát, a második a modult, ami a riasztást küldte, a harmadik a riasztás erősségét, a negyedik pedig riasztás időpontját. Ebben a környezetben mintára lehet példa az alábbi:

$$p(x,y) = x.típus=2356 \wedge y.típus=7401 \wedge x.time \leq y.time \wedge x.modul=y.modul$$

ami azt jelenti, hogy egy 2356 és egy 7401 típusú riasztás érkezett ebben a sorrendben ugyanabból a modulból. Bevezethetjük például a *szomszédja* – modul attribútumra vonatkozó – kétváltozós predikátumot, ha úgy gondoljuk hogy fontos lehet ennek vizsgálata. Ekkor a

$$p'(x,y) = x.típus=2356 \wedge y.típus=7401 \wedge szomszédja(x.modul=y.modul)$$

azt fejezi ki hogy a 2356 és 7401 típusú riasztások szomszédos modulból érkeztek.

A  $p(x_1, x_2, \dots, x_m)$   $m$  változós minta *illeszkedik* az  $\langle S_1, S_2, \dots, S_v \rangle$  sorozatra, ha léteznek  $i_1, i_2, \dots, i_m$  egészek úgy, hogy  $p(S_{i_1}, S_{i_2}, \dots, S_{i_m})$  igaz értéket ad.

## 11.3. Gyakori epizódok

Az eddig részekben sok elemhalmaz, sorozat volt adva, és kerestük a gyakori mintákat. Ezek a minták általánosan érvényes információt adtak: az adott vásárlói minta sok vásárlóra

„Kutatási eredmények igazolják, hogy a csoportban működőknek teljesebb szülésélményben van részük, körükben alacsonyabb a koraszülések száma, és a babák súlya is nagyobb az egyéni felkészülésben részesülőknél.”  
Forrás: Baba Patika X. évfolyam 10. szám, 56. oldal 2007. október

jellemző. Ha a sok sorozatból kiválasztunk egyet és azt elemezzük, akkor az adott sorozatra jellemző információt nyerünk ki. Megtudhatjuk például, mi jellemző az adott ügyfélre, amit felhasználhatunk akkor, amikor személyre szabott ajánlatot szeretnénk tenni (például azért mert az ügyfél elégedetlen szolgáltatásainkkal, és vissza akarjuk szerezni bizalmát).

Epizódkutatásról beszélünk, ha egyetlen sorozat van adva, és ebben keressük a gyakran előforduló mintákat[87, 88]. Az epizódkutatásnak egyik fontos területe a telekommunikációs rendszerek vizsgálata. Az olyan epizódok feltárása, amelyben riasztás is előfordul, alkalmas lehet a riasztás okának felderítésére, vagy előrejelzésére.

Nem vezetünk be új típusú mintát, tehát most is elemhalmazokat, sorozatokat keresünk, de a formalizmus könnyen általánosítható elemhalmazokat tartalmazó sorozatokra, vagy általános mintára is. A támogatottsági függvény lesz új, ami abból fakad, hogy egyetlen bemeneti sorozat van adva.

### 11.3.1. A támogatottság definíciója

Legyen  $\mathcal{J}$  elemek (items) halmaza. A bemenet az  $\mathcal{J}$  felett értelmezett sorozat.

$$\text{bemenet} : \mathcal{S} = \langle i_1, i_2, \dots, i_n \rangle,$$

ahol  $i_k \in \mathcal{J}$  minden  $k$ -re,

**11.7. definíció.** Az  $\mathcal{S} = \langle i_1, i_2, \dots, i_n \rangle$  sorozatnak a  $\langle i_j, i_{j+1}, \dots, i_{j+w-1} \rangle$  sorozat egy  $w$  elem széles összefüggő részsorozata, ha  $1 \leq j \leq n+1-w$ .

Ha  $w < n$ , akkor valódi összefüggő részsorozatról beszélünk.

Legyen adva  $\mathcal{MK}$  mintakörnyezet, és értelmezzük valahogy a  $\tau$  anti-monoton illeszkedési predikátumot.  $\tau_{\mathcal{S}}(m)$  igaz értéket ad, ha az  $m$  minta illeszkedik az  $\mathcal{S}$  sorozatra.

**11.8. definíció.** A  $m$  minta minimálisan illeszkedik az  $\mathcal{S}$  sorozatra, ha  $\mathcal{S}$ -nek nincsen olyan valódi összefüggő részsorozata, amelyre illeszkedik  $m$ .

Ha például a mintatér elemei  $\mathcal{J}$  részhalmazai, akkor a  $\mathcal{S} = \langle i_1, i_2, \dots, i_n \rangle$  sorozatra illeszkedik az  $I$  halmaz, amennyiben minden  $i \in I$ -hez létezik  $1 \leq j \leq n$ , amelyre  $i = i_j$ . Elemssorozat típusú minta esetén  $S$  akkor illeszkedik az  $\mathcal{S}$  sorozatra, ha  $S$  részsorozata  $\mathcal{S}$ -nek, ahol a részsorozat definíciója megegyezik a 11.1 részben megadottal.

Két különböző támogatottsági definíció terjedt el.

**11.9. definíció.** Legyen  $\mathcal{S}$  bemeneti sorozat,  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet és  $\tau$  anti-monoton illeszkedési predikátum. Az  $m \in \mathcal{M}$  minta támogatottsága megegyezik

- I.  $S$  azon összefüggő részsorozatainak számával, amelyekre  $m$  minimálisan illeszkedik.
- II.  $S$  azon  $w$  széles részsorozatainak számával, amelyekre  $m$  illeszkedik. Itt  $w$  előre megadott konstans.

Ha a támogatottság így van definiálva, akkor a mintatér elemeit *epizódoknak* nevezzük. Egy epizód *gyakori*, ha támogatottsága nem kisebb egy előre megadott korlátnál, amit általában *min\_supp*-al jelölünk.

Epizódkutatásnál adott  $\mathcal{S}$  bemeneti sorozat  $\mathcal{MK} = (\mathcal{M}, \preceq)$  mintakörnyezet (esetleg  $w$ ) és  $\tau$  illeszkedési predikátum, célunk megtalálni a gyakori epizódokat.

### 11.3.2. APRIORI

Az illeszkedési predikátum anti-monoton tulajdonságából következik a támogatottság anti-monotonitása, amiből jön, hogy gyakori epizód minden részepizódja gyakori. Mi sem természetesebb, hogy a gyakori epizódok kinyeréséhez az APRIORI algoritmust használjuk. Az jelöltek-előállítás és a gyakori epizódok kiválogatása ugyanaz, mint a támogatottságot a régi módszerrel definiálnánk (lásd 4.2 11.1.2 rész). Egyedül a támogatottság meghatározásán kell változtatnunk. A következőkben feltesszük, hogy a támogatottságot a második definíció szerint értjük ( $w$  széles ablakok száma).

A támogatottság meghatározásának egy butuska módszere lenne, ha az eseménysorozaton egyszerűen végigmasírozva minden összefüggő részsorozatnál meghatároznánk, hogy tartalmazza-e az egyes jelölt epizódokat. Hatékonyabb algoritmushoz juthatunk, ha felhasználjuk azt, hogy szomszédos sorozatok között pontosan két elem eltérés van. Vizsgáljuk meg az első sorozatot, majd nézzük az eggyel utána következőt, és így tovább addig, amíg el nem érjük az utolsót. Mintha egy ablakot tolnánk végig a sorozaton.

Vezetjük be a következő változókat. Minden  $i$  elemhez tartozik:

- $i$ .számláló, ami megadja, hogy a jelenlegi összefüggő részsorozatba hányszor fordul elő az  $i$  elem.
- $i$ .epizódjai lista, amelyben az  $i$  elemet tartalmazó epizódok találhatók.

Epizódjelöltekhez pedig a következőkre lesz szükségünk:

- $j$ .kezdeti\_index: annak a legkorábbi elemnek az indexe, amely után minden részsorozatban előfordult az epizód egészen a jelenlegi részsorozatig.
- $j$ .számláló, ami megadja, hogy hány kezdeti\_index előtti összefüggő részsorozatban fordult elő  $j$  jelölt. A bemenet feldolgozása után  $e$  változó fogja tartalmazni a jelölt támogatottságát.
- $j$ .hiányzás egész szám adja meg, hogy  $j$  elemei közül hány nem található a jelenlegi összefüggő részsorozatban. Nyilvánvaló, hogy ha  $\varphi$  előfordul a jelenlegi részsorozatban, akkor  $j$ .hiányzás=0.

„Nemzetközi tanulmányok alapján elmondhatjuk, hogy a magzati fejlődési rendellenességek ( az agykoponya hiánya, nyitott hátgerinc), továbbá a szív és a vese rendellenességei megelőzhetőek, ha a terhes kismama a fogamzást megelőzően legalább négy hétig, majd a terhesség első három hónapjában folsav tartalmú készítményt szed.” Forrás: Baba Patika X. évfolyam 10. szám, 48. oldal, 2007. október:

#### Elemhalmazok támogatottságának meghatározása

Amikor lépünk a következő részsorozatra, akkor egy új elem kerül bele az ablakba, amit jelöljünk  $i_{új}$ -al, ugyanakkor egy elem eltűnik a sorozatból, ezt pedig jelöljük  $i_{rég}$ -vel.

Egy elem kilépésének következtében epizódok is kiléphetnek.  $i_{rég}$ .számláló segítségével megállapíthatjuk, hogy maradt-e még ilyen elem az ablakban, mert ha igen, akkor az eddig tartalmazott epizódokat az új ablak is tartalmazza. Ha nem maradt, akkor  $i$ .epizódjai és epizódok hiányzás számlálója alapján megkaphatjuk azon epizódokat, amelyek kiléptek a sorozatból. Ezek előfordulásának értékét kell növelni. Ebben segítségünkre van a kezdeti\_index érték, ami

```

irégi.számláló ← irégi.számláló-1;
if( irégi.számláló = 0)
  forall j in irégi.epizódjai
  {
    j.hiányzás ← j.hiányzás+1;
    if( j.hiányzás = 1) then
      j.számláló ← j.számláló + j.kezdeti_index-jelenlegi_index;
  }

```

11.2. ábra. régi elem kilépése

megadja, hogy mióta van jelen az epizód a sorozatokban. Az algoritmus pseudokódja az alábbi ábrán látható.

Könnyű kitalálni ezek alapján, hogy mit kell tenni egy új elem belépésénél. Ha az új elem még nem szerepelt az ablakban, akkor végig kell nézni az új elemet tartalmazó epizódokat. Azon epizód kezdeti indexét kell a jelenlegi indexre beállítani, amelyekből csak ez az egyetlen elem hiányzott (11.3 ábra).

```

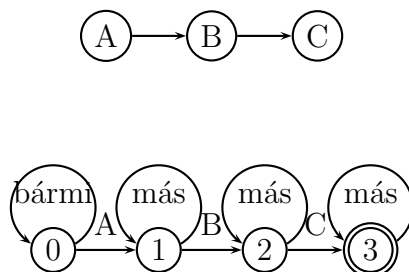
iúj.számláló ← iúj.számláló+1;
if( eúj.számláló = 1 )
  forall j in iúj.epizódjai
  {
    j.hiányzás ← j.hiányzás-1;
    if j.hiányzás=0 then
      j.kezdeti_index ← jelenlegi_index;
  }

```

11.3. ábra. új elem belépése

### Elem sorozatok támogatottságának meghatározása

Az elem sorozatok felismerése determinisztikus véges automatákkal történik, amelyek az egyes elem sorozatokat fogadják el. Az epizód alapján az automata előállítás egyszerű, az alábbi ábra erre mutat példát.



A teljes elemsorozatot egyesével olvassuk végig az első elemtől kezdve. Ha valamely epizód első eleme megegyezik az éppen olvasott elemmel, akkor új automatát hozunk létre. Ha ez az elem elhagyja az ablakot, akkor töröljük az automatát. Amikor egy automata elfogadó állapotba lép (jelezve, hogy az epizód megtalálható az ablakban), és nincs ehhez az epizódhoz tartozó másik – szintén elfogadó állapotban lévő – automata, akkor *kezdeti\_index* felveszi az aktuális elem indexét. Amennyiben egy elfogadó állapotban lévő automatát törölünk, és nincs más, ugyanahhoz az epizódhoz tartozó elfogadó állapotú automata, akkor a *kezdeti\_index* alapján növeljük az epizód *számlálóját*, hiszen tudjuk, hogy az epizód a kezdeti idő utáni összes részsorozatban megtalálható volt egészen az aktuális részsorozat előtti részsorozatig.

Vegyük észre, hogy felesleges adott epizódhoz tartozó, ugyanabban az állapotban lévő automatákat többszörösen tárolni: elég azt ismernem, amelyik utoljára lépett be ebbe az állapotba, hiszen ez fog utoljára távozni. Emiatt  $j$  jelölthöz maximum  $j$  darab automatára van szükség.

Egy új elem vizsgálatakor nem kell az összes automatánál megnéznünk, hogy új állapotba léphetnek-e, mert az elem *epizódjai* listájában megtalálható az őt tartalmazó összes epizód.

Az előzőekben ismertetett epizódkutatási algoritmus olyan adatbányászati problémára adott megoldást, ami az ipari életben merült fel, és hagyományos eszközök nem tudták kezelni. Az algoritmus telekommunikációs hálózatok riasztásáról eddig nem ismert, az adatokban rejlő információt adott a rendszert üzemeltető szakembereknek. Erről bővebben a [76][81] [83][82][58] cikkekben olvashatunk.

## 12. fejezet

# Gyakori fák és feszített részgráfok

Amikor gyakori elemhalmazokat kerestünk, akkor azt néztük, hogy mely elemek fordulnak elő együtt gyakran. Sorozatok keresésénél ennél továbbléptünk, és azt is néztük, hogy milyen sorrendben fordulnak elő az elemek, azaz melyek elemek előznek meg más elemeket. Ez már egy bonyolultabb kapcsolat. Még általánosabb kapcsolatok leírására szolgálnak a gráfok: a felhasználási terület entitásainak felelnek meg a gráf csúcsai vagy a csúcsainak címkéi, amelyeket él köt össze, amennyiben van közöttük kapcsolat. A kapcsolat típusát, sőt az entitások jellemzőit is kezelni tudjuk, amennyiben a gráf csúcsai és élei címkézettek.

Ezt a fejezetet először a gráf egy speciális esetével a gyökeres fák vizsgálatával kezdjük, majd rátérünk a gyakori általános gráfok keresésére. Ellentétben az elemhalmazokkal vagy a sorozatokkal a támogatottságot megadó illeszkedési predikátumot a gráfnál többféleképpen definiálhatjuk: részgráf, feszített részgráf, topologikus részgráf. Ez tovább bővíti a megoldandó feladatok körét.

### 12.1. Az izomorfia problémája

Ha gráfokra gondolunk, akkor szemünk előtt vonalakkal – irányított gráfok esetében nyilakkal – összekötött pontok jelennek meg. Címkézett gráfnál a pontokon és/vagy az éleken címkék, általában számok szerepelnek. Különböző pontoknak lehetnek azonos címkéi. Egy ilyen pontokat és vonalakat tartalmazó rajz a gráf egy lehetséges ábrázolása. Matematikailag egy gráf egy páros, amelynek első eleme egy alaphalmaz, a második eleme ezen alaphalmazon értelmezett bináris reláció.

Különböző gráfoknak lehet azonos a rajzuk. Például a  $G_1 = (\{a, b\}, \{a, b\})$  és a  $G_2 = (\{a, b\}, \{b, a\})$  gráfok rajza ugyanaz lesz: az egyik pontból egy nyíl indul a másik pontba. Ugyanúgy azonos ábrát készítenénk, ha az egyetlen élnek címkéje lenne, vagy a két pontnak ugyanaz lenne a címkéje. Az alkalmazások többségében a gráf rajza, topológiája továbbá a címkék az érdekesek és nem az, hogy a pontokat hogyan azonosítjuk annak érdekében, hogy a bináris relációt fel tudjuk írni. Ezen alkalmazásokban nem akarjuk megkülönböztetni az *izomorf* gráfokat (pontos definíciót lásd alapfogalmak gráfelmélet részében). Ez a helyzet áll fenn, például amikor kémiai vegyületeket vizsgálunk. Itt a gráf címkéi jellemzik az atomot (esetleg még további információt, pl. töltést) az él a kötést, az él címkéi pedig a kötés típusát (egyszeres kötés, kétszeres kötés, aromás kötés). Amikor gyakori gráfokat keresünk, akkor mindenképpen el kell döntenünk, hogy az izomorf gráfokat megkülönböztetjük, vagy nem. Mielőtt

rátérünk a gyakori gráfok keresésére járjuk egy kicsit körül az izomorfia kérdését.

Két gráf izomorfijának eldöntésére nem ismerünk polinom idejű algoritmust, sőt azt sem tudjuk, hogy a feladat NP-teljes-e. Hasonló feladat a *részgráf izomorfia* kérdése, ahol azt kell eldönteni, hogy egy adott gráf izomorf-e egy másik gráf valamely részgráfjával. Ez a feladat NP-teljes. Ha ugyanis az egyik gráf egy  $k$ -csúcsú teljes gráf, akkor a feladat az, hogy keressünk egy gráfban  $k$ -csúcsú klikket, ami bizonyítottan NP-teljes. Szerencsére kisebb méretű gráfok esetében az izomorfia eldöntése egyszerűbb algoritmusokkal is megoldható elfogadható időn. A két legismertebb részgráf izomorfia eldöntő algoritmus Ullmanntól a backtracking [138] és B.D.McKaytól a Nauty [91].

A gráf izomorfia eldöntő módszerek a csúcsok *invariánsait* használják. Az invariáns tulajdonképpen egy tulajdonság. Például invariáns a csúcs címkéje, fokszáma, illetve irányított gráfok esetében a befok és a kifok is két invariáns. Amennyiben a  $G_1, G_2$  gráfok a  $\phi$  bijekció alapján izomorfak, akkor az  $u$  csúcs minden invariánsa megegyezik a  $\phi(u)$  csúcs megfelelő invariánsaival a  $G_1$  minden  $u$  csúcsára. Ez tehát egy szükséges feltétel: az  $u$  csúcsához csak azt a csúcsot rendelheti a bijekció, amelynek invariánsai páronként azonosak az  $u$  invariánsaival.

Az izomorfia eldöntésének naív módszere az lenne, ha az összes bijekciót megvizsgálnánk egyesével. Egy bijekció a csúcsoknak egy permutációja, így  $n$  csúcsú gráfok esetében  $n!$  bijekció létezik. Csökkenthetjük ezt a számot az invariánsok segítségével. Osszuk részekre a csúcsokat. Egy csoportba azon csúcsok kerüljenek, amelyeknek páronként minden invariánsuk azonos. Nyilvánvaló, hogy az olyan bijekciókat kell megvizsgálni, amelyek csak ugyanazon invariánsok által leírt csoportba tartoznak. Ha az invariánsokkal a  $V$  csúcsokat szétosztottuk a  $V_1, \dots, V_k$  csoportokba, akkor a szóba jövő bijekciók száma  $\prod_{i=1}^k |V_i|!$ -re csökken. Minél több csoportot hoznak létre az invariánsok annál többet nyerünk ezzel az egyszerű trükkel. Az invariánsok nem csökkentik asszimptotikusan a számítás komplexitását. Ha például a gráf reguláris és a csúcsoknak nincsenek címkéjük, akkor minden csúcs azonos csoportba kerül, azaz nem nyerünk a trükkel semmit.

Eddigi ismereteink alapján elmondhatjuk, hogy minél bonyolultabb gyakori mintát keresünk, annál nehezebb a feladat és annál erőforrás-igényesebbek a megoldó algoritmusok. A címke nélküli gráfok egy általánosítása a címkézett gráfok, így azt várjuk, hogy címkézett gráfokhoz még több számítást kell majd végezni. Az előbb bemutatott módszer szerencsére az ellenezőjét állítja, hiszen a címke egy invariáns, ami újabb csoportokat hozhat létre. Sőt minél több a címke, annál több a csoport és annál gyorsabban döntjük el, hogy két gráf izomorf-e.

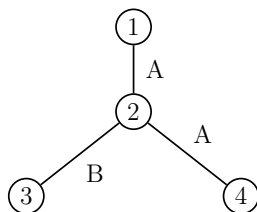
A gráf izomorfijából született probléma a gráfok *kanonikus kódolásának* problémája.

**12.1. definíció.** *A gráfok kanonikus kódolása (vagy kanonikus címkézése) egy olyan kódolás, amely az izomorf gráfokhoz és csak azokhoz azonos kódsorozatot rendel.*

Nyilvánvaló, hogy egy kanonikus kódolás előállítása ugyanolyan nehéz feladat, mint két gráf izomorfijának eldöntése, hiszen két gráf izomorf, ha kanonikus kódjaik megegyeznek. Például egy egyszerű kanonikus kód az, amit úgy kapunk, hogy a gráf szomszédossági mátrix oszlopai permutálásai közül kiválasztjuk azt, amely elemeit valamely rögzített sorrendben egymás után írva a legkisebbet kapjuk egy előre definiált lexikografikus rendezés szerint.

„A legújabb kutatások szerint bizonyos vitaminok képesek a hibás gének okozta fejlődési rendellenességek kivédésére.” Forrás: Baba Patika X. évfolyam 10. szám, 44. oldal, 2007. október

A szomszédossági mátrix alapú kanonikus kód előállításához szintén az invariánsokat célszerű használni. Ezáltal az oszlopok összes permutációjához tartozó kódok kiértékelése helyett egy oszlopot csak a saját csoportján belüli oszlopokkal kell permutálni.



12.1. ábra. Példa kanonikus kódolásra

Nézzük példaként a 12.1 ábrán látható csúcs- és élcímkézett gráfot (a csúcsokban szereplő számok a csúcsok azonosítói). Legyen  $cimke(1) = e$ ,  $cimke(2) = e$ ,  $cimke(3) = e$ ,  $cimke(4) = f$ . A csúcsok címkéi szerint két csoportot hozunk létre. Ha figyelembe vesszük a foksámot is, akkor a nagyobb csoportot két részre osztjuk ( $\{1,3\}$ ,  $\{2\}$ ,  $\{4\}$ ). A  $4!=24$  kombináció helyett csak  $2!=2$  permutációt kell kiértékelnünk, ami alapján megkapjuk a kanonikus kódot:  $\langle e000A0e0A00fBAABe \rangle$  lesz, ha a címkéken az abc szerinti rendezést vesszük és a 0 minden betűt megelőző.

## 12.2. A gyakori gráf fogalma

Annak alapján, hogy az izomorf gráfokat megkülönböztetjük, vagy nem a gyakori gráfok kinyerésének feladatát két csoportra osztjuk. Legyen  $\mathcal{V} = \{v_1, v_2, \dots, v_m\}$  csúcsok halmaza. A mintakörnyezet ekkor az  $MK = (\{G_1 = (V_1, E_1), G_2 = (V_2, E_2), \dots\}, \preceq)$  pár, ahol  $V_i \subseteq \mathcal{V}$ , minden gráf összefüggő és  $G_i \preceq G_j$ , amennyiben  $G_i$  a  $G_j$ -nek részgráfja. A bemenet szintén olyan gráfok sorozata, amelyek csúcshalmaza  $\mathcal{V}$ -nek részalmazai. A gráfok csúcsainak és/vagy éleinek lehetnek címkéi. A továbbiakban az élek és csúcsok címkéjét a  $c_E$  és  $c_V$  függvények adják meg. Az általánosság megsértése nélkül feltehetjük, hogy a címkék pozitív egész számok.

A támogatottságot illeszkedési predikátum alapján definiáljuk. Attól függően, hogy a csúcsok értéke fontos, vagy csa a címkéjük, az illeszkedést kétféleképpen definiálhatjuk:  $G'$  gráf illeszkedik a  $G$  bemeneti gráfra, ha

- $G'$  részgráfja/feszített részgráfja/topologikus részgráfja  $G$ -nek,
- létezik  $G$ -nek olyan részgráfja/feszített részgráfja/topologikus részgráfja, amely izomorf  $G'$ -vel.

A fenti lehetőségek közül az alkalmazási terület ismerete alapján választhatunk.

A topologikus részgráf fogalma nem tartozik az alapfogalmak közé, így ennek jelentését meg kell adnunk.

**12.2. definíció.** A  $G' = (V', E')$  gráf a  $G = (V, E)$  gráf topologikus részgráfja, ha  $V' \subseteq V$  és  $(u, v) \in E'$  akkor és csak akkor, ha  $u$ -ból vezet út  $v$ -be a  $G$  gráfban.

Gráfok esetében használt fogalom a *súlyozott támogatottság*, melynek kiszámításához illeszkedési predikátum helyett illeszkedési függvényt használunk. Az illeszkedési függvény megadja a bemeneti gráf különböző részgráfjainak/feszített részgráfjainak/topologikus részgráfjainak



számát, amely azonosak/izomorfak a mintagráffal. A  $G$  gráf súlyozott támogatottsága a bemeneti elemeken vett illeszkedési függvény összege.

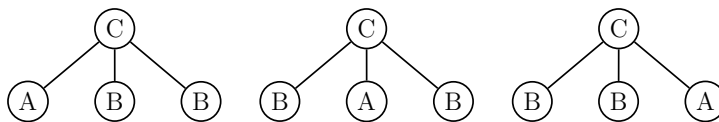
Mielőtt rátérnénk az általános eset tárgyalására nézzük meg, hogyan lehet kinyerni a gyakori címkézett fákat.

### 12.3. gyakori gyökeres fák

Ebben a részben feltesszük, hogy a mintatér és a bemeneti sorozat elemei csúscímkézett gyökeres fák. Egy fa mérete a csúcsainak számát adja meg. Csak a címkék fontosak, ezért az illeszkedési predikátumnak a második fajtáját használjuk: akkor illeszkedik egy mintafa egy bementi fára, ha annak létezik olyan topologikus részgráfja, amellyel a mintafa izomorf.

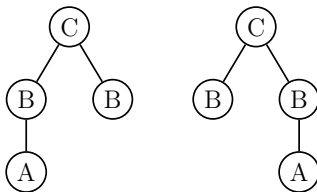
A gyakori fák kinyerése hasznos a bioinformatikában, a webelemzésnél, a félig strukturált adatok vizsgálatánál stb. Az egyik legszemléletesebb felhasználási terület a webes szokások elemzése. Gyakori elemhalmaz-kinyerő algoritmussal csak azt tudnánk megállapítani, hogy melyek a gyakran látogatott oldalak. Ha gyakori szekvenciákat keresünk, akkor megtudhatjuk, hogy az emberek milyen sorrendben látogatnak el az oldalakra leggyakrabban. Sokkal élethűbb és hasznosabb információt kapunk, ha a weboldalakból felépített gyakori fákat (vagy erdőket) keresünk. Egy internetező viselkedését egy fa jobban reprezentálja, mint egy sorozat.

*Rendezett* gyökeres fáknál további feltétel, hogy az egy csúcsból kiinduló élek a gyerek csúcs címkéje szerint rendezve legyenek. Ez tulajdonképpen egy átmenet afelé, hogy az izomorf gráfokat ne különböztessük meg, vagy másként szólva a mintatérben ne legyenek izomorf gráfok. Ha a címkék rendezése abc szerint történik, akkor például a következő 3 fa közül csak az első tartozik a mintatér elemei közé.



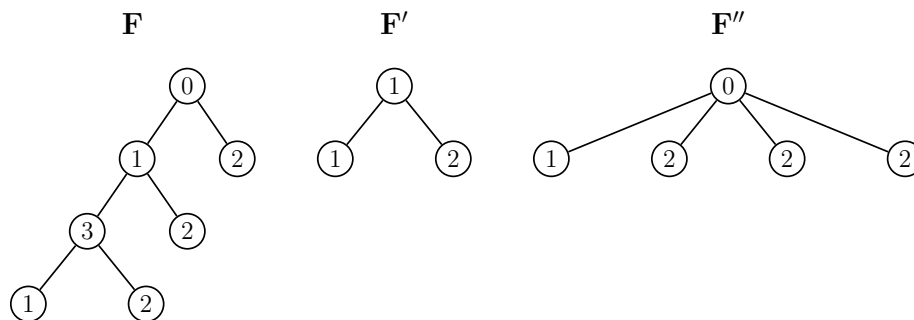
12.2. ábra. Példa: rendezés nélküli, címkézett, gyökeres fák

A rendezettség nem biztosítja azt, hogy a mintatérben ne legyenek izomorf fák. Például a következő ábrán látható két rendezett fa izomorf egymással, és mindkettő a mintatérnek elemei.



12.3. ábra. Példa: izomorf rendezett, gyökeres fák

Mivel az illeszkedés során izomorf részfákat keresünk, ezért feltehetjük, hogy a fa csúcsai természetes számok, és az  $i$  csúcs azt jelenti, hogy a csúcsot az  $i$ -edik lépésben látogatjuk meg a gráf preorder, mélységi bejárása során. Legyen a gyöker csúcs a 0. Az  $F$  fa  $i$  csúcsjának címkéjét  $c_F(i)$ -vel a szülőjét pedig  $szulo_F(i)$ -vel jelöljük. Elhagyjuk az  $F$  alsó indexet azokban az esetekben, ahol ez nem okozhat félreértést.



12.4. ábra. Példa: gyökeres részfák tartalmazására

A 12.4 ábrán egy példát láthatunk illeszkedésre (topologikus részfára). A fák csúcsaiba írt számok a csúcsok címkéit jelölik. Az  $F'$  és  $F''$  is illeszkedik a  $F$  fára.

Amennyiben egy gráf ritka (kevés élt tartalmaz), akkor azt szomszédossági listával (lásd alapfogalmak 2.3 rész) célszerű leírni. Fák esetében a *címkeláncok* még kevesebb helyet igényelnek a memóriából. A címkeláncot úgy kapjuk meg, hogy bejárjuk a fát preorder, mélységi bejárás szerint, és amikor új csúcsba lépünk akkor hozzáírjuk az eddigi címkelánchoz az új csúcs címkéjét. Amikor visszalépünk, akkor egy speciális címkét (\*) írunk. Például az előző ábrán  $F$  címkelánca:  $\mathcal{F}=0,1,3,1,*,2,*,*,2,*,*,2,*$  és  $\mathcal{F}'=1,1,*,2,*$ . Címkesorozatnak hívjuk és  $l(F)$ -vel jelöljük azt a sorozatot, amit a  $F$  gráf címkeláncából kapunk meg, ha elhagyjuk a \* szimbólumot. Nyilvánvaló, hogy a címkesorozat – a címkelánccal ellentétben – nem őrzi meg a fa topológiáját.

Hasonlóan a gyakori elemhalmazok kereséséhez most is megkülönböztetünk horizontális és vertikális adatábrázolási módot. Horizontális ábrázolásnál a bemenet gráfok leírásának (például címkelánc) sorozata. Vertikális tárolásnál minden címkéhez tartozik egy párokból álló sorozat. Az  $i$  címkéhez tartozó sorozatban a  $(j, k)$  pár azt jelenti, hogy a  $j$ -edik bemeneti gráf preorder bejárás szerinti  $k$ -edik csúcs címkéje  $i$ .

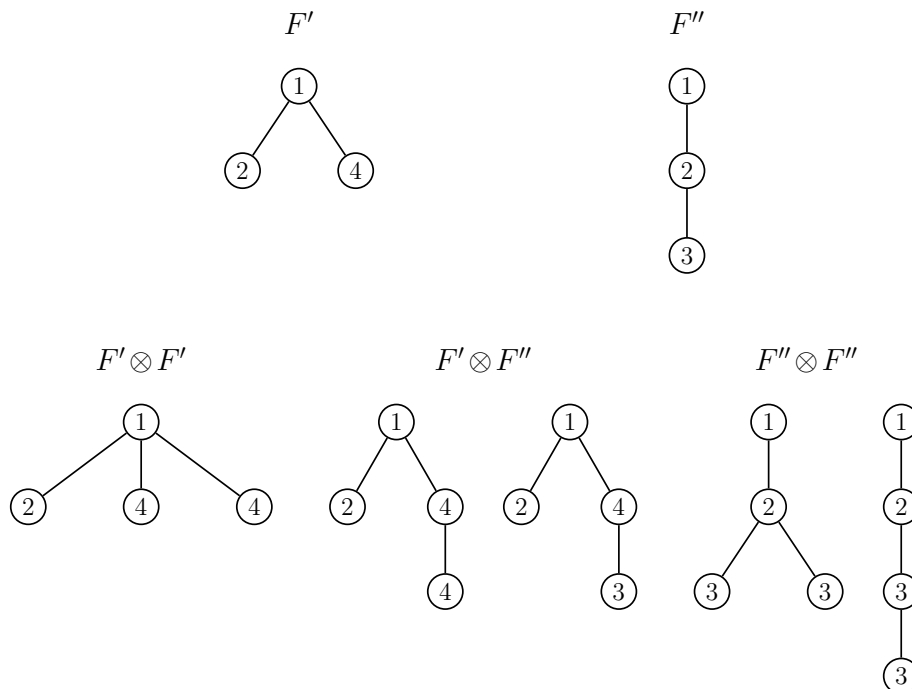
### 12.3.1. TreeMinerH

A TreeMinerH [145] az APRIORI sémára épül (annak ellenére, hogy Zaki publikálta). Nézzük meg, hogyan állítjuk elő a jelölteket és hogyan határozzuk meg a támogatottságukat.

#### Jelöltek előállítás

Egy  $\ell$ -elemű jelöltet két  $(\ell - 1)$ -elemű gyakori fa ( $F'$  és  $F''$ ) illesztésével (jelölésben:  $\otimes$ ) kapjuk meg. Hasonlóan az eddigiekhez a két  $(\ell - 1)$ -elemű fa csak a legnagyobb elemükben különböznek, amely esetünkben azt jelenti, hogy ha elhagynánk a legnagyobb csúcsot (és a hozzá tartozó élt), akkor ugyanazt a fát kapnánk. Az általánosság megsértése nélkül feltehetjük, hogy  $szulo_{F'}(\ell - 1) \leq szulo_{F''}(\ell - 1)$ . A potenciális jelölt a  $G'$  gráf egy csúccsal való bővítése lesz, ahol az új csúcs címkéje a  $c_{F''}(\ell - 1)$  lesz.

Kételemű (egy élt tartalmazó) jelöltek előállításánál nincs sok választás: az új élt egyetlen helyre illeszthetjük. Ha  $\ell > 2$ , akkor két esetet kell megkülönböztetnünk. Az első esetben  $szulo_{F'}(\ell - 1) = szulo_{F''}(\ell - 1)$ . Ekkor két jelöltet állítunk elő. Az elsőben az új élt a  $szulo(\ell - 1)$  csúcsához, a másodikban a  $szulo(\ell - 1) + 1$  csúcsához kapcsoljuk. Ha  $szulo_{F'}(\ell - 1) < szulo_{F''}(\ell - 1)$ , akkor az új élt a  $szulo_{F''}(\ell - 1)$ -hez csatoljuk. Jelölt-előállításra mutat példát a következő ábra:



12.5. ábra. Példa jelöltek előállítására

Szokásos módon a jelöltek előállításának második lépésében minden  $\ell - 1$  elemű részfat ellenőrizni kell, hogy gyakori-e.

### Támogatottság meghatározása

Az egy- és kételemű fák támogatottságát vektorral, illetve tömbbel célszerű meghatározni. A vektor  $i$ -edik eleme tárolja a támogatottságát az  $i$ -edik címkének. A tömb  $i$ -edik sorának  $j$ -edik eleme tárolja a támogatottságát annak a kételemű fának, amelyben a gyökér címkéje az  $i$ -edik gyakori címke, a másik csúcs címkéje a  $j$ -edik gyakori címke.

A kettőnél nagyobb elemszámú fák támogatottságának meghatározásánál szöfa jellegű adatstruktúrát javasoltak. A fát a fák címsorozatai alapján építjük fel, de a levelekben a címkeláncot tároljuk. Egy levélben több jelöltfa is lehet, hiszen különböző fáknek lehet azonos a címkeláncuk. Amikor egy bemeneti fára illeszkedő jelölteket kell meghatároznunk, akkor a bemenet címsorozata alapján eljutunk azokhoz a jelöltekhez, amelyek illeszkedhetnek a bemeneti fára. Egy jelölt címsorozatának illeszkedése szükséges feltétele annak, hogy maga a jelölt is illeszkedjen a bemeneti fára. Ha eljutunk egy levélbe, akkor az ott található címsorozatok mindegyikét megvizsgáljuk egyesével, hogy topologikusan illeszkedik-e a bemenet címkeláncra. Ennek részleteit nem ismertetjük.

### 12.3.2. TreeMinerV

A TreeMiner algoritmus [146] Zaki módszerét használja. A vertikális adatbázisból kiindulva előállítja a egyelemű fák illeszkedési listáit és a továbbiakban már csak ezen listákkal dolgozik.

Zaki módszerét ismertettük a 10.5.2 részben, a jelölt-előállítás pedig megegyezik a TreeMinerH jelölt-előállításának első lépésével. Csak azt kell tisztáznunk, hogyan határozzuk meg a

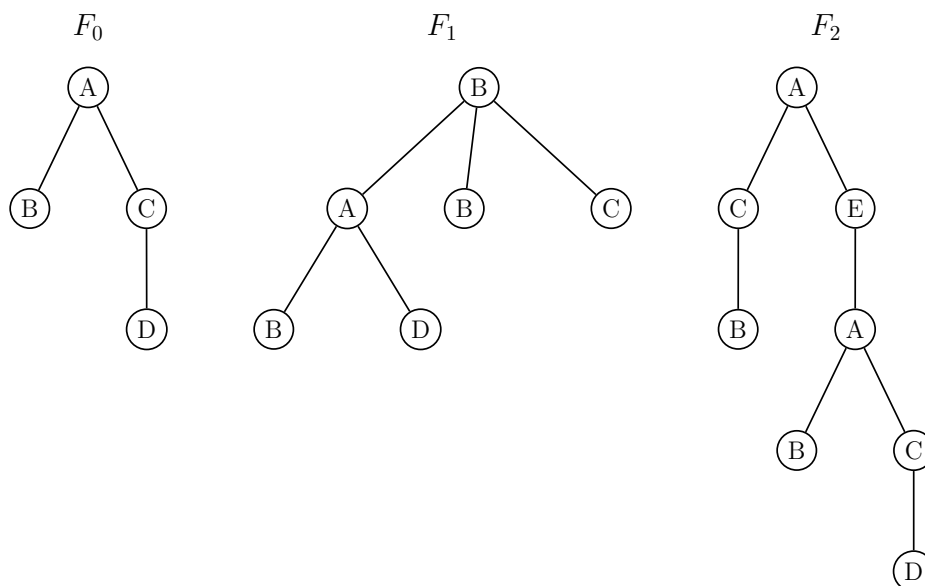
jelöltek támogatottságát.

### Jelöltek támogatottságának meghatározása

Az egyelemű jelöltek meghatározásához ismét elegendő egy lista, a kételemű jelöltekhez pedig egy tömb. A kettőnél nagyobb méretű jelöltek meghatározásának kulcsa az illeszkedési lista. A kiinduló illeszkedési listákat (amelyek az egyelemű gyakori fákhoz tartoznak) kételemű jelöltek meghatározása közben építjük fel.

Az a fő kérdés, hogy miként kell definiálni az illeszkedési listákat címkézett gyökeres fák esetében ahhoz, hogy teljesüljön a két elvárás (emlékeztetőül: a támogatottság egyértelműen meghatározható legyen belőle, és a jelöltek illeszkedési listáit a generátoraiból elő tudjuk állítani).

A fogalmak szemléltetésére a 12.6 ábrán található fákat fogjuk használni. Jelöljük egy



12.6. ábra. Példaadatbázis: címkézett gyökeres fák

tetszőleges  $F$  gyökeres fa  $j$  csúcsából induló részfa legnagyobb sorszámát  $MAX_F(j)$ -vel. Például  $MAX_{F_0}(0) = 3$ ,  $MAX_{F_2}(4) = 7$ ,  $MAX_{F_2}(1) = 2$ .

Az  $\ell$ -elemű  $F$  fa illeszkedési listájának minden eleme  $F$ -nek egy előfordulását rögzíti. Tegyük fel, hogy a  $F$  fa része az  $i$ -edik bementi fának ( $F_i$ -nek) és a tartalmazás injektív függvényét  $f$ -el jelöljük. Ekkor az illeszkedési lista ezen illeszkedését leíró eleme a következő 4-es:  $(i, \langle f(0), f(1), \dots, f(\ell-1) \rangle, f(\ell), MAX_{F_i}(f(\ell)))$ . Nyilvánvaló, hogy az illeszkedési listából a támogatottság és a súlyozott támogatottság is könnyűszerrel meghatározható. Már csak azt kell megnéznünk, hogy állítjuk elő a jelölt illeszkedési listáját, azaz hogyan illesztünk két illeszkedési listát.

Két illeszkedési lista illesztésének alapfeltétele, hogy a 4-esek első két tagjai megegyezzenek, hiszen azonos gráfban lévő illeszkedéseket keresünk (1. tag) és a generátorok prefixei azonosak (2. tag). Emlékszünk, hogy a  $F'$ ,  $F''$  illesztésénél két esetet különböztettünk meg, attól függően, hogy az új csúcsot  $F'$  legnagyobb sorszámú elemének testvére lesz vagy gyereke. Jelöljük a két fa illeszkedési listáinak 3. és 4. tagját  $f(\ell)$ ,  $MAX_{F_i}(f(\ell))$  és  $f'(\ell)$ ,  $MAX_{F_i}(f'(\ell))$ -el.

Az első típusú illesztés feltétele, hogy  $MAX_{F_i}(f(\ell)) < f'(\ell)$ . Ekkor a jelölt illeszkedési listája:  $i, \langle f(0), f(1), \dots, f(\ell-1), f(\ell) \rangle, f'(\ell), MAX_{F_i}(f'(\ell))$ .

A második típusú illesztés csak akkor jöhet létre, ha  $f(\ell) \leq f'(\ell)$  és  $MAX_{F_i}(f(\ell)) \geq MAX_{F_i}(f'(\ell))$ . A kapott jelölt illeszkedési listája az alábbi lesz:  $i, \langle f(0), f(1), \dots, f(\ell-1), f(\ell) \rangle, f'(\ell), MAX_{F_i}(f(\ell))$

Nézzünk példákat. Illesszünk az  $A \rightarrow C$  fát az  $A \rightarrow D$  fával. gráffal. Legyen az első fa illeszkedési listája:  $\langle (0, \langle 0 \rangle, 2, 3), (2, \langle 0 \rangle, 6, 7), (2, \langle 4 \rangle, 6, 7) \rangle$ , a másodiké:  $\langle (0, \langle 0 \rangle, 3, 3), (1, \langle 1 \rangle, 3, 3), (2, \langle 0 \rangle, 7, 7), (2, \langle 4 \rangle, 7, 7) \rangle$ . A generátorokból két jelöltet állítunk elő. Az elsőben a  $D$  címkéjű csúcs a  $C$  címkéjű csúcs testvére lesz, a másodikban a gyereke. Az első jelölt illeszkedési listája üres lesz, hiszen a 4 tagok egyenlők az azonos bemeneti fákhhoz tartozó listákban. A második jelölt illeszkedési listája:  $\langle (0, \langle 0, 2 \rangle, 3, 3), (2, \langle 0, 6 \rangle, 7, 7) \rangle$ , amiből azonnal meg tudjuk állapítani, hogy a jelölt támogatottsága 2.

Amennyiben egy prefix csak egyszer fordul elő egy bemeneti fában, akkor két generátor illesztésénél elég csak a 4-esek első elemének egyenlőségét vizsgálni. Ha ezek egyeznek, akkor nyilvánvaló, hogy a prefixek ugyanott fordulnak elő. Memóriefogyasztás csökkentése érdekében ezért a második tagot csak akkor tároljuk, ha a prefix a bemeneti fában többször előfordul.

## 12.4. Gyakori részfák

FOLYT. KÖV.

Fák esetében a részfa izomorfia eldöntésére létezik polinom idejű (pontosabban  $O(n \frac{k^{1.5}}{\log k})$ , ahol  $k$  a mintafa,  $n$  a másik fa csúcsainak száma) algoritmus [124].

FOLYT. KÖV.

## 12.5. A gyakori feszített részgráfok

Ebben a részben bemutatjuk a legismerteb gyakori feszített részgráfokat kinyő algoritmust. A  $\mathcal{MK} = (\mathcal{G}, \preceq)$ -ben mintatér elemei címkézett egymással nem izomorf gráfok és  $G' \preceq G$ , ha  $G'$  a  $G$ -nek feszített részgráfja. A gráf méretét a csúcsainak száma adja meg. A bemenet címkézett gráfok sorozata. A  $G$  gráf támogatottságán azon bemeneti elemek a számát értjük, amelyeknek létezik  $G$ -vel izomorf feszített részgráfjuk (feszített részgráf fogalma lásd alapfogalmak 2.3 rész).

### 12.5.1. Az AcGM algoritmus

Az AcGM algoritmus [64] – ami az AGM javított változata [63] – a gyakori feszített részgráfokat nyeri ki. Az algoritmus az APRIORI sémát követi. Ahhoz, hogy az összes összefüggő feszített részgráfot megtalálja előállítja a *félíg összefüggő* feszített részgráfokat is. Egy gráf félíg összefüggő, ha összefüggő, vagy két összefüggő komponensből áll, ahol az egyik komponens egyetlen csúcsot tartalmaz.

Az egész algoritmus során a gráfok szomszédsági mátrixzaival dolgozunk. A szomszédsági mátrix eredeti definíciója alapján nem tárolja a címkéket, ezért ebben a részben a  $G = (V, E, c_V, c_E)$  gráf  $f$  bijekciójához tartozó  $A_{G,f}$  szomszédsági mátrixának elemei  $(a_{ij}$  a mátrix

$i$ -edik sorának  $j$ -edik elemét jelöli):

$$a_{i,j} = \begin{cases} c(e_{ij}) & , \text{ ha } i \neq j \text{ és } (f^{-1}(i), f^{-1}(j)) \in E, \\ c(f^{-1}(i)) & , \text{ ha } i = j, \\ 0 & , \text{ különben} \end{cases}$$

Az  $A_{G,f}$  elemeiből és a csúcsok címkéiből egy kódot rendelhetünk a  $G$  gráfhoz:

$$CODE(A_{G,f}) = a_{1,1}, a_{2,2}, \dots, a_{k,k}, c_V(f^{-1}(k))a_{1,2}, a_{1,3}, a_{2,3}, a_{1,4}, \dots, a_{k-2,k}, a_{k-1,k},$$

azaz először felsoroljuk a csúcsok címkéit, majd a szomszédossági mátrix felső háromszög-mátrixának elemeit.

Különböző bijekciók különböző szomszédossági mátrixot, és így különböző kódokat eredményeznek. Amennyiben a címkéken tudunk egy rendezést definiálni, akkor a kódokat is tudjuk rendezni. Legyen a  $G$  gráf kanonikus kódolása az a kód, amelyik a legnagyobb ezen rendezés szerint. A kanonikus kódhoz tartozó szomszédossági mátrixot *kanonikus szomszédossági mátrixnak* hívjuk.

Az eddigiekhez hasonlóan most is azt kell tisztáznunk, hogy miként állítjuk elő a jelölteket és hogyan határozzuk meg a támogatottságukat.

### Jelöltek előállítás

Az  $X = A_{G',f}$  és  $Y = A_{G'',g}$   $\ell \times \ell$  méretű szomszédossági mátrixokat, ahol  $G'$  összefüggő,  $G''$  pedig félig összefüggő gráf, akkor illesztjük, ha teljesül három feltétel:

- Ha az  $X$  és  $Y$ -ből töröljük az utolsó sort és oszlopot, akkor azonos ( $T$ ) mátrixot kapunk, és a csúcsok címkéi is rendre megegyeznek:

$$X_\ell = \begin{pmatrix} T & x_1 \\ x_2^T & x_{l,l} \end{pmatrix}, Y_\ell = \begin{pmatrix} T & y_1 \\ y_2^T & y_{l,l} \end{pmatrix}.$$

- $T$  egy kanonikus szomszédossági gráf.
- ha  $x_{l,l} = y_{l,l}$ , akkor legyen  $code(X) < code(Y)$ , ellenkező esetben  $x_{l,l} < y_{l,l}$  vagy  $G'$  ne legyen összefüggő.

A potenciális jelölt szomszédossági mátrixa a következő lesz:

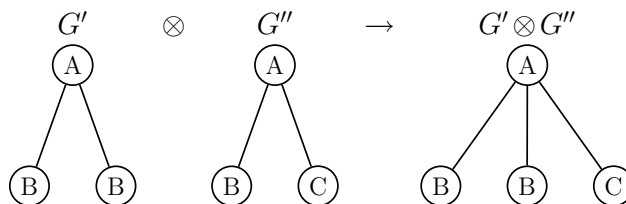
$$Z_{\ell+1} = \begin{pmatrix} T & x_1 & y_1 \\ x_2^T & x_{l,l} & z_{\ell,\ell+1} \\ y_2^T & z_{\ell+1,\ell} & y_{l,l} \end{pmatrix},$$

ahol  $z_{\ell,\ell+1}$  és  $z_{\ell+1,\ell}$  0-át és az összes lehetséges élcímke értékét felvehetik. Irányítatlan gráfok esetében a két értéknek meg kell egyeznie. Az ilyen módon létrehozott szomszédossági mátrixot a szerzők *normál formájú* szomszédossági mátrixnak nevezik.

Az első feltétel szerint nem csak azt várjuk el, hogy a két illesztendő mintának legyen  $(\ell - 1)$ -elemű közös részmintája, hanem még azt is, hogy ez a rész minta mindkét generátor prefixe is legyen. Tulajdonképpen ez biztosítja azt, hogy az illesztésként kapott jelölt mérete  $\ell + 1$

legyen. Ha a második és harmadik feltételnek nem kellene teljesülnie, akkor sokszor ugyanazt a potenciális jelöltet hoznánk létre. Az algoritmus nem lenne teljes, amennyiben csak összefüggő gráfok lehetnének a generátorok. Az  $\textcircled{A}-\textcircled{B}-\textcircled{C}-\textcircled{D}$  gráfot például a fenti jelölt előállításával nem lehetne kinyerni.

Nézzünk egy példát. A következő ábrán két gyakori 3 csúszú gráfot láthatunk, amelyből a jelölt előállítás során a jobb oldalon látható gráfot hozzuk létre. Az első gráf szomszédossági mátrixa  $\begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ , a másodiké  $\begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{pmatrix}$ , az illesztés során kapott szomszédossági mátrix pedig  $\begin{pmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 1 \\ 0 & 1 & 0 & z \\ 0 & 1 & z & 0 \end{pmatrix}$ .



12.7. ábra. Példa jelöltek előállítására

A jelölt-előállítás második fázisában minden  $\ell$  elemű feszített részgráfról el kell dönteni, hogy gyakori-e. Amennyiben az összes részgráf gyakori, akkor a potenciális jelölt valódi jelölt lesz, ami azt jelenti, hogy meg kell határozni a támogatottságát.

Sajnos ez a második lépés nem annyira egyszerű, mint elemhalmazok, sorozatok, gyökeres fák esetében. A feszített részgráf egy szomszédossági mátrixát megkaphatjuk, ha töröljük a mátrix adott indexű sorát és oszlopát. A problémát az okozza, hogy az így kapott mátrix nem biztos, hogy normál formájú lesz. Az AcGM a következő módszerrel alakítja át a részmatricot normál formájúvá.

FOLYT. KOV.

### támogatottságok meghatározása

A jelöltek előállítása után rendelkezésünkre fog állni egy nagy halom normál formájú szomszédossági mátrix. Ugyanannak a gráfnak több normál formájú szomszédossági mátrixa létezik ezért minden mátrixhoz hozzá kell rendelni az általa reprezentált gráf kanonikus kódját.

FOLYT. KOV.

Ha az azonos gráfot reprezentáló normál formájú szomszédossági mátrixok közül ki tudtuk választani a normál formájú szomszédossági mátrixot, akkor a továbbiakban már csak ezekkel dolgozunk, tehát csak ezekhez rendelünk – kezdetben 0 értékű – számlálókat.

A bemeneti gráfokat egyesével vesszük és minden jelöltet megvizsgálunk, hogy izomorf-e a bemeneti gráf valamely feszített részgrádjával. Feltételezzük, hogy a bemeneti mátrix kanonikus szomszédossági mátrixa rendelkezésünkre áll.

Ez a részfeladat tulajdonképpen a részgráf izomorfia feladata, amiről tudjuk, hogy NP-teljes. A feladatot azonban gyorsan megoldhatjuk, ha tudjuk, hogy a jelölt  $\ell$ -elemű feszített részgráfa a bemeneti gráf melyik feszített részgrádjával volt izomorf. Nem kell mást tennünk, mint megvizsgálni, hogy az új csúcs és a hozzá tartozó él illeszkedik-e a bemeneti gráf részgrádjára.

## 12.6. A gyakori részgráfok keresése

Ebben a részben feltesszük, hogy a mintatér elemei összefüggő gráfok és  $G' \preceq G$ , ha  $G'$  a  $G$  gráfnak részgráfja. Eben a mintakörnyezetben egy gráf méretét az éleinek száma adja meg. A bemenet címkézett gráfok sorozata. A  $G$  gráf támogatottságán azon bemeneti elemeknek a számát értjük, amelyeknek létezik  $G$ -vel izomorf részgráfja. Bemutatjuk a két legismertebb algoritmust az FSG-t és a gSpan-t.

### 12.6.1. Az FSG algoritmus

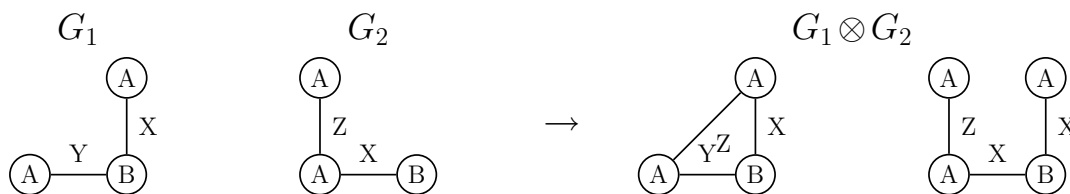
Az FSG algoritmus [78] az APRIORI sémára épül. A gráfok tárolásához szomszédossági listát használ. Amikor egy gráfnak elő kell állítani a kanonikus kódját, akkor a szomszédossági listát szomszédossági mátrixá alakítja. Amennyiben a gráfok ritkák, a szomszédossági listák kevesebb helyet igényelnek, mint a mátrixok.

Megszokhattuk már, hogy a fő lépés a jelöltek előállítás.

#### Jelöltek előállítása

Két  $\ell$ -elemű  $G_1 = (V_1, E_1), G_2$  gráfot akkor illesztünk, ha van  $(\ell - 1)$ -elemű közös részgráfjuk (ezt hívtuk magnak), és az  $G_1$  kanonikus kódja nem nagyobb  $G_2$  kanonikus kódjánál. Ez azt jelenti, hogy minden gráfot önmagával is illesztünk. Két gráf illesztésénél – akárcsak két elemsorozatok esetében – több gráf jön létre. Jelöljük a  $G_2$ -nek a magba nem tartozó élét  $e = (u, v)$ -vel. Az előállított gráfok a  $G_1$  bővítése lesz egy olyan  $e' = (u', v')$  éllel, amelyre  $u' \in V_1, e' \notin E_1, c_E(e) = c_E(e'), c_V(u) = c_V(u')$  és  $c_V(v) = c_V(v')$ . Tehát egy megfelelően címkézett él helyezünk be a  $G_1$  gráfba. Ezt többféleképpen tehetjük, így több potenciális jelöltet hozunk létre.

Lehet, hogy az új él új csúcsot is fog eredményezni, de az is lehet, hogy csak két meglévő pont között húzunk be egy új élt. Ezt szemlélteti a 12.8 ábra.



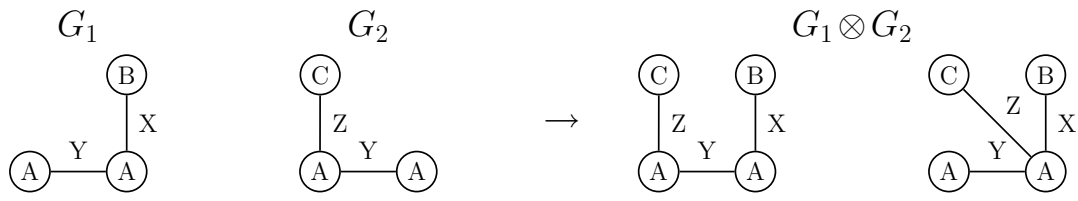
12.8. ábra. Példa: gráf illesztése

Az előállított potenciális jelöltek számát növeli az a tény is, hogy a magnak több automorfizmusa lehet, így az új élt több csúcsokhoz is illeszthetjük. Erre mutat példát a következő ábra.

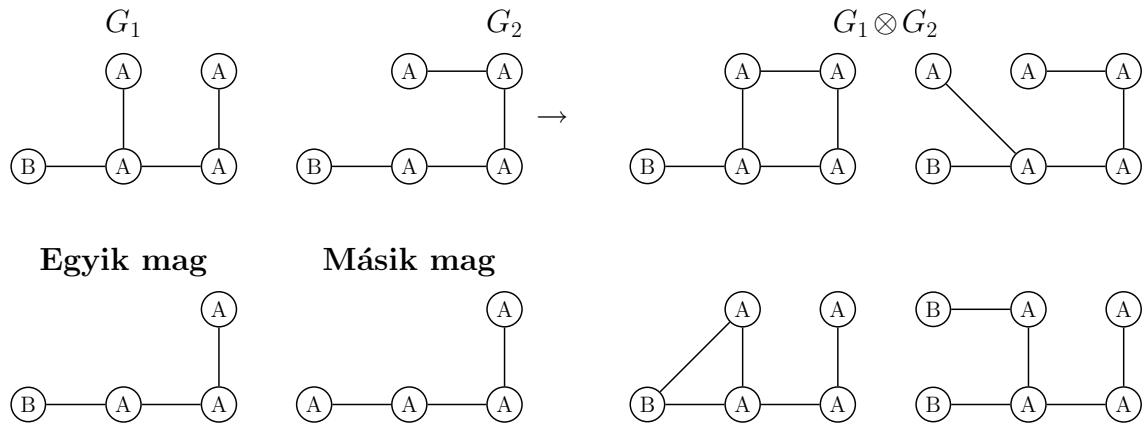
A harmadik ok, amiért két gráf több potenciális jelöltet állíthat elő az, hogy két gráfnak több közös részgráfja (magja) is lehet. Egy ilyen eset látható a 12.10 ábrán.

Miután előállítottuk a potenciális jelölteket, minden potenciális jelölt  $(\ell - 1)$ -elemű részgráfját ellenőrizzük, hogy gyakori-e. Azok a potenciális jelöltek lesznek jelöltek, amelyek minden valódi részhalmaza gyakori és még nem vettük fel a jelöltek közé. Ez utóbbi feltétel már sejteti, hogy a fenti jelölt-előállítás nem ismétlés nélküli. Az algoritmus a gráfok kanonikus





12.9. ábra. Példa: gráf illesztése - mag automorfizmusok



12.10. ábra. Példa: gráf illesztése - több közös mag

kódolását használja annak eldöntésére, hogy egy potenciális jelölt adott részgráfja gyakori-e, illetve a jelölt szerepel-e már a jelöltek között.

A jelöl-előállításnak tehát három fő lépése van: mag azonosítás (ha létezik egyáltalán), él-illesztés és a részgráfok ellenőrzése. Az első lépést gyorsíthatjuk, ha minden gyakori gráfnak egy listában tároljuk az  $(\ell - 1)$ -elemű részgráfjainak kanonikus kódjait. Ekkor a közös mag meghatározása tulajdonképpen két lista metszetének meghatározását jelenti.

### támogatottság meghatározása

A bemeneti gráfokat egyesével vizsgálva meg kell határozni, hogy melyek azok a jelöltek, amelyek izomorfak a bemeneti gráf valamely részgráfjában. A részgráf izomorfia eldöntése NP-teljes, de ezen feladat eldöntésére használt algoritmusok számát csökkenthetjük, ha minden részgráfnak rendelkezésünkre áll a TID-hamaza, azon bemeneti gráfok sorszámai, amelyek tartalmazzák a részgráfot. Egy jelölt vizsgálatánál csak azon bemeneti elemeket kell megvizsgálnunk (ha ezek száma nagyobb  $min\_supp$ -nál), amely sorszámja minden részgráf TID-halmazában szerepel.

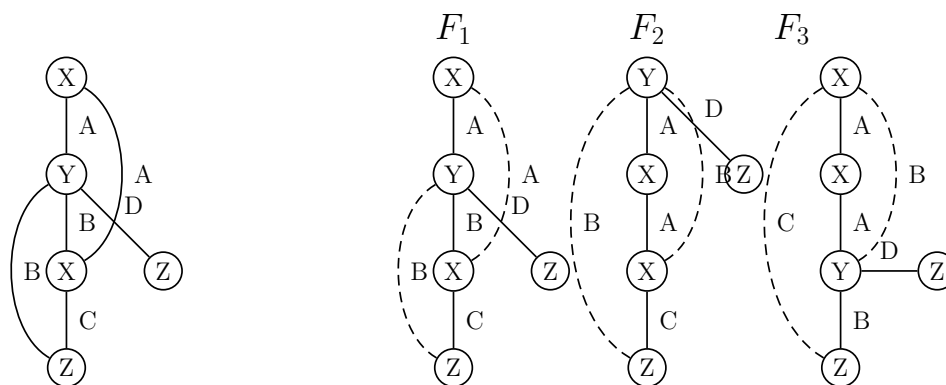
### 12.6.2. gSpan

A gráfok ábrázolására a gSpan a *DFS-kódokat*, illetve az abból előállított kanonikus kódolást használja. A mélységi kód előállításához ki kell választanunk egy gyöker csúcsot, majd ebből a csúcsból indulva bejárni a gráfot, mintha egy gyökeres fát járnánk be mélységi bejárás szerint. A bejárás szerint minden csúcshoz időcímkét rendelhetünk, amely megadja, hogy hanyadik lépés során látogattunk meg egy csúcsot. Mivel a gráf tartalmazhat köröket is, ezért előfordulhat,

hogy egy csúcsot többször meglátogatunk. Ilyen esetben a csúcs időcímkéjét ne írjuk felül (és az idő számlálóját se növeljük). Hívjuk *előreélnek* azokat az éleket, amelyek még nem látogatt csúcsba vezetnek, a többit élt pedig *visszaélnek*.

A gráf bejárása során minden lépésnek egy elem felel meg a DFS-kódban, azaz a kód hossza megegyezik a gráf éleinek a számával. Minden elem egy ötös, amelynek első két eleme az indulási és az érkezési csúcsok időbélyegét adja, a harmadik és ötödik elem ezen csúcsok címkéit és a negyedik elem az él címkéjét tárolja.

Természetesen egy adott gráfnak több DFS-kódja is lehet attól függően, hogy melyik csúcsot választjuk gyökének és milyen sorrendben vesszük egy csúcs gyermekeit. A 12.11 ábrán egy példagráfot, három különböző mélységi bejárást és az azokhoz tartozó DFS-kódokat láthatjuk. A visszaéleket szagatott vonallal jelöltük.



él	$F_1$	$F_2$	$F_3$
0	(0,1,X,a,Y)	(0,1,Y,a,X)	(0,1,X,a,X)
1	(1,2,Y,b,X)	(1,2,X,a,X)	(1,2,X,a,Y)
2	(2,0,X,a,X)	(2,0,X,b,Y)	(2,0,Y,b,X)
3	(2,3,X,c,Z)	(2,3,X,c,Z)	(2,3,Y,b,Z)
4	(3,1,Z,b,Y)	(3,0,Z,b,Y)	(3,0,Z,c,X)
5	(1,4,Y,d,Z)	(0,4,Y,d,Z)	(2,4,Y,d,Z)

12.11. ábra. Példa: mélységi fák és mélységi kódok

A címkéken tudunk egy rendezést definiálni, ami alapján az ötösöket is rendezni tudjuk. Ezen rendezés szerint lexikografikusan rendezni tudjuk a kódokat is. Egy gráf kanonikus kódja legyen az a DFS-kódja, amely ezen rendezés szerint a legkisebb.

Legyen  $\alpha = \langle a_0, a_1, \dots, a_m \rangle$  egy DFS-kód. Ekkor a  $\beta = \langle a_0, a_1, \dots, a_m, b \rangle$ -t az  $\alpha$  *gyermekének* hívjuk,  $\alpha$ -t pedig a  $\beta$  *szülőjének*. Ahhoz, hogy a  $\beta$  tényleg DFS-kód legyen a  $b$  címkéjű élnek az  $\alpha$  által kódolt mélységi fa legjobboldali ágán kell elhelyezkednie. Erre a DFS-kódnövelésre láthatunk példát a következő ábrán.

Könnyű belátni, hogy amennyiben az új él visszaél, akkor csak a legjobboldalibb csúcsból indulhat.

A szülő-gyerek reláció megadásával definiálhatunk a DFS-kódfa fogalmát. A DFS-kódfa egy olyan fa, amelynek csúcsaiban DFS-kódok ülnek és minden szülő-gyerek csúcs által reprezentált DFS-kódokra teljesül a fenti szülő-gyerek kapcsolat és a fa rendezett, azaz minden csúcs gyermeke

## 12.12. ábra. Példa: mélységi kód

a DFS-kód szerint növevő sorrendbe van rendezve.

Amennyiben egy kanonikus kódhoz hozzáveszek egy új élt úgy, hogy ez DFS-kódot eredményezzek, az nem jelenti azt, hogy ez a kód kanonikus kód lesz. A DFS-kódfában minden kanonikus kód megtalálható, de emellett számos nem kanonikus kód is szerepel. A rendezés azonban garantálja, hogy ha pre-ordes bejárás szerint bejárnánk a fát, akkor tetszőleges gráf első DFS-kódja egybe kanonikus kód is. A DFS-kódfát ezek szerint egyszerűsíthetjük, hogy kimetszük azon részfákat, amelyek csúcsai nem kanonikus kódokat tartalmaznak. A gSpan algoritmus tulajdonképpen ezt a gyakori gráfokat tároló egyszerűsített DFS-kódfát állítja elő. Mihelyt egy olyan DFS-kódot állít elő, amely nem minimális, a fát nem növeszti tovább ezen az ágon.

Könnyű belátni, hogy a  $G=(V, E)$  gráfnak nem kell  $|V| \cdot |E|$ -nél többször törölni nem kanonikus DFS-kódját. A  $G$  gráfot csak  $(|E|-1)$ -elemű részgráfjából származtathatjuk, amelyek száma legfeljebb  $|E|$ . A részgráf által nem tartalmazott élt, amennyiben az előreél  $|V|-1$  féleképpen illeszthetjük a legjobboldalibb ághoz. Visszaél esetében pedig a legjobboldalibb csúcshoz kell tennünk, ezen lehetőségek száma pedig  $|V|-2$ . Ez a korlát elég gyenge, hiszen csak annyit tett fel, hogy a legjobboldali úton található csúcsok száma kisebb  $|V|$ -nél. Az esetek többségében ez az út az élszámnál jóval kisebb, így a nemkanonikus kódok törlésének száma jóval kevesebb.

FOLYT KÖV.

## 13. fejezet

# Adatbányászat a gyakorlatban

Az eddigi fejezetekben matematikai modellekről, megoldandó feladatokról és algoritmusokról beszéltünk. E fejezet nem lesz ennyire tudományos: az adatbányászat legtipikusabb felhasználási területeit fogjuk átnézni. Azt vizsgáljuk, hogy milyen jellegű összefüggések után érdemes kutatni, és hogy ezen összefüggések felderítése milyen előnyökkel jár.

Az adatbányászat napjaink egyik legnépszerűbb területe. Nem meglepő, hogy napról napra új adatbányászati szoftver jelenik meg, hirdetve magáról azt, hogy a piac legjobb terméke. A fejezet második részében összefoglaljuk a jelenleg kapható adatbányászati szoftvereket, majd kitérünk arra, hogy milyen szempontokat vegyen figyelembe egy cég a megfelelő szoftver kiválasztásánál.

### 13.1. Felhasználási területek

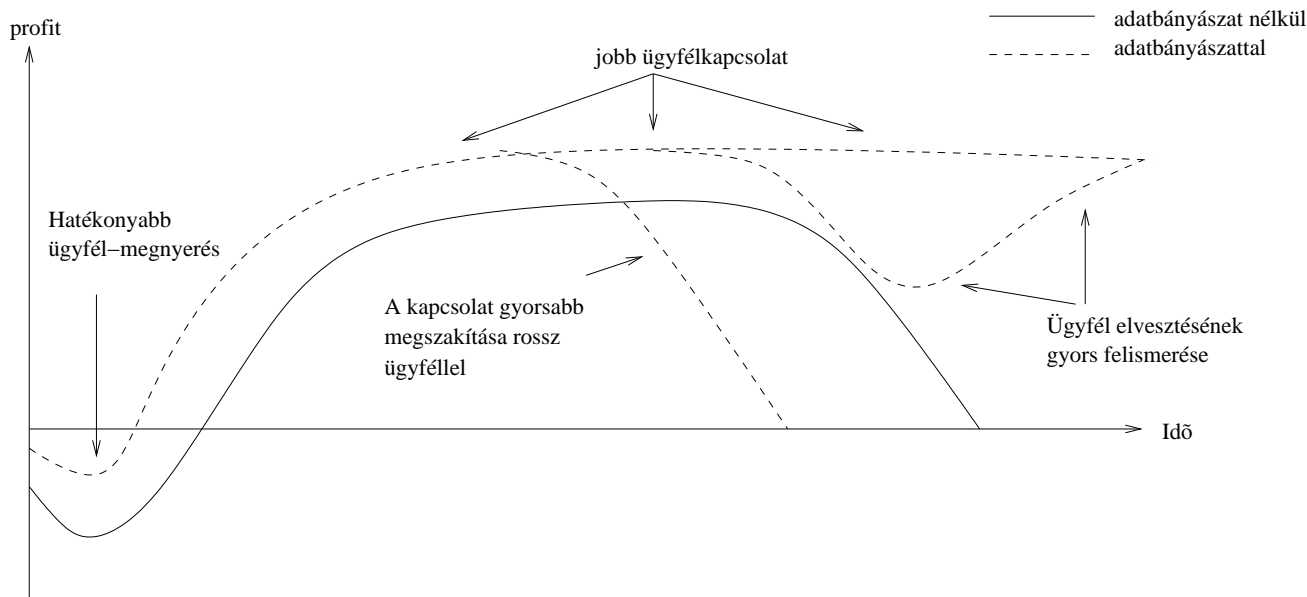
A sikeres alkalmazások hatására az adatbányászat egyre elfogadottabb tudományággá vált. Már szinte mindenhol fontos az adatok tárolása mellett azok feldolgozása és elemzése. A kinyert információ új tételek, törvényszerűségek felfedezését segítheti elő, vagy éppen fordítva: meglévő hipotéziseket cáfolhat meg. Ebben a részben 3 olyan területről szólunk, ahol az adatbányászat már mélyen gyökeret vert és az egyik legfontosabb eszközzé vált. Ezek a területek pedig: (1.) kereskedelem, (2.) pénzügy, (3.) biológia és orvostudomány. Az itt leírtakon túl számos esettanulmányt sikeres alkalmazásról szóló hírt lehet találni például a magyar adatbányászok honlapján (<http://www.datamining.hu>).

#### 13.1.1. Az ügyfél életciklusa

A kereskedelemben és a pénzügyben is a profitot az ügyfelek termelik. A következő ábrán láthatjuk, hogy milyen változásokat képes hozni az adatbányászat az ügyfelek életciklusában.

Az adatbányászat segítségével hatékonyabban fel tudjuk ismerni a potenciális ügyfelek körét. Így kevesebbet költünk azon ügyfelekre, amelyek nagy valószínűséggel nem lesznek ügyfeleink, azaz faragunk a költségeken.

Meg tudjuk különböztetni a jó és a rossz ügyfeleket. A rossz ügyfelektől hamarabb megszabadulhatunk, és az ezáltal megtakarított összeget a jó ügyfelekre fordíthatjuk, ami még gyümölcsözőbb kapcsolatot eredményezhet. Azt is idejében észrevehetjük, ha egy jó ügyfelünknek növekszik az elégedetlensége velünk szemben.



13.1. ábra. Az ügyfél életciklusa

### 13.1.2. Kereskedelem

Többször hoztunk fel példát az adatbányászat kereskedelmi felhasználásáról. Magát az asszociációs szabályokat is egy kereskedelmi példán keresztül vezettük be. Ez nem véletlen, hiszen a vásárlói kosarak elemzésének igénye keltette életre ezt a területet.

A kereskedelemben ma már minden üzletben megtalálhatóak a működést segítő számlázó, raktárkészlet-kezelő, programok. Ezek egyre összetettebbek, a pusztán vásárlások felsorolásánál és visszakeresésénél jóval többet tudnak: ha rendelkezésünkre áll valamilyen vevőazonosító, akkor a vevők teljes vásárlói történetét megkaphatjuk, de ezenfelül hitelekéről, beszerzésekről, szállításokról is rögzíthetünk adatokat. A nagy multik ma már tudják, hogy a törzsvásárlói kártyák növelik a vevő hűségkedvét. Ezek a kártyák újabb adatokat így újabb hasznos elemzéseket tesznek lehetővé. A szervizekből érkező visszacsatolásoknak is fontos szerepe lehet egy termék sikerénél.

Az on-line áruházak elterjedésével a vevőkről begyűjtött adatok minősége tovább javul. Az elektronikus kereskedelemről külön részben szólnunk bővebben.

Az adatbányászatnak a kereskedelem területén a következő céljai lehetnek.

- Vásárlói szokások elemzése az asszociációs szabályok, illetve az epizódok kinyeréséhez vezetnek. A szabályokat felhasználhatjuk eladást ösztönző akciók szervezésénél, áruházak térképének kialakításánál, prospektusok tervezésénél, eladáshelyi reklámeszközök kialakításánál, termékfejlesztésnél, ....
- Klaszterezés és osztályozás segítségével vásárlói csoportokat hozhatunk létre. A célcsoportok pontosabb behatárolásával irányíthatóbb, emberközelibb, interaktív, egyedi és hatékonyabb reklám- és marketingstratégiát készíthetünk.
- A személyre szabott ügyfélszolgálat nagyban fokozza a vásárlók elégedettségét.

- A vásárlói szokások jobb megismerésével pontosabban tudjuk megjósolni az egyes termékek értékesítési adatait. Ha egy üzletnek pontos képe van a raktárkészletről, a fogyás üteméről és az igényekről, akkor hatékonyabban tudja megszervezni a beszerzéseket, a disztribúciós csatornák típusát, nagyságát, a raktározás módját (pl.: just-in time). A hatékonyság növelésével csökkenteni tudjuk a költséget és jobban elosztani az erőforrásokat.
- Az új vevők toborzása mellett a régiók megtartása egyre fontosabb (CRM - Customer Relation Management). A vevők vásárlói sorozatainak elemzésével képet kaphatunk arról, hogy kinek csökkent a vásárlói kedve, így még azelőtt tehetünk ellene valamit (pl. hűségakció), hogy végképp elpártolna tőlünk.

### 13.1.3. Pénzügy

A bankok szolgáltatásai közül kiemelten fontosak a számlák, lekötött betétek vezetése, a hitelek nyújtása és egyéb pénzügyi tranzakciók lebonyolítása. Ezeken a területeken mára elismertté vált az adatbányászat szerepe.

- A bankok eredetileg azért jöttek létre, hogy mások értékeit megőrizték. Az ügyfelek így biztonságban tudták a pénzüket, ami a kamatok miatt gyarapodott is, a bankok pedig nagy tőkékhez jutottak, amit be tudtak fektetni. Mára a bankok az ügyfeleket eltérően kezelik (lakossági, vállalati ügyfelek, számlaforgalomtól függően átlagos, fontos, kiemelt ügyfelek ...). Az ügyfelek és a tranzakciók nagy száma miatt az ügyfélcsoportok manuális kialakítása lehetetlen feladat. A klaszterezés és az osztályozás ezért ezen a területen kiemelten fontos eszköz.
- Hitelek nyújtása a kamatok és a rendszeres jövedelemforrás miatt jó befektetés a banknak. A kérelmező körülményeinek megvizsgálása nélkül osztogatni a hiteleket azonban kockázatos, mert lehet, hogy az ügyfél nem tudja visszafizetni. Ha az ügyfelekről sok adat áll rendelkezésre (nettó jövedelem, beosztás, családi állapot, korábbi banki tranzakciói ...), akkor az osztályozást felhasználva olyan döntési fákot lehet létrehozni, amelyek nagy bizonyossággal megállapítják adott ügyfélről, hogy megbízható hitel szempontjából vagy nem. Ezt a módszert elsősorban olyan országokban alkalmazzák, ahol a hitel odaítélését nem kötik nagyon szigorú feltételekhez. Amerikában például elterjedt szokás, hogy Karácsony előtt a bankok előzetes megrendelés nélkül hitelkártyákat küldenek szét, amit a címzett nem köteles használni, de ha fizet vele, akkor a hitelt néhány hónapon belül vissza kell fizetnie. Nyilvánvalóan a megnövekedett vásárlói kedv ily módon való ösztönzése kiemelt jelentőségű lehet egy bank számára partnerkörének kibővítésénél, megtartásánál és természetesen a plusz generált pénzforgalom figyelembe vételénél, de ezen marketingakció kockázata igen magas, ha a hitelkártyát használó a későbbiekben nem fizeti vissza a bank pénzét.
- A bank/hitelkártyával történő fizetés és készpénzfelvétel a civilizáció nélkülözhetetlen eleme. A rengeteg biztonsági intézkedés ellenére a kártyás csalások még mindig sok kárt okoznak. Mivel a tranzakciók száma óriási, ezért manuális eszközökkel ebben az esetben is lehetetlen feladat kiszűrni a szokatlan viselkedést, ami a csalókra, kártyatolvajokra jellemző. Az eltéréselemzés az adatbányászat azon területe, ahol a szokásostól eltérő viselkedés, mintázat felfedezése a cél.

### 13.1.4. Biológia és Orvostudomány

A biológiában és az orvostudományban az adatok elemzéséből kapott törvényszerűségek értéke felbecsülhetetlen. Adatbányászat segítségével fejlesztenek új gyógyszereket, segít a rák elleni terápia hatékonyabb kialakításában, különböző betegségek tüneteinek meghatározásában. . . . Ebben a részben két fontos alkalmazásról szólunk: a DNS láncok elemzéséről és a cukorbetegség kezelésének segítségéről.

#### DNS láncok elemzése

Az orvostudomány talán legnagyobb megoldatlan feladata a DNS láncok teljes megfejtése. Tudjuk, hogy minden élőlényt egyértelműen azonosít a DNS lánc, mintha egy genetikai kódot kaptunk volna a természettől! A DNS láncok a felelősök többek között a betegségekért, bizonyos emberi tulajdonságokért, hajlamokért, allergiákért. . . . Éppen ezért a kiemelt szerepért a DNS láncok fontosságát aligha lehet túlbecsülni.

Minden DNS lánc 4 építőközből épül fel, ezek a nucleotidok: adenin, cytosin, guanin és thymin. Ez a négy nucleotid alkot egy hosszú láncot, ami leginkább egy spirál alakú létrára emlékeztet. Egy ember kb. 100 ezer génnel rendelkezik, egy gén pedig általában többszáz nucleotidból épül fel, ahol a nucleoidok sorrendjének fontos szerepe van. A DNS láncok elemzése nem pusztán epizód kutatásról szól. A tudás kinyeréséhez ötvözni kell a különböző adatbányászati technikákat!

- DNS láncokat gyakran kell összehasonlítani, ezért sorozatok hasonlóságának elemzése fontos módszer. Beteg és egészséges szövetekből vett minták összevetéséből megállapíthatjuk a kritikus eltéréseket. Először a két mintát külön vizsgálják és nyerik ki a gyakran előforduló mintázatokat. Később már csak ezeket a mintázatokat vetik össze. A beteg szövetben jóval gyakrabban előforduló mintázatok lehetnek a betegség genetikai tényezői. Vagy fordítva, az egészséges szövetben gyakrabban előforduló mintázatok adhatnak alapot a gyógyszer elkészítéséhez.
- A betegségekért általában nem csak egy gén felelős, hanem a gének egy kombinációja. Az asszociációs szabálykeresésnél megismert módszerekkel lehet feltárni a gyakran együtt előforduló esetleg felelős géneket beteg egyedek egy adott csoportjában.
- A gének és betegségek világát tovább bonyolítja, hogy a betegség különböző fázisaiban esetleg más-más gének aktívak. Ha egy adott betegségnél sikerülne ezt feltérképezni, akkor a különböző fázisokhoz elkészített gyógyszerek kifejlesztésével növelni lehetne a kezelés hatékonyságát.

#### Cukorbetegség

A cukorbetegség egy elterjedt és nem megfelelő kezelés esetén halált okozó betegség. Kezelésére a betegnek inzulint kell a szervezetébe juttatnia. Amennyiben a beteg túl sok inzulint kap, akkor szervezete tovább csökkenti a cukor termelését, és betegség tovább súlyosbodik. Ha viszont a beteg kevés inzulint kap, akkor szervezetében cukorhiány mutatkozik, aminek hatásaként szédülés, ájulás, sőt akár bénulás is előállhat.

Az inzulin megfelelő adagolása ezért kiemelten fontos a cukorbetegség kezelésében. A tudomány jelenlegi állása szerint nincs pontos képlet arra nézve, hogy egy adott paraméterekkel

rendelkező beteg mekkora adagot kapjon. Habár egyre több eszköz jön létre a minél gyorsabb visszajelzésre, az adagok meghatározása még mindig ösztönszerűen, az orvos le nem írt, konkrétan meg nem fogalmazott tapasztalatai alapján történik.

Az inzulin megfelelő adagjának kiválasztása rengeteg paramétertől függ (testsúly, kor, nem, betegségre jellemző adatok stb.). A különböző mérő- és figyelőeszközök piacra kerülésével egyre több adat gyűlik össze, így lehetővé válik ezek elemzése. Az osztályozás, korrelációanalízis, asszociációkutatás mind fontos eszközök a cukorbetegség kutatásában.

## 13.2. Az adatbányászat bölcsője: az elektronikus kereskedelem (e-commerce)

A bevezetőben szó volt arról, milyen feltételei vannak a sikeres adatbányászatnak (lásd 22.oldal). Idézzük fel ezeket a feltételeket, és nézzük meg, hogyan teljesülnek az elektronikus kereskedelemben.

**sok adat:** Közismert weboldalnak nagy a látogatottsága.

**sok attribútum:** Az on-line áruházaknál lehetőség van a vásárló fontosabb adatainak tárolására, de ezenfelül több más információt is megtudhatunk róla, pl. hogy mi iránt érdeklődik gyakran a látogató, milyen reklámokat néz meg, ...

**tiszta adat:** Az adatok az emberi rögzítés hibájától mentesek. A weboldal készítője határozhatja meg, milyen típusú adatok legyenek tárolva, illetve, mely mezőket kell kötelezően kitölteni.

**akcióképesség:** A kinyert tudás birtokában megváltoztathatjuk a weboldalt (akár a teljes designt, akár csak a linkeket), személyre szabott oldalakat készíthetünk, e-maileket küldhetünk ki, ...

**befektetés megtérülése:** Mivel minden elektronikusan zajlik a bevételnövekedés kiszámítása alapfeladat. Sőt még a célzott marketing hatékonyságát is könnyedén megállapíthatjuk, hiszen rögzíteni tudjuk, ha valaki egy e-mailen keresztül jutott az oldalunkra.

A fentiek ellenére az adatbányászat alkalmazása az elektronikus kereskedelemben korántsem akadálytalan [77]. A problémák forrása az, hogy az adatokat a webszerverek mentik el. A webszerver adatainak bányászata kézenfekvőnek tűnik, hiszen a webszerverek szinte mindent rögzítenek. A naplófájlokat azonban eredetileg a webszerverek debuggolására találták ki, nem pedig az adatbányászat támogatására.

A legfőbb problémák az alábbiak:

- Nem lehet egyértelműen azonosítani a felhasználót. Szemben az adatbányászattal, a webszerverek számára ez ugyanis nem fontos információ. Próbálkoznak a felhasználók cookie, IP, vagy böngésző szerinti azonosításával [31], de ezek közül egy sem nyújtja a tökéletes megoldást [14].
- A webszerverek nem tárolnak minden fontos adatot. A naplófájlokban nincs nyoma például a „berakom a kosárba”, „mennyiség megváltoztatása”, „termék törlése” műveleteknek.



- A form-ok adatai nincsenek tárolva. Pedig gondoljuk meg, hogy például a keresési form-ok éppen a vásárlások érdeklődését tükrözik.
- A naplófájlokban URL-ek szerepelnek, nem pedig az oldal tartalma. Nem mindig könnyű meghatározni, hogy adott termék melyik oldalhoz tartozik. A helyzetet tovább bonyolítja, hogy gyakran ugyanaz az információ több nyelven is elérhető.
- A dinamikus oldalak tartalmát sem lehet egyértelműen meghatározni. Melyik termék érdekelte a látogatót, ha az összes terméket a `termek.jsp` oldal mutatja? Vagy csak egy reklám volt a felbukkanó ablak? Esetleg egy „Nincs raktáron!” üzenet? Sikeres volt a keresés vagy nem hozott eredményt? Ezek a kérdések legtöbbször dinamikus oldalakhoz kötődnek és megválaszolásuk a naplófájlok alapján lehetetlen feladat.
- Az igazán nagy oldalaknak több webszerverük van, amelyek különböző helyeken helyezkednek el. Ezek mind saját naplófájllal dolgoznak, egyesítésüket nehezíti, hogy különböző időzónákban lehetnek.

A fenti problémákra megoldás nyújt, ha a vásárlásokkal kapcsolatos információk tárolását a vásárlásokat kiszolgáló program végzi, azaz az adatok tárolását az alkalmazási rétegre bízjuk. A valóságot tükröző adatok előállításának nagy ellenségei az Internetes robotok. Ezek olyan lekérdezéseket, oldalletöltéseket generálnak, amelyek nem tükröznek valóságos emberi érdeklődést. Intenzív kutatás tárgyát képezi a robotok által generált hamis adatok kiszűrése.

### 13.3. Adatbányász szoftverek

A továbbiakban egy rövid összefoglalót adunk a ma kapható legfontosabb adatbányász szoftverekről. A lista korántsem teljes. Ennek oka egyrészt a terjedelmi korlát, másrészt a nap mint nap változó piac.

**weka** (<http://www.cs.waikato.ac.nz/ml/weka/>) Az új-zélandi Waikato Egyetem fejleszti a szabad forráskódú a WEKA nevű adatbányászati programcsomagot. Szimbólikus elnevezése az ország nemzeti madaráról, a kiviról származik: az adatbányász "rejtett tudást" keres, a kivimadár (weka) pedig fejét vízbe dugva kutat a "rejtett" táplálék után.

A különböző adatbányászati algoritmusok igen széles körét találjuk meg a szoftvercsomagban. Az implementált eljárások száma nemcsak abszolútértékben, hanem az igen drága kereskedelmi termékhez viszonyítva is magas.

A WEKA-t JAVA nyelven fejlesztik, az egyes osztályok forráskódja mellett azok dokumentációi is hozzáférhetők az interneten, mely remek lehetőséget kínál a kutatóknak, diákoknak, adatbányászati iránt érdeklődőknek.

A WEKA felhasználóbarát, logikus, jól áttekinthető grafikus felülete vezeti végig a felhasználót az adatbányászati lépésein. Oktatási és demonstrációs célra is kiváló. A WEKA adatforrások széles körét támogatja. Az elemzendő adatok származhatnak például JDBC-n keresztül elérhető adatbázisoktól vagy fájlokból. Előnyös tulajdonságainak köszönhetően világszerte ismert és elismert szoftver.

**Enterprise Miner** (<http://www.sas.com/products/miner/index.html>) A SAS Institute, Inc. fejlesztette ezt a programcsomagot. A cég komoly múltra tekint vissza statisztikai elemzések terén. Az Enterprise Miner is számos statisztikai eszközt kínál fel, de már megtalálható az összes többi adatbányászati feladatra megoldás, csak úgy mint döntési fák, neurális hálózatok, regresszió, klaszterezés, sorozat-elemzés, asszociációbányászat.

**Clementine** (<http://www.spss.com/spssbi/clementine>) A Clementine az SPSS Inc. terméke. Integrált adatbányászati környezetet biztosít végfelhasználók és fejlesztők részére. Adatbányászati eszközök közül megtalálható a neurális hálózatok, osztályozás, sorozat-elemzés stb. A Clementine szoftverében egyedi az az objektum-orientált interfész, amin keresztül a felhasználó saját algoritmusokat és funkciókat adhat meg.

**Intelligent Miner** (<http://www-3.ibm.com/software/data/iminer/foridata>) Az IBM terméke talán a legismertebb és a legelterjedtebb adatbányászati eszköz. Emellett fontos érv szól mellette: az IBM kutatóintézetében született jónéhány neves publikáció, tehát e szoftver mögött áll a legfelkészültebb kutatógárda. A programmal lehet bányászni asszociációkra, epizódokra, alkalmas osztályozási, klaszterezési feladatok ellátására, de ezenkívül lehet regressziót számolni és eltérést keresni. A fejlett adatmegjelenítés mellett képes statisztikai elemzésre és neurális hálózatokon alapuló algoritmusok futtatására. Az Intelligent Miner használatához IBM DB2 relációs adatbáziskezelő rendszernek is futnia kell.

**DBMiner** (<http://www.dbminer.com>) A DBMinert a Simon Fraser University által elkészített programból fejlesztette tovább a DBMiner Technology Inc.. Adatbányászati funkciók közül megtalálhatók a asszociációbányászat, karakterizáció, osztályozás, klaszterezés és jóslás. Ennek a programnak legszorosabb, legintegráltabb a kapcsolata az OLAP-pal. A szoros kapcsolat miatt itt már OLAM-ról (On-Line Analytical Mining) beszélünk. A program egy interaktív környezetet kínál a felhasználónak, aki dinamikusan váltogathat OLAP operációk és adatbányászati funkciók között.

**MineSet** A MineSet legnagyobb erőssége a fejlett vizualizációs képessége. Ez nem meglepő, hiszen a szoftvert a Silicon Graphics fejlesztette, amely cég mindig is a legjobbak közé tartozott a grafikában. A MineSetben megtalálható szinte az összes ismert adatbányászati funkció. További előny, hogy a MineSet egyben fejlesztői környezetet biztosít új algoritmusok implementálásához, így ha valamely feladatra nincs kész megoldás, akkor megírhatjuk magunk, majd az eredmény megtekintéséhez használhatjuk a MineSet vizualizációs eszközeit.

A fenti öt óriásszoftver mellett felsorolás szinten szólnunk kell még az alábbi programokról: 4Thought, Alice, Darwin, Datascope, Scenario, Data Surveyor & Expert Surveyor.

### 13.3.1. Adatbányászati rendszerek tulajdonságai

Az előzőekben felsoroltunk néhány adatbányászati szoftvert. A felsoroltakon kívül léteznek még további szoftverek, amelyek bizonyos tekintetben akár jobbak is lehetnek a fentieknél. Ekkora választékban hogyan tudjuk megtalálni a nekünk megfelelő szoftvert, mik azok a tulajdonságok, amit mindeképpen meg kell vizsgálnunk egy ilyen beruházás előtt.

**Adatbányászati funkciók.** Egy cég azért vásárol adatbányászati szoftvert, mert összefüggést akar kinyerni az adataiból. Már a szoftv vásárlás előtt hasznos, ha pontos elképzelése van arról, hogy milyen típusú összefüggéseket fognak keresni (asszociációs szabályok, epizódok, klaszterek stb.). A legfontosabb, hogy a szoftver funkciói között megtalálhatók legyenek az ilyen típusú összefüggések kinyerésének lehetősége.

Nem biztos, hogy a nekünk megfelelő szoftver lesz a legtöbb adatbányászati feladat megoldását támogató. Egyre több szoftver jelenik meg, amely egy adott feladatra szakosodik (pl.: weblog elemző szoftver), ugyanakkor az átfogó képességgel rendelkezők mellett szól, hogy a jövőre is célszerű gondolni: milyen típusú összefüggéseket keresünk esetleg később.

**Adattípus.** A legtöbb szoftver a relációs adatbázisokban található adatokat tudja feldolgozni, de ezenkívül a sima szövegfájlt, munklapokat, ismertebb formátumú fájlokat is kezelik. Fontos tehát ellenőrizni, hogy pontosan milyen formátumú adatokon dolgozik. Ma már léteznek szoftverek, amelyek speciális adatformátumokat is kezelni tudnak, mint például földrajzi, multimédiás, web logok, DNS adatbázisok.

**Adatforrás.** Vannak adatbányász szoftverek, amelyeket fel kell tölteni az adatokkal mielőtt dolgozni lehet velük. Hasznosabb azonban, ha a szoftver a más adatbázisokban található adatokat is kezelni tudja. Fontos, hogy a rendszer támogassa az ODBC kapcsolatot vagy az OLE DB for ODBC-t. Ez lehetővé teszi a hozzáférést sok más relációs adatbázishoz (DB2, Informix, Microsoft SQL Server, Microsoft Access, Excel, Oracle stb.).

**Adatméret, skálázhatóság.** Tudnunk kell, hogy a szoftver mekkora adattal képes megbirkózni továbbá, hogy az adatbázis növelésével hogyan romlik a futási idő. Skálázhatóság szempontjából megkülönböztetünk *sor szerint skálázható* és *oszlop szerint skálázható* szoftvereket. Az első azt jeleti, hogy ha megduplázom a sorok számát, akkor nem nő duplájára a futási idő/memória igény. Az oszlop szerint skálázhatóság szerint a futási idő/memória igény az oszlopok számával lineárisnál nem rosszabb. Ez utóbbi feltétel teljesüléséhez kifinomultabb algoritmusokra van szükség.

**Megjelenítési eszközök.** A vizualizáció egy külön szakma. Az adatbányászati algoritmusok eredményeinek áttekinthető, szemléletes megjelenítése sokat segít az értelmezésben. A 3D ábrák, grafikonok, táblázatok nagyon hasznosak és sokat segítenek az adatbányászat használhatóságában és az eredmények interpretálhatóságában.

Az adatbányászat nagyon fiatal tudományág, így a szoftverek sem tekinthetnek vissza nagy múltra. A szoftverek szinte minden tekintetben különböznek egymástól. A megjelenítéssel, adatbányászati funkciókkal, terminológiával kapcsolatos egységes koncepció kialakulásáig még várnunk kell.

### 13.3.2. Esettanulmányok röviden

A következőkben vázolunk néhány sikeres adatbányászati projektet [95]<sup>1</sup>.

---

<sup>1</sup>Egyes részeket a BME diákjai fordították.

## Szlovén médiaszokások feltárása

Ma a médiumok kezében óriási hatalom van mind politikai, mind üzleti értelemben. Az egyes újságok, tv műsorok „fogyasztóinak” megismerésével közvetlenül elérhetik az egyes cégek a célközönségeiket.

A szlovén Mediana mintegy 8000 (20 oldalas!) kérdőív adatait elemeztette adatbányászati módszerekkel. Az adatok kitűnő minőségűek voltak és rengeteg attribútumot tartalmaztak többek között az egyes személyek különböző médiumokhoz fűződő viszonyát, a személyek érdeklődési körét, életstílusát, anyagi helyzetét, demográfia adatait (lakásának, munkahelyének fekvése). Az elemzések során a következő kérdésekre keresték a választ:

- mely más újságot/magazint olvasnak még szívesen bizonyos nyomtatott médiumok olvasói,
- mi jellemző az olvasóira/hallgatóira/nézőire az egyes médiumoknak,
- milyen tulajdonságok különböztetik meg a különböző újságok olvasóit,
- az ügyfeleiket tekintve mely médiumok hasonlóak?

A kérdések megválaszolásához számos adatbányászati módszert használtak fel, csak úgy mint korreláció-elemzés, klaszterezés, döntési fák, asszociáció szabályok, Kohonen hálók. Például döntési fák segítségével próbálták megtudni, hogy jellemzően kik olvassák a 'Delo' és a 'Slovenske Novice' újságokat. A kinyert szabályokból két példa: „A Delo tipikus olvasója egy héten több alkalommal olvas újságot, az átlagnál magasabb az iskolai végzettsége, ismeri a különböző magazinokat, autó és sörmárkákat, szeret tv-zni, stb.” ezzel szemben a „Slovenske Novice olvasói szertenc kávézóknak és bárókban időzni, kevésbé tájékozottak márkák ismeretében, mint a Delo olvasói, és általában olvassák még a Slovenski Delnicar, Jana, stb. magazinokat is.”

Klaszterezéssel profilszortokat hoztak létre. Kohonen háló segítségével megállapították a csoportok számát, majd a  $k$ -közép algoritmussal létrehoztak négy klasztert. A kapott klaszterek sajátosságait ezután döntési fák segítségével próbálták felderíteni. Az egyes csoportokat a jellemzők meghatározása után a „elhivatott fiatalok”, „inaktív idősek”, „ambíciózus emberek” és „aktív idősek” jelzőkkel illették. Például az „inaktív idősek” csoportját jellemzi, hogy nem szeretik a kihívásokat, nem érdekli őket a szórakoztatóipar, tudomány és technika, a fő örömforrásuk a család és nem szeretik a változásokat.

## Az Egyesült Királyság baleseteinek elemzése

A baleseteket leíró adatbázisokból kinyert hasznos információk életet menthetnek meg. Az okozatok, kapcsolatok feltárása olyan közúti vagy közlekedési szabályokat érintő módosításokhoz vezethet, amelyek segítségével megelőzhetőek a balesetek anélkül, hogy az agyonszabályozások következtében megnehezítenék az autósok életét.

Az Egyesült Királyság adatbázisának elemzése egy nagyon sikeres adatbányászati projekt volt. Az adatbázis az 1979-1999-ig terjedő intervallum adatait, mintegy ötmillió rekordot tárolt. Minden rekordhoz tárolták a baleset körülményeit, az autó, ill. a vezető továbbá az esetleges sérülések adatait.

Az elemzéshez vizualizációs eszközöket, számos klasszikus statisztikai (pl. regresszió) és adatbányász módszert alkalmaztak. A balesetek körülményei szöveggel voltak megadva, ezért ezek feldolgozásában kiemelt szerepet kaptak a szövegbányász módszerek.

Egy érdekes és újszerű megoldás volt a földrajzi helyek klaszterezése, melynek során a hasonló baleseti dinamikával rendelkező helyek kerültek egy csoportba. Különböző időfelbontásokhoz (évi/heti/napi balesetszám alakulás) különböző klaszterezést készítettek. Észrevették például, hogy az olyan esetek amelyek a havi szinten emelkedő baleseti számmal rendelkeznek és a balesetek száma nyáron éri el a maximumot megegyeznek a közkedvelt turisztikai helyekkel. Ezeken a helyeken tehát csak a főszezonban szükséges emelni a biztonsági szintet. További fontos csoportot jellemezett az a görbe, amely napközben és hétvégén alacsony szinten volt, de a munkaidő utáni időszakban megugrott. Ez a görbe az „ipari” területekre volt jellemző.

Vegyük észre, hogy a közlekedési balesetekre nagyon jellemző a lokalitás vagy más szavakkal az ideiglenes gyakoriság. Ez azt jelenti, hogy összességében kevés baleset történik, viszont a kevés baleset nagy része ugyanabban az időben (például hóesésben, vagy munkaidő után) esik meg. Ezeket az eseteket naív módon, gyakori mintákat kinyerő algoritmusokkal nem lehetne felderíteni, hiszen az összes baleset száma csekély.

A balesetek súlyosságának megállapítására felállított döntési fa is számos értékes összefüggéssel szolgált. Megmutatta például, hogy az 20 óra után történt motorkerékpárossal történt balesetekben a súlyos sérülések aránya jóval magasabb, mint általában.

### Portugál Statisztikai Hivatal weblapjának elemzése

Az internet rohamos fejlődésével egyre bővül az elérhető információ mennyisége, így folyamatosan nagyobb szerephez jut a megfelelő adatok keresése és kiválasztása. Ezen szempontok figyelembe vételével fejlesztett weblapok nagyban megkönnyítik a felhasználók dolgát, ezért döntött a Portugál Statisztikai Hivatal is az oldalát látogatók szokásainak elemzése mellett. Három fő célt jelöltek meg: ajánlattevő rendszer fejlesztése, felhasználói profilok kialakítása, weblap vizualizáció.

A log fájlban tárolt adatok (3GB) jelentős szűrésen estek át, mivel csak regisztrált felhasználók azonosíthatóak egyértelműen, a további információk megtévesztőek lehetnek. Az ajánlattevő rendszer fejlesztése során a rendelkezésre álló adatokból felhasználó,oldalj párokat hoztak létre, és ezekből vezettek le asszociációs szabályokat, aminek segítségével minden oldalhoz meghatározták a legjobban hasonlító  $N$  oldalt. A honlap architektúráját tekintve három rétegű volt: téma, altéma, fejezet. Ezekre külön modelleket hoztak létre, és külön tesztelték őket. A teszt során egy felhasználó által látogatott oldalak közül egyet kivéve vizsgálták, hogy a rendszer milyen arányban ajánlja a hiányzó oldalt. Az eredmények bizonyították az ajánlattevő rendszer használhatóságát, különösen kis  $N$ -ekre értek el jelentős javulást az adatbányászati módszereket nem alkalmazó rendszerekhez képest ( $N=1$  recall/recall.default=3).

A felhasználói profilok kialakításának alapötlete, hogy hasonló érdeklődési körű felhasználók nagyjából hasonló oldalakat látogatnak. Ha minden felhasználóhoz hozzárendelünk egy URL-vektort, ami az általa felkeresett oldalak címét tartalmazza, akkor ezek klaszterezésével felhasználói csoportok alakíthatók ki. K-means algoritmus segítségével 10 csoportot különítették el, amelyeket a rájuk legjobban jellemző oldalakkal írtak le. Ezekkel az eredményekkel új oldalról közelíthető meg az ajánlattevő rendszer, ugyanis két oldal "jobban" hasonlít, ha azonos csoportba tartozó felhasználók látogatják őket, a módszer neve collaborative-filtering.

A honlap szerkezetének vizualizációjához magukat az oldalakat kellett csoportosítani, tartalom alapján klaszterezni. Kétféle megoldás készült: egy gráf alapú, és egy hierarchikus. Gráf megjelenítés esetén a csomópontok jelölik a klasztereket, kulcsszavakkal jellemezve, míg a kapcsolatokra kerülnek az adott csoportok hasonlóság értékei. Ezeket a központi vektorok cosinus-távolságával számolták ki. A hierarchikus klaszterezés csoportokat képez a létrejött 20 klaszterből. A módszer során változó méretű klaszterek jönnek létre (különböző számú oldalt tartalmaznak), ezeket téglalapokkal jelölik, a hasonlóságot pedig távolságuk adja, ami hasonló módon számolható, mint az előző esetben.

### Döntéstámogató rendszerek alkalmazásai

A következő 5 döntéstámogató rendszer Szlovéniában került alkalmazásra, mindegyik más-más jellemvonásokkal rendelkezik. Lakástámogató program: A feladat bankok megbízása államilag támogatott hitelek nyújtására. A konstrukció rendkívül kedvező a bankok számára, minél több szerződésre szeretnének jogot szerezni. A projekt nagy anyagi kereteit figyelembe véve mindenképp egy átlátható modell szükséges, a döntést követő támadási felületek csökkentésére. A nehézséget a mindössze egy hónapos határidő jelentette. A modell létrehozása során meghatározták a bankoktól bekérendő adatokat (hard data, magyarázó attribútumok), majd ezeket csoportosítva új, diszkrét tulajdonságokat hoztak létre (magyarázandó attribútumok). A diszkrét értékeket meghatározó függvény előállítása szakértők bevonásával történt. Ezek alapján már hozzárendelhető minden bankhoz egy prioritás, amit a bank méretével súlyozva alakulnak ki a kiosztott szerződésszámok.

Lakásfelújítási program: Lakótelepek renoválására írtak ki pályázatot, melyek elbírálásához kértek döntéstámogató rendszereket nyújtotta segítséget. A bekért adatokból a következő aggregált tulajdonságokat hozták létre: az épület állapota, a jelentkező adatai, a jelentkező státusza. Utóbbi esetén külön kezelték a tulajdonos által lakott épületeket és bérbe adottakat. A modell segítségével két lépésben összesen 250 jelentkezést fogadtak el.

Betegek állapotelemzése: Cukros betegek állapotának felmérése után döntéstámogató módszerek segítségével határozták meg az új betegek rizikófaktorait és javasolt kezelési módszerét. A projekt időtartama 3 év, a modell kialakításakor 3500 beteg adatait dolgozták fel orvosszakértők bevonásával. A főbb tulajdonságok a kórtörténet, a jelenlegi státusz és a teszteredmények voltak, ezek kiértékelési függvényét adatbányászati módszerekkel határozták meg a kiindulási adatokból. Az új betegek állapotának alakulásával párhuzamosan frissítették a modellt a nagyobb hatékonyság elérése érdekében.

Minőségelemzés, ajánlat kiválasztás: A szlovén informatikai hivatal két folyamatának automatizálását tűzte ki célul: beszállítók ajánlatainak összehasonlítása, a megvalósítási technikák összehasonlítása. A lehetőségek összevetése nagy szakmai tapasztalatot igényelt, sok informatikai szakértő bevonására volt szükség. A létrehozott modellek közös tulajdonsága volt a nagyon sok attribútum (18-19 alap és 10-12 aggregált). A tesztelés után éles alkalmazásra nem került sor, a modellek későbbi felhasználásra készültek.

Kulturális pályázatok elbírálása: Egyszeri pályázati támogatások kiosztására hoztak létre döntéstámogató rendszert. A nehézséget az adatok jellege jelentette, ugyanis a pályázatok szöveges formában (soft data) kerültek beadásra. Az elemzést két független szakértő végezte minden pályázat esetében, majd ezeket összevetve állapították meg a döntési modell alapulajdonságainak értékeit. A létrehozott modell szintén sok attribútumot tartalmazott (15 alap, 9 aggregált). A rendszer második fázisa segít kiküszöbölni a szubjektivitást, de a szakértői

vélemények támadhatók, ami a fellebbezések nagy számával járt.

### Terhelés előrejelzés

A villamosenergia-szolgáltató iparban rendkívül fontos szerepet játszik a jövőbeli igények minél korábbi felmérése. Minden év minden egyes szakának, hónapjának, napjának és órájának minimum és maximum terhelésére vonatkozó pontos becsléseinek birtokában, a közművek jelentős megtakarításokra tehetnek szert a működési tartalék beállítása, a karbantartás-ütemezés és a tüzelőanyag-készlet management területein.

Az egyik nagy szolgáltatónál az elmúlt évtizedben egy automatizált terhelés-előrejelző modul működött a következő két nap órákra lebontott terhelésének becslésére.

A szolgáltató első lépése a megelőző tizenöt év összegyűjtött adataiból egy szofisztikált terhelésmodell manuális megtervezése volt. A modell három komponensből épült fel: éves alapterhelés, évközi periodikus terhelés, és az ünnepnapok extraterhelése. Az alapterheléshez való optimalizációra az adatokat az óránkénti aktuális terhelésből az éves átlagterhelést levonva, majd a kapott értéket az éves szórással leosztva standardizálták. Az elektromos terhelés három ciklus szerint mutat periodicitást: napi (itt a terhelés minimuma reggelre, maximuma pedig délre és délutánra esik), heti (hétvégente mutat alacsony értéket) és évszakonként (a nyaranta és telente megnövekvő fűtési illetve hűtési igény miatt). A nagyobb ünnepnapok, mint a Hálaadás napja, Karácsony vagy az Újév napja jelentős eltérést mutatnak az egyébként szokásos terhelésadatoktól, így ezeket külön-külön modellezték, a megelőző tizenöt év adott napján és órájában mért átlagokkal. A kisebb állami ünnepnapokat, például a Kolombusz napját egy kalap alá veszik az iskolai szünnapokkal és a normál napi minta járulékos terheléseként kezelik. Mindezek a hatások megjelennek az évet tipikus napok sorozataként rekonstruálva, az ünnepnapokat helyükre illesztve, denormalizálva a teljes növekményt kitevő terhelést.

Mindeddig a terhelésmodell statikus volt, amelyet a korábbi adatok felhasználásával, manuálisan szerkesztettek meg, és implicite magában hordozta azt a feltevést, hogy a klímaviszonyok nem térnek el a „normálistól” az év során. Az utolsó lépés az időjárás, mint külső faktor modellbe illesztése volt. A módszer a mentett adatok közt az aktuálishoz hasonló időjárási körülmények után kutat, és a megtalált nap értékeit használja a napi terhelés előrejelzéséhez. Ebben az esetben a rendszer a becsléssel korrigálja a statikus terhelésmodellt. Az extrém értékek elleni védelemre a rendszer a nyolc leginkább hasonlatos nap eltéréseinek átlagával korrigál. Óránkénti felbontású adatbázist hoztak létre három helyi meteorológiai állomás 15 évre visszamenő hőmérséklet, páratartalom, szélsősebesség és felhőtakaróra vonatkozó adataiból, az aktuális terhelés és a statikus modell által jósolt terhelés különbségével kiegészítve. Az iménti paraméterek terhelésre gyakorolt hatását lineáris regresszióanalízissel számították ki, majd az együtthatókkal súlyozták a hasonló napok keresésére használt távolságfüggvényt. Az így kapott rendszer ugyanolyan teljesítményt produkált sokkal gyorsabban, mint a képzett meteorológusok, órák helyett percek alatt készítve el a napi előrejelzést. A rendszer operátorai tesztelhetik az előrejelzés érzékenységét az időjárás szimulált változásainak hatására, és megvizsgálhatják az aktuálishoz leghasonlóbb napokat, amelyeket az algoritmus a finomhangoláshoz felhasznált.

## Diagnosztika

A szakértői rendszerek legfontosabb alkalmazási területe a diagnosztika. Bár sokszor a kézzel beállított szabályok is jól teljesítenek a szakértői rendszerekben, a gépi tanulás lehetősége számos lehet olyan esetekben, amikor a manuális szabályalkotás túlságosan munkaigényes.

Az elektromechanikus eszközök (például motorok és generátorok) megelőző szervizelését vehet az ipari folyamatokat megzavaró hibák kialakulásának. A technikusok rendszeresen felülvizsgálják minden eszközt, különböző pontokon mérve a rezonanciát, hogy felmérjék mely eszközök szorulnak szervizelésre. Tipikus meghibásodások közé tartoznak: a tengely elállítódása, mechanikai kilazulások, hibás csapágyak és kiegyenlítetlen szivattyúk. Az egyik vegyi üzem több mint 1000 különböző eszközt használ, a kisebb szivattyúktól egészen a hatalmas turbógenerátorokig, amelyek mindegyikét egészen a közelmúltig egy 20 éves tapasztalattal rendelkező szakember szervizelte. Az eszköz talapzatán, különböző pontokon végeznek vibrációmérést, és az energiaszintet vizsgálják Fourier-analízis segítségével az alap forgási sebesség minden felharmonikusának mindhárom irányában a hibák detektálására. Ennek, az – a mérési és rögzítési folyamat korlátoltsága miatt rendkívül nagy hibaarányú – információnak a tanulmányozását a diagnosztizáló szakember végzi. Néhány szituációra ugyan manuálisan kifejlesztettek szabályrendszereket, de a fejlesztési folyamatot számos különböző berendezésre külön-külön meg kellett volna ismételni, így került látóterbe a gépi tanulás.

Hatszáz hiba, mindegyik mérési adatokkal és a szakértő diagnózisával rendelkezésre állt, a problémakör húszévnnyi tapasztalataként. Az adatok körülbelül fele különböző okok miatt elégtelen volt, így figyelmen kívül hagyták, a maradékot pedig tanítóhalmaznak használták. A cél nem a hiba detektálása, hanem annak kategorizálása, diagnosztizálása volt. Így nem volt szükség hibamentes esetek adataival bővíteni a tanítóhalmazt. Mivel a mért attribútumok meglehetősen alacsony szintűek voltak, így ki kellett bővíteni azokat a származtatott – terület-specifikus információval bíró – fogalmakkal, például az alap tulajdonságok funkcióival, amelyeket a szakértő bevonásával definiáltak. A származtatott attribútumokra lefuttattak egy indukciós algoritmust, hogy előállítsák a diagnosztizáló szabálykészletet. Kezdetben a szakértő nem volt elégedett a szabályokkal, mert nem tudta azokat a saját tudásához és tapasztalatahoz viszonyítani. Számára a pusztán statisztikai megalapozottság önmagában nem volt elegendő bizonyíték. További háttértudást kellett felhasználni mielőtt elkészülhetett a megfelelő szabályrendszer. Ugyan az eredményül kapott szabályok meglehetősen bonyolultultra sikerültek, a szakértőnek tetszettek, mert mechanikai tudása és tapasztalata alapján ellenőrizni tudta azokat. Örült, hogy a szabályok harmada egybeesett az általa is használtakkal, a maradék egy része pedig szélesítette rálátását a rendszerre.

A teljesítménylesztek azt mutatták ki, hogy az újonnan kinyert szabályok alig voltak jobbak a korábban manuálisan felállítottaknál. Az eredményt később a vegyi üzem is megerősítette. Ugyanakkor érdemes megjegyezni, hogy a szabályrendszert nem a jó teljesítménye miatt állították üzembe, hanem mert a terület szakértője elismerően vélekedett róla.



# Függelék

## Függelék A

**.1. tétel.** *A Gyűjtőlapok és Tekintélyek során alkalmazott iteráció során  $t^{(i)}$ , illetve  $g^{(i)}$  sorozatok konvergálnak nemnegatív értékű vektorokhoz. Tehát lássuk be, hogy amennyiben  $A$  egy tetszőleges gráf adjacencia mátrixa és  $v^{(0)} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} = j^t$ , akkor a*

$$v^{(i)} = \frac{AA^T v^{(i-1)}}{[AA^T v^{(i-1)}]}$$

*iteráció által kapott sorozat konvergál.*

Megjegyzés 1: Az iterációs lépésből közvetlenül adódik, hogy  $v^{(i)}$  az  $(AA^T)^i j^t$  irányú egységvektor.

Megjegyzés 2:  $g^{(i)}$  konvergenciájából  $t^{(i)}$  konvergenciája is következik  $A$  és  $A^T$  felcserélésével. A tétel bizonyításához szükségünk van néhány segédtétele.

**.2. lemma.** *Legyen  $A \in \mathbb{R}(n \times n)$ . Ekkor  $AA^T$  (és hasonlóan  $A^T A$  is) pozitív szemidefinit szimmetrikus mátrix.*

*Bizonyítás:* A szimmetrikusság a mátrixszorzás szabályából közvetlenül adódik. Felhasználva a  $vA = (A^T v^T)^T$  azonosságot

$$vAA^T v^T = (A^T v^T)^T (A^T v^T) = w^T w \geq 0$$

adódik, ami bizonyítja, hogy  $AA^T$  pozitív szemidefinit. ■

**.3. lemma.** *Ha  $M$  mátrix pozitív szemidefinit és szimmetrikus, akkor sajátértékei valósak és nemnegatívak.*

**.4. tétel (Perron-Frobenius).** *Ha egy mátrix aperiodikus, irreducibilis és nemnegatív elemű, akkor legnagyobb abszolútértékű sajátértékhez tartozó sajátvektor nemnegatív koordinatájú, és nincs más, ilyen abszolút értékű, sajátérték.*

**.5. lemma.**  *$M$  mátrix pozitív szemidefinit szimmetrikus,  $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_k \geq 0$ , ( $k < n$ ) sajátértékekkel. Ekkor tetszőleges  $v \in \mathbb{R}^n$  felírható  $v = \sum_{i=1}^k \alpha_i w^{(i)}$  alakban, ahol  $\|w^{(i)}\| = 1$ ,  $w^{(i)} w^{(j)} = 0$  ha  $i \neq j$  és  $Mw^{(i)} = \lambda_i w^{(i)}$ .*

Térjünk vissza az .1-ös tétel bizonyításához.

*Bizonyítás:*

Jelöljük  $AA^T$  mátrixot  $M$ -el. Feltehetjük, hogy  $M$  aperiodikus, hiszen  $m_{ii}$  az  $i$ -edik pontból más pontba mutató élszám négyzetének összegét adja meg ( $\sum_k m_{ik}^2$ ), ami csak akkor lehet 0, ha  $i$ -edik pontból nem indul él. Ez a pont a konvergencia tényét nem befolyásolja, mert  $M$  minden hatványának megfelelő sora és oszlopa csupa 0 elemből fog állni, tehát jogos a feltételezés. Azt is feltehetjük, hogy  $M$  irreducibilis, mert ha nem az, akkor mátrixot irreducibilis blokkmátrixokra bonthatjuk, és a hatványozást blokkonként végezhetjük.

Tudjuk tehát, hogy  $M$  nemnegatív elemű, aperiodikus, irreducibilis, pozitív szemidefinit szimmetrikus mátrix, ami miatt minden sajátérték nemnegatív, a legnagyobb sajátértéke egyszerű, továbbá az ehhez tartozó sajátvektor nemnegatív elemű. Legyen  $v \in \mathbb{R}^n$  tetszőleges vektor. .5 alapján  $v = \sum_{i=1}^k \alpha_i w^{(i)}$  és  $w^{(1)}$  egyértelmű, nemnegatív elemű vektor. A  $\frac{M^j v}{\|M^j v\|}$  kifejezés  $w^{(1)}$ -hez tart ha  $j \rightarrow \infty$ , mert

$$\frac{M^j v}{\|M^j v\|} = \frac{\sum_{i=1}^k \alpha_i M^j w^{(i)}}{\|\sum_{i=1}^k \alpha_i M^j w^{(i)}\|} = \frac{\sum_{i=1}^k \alpha_i \lambda_i^j w^{(i)}}{\sqrt{\sum_{i=1}^k (\alpha_i \lambda_i^j)^2}}$$

$$\frac{\alpha_1 \lambda_1^j w^{(1)} + \sum_{i=2}^k \alpha_i \lambda_i^j w^{(i)}}{\sqrt{(\alpha_1 \lambda_1^j)^2 + \sum_{i=2}^k (\alpha_i \lambda_i^j)^2}} \cdot \frac{1}{\lambda_1^j} = \frac{\alpha_1 w^{(1)} + \sum_{i=2}^k \alpha_i \left(\frac{\lambda_i}{\lambda_1}\right)^j w^{(i)}}{\sqrt{\alpha_1^2 + \sum_{i=2}^k \left(\alpha_i \left(\frac{\lambda_i}{\lambda_1}\right)^j\right)^2}} \rightarrow w^{(1)}$$

A normálás során felhasználtuk, hogy a  $w^{(i)}$  vektorok merőlegesek egymásra, és egységnyi hosszúak, a határérték meghatározásakor pedig azt, hogy  $\lambda_1$  a legnagyobb sajátérték, tehát  $\frac{\lambda_i}{\lambda_1} < 1, i = 2, \dots, k$ -ra. Tehát ha  $v$  nem merőleges  $w^{(1)}$ -re, akkor  $\frac{M^j v}{\|M^j v\|}$  vektor  $w^{(1)}$ -hez konvergál. Ez azonban nem áll fenn, lévén  $j w^{(1)} > 0$ , mert  $w^{(1)}$  nemnegatív elemű vektor.

■

ANGOL	MAGYAR
antecedent	feltételrész
approximate dependency	közelítő függőség
association rule	asszociációs szabály
authority	tekintélylap
basket	kosár
candidate	jelölt
classification	osztályozás
confusion matrix	keveredési mátrix
consequent	következményrész
clustering	klaszterezés
confidence	bizonyosság
conviction	meggyőződés
data mining	adatbányászat
dead end problem	zsákutca probléma
decision rule	döntési szabály
decision tree	döntési fa
dense	sűrű
episode	epizód
false-positive	hamis jelölt
false-negative	hiányzó elem
frequent	gyakori
gain ratio	nyereségarány
goodness-of-split	vágás jósága
hash-tree	hash-fa
hub	gyűjtőlap
impurity-based criteria	
item	elem
knowledge retrieval	tudásfeltárás
kurtosis	lapultság
levelwise	szintenként haladó
lift	függetlenségi mutató
locality-sensitive hashing (LSH)	hely-érzékeny hashelés (HÉH)

1. táblázat. Idegen kifejezések fordítása (a-l)

ANGOL	MAGYAR
market-basket problem	piaci-kosár probléma
mode	módusz
negative border	esélyes jelölt
oblivious decision tree	hanyag döntési fák
outlier analysis	eltérés elemzés
pattern	minta
potpruning	utóőnyesés
power divergence function	erő divergencia függvény
prepruning	előnyesés
principal component analysis	főkomponens analízis
product	termék
ranking	rangsorolás
replicated subtree problem	ismétlődő részfa probléma
sequence matching	sorozatillesztés
signature	lenyomat
singular value decomposition	szinguláris felbontás
skewness	ferdeség
sparse	ritka
spider trap problem	pókháló probléma
stripped partition	redukált partíció
support	támogatottság
threshold	küszöb
transaction	tranzakció
valid	érvényes
z-score normalization	standard normalizálás

2. táblázat. Idegen kifejezések fordítása (m-z)

# Tárgymutató

- $\chi^2$  próba, 29
- átlagos négyzetes hiba, 141
- átlagos négyzetes hibagyök, 141
- min\_supp*, 60
- APRIORI, 63
- ECLAT algoritmus, 76
- FP-GROWTH algoritmus, 79
- 11 pontos átlagos pontosság, 176
  
- A minta nagysága, 84
- abszolút hibaátlag, 141
- AdaBoost, 175
- adat
  - strukturálatlan, 165
  - strukturált, 165
  - tanuló, 171
  - teszt, 171
  - validációs, 171
- adatbázis
  - horizontális, 61
  - vertikális, 61
- algoritmus
  - helyesen működő, 221
  - mohó, 177
  - teljes, 221
- anti-monoton, 216
- APRIORI
  - módszer, 221
- apriori algoritmus, 197
- APRIORI-CLOSE, 224
- asszociációs szabály, 90, 91
  - érdekessége, 93
  - érvényes, 91
  - bizonyossága, 91
  - egzakt, 91
  - hierarchikus, 101, 102
  - támogatottsága, 91
- average linkage, 159
  
- Bayes-módszer
  - naiv, 171, 177
- bemeneti sorozat, 60
- bizottság
  - osztályozóké, 174
  - tagok, 174
- boosting eljárások
  - AdaBoost, 175
  
- centroid kapcsolódás, 183
- centroid-egyszerű kapcsolódás, 183
- $\chi^2$ -statisztika, 169
- complet linkage, 159
- csoportosítás
  - szövegeké, 180
    - hierarchikus klaszterezők, 183
    - jellegzetességek, 181
    - k-átlag módszerek, 183
  
- dimenzió csökkentése, 168
  - kategorizálásnál, 169
- dokumentum
  - ábrázolása, 167
  - előfeldolgozása, 166
  - reprezentációja, 166
    - bináris, 168
    - csoportosításnál, 181
- dokumentum frekvencia küszöbölő, 169
- dokumentumgyűjtemény
  - reprezentálása, 168
- dokumentumok csoportosítása, 180
- dokumentumok előfeldolgozása, 166
- döntési fa
  - szövegosztályozó, 172
- Duquenne–Guigues-bázis, 92
  
- egyensúlyi pont
  - felidézése és pontosság, 176
- ekvivalencia-reláció, 23

- elemhalmaz, 237
  - fedés, 60
  - gyakori, 60
  - gyakorisága, 61
- eloszlás
  - $\chi^2$ , 27
  - binomiális, 26
  - hipergeometrikus, 27
  - normális, 27
  - Poisson, 26
- entrópia, 28
- entrópia, 182
- Euklideszi-norma, 40
- függetlenségvizsgálat, 30
- főkomponens analízis, 54
- felügyelet nélküli tanulás, 145
- felidézés, 143
- felidézés, 175, 187
  - szintenkénti, 178
- felügyelet nélküli tanulás, 181
- felügyelt tanulás, 171
- ferdeség, 27
- F-mérték, 176
  - csoportosításnál, 182
  - szintenkénti, 178
- fogalomtársítás, 195
- fontosság, 202
- FP-fa
  - vetített, 81
- funkció szavak, 169
- funkció szavak
  - elhagyása, 169
- funkcionális függőség, 263
- FUP algoritmus, 234
- Galois-kapcsolat, 86
- Galois-lezárás operátor, 86
- GSP, 240
- gyűjtőlap, 206
- gyakorisági küszöb, 225
- gyakorisági küszöböt, 61
- gyakoriság, 168
- halmaz, 23
  - lokálisan véges, 216
  - rangsámozott, 216
- halmazcsalád, 73
- hatékonyság mérése
  - szövegbányászatnál általában, 170
  - szövegek csoportosításánál, 182
- szövegosztályozás
  - egyszerű, 175
  - hierarchikus, 178
- hiba
  - szövegosztályozásnál, 175
- hierarchikus asszociációs szabály
  - érdekessége, 103
- hierarchikus klaszterező, 183
  - egyesítő, 183
  - felosztó, 183
- ierarchikus klaszterező
  - UPGMA, 183
- hierarchikus szövegosztályozás, 177
- HITEC, 178
- Hoeffding-korlát, 27
- információ nyereség módszer, 169
- információkinyerés, 195
- invariáns hasonlóság, 39
- inverz dokumentum frekvencia, 168
- Ismételt mintavételezés, 138
- Jaccard-koefficiens, 40
- jelölt, 63, 222
  - hamis, 222
- jelölt-előállítás
  - ismétlés nélküli, 63
- jellemzők kiválasztása, 169
- k-legközelebbi szomszéd gráf, 149
- közelítő függőség, 263
- kényszer
  - erősen átalakítható, 219
- kanonikus reprezentáció, 73
- kappa statisztika, 143
- kategóriaösvény, 178
- kategóriarendszer, 170
- kategorizálás *lásd* osztályozás 170
- k-átlag eljárás
  - kettészeli, 183
- kérdés-megválaszoló rendszerek, 196
- Kereszt-validáció, 139
- kettészeli k-átlag eljárás, 183

- keveredési mátrix, 142
- kivonatolás, 185
  - csoportosítás alapú módszerek, 190
  - definíció, 185
  - hatékonyságának mérése, 187
  - jellemzők, 187
  - klasszikus módszer, 188
  - MEAD módszer, 191
  - MMR módszer, 190
  - mondatkiválasztással, 187
  - TF-IDF alapú módszer, 189
  - weboldalake, 193
- Klaszterezés, 144
- klaszterezés *lásd* csoportosítás 180
- k-NN *lásd* legközelebbi szomszédok 172
- kontingencia-táblázat, 30
- korrelációs együttható, 141
- koszinusz-mérték, 42
- kulcs, 265
  
- Laplace estimation, 135
- lapultság, 27
- látens szemantikus indexelés (LSI), 169
- leave-one-out, 139
- legközelebbi szomszédok
  - szövegosztályozó, 172
- lexikografikus rendezés, 24
- lexikon *lásd* szótár 167
- lineáris kiterjesztés, 217
- lineárisan szeparálható osztályok, 110
- LSI *lásd* látens szemantikus indexelés 169
- lusta tanuló, 172
  
- Manhattan-norma, 40
- min\_freq, 61, 225
- Minkowski-norma, 40
- minta, 216
  - üres, 216
  - elhanyagolt, 222
  - gyakori, 216
  - gyakorisága, 225
  - jelölt, 222
  - mérete, 216
  - nem bővíthető, 218
  - ritka, 216
  - támogatottsága, 216
  - zárt, 218
- mintafelismerés, 169
- mintahalmaz, 216
- mintatér, 216
- mohó algoritmus, 177
  
- névelem, 197
- naiv Bayes-módszer, 171
  - hierarchikus osztályozás, 177
- Naiv mintavételező algoritmus, 231
- neurális hálózat, 173
  
- oldalak rangsorolása, 201
- „oszd meg és uralkodj” stratégia, 173
- osztályozás
  - egyszerű, 170
  - hierarchikus, 170
  - szövegeké, 170
- összegzőkészítés, 185
  - általános, 186
  - indikatív, 186
  - informatív, 186
  - kérdés-vezérelt, 186
  
- pókháló probléma, 204
- Page Rank, 202, 205
- partíció, 264
- partíció finomítása, 265
- partíciós algoritmus, 231
- PATRICIA fa, 33
- perceptron, 173
- pontosság, 175, 187
  - szintenkénti, 178
- Porter-algoritmus, 197
- próba
  - Student t-próba, 31
- predikátum
  - anti-monoton, 219
  - monoton, 219
  - prefix anti-monoton, 219
  - prefix monoton, 219
  - triviális, 219
- prefix, 217
- pszeudo-zárt elemhalmaz, 92
  
- részben rendezés, 23
- rész minta, 216

- valódi, 216
- rétegzett particionálás, 138
- rang-vektor, 202
- redukált partíció, 266
- relatív abszolút hiba, 141
- relatív négyzetes hiba, 141
- relatív négyzetes hibagyök, 141
- Reuters-gyűjtemény, 177
- Rocchio-eljárás, 112
  
- shrinkage, 177
- single linkage eljárás, 158
- sorozat, 24
- stopwords*, 169
- strukturálatlan adat, 165
- strukturált adat, 165
- súlybeállítás
  - additív, 173
  - multiplikatív, 173
- súlyozás
  - bináris, 168
  - TF, 168
  - TFIDF, 168, 189
- SVD *lásd* szinguláris értékbontás 169
- SVM, 174
- szófa, 31, 66
  - láncolt listás implementáció, 32
  - nyesett, 33
  - táblázatos implementáció, 32
- szabatosság, 175
- szavazásos osztályozás, 174
- szerkesztési távolság, 42
- szerkesztési elv, 190
- szeszélyes sztochasztikus szörfölő, 205
- szinguláris felbontás, 54
- szinguláris értékbontás (SVD), 169, 192
- szó–dokumentum mátrix, 168
- szófajcímkéző, 197
- szótár, 167
  - mérete, 168
  - méretének csökkentése, 168
    - kategorizálásnál, 169
- szótövező, 167, 197
- szövegbányászat
  - általános modellje, 166
  - definíció, 165
- szövegek kategorizálása, 170
- szöveges információk vizualizálása, 196
- szövegosztályozás, 170
  - hierarchikus, 177
- szövegosztályozó
  - bizottság, 174
  - döntési fa alapú, 172
  - HITEC, 178
  - legközelebbi szomszédokon alapuló, 172
  - naiv Bayes-módszer, 171, 177
  - neurális hálózat alapú, 173
  - SVM, 174
  - szavazásos, 174
- sztochasztikus szörfölő, 203
- szuperkulcs, 265
  
- támogatottsági függvény, 216
- támogatottsági küszöb, 60, 216
- TANE, 264
- tanulás
  - felügyelet nélküli, 181
  - felügyelt, 171
- tanulási ráta, 173
- tanulóhalmaz, 171
- taxonómia, 101
- taxonómia, 170, 177, 181
- tekintélylapok, 206
- teljes rendezés, 23
- témakövetés, 195
- teszthalmaz, 171
- tesztkorpuszok
  - szövegszo2vegklaszterezeshoz, 184
  - szövegosztályozáshoz, 198
- TID-halmaz, 76
- token, 167
- tranzakció, 60
  
- újraparametrizálás, 169
- univerzálisan népszerű lapok, 207
- UPGMA módszer, 183
  
- validációs halmaz, 171
- variáns hasonlóság, 39
- vektortér-modell, 167
  
- Ward módszer, 159
- Webes adatbányászat, 201



- weka
  - Associate fül, 92
  - Classify fül, 108
  - Arff formátum, 38
  - sparse arff formátum, 38
  - weka.associations.Apriori, 64, 104
- weka.associations
  - Apriori, 100
  - Conviction, 100
  - Leverage, 100
- weka.classifiers
  - Classifier evaluation options, 142, 144
  - Classifier output, 143
  - functions.LinearRegression, 131
  - functions.MultilayerPerceptron, 134
  - functions.SimpleLinearRegression, 131
  - functions.Winnow, 112
  - lazy.IB1, 116
  - lazy.IBk, 116
  - Result list panel, 122
  - rules.OneR, 120
  - rules.Prism, 121
  - rules.ZeroR, 120
  - Test options panel, 139, 142
  - trees csomag, 122
  - trees.Id3, 126
  - trees.J48, 129
  - trees.UserClassifier, 125
- weka.clusterers
  - DBScan, 163
  - SimpleKMeans, 156
- weka.filters.supervised
  - attribute.Discretize, 47
  - instance.Resample, 53
  - instance.SpreadSubsample, 53
  - instance.StratifiedRemoveFolds, 53
- weka.filters.unsupervised
  - attribute.Add, 44
  - attribute.AddExpression, 44
  - attribute.AddID, 44
  - attribute.AddNoise, 46
  - attribute.Center, 48
  - attribute.ChangeDateFormat, 38
  - attribute.Copy, 44
  - attribute.Discretize, 46
  - attribute.FirstOrder, 44
  - attribute.InterquartileRange, 45
  - attribute.MathExpression, 44
  - attribute.MergeTwoValues, 38
  - attribute.NominalToBinary, 38, 104
  - attribute.Normalize, 48
  - attribute.NumericCleaner, 45
  - attribute.NumericToNominal, 38
  - attribute.NumericTransform, 44
  - attribute.Obfuscate, 46
  - attribute.PKIDiscretize, 46
  - attribute.PrincipalComponents, 57
  - attribute.Remove, 45
  - attribute.RemoveType, 45
  - attribute.RemoveUseless, 45
  - attribute.ReplaceMissingValues, 43
  - attribute.Standardize, 48
  - instance.RemoveFolds, 53
  - instance.RemoveMisclassified, 45
  - instance.RemovePercentage, 53
  - instance.RemoveWithValues, 45
  - instance.Resample, 53
  - instance.ReservoirSample, 53
- Winnow, 173
  - kiegyensúlyozott, 173
- zárt elemhalmaz, 86
- zsákutca probléma, 204

# Irodalomjegyzék

- [1] Pieter Adriaans–Dolf Zantinge: *Adatbányászat*. Budapest, 2002, Panem Kiadó.
- [2] Ramesh C. Agarwal–Charu C. Aggarwal–V. V. V. Prasad: A tree projection algorithm for generation of frequent item sets. *Journal of Parallel and Distributed Computing*, 61. évf. (2001) 3. sz. URL <http://citeseer.nj.nec.com/agarwal99tree.html>.
- [3] Rakesh Agrawal–Tomasz Imielinski–Arun N. Swami: Mining association rules between sets of items in large databases. In Peter Buneman–Sushil Jajodia (szerk.): *Proceedings of the 1993 ACM SIGMOD International Conference on Management of Data* (konferenciaanyag). Washington, D.C., 1993. 26-28, 207–216. p. URL <http://citeseer.nj.nec.com/agrawal93mining.html>.
- [4] Rakesh Agrawal–Heikki Mannila–Ramakrishnan Srikant–Hannu Toivonen–A. Inkeri Verkamo: Fast discovery of association rules. In *Advances in Knowledge Discovery and Data Mining* (konferenciaanyag). 1996, 307–328. p.
- [5] Rakesh Agrawal–Ramakrishnan Srikant: Fast algorithms for mining association rules. In Jorge B. Bocca–Matthias Jarke–Carlo Zaniolo (szerk.): *Proceedings of the 20th International Conference Very Large Data Bases, VLDB* (konferenciaanyag). 1994. 12-15, Morgan Kaufmann, 487–499. p. ISBN 1-55860-153-8. URL <http://citeseer.nj.nec.com/agrawal94fast.html>.
- [6] Rakesh Agrawal–Ramakrishnan Srikant: Mining sequential patterns. In Philip S. Yu–Arbee L. P. Chen (szerk.): *Proceedings of the 11th International Conference on Data Engineering, ICDE* (konferenciaanyag). 1995. 6-10, IEEE Computer Society, 3–14. p. ISBN 0-8186-6910-1. URL <http://citeseer.nj.nec.com/agrawal95mining.html>.
- [7] Rényi Alfréd: *Valószínűségszámítás*. 1968, Tankönyvkiadó.
- [8] Brian Amento–Loren G. Terveen–William C. Hill: Does „authority” mean quality? predicting expert quality ratings of web documents. In *Research and Development in Information Retrieval* (konferenciaanyag). 2000, 296–303. p. URL <http://citeseer.nj.nec.com/417258.html>.
- [9] Amihood Amir–Ronen Feldman–Reuven Kashi: A new and versatile method for association generation. In *Principles of Data Mining and Knowledge Discovery* (konferenciaanyag). 1997, 221–231. p. URL <http://citeseer.nj.nec.com/amir97new.html>.
- [10] Franz Aurenhammer: Voronoi diagrams—a survey of a fundamental geometric data structure. *ACM Comput. Surv.*, 23. évf. (1991) 3. sz. ISSN 0360-0300.

- [11] Necip Fazil Ayan – Abdullah Uz Tansel – M. Erol Arkun: An efficient algorithm to update large itemsets with early pruning. In *Knowledge Discovery and Data Mining* (konferenciaanyag). 1999, 287–291. p. URL <http://citeseer.nj.nec.com/ayan99efficient.html>.
- [12] Yves Bastide – Rafik Taouil – Nicolas Pasquier – Gerd Stumme – Lotfi Lakhal: Mining frequent patterns with counting inference. *SIGKDD Explor. Newsl.*, 2. évf. (2000) 2. sz.
- [13] Jon Louis Bentley: Multidimensional binary search trees used for associative searching. *Commun. ACM*, 18. évf. (1975) 9. sz. ISSN 0001-0782.
- [14] B. Berendt – B. Mobasher – M. Spiliopoulou – J. Wiltshire.: Measuring the accuracy of sessionizers for web usage analysis, 2001.  
URL <http://citeseer.nj.nec.com/berendt01measuring.html>.
- [15] Alina Beygelzimer – Sham Kakade – John Langford: Cover trees for nearest neighbor. In *ICML '06: Proceedings of the 23rd international conference on Machine learning* (konferenciaanyag). New York, NY, USA, 2006, ACM, 97–104. p. ISBN 1-59593-383-2.
- [16] Krishna Bharat – Monika Rauch Henzinger: Improved algorithms for topic distillation in a hyperlinked environment. In *Research and Development in Information Retrieval* (konferenciaanyag). 1998, 104–111. p.  
URL <http://citeseer.nj.nec.com/bharat98improved.html>.
- [17] Ferenc Bodon: A fast apriori implementation. In Bart Goethals – Mohammed J. Zaki (szerk.): *Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*, CEUR Workshop Proceedings konferenciasorozat, 90. köt. Melbourne, Florida, USA, 2003. November 19..
- [18] Richard J. Bolton – David J. Hand: Significance tests for patterns in continuous data. In *Proceedings of the 2001 IEEE International Conference on Data Mining (ICDE)* (konferenciaanyag). 2001.
- [19] Christian Borgelt: Efficient implementations of apriori and eclat. In Bart Goethals – Mohammed J. Zaki (szerk.): *Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*, CEUR Workshop Proceedings konferenciasorozat, 90. köt. Melbourne, Florida, USA, 2003.
- [20] Christian Borgelt – Rudolf Kruse: Induction of association rules: Apriori implementation. In *Proceedings of the 15th Conference on Computational Statistics (Compstat 2002, Berlin, Germany)* (konferenciaanyag). Heidelberg, Germany, 2002, Physika Verlag.
- [21] Leo Breiman – Jerome Friedman – Charles J. Stone – R. A. Olshen: *Classification and Regression Trees*. 1984. January, Chapman & Hall/CRC. ISBN 0412048418.
- [22] Sergey Brin – Rajeev Motwani – Jeffrey D. Ullman – Shalom Tsur: Dynamic itemset counting and implication rules for market basket data. *SIGMOD Record (ACM Special Interest Group on Management of Data)*, 26(2):255, 1997.

- [23] Sergey Brin–Lawrence Page: The anatomy of a large-scale hypertextual Web search engine. *Computer Networks and ISDN Systems*, 30. évf. (1998) 1–7. sz.  
URL <http://citeseer.nj.nec.com/brin98anatomy.html>.
- [24] Douglas Burdick–Manuel Calimlim–Johannes Gehrke: Mafia: A maximal frequent itemset algorithm for transactional databases. In *Proceedings of the 17th International Conference on Data Engineering* (konferenciaanyag). Heidelberg, Germany, 2001, IEEE Computer Society, 443–452. p. ISBN 0-7695-1001-9.
- [25] Krisztián Antal Búza: „Egyszerű asszociációs szabályok jelenséghálózatokkal támogatott keresése”. Doktori értekezés (Budapesti Műszaki és Gazdaságtudományi Egyetem, Hungary). 2007.
- [26] Jadzia Cendrowska: Prism: An algorithm for inducing modular rules. *International Journal of Man-Machine Studies*, 27. évf. (1987) 4. sz.
- [27] Soumen Chakrabarti–Byron Dom–Prabhakar Raghavan–Sridhar Rajagopalan–David Gibson–Jon Kleinberg: Automatic resource compilation by analyzing hyperlink structure and associated text. *Computer Networks and ISDN Systems*, 30. évf. (1998) 1–7. sz.  
URL <http://citeseer.nj.nec.com/chakrabarti98automatic.html>.
- [28] Pete Chapman–Julian Clinton–Randy Kerber–Thomas Khabaza Thomas Reinartz–Colin Shearer–Rüdiger Wirth: Cross industry standard process for data mining (crisp-dm) – step by step data mining guide. Jelentés, 1999.
- [29] David Wai-Lok Cheung–Jiawei Han–Vincent Ng–C. Y. Wong: Maintenance of discovered association rules in large databases: An incremental updating technique. In *ICDE* (konferenciaanyag). 1996, 106–114. p.  
URL <http://citeseer.nj.nec.com/cheung96maintenance.html>.
- [30] David Wai-Lok Cheung–Sau Dan Lee–Ben Kao: A general incremental technique for maintaining discovered association rules. In *Database Systems for Advanced Applications* (konferenciaanyag). 1997, 185–194. p.  
URL <http://citeseer.nj.nec.com/cheung97general.html>.
- [31] Robert Cooley–Bamshad Mobasher–Jaideep Srivastava: Data preparation for mining world wide web browsing patterns. *Knowledge and Information Systems*, 1. évf. (1999) 1. sz. URL <http://citeseer.nj.nec.com/cooley99data.html>.
- [32] Thomas M. Cover–Joy A. Thomas: *Elements of Information Theory*. Wiley Series in Telecommunications sorozat. 1991, John Wiley & Sons, Inc.
- [33] R. de la Briandais: File searching using variable-length keys. In *Western Joint Computer Conference* (konferenciaanyag). 1959. March, 295–298. p.
- [34] T. G. Dietterich–M. Kearns–Y. Mansour: Applying the Weak Learning Framework to Understand and Improve C4.5. In L. Saitta (szerk.): *Proceedings of the 13th International Conference on Machine Learning, ICML'96* (konferenciaanyag). San Francisco, CA, 1996, Morgan Kaufmann, 96–104. p.

- [35] Margaret H. Dunham: *Data Mining: Introductory and Advanced Topics*. Upper Saddle River, NJ, USA, 2002, Prentice Hall PTR. ISBN 0130888923.
- [36] Herb Edelstein: Mining large databases – a case study. Jelentés, 1999, Two Crows Corporation.
- [37] M. Ester–H.-P. Kriegel–X. Xu.: A database interface for clustering in large spatial databases. In *Proceedings of the Knowledge Discovery and Data Mining Conference, Montreal, Canada* (konferenciaanyag). 1995, 94–99. p.
- [38] Martin Ester–Hans-Peter Kriegel–Jorg Sander–Xiaowei Xu: A density-based algorithm for discovering clusters in large spatial databases with noise. In Evangelos Simoudis–Jiawei Han–Usama Fayyad (szerk.): *Second International Conference on Knowledge Discovery and Data Mining* (konferenciaanyag). Portland, Oregon, 1996, AAAI Press, 226–231. p. URL <http://citeseer.nj.nec.com/chu02incremental.html>.
- [39] Usama M. Fayyad–Gregory Piatetsky-Shapiro–Padhraic Smyth: From data mining to knowledge discovery: An overview. In *Advances in Knowledge Discovery and Data Mining*. 1996, AAAI Press/The MIT Pres, 1–34. p.
- [40] William Feller: *Bevezetés a Valószínűségszámításba és Alkalmazásaiba*. 1978, Műszaki Könyvkiadó.
- [41] Bodon Ferenc: Hash-fák és szófák az adatbányászatban. *Alkalmazott Matematikai Lapok*, 21. évf. (2003).
- [42] E. W. Forgy: Cluster analysis of multivariate data: Efficiency versus interpretability of classifications. *Biometric Soc. Meetings, Riverside, California*, 21. évf. (1965).
- [43] Scott Fortin–Ling Liu: An object-oriented approach to multi-level association rule mining. In *CIKM* (konferenciaanyag). 1996, 65–72. p.
- [44] Edward Fredkin: Trie memory. *Communications of the ACM*, 3. évf. (1960) 9. sz. ISSN 0001-0782.
- [45] Y. Fu.: Discovery of multiple-level rules from large databases, 1996. URL <http://citeseer.nj.nec.com/fu96discovery.html>.
- [46] Iván Futó (szerk.): *Mesterséges Intelligencia*. Budapest, 1999, Aula Kiadó.
- [47] S.B. Gelfand–C.S. Ravishankar–E.J. Delp: An iterative growing and pruning algorithm for classification tree design. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 13. évf. (1991) 2. sz. ISSN 0162-8828.
- [48] Bart Goethals: Survey on frequent pattern mining. 2002. Manuskript.
- [49] Bart Goethals–Mohammed J. Zaki: Advances in frequent itemset mining implementations: Introduction to fimi03. In Bart Goethals–Mohammed J. Zaki (szerk.): *Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*, CEUR Workshop Proceedings konferenciasorozat, 90. köt. Melbourne, Florida, USA, 2003. November 19..

- [50] Gosta Grahne–Jianfei Zhu: Efficiently using prefix-trees in mining frequent itemsets. In Bart Goethals–Mohammed J. Zaki (szerk.): *Proceedings of the IEEE ICDM Workshop on Frequent Itemset Mining Implementations (FIMI'03)*, CEUR Workshop Proceedings konferenciasorozat, 90. köt. Melbourne, Florida, USA, 2003. November 19..
- [51] Sudipto Guha–Rajeev Rastogi–Kyuseok Shim: CURE: an efficient clustering algorithm for large databases. In *ACM SIGMOD International Conference on Management of Data* (konferenciaanyag). 1998. June, 73–84. p.  
URL <http://citeseer.nj.nec.com/article/guha98cure.html>.
- [52] Torben Hagerup–C. Rüb: A guided tour of chernoff bounds. *Inf. Process. Lett.*, 33. évf. (1990) 6. sz. ISSN 0020-0190.
- [53] J. Han–Y. Fu: Discovery of multiple-level association rules from large databases. *Proceedings of the 21st International Conference on Very Large Databases (VLDB), Zurich, Switzerland*, 1995.
- [54] Jiawei Han–Micheline Kamber: *ADATBÁNYÁSZAT - Koncepciók és technikák*. 2004, Panem Könyvkiadó.
- [55] Jiawei Han–Micheline Kamber: *Data mining: concepts and techniques (Second Edition)*. 2006, Morgan Kaufmann Publisher.
- [56] Jiawei Han–Jian Pei–Yiwen Yin: Mining frequent patterns without candidate generation. In Weidong Chen–Jeffrey Naughton–Philip A. Bernstein (szerk.): *2000 ACM SIGMOD International Conference on Management of Data* (konferenciaanyag). 2000. 05, ACM Press, 1–12. p. ISBN 1-58113-218-2.  
URL <http://citeseer.nj.nec.com/han99mining.html>.
- [57] Trevor Hastie–Robert Tibshirani–Jerome Friedman: *The Elements of Statistical Learning: Data Mining, Inference and Prediction*. 2001, Springer-Verlag.
- [58] K. Hatonen–Mika Klemettinen–Heikki Mannila–P. Ronkainen–Hannu Toivonen: Knowledge discovery from telecommunication network alarm databases. In Stanley Y. W. Su (szerk.): *Proceedings of the twelfth International Conference on Data Engineering, February 26–March 1, 1996, New Orleans, Louisiana* (konferenciaanyag). 1109 Spring Street, Suite 300, Silver Spring, MD 20910, USA, 1996, IEEE Computer Society Press, 115–122. p. URL <http://citeseer.nj.nec.com/hatonen96knowledge.html>.
- [59] Robert C. Holte: Very simple classification rules perform well on most commonly used datasets. *Mach. Learn.*, 11. évf. (1993) 1. sz. ISSN 0885-6125.
- [60] Maurice Houtsma–Arun Swami.: Set-oriented mining of association rules, 1993.
- [61] D. A. Hull: Improving text retrieval for the routing problem using latent semantic indexing. In *Proc. of SIGIR-94, 17th ACM Int. Conf. on Research and Development in Information Retrieval* (konferenciaanyag). Dublin, Ireland, 1994, 282–289. p.
- [62] Index.hu.: Rákkeltő anyagok a mcdonaldsban és burger kingben.  
URL <http://index.hu/gazdasag/vilag/mcrak060929>.

- [63] Akihiro Inokuchi – Takashi Washio – Hiroshi Motoda: An apriori-based algorithm for mining frequent substructures from graph data. In *Proceedings of the 4th European Conference on Principles of Data Mining and Knowledge Discovery* (konferenciaanyag). 2000, Springer-Verlag, 13–23. p. ISBN 3-540-41066-X.
- [64] Akihiro Inokuchi – Takashi Washio – Nishimura Yoshio – Hiroshi Motoda: A fast algorithm for mining frequent connected graphs,. Jelentés, 2002, IBM research, Tokyo Research Laboratory.
- [65] Fazekas István: *Bevezetés a matematikai statisztikába*. 2000, Debreceni Egyetem Kossuth Egyetemi Kiadója.
- [66] R. C. Jancey: Multidimensional group analysis. *Austral. J. Botany*, 14. évf. (1966).
- [67] Dr. Abonyi János: *Adatbányászat a hatékonyság eszköze*. Budapest, 2006, Computerbooks.
- [68] Richard A. Johnson – Dean W. Wichern: *Applied Multivariate Statistical Analysis*. Fifth. kiad. Upper Saddle River, NJ, 2002, Prentice-Hall.
- [69] Ravi Kannan – Santosh Vempala – Adrian Vetta: On clusterings: Good, bad and spectral. In *Proceedings of the 41th Annual Symposium on Foundations of Computer Science* (konferenciaanyag). 2000. URL <http://citeseer.nj.nec.com/495691.html>.
- [70] O. Kariv – S.L.Hakimi: An algorithmic approach to network location problems, part ii: p-medians. *SIAM J. Appl. Math.*, 37. évf. (1979).
- [71] L. Kaufman – P.J. Rousseeuw: *Finding Groups in Data: an Introduction to Cluster Analysis*. 1990, John Wiley & Sons.
- [72] Michael Kearns – Yishay Mansour: On the boosting ability of top-down decision tree learning algorithms. In *STOC '96: Proceedings of the twenty-eighth annual ACM symposium on Theory of computing* (konferenciaanyag). New York, NY, USA, 1996, ACM Press, 459–468. p. ISBN 0-89791-785-5.
- [73] Ashraf M. Kibriya – Eibe Frank: An empirical comparison of exact nearest neighbour algorithms. In *Proc 11th European Conference on Principles and Practice of Knowledge Discovery in Databases*, Warsaw, Poland konferenciasorozat. 2007, Springer, 140–151. p.
- [74] Jon Kleinberg: An impossibility theorem for clustering. *Advances in Neural Information Processing Systems (NIPS) 15*, 2002. URL <http://citeseer.nj.nec.com/561287.html>.
- [75] Jon M. Kleinberg: Authoritative sources in a hyperlinked environment. *Journal of the ACM*, 46. évf. (1999) 5. sz.  
URL <http://citeseer.nj.nec.com/kleinberg97authoritative.html>.
- [76] Mika Klemettinen: A knowledge discovery methodology for telecommunication network alarm databases, 1999.  
URL <http://citeseer.nj.nec.com/klemettinen99knowledge.html>.

- [77] Ron Kohavi: Mining e-commerce data: The good, the bad, and the ugly. In Foster Provost–Ramakrishnan Srikant (szerk.): *Proceedings of the Seventh ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (konferenciaanyag). 2001, 8–13. p. URL <http://citeseer.nj.nec.com/kohavi01mining.html>.
- [78] Michihiro Kuramochi–George Karypis: Frequent subgraph discovery. In *Proceedings of the 2001 IEEE International Conference on Data Mining* (konferenciaanyag). 2001, IEEE Computer Society, 313–320. p. ISBN 0-7695-1119-8.
- [79] Rónyai Lajos–Ivanyos Gábor–Szabó Réka: *Algoritmusok*. 1998, Typotex Kiadó.
- [80] Nada Lavrac–Dragan Gamberger–Hendrik Blockeel–Ljupco Todorovski (szerk.). *ExAnte: Anticipated Data Reduction in Constrained Pattern Mining*, Lecture Notes in Computer Science konferenciasorozat, 2838. köt. Springer, 2003. ISBN 3-540-20085-1.
- [81] Wenke Lee–Salvatore Stolfo: Data mining approaches for intrusion detection. In *Proceedings of the 7th USENIX Security Symposium* (konferenciaanyag). San Antonio, TX, 1998. URL <http://citeseer.nj.nec.com/article/lee00data.html>.
- [82] Wenke Lee–Salvatore J. Stolfo: A framework for constructing features and models for intrusion detection systems. *ACM Transactions on Information and System Security*, 3. évf. (2000) 4. sz. URL <http://citeseer.nj.nec.com/article/lee00framework.html>.
- [83] Wenke Lee–Salvatore J. Stolfo–Kui W. Mok: A data mining framework for building intrusion detection models. In *IEEE Symposium on Security and Privacy* (konferenciaanyag). 1999, 120–132. p. URL <http://citeseer.nj.nec.com/article/lee99data.html>.
- [84] R. Lempel–S. Moran: The stochastic approach for link-structure analysis (SALSA) and the TKC effect. In *WWW9* (konferenciaanyag). 2000. URL <http://citeseer.nj.nec.com/346353.html>.
- [85] Bodrogi Lilla.: Legitimebb titkunk: a saját genetikai állományunk, 2007. URL <http://www.origo.hu/tudomany/20070919-egyedi-emberi-genetikai-allomany-genom-elemzese-eloszor.html>.
- [86] Heikki Mannila–Hannu Toivonen: Discovering generalized episodes using minimal occurrences. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining (KDD'96)* (konferenciaanyag). 1996. August, AAAI Press, 146–151. p. URL <http://citeseer.nj.nec.com/mannila96discovering.html>.
- [87] Heikki Mannila–Hannu Toivonen–A. Inkeri Verkamo: Discovering frequent episodes in sequences. In *Proceedings of the First International Conference on Knowledge Discovery and Data Mining (KDD'95)* (konferenciaanyag). 1995. August, AAAI Press, 210–215. p.
- [88] Heikki Mannila–Hannu Toivonen–A. Inkeri Verkamo: Discovery of frequent episodes in event sequences. *Data Mining and Knowledge Discovery*, 1. évf. (1997) 3. sz. ISSN 1384-5810. URL <http://citeseer.nj.nec.com/mannila97discovery.html>.



- [89] Heikki Mannila–Hannu Toivonen–A. Inkeri Verkamo: Efficient algorithms for discovering association rules. In Usama M. Fayyad–Ramasamy Uthurusamy (szerk.): *AAAI Workshop on Knowledge Discovery in Databases(KDD-94)* (konferenciaanyag). Seattle, Washington, 1994, AAAI Press, 181–192. p.  
URL <http://citeseer.nj.nec.com/mannila94efficient.html>.
- [90] R. López De Mántaras: A distance-based attribute selection measure for decision tree induction. *Mach. Learn.*, 6. évf. (1991) 1. sz. ISSN 0885-6125.
- [91] Brendan D. McKay: Practical graph isomorphism. *Congressus Numerantium*, 30. évf. (1981). URL <http://cs.anu.edu.au/people/bdm/nauty/>.
- [92] N. Megiddo–K. Supowit: On the complexity of some common geometric location problems. *SIAM J. Comput.*, 1984.
- [93] Jesus Mena: *Data Mining und E-Commerce*. Düsseldorf, 2000, Symposion Publishing.  
URL <http://www.symposion.de/datamining>.
- [94] Ulrich Meyer–Peter Sanders–Jop F. Sibeyn (szerk.). *Algorithms for Memory Hierarchies, Advanced Lectures [Dagstuhl Research Seminar, March 10-14, 2002]*, Lecture Notes in Computer Science konferenciasorozat, 2625. köt. Springer, 2003. ISBN 3-540-00883-7.
- [95] Dunja Mladenic–NADA Lavrac–Marko Bohanec–Steve Moyle: *Data Mining and Decision Support: Integration and Collaboration*. 2003, Kluwer Academic Publishers.
- [96] Andreas Mueller: Fast sequential and parallel algorithms for association rule mining: A comparison. CS-TR-3515. Jelentés, College Park, MD, 1995, Department of Computer Science, University of Maryland.  
URL <http://citeseer.nj.nec.com/mueller95fast.html>.
- [97] Raymond T. Ng–Jiawei Han: Efficient and effective clustering methods for spatial data mining. In Jorge B. Bocca–Matthias Jarke–Carlo Zaniolo (szerk.): *Proceedings of the 20th International Conference Very Large Data Bases, VLDB* (konferenciaanyag). 1994. 12-15, Morgan Kaufmann, 144–155. p. ISBN 1-55860-153-8.  
URL <http://citeseer.nj.nec.com/571734.html>.
- [98] Edward Omiecinski–Ashoka Savasere: Efficient mining of association rules in large dynamic databases. In *British National Conference on Databases* (konferenciaanyag). 1998, 49–63. p.
- [99] Stephen M. Omohundro: Five balltree construction algorithms. Jelentés, 1989. December, International Computer Science Institute.
- [100] Stifán Orsolya: Adatbányászat és adatvédelem. In Dr. Székely Iván–Dr. Szabó Máté (szerk.): *Szabad adatok, védett adatok*. ALMA MATER sorozat, 10. köt. Budapest, 2005, BME GTK ITM, 169–196. p.
- [101] Banu Ozden–Sridhar Ramaswamy–Abraham Silberschatz: Cyclic association rules. In *ICDE* (konferenciaanyag). 1998, 412–421. p.  
URL <http://citeseer.nj.nec.com/ozden98cyclic.html>.

- [102] Lawrence Page–Sergey Brin–Rajeev Motwani–Terry Winograd: The pagerank citation ranking: Bringing order to the web. Jelentés, 1998, Stanford Digital Library Technologies Project. URL <http://citeseer.nj.nec.com/page98pagerank.html>.
- [103] Jong Soo Park–Ming-Syan Chen–Philip S. Yu: An effective hash based algorithm for mining association rules. In Michael J. Carey–Donovan A. Schneider (szerk.): *Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data* (konferenciaanyag). San Jose, California, 1995. 22-25, 175–186. p. URL <http://citeseer.nj.nec.com/park95effective.html>.
- [104] N. Pasquier–Y. Bastide–R. Taouil–L. Lakhal: Pruning closed itemset lattices for association rules. In *Proceedings of the BDA French Conference on Advanced Databases* (konferenciaanyag). 1998. October. URL <http://citeseer.nj.nec.com/pasquier98pruning.html>.
- [105] N. Pasquier–Y. Bastide–R. Taouil–L. Lakhal: Efficient mining of association rules using closed itemset lattices. In *Journal of Information systems* (konferenciaanyag). 1999, 25–46. p.
- [106] Nicolas Pasquier–Yves Bastide–Rafik Taouil–Lotfi Lakhal: Discovering frequent closed itemsets for association rules. In *ICDT* (konferenciaanyag). 1999, 398–416. p. URL <http://citeseer.nj.nec.com/pasquier99discovering.html>.
- [107] Jian Pei–Jiawei Han–Laks V. S. Lakshmanan: Mining frequent item sets with convertible constraints. In *ICDE* (konferenciaanyag). 2001, 433–442. p. URL <http://citeseer.ist.psu.edu/383962.html>.
- [108] Jian Pei–Jiawei Han–Runying Mao: CLOSET: An efficient algorithm for mining frequent closed itemsets. In *ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery* (konferenciaanyag). 2000, 21–30. p. URL <http://citeseer.nj.nec.com/pei00closet.html>.
- [109] Wim Pijls–Jan C. Bioch: Mining frequent itemsets in memory-resident databases. In *Proceedings of the Eleventh Belgium /Netherlands Artificial Intelligence Conference (BNAIC 99)* (konferenciaanyag). 1999, 75–82. p. URL <http://citeseer.nj.nec.com/pijls99mining.html>.
- [110] Jim Porter: Disk/trend report. In *Proceedings of the 100th Anniversary Conference on Magnetic Recording and Information Storage*. Santa Clara University, 1998.
- [111] J. R. Quinlan: Induction of decision trees. *Mach. Learn.*, 1. évf. 1. sz. ISSN 0885-6125.
- [112] J. R. Quinlan: Simplifying decision trees. *Int. J. Man-Mach. Stud.*, 27. évf. (1987) 3. sz. ISSN 0020-7373.
- [113] J. Ross Quinlan: *C4.5: programs for machine learning*. San Francisco, CA, USA, 1993, Morgan Kaufmann Publishers Inc. ISBN 1-55860-238-0.
- [114] T. R. C. Read–N. A. C. Cressie: *Goodness-of-Fit Statistics for Discrete Multivariate Data*. Springer Series in Statistics sorozat. New York, 1988, Springer-Verlag.

- [115] Pál Rózsa: *Lineáris algebra és alkalmazásai*. 1991, Tankönyvkiadó, Budapest.
- [116] S. Sahni–T. Gonzales: P-complete approximation problems. *JACM*, 23. évf. (1976).
- [117] Nandlal L. Sarda–N. V. Srinivas: An adaptive algorithm for incremental mining of association rules. In *DEXA Workshop* (konferenciaanyag). 1998, 240–245. p.
- [118] Ashoka Savasere–Edward Omiecinski–Shamkant B. Navathe: An efficient algorithm for mining association rules in large databases. In *The VLDB Journal* (konferenciaanyag). 1995, 432–444. p. URL <http://citeseer.nj.nec.com/sarasere95efficient.html>.
- [119] R. E. Schapire–Y. Singer–A. Singhal: Boosting and Rocchio applied to text filtering. In *Proc. of SIGIR-98, 21st ACM International Conference on Research and Development in Information Retrieval* (konferenciaanyag). Melbourne, Australia, 1998, 215–223. p.
- [120] Matthew G. Schultz–Eleazar Eskin–Salvatore J. Stolfo.: Mef: Malicious email filter - a unix mail filter that detects malicious windows executables.  
URL <http://citeseer.nj.nec.com/417909.html>.
- [121] Matthew G. Schultz–Eleazar Eskin–Erez Zadok–Salvatore J. Stolfo.: Data mining methods for detection of new malicious executables.  
URL <http://citeseer.nj.nec.com/417492.html>.
- [122] F. Sebastiani: Machine learning in automated text categorization. *ACM Computing Surveys*, 34. évf. (2002. March) 1. sz.
- [123] Dennis G. Severance: Identifier search mechanisms: A survey and generalized model. *ACM Comput. Surv.*, 6. évf. (1974) 3. sz. ISSN 0360-0300.
- [124] Ron Shamir–Dekel Tsur: Faster subtree isomorphism. *Journal of Algorithms*, 33. évf. (1999) 2. sz. ISSN 0196-6774.
- [125] Li Shen–Hong Shen: Mining flexible multiple-level association rules in all concept hierarchies (extended abstract). In *Database and Expert Systems Applications* (konferenciaanyag). 1998, 786–795. p.
- [126] Y.-S. Shih: Families of splitting criteria for classification trees. *Statistics and Computing*, 9. évf. (1999) 4. sz. ISSN 0960-3174.
- [127] Abraham Silberschatz–Alexander Tuzhilin: On subjective measures of interestingness in knowledge discovery. In *Knowledge Discovery and Data Mining* (konferenciaanyag). 1995, 275–281. p. URL <http://citeseer.nj.nec.com/silberschatz95subjective.html>.
- [128] Spencer: The probabilistic method. In *SODA: ACM-SIAM Symposium on Discrete Algorithms (A Conference on Theoretical and Experimental Analysis of Discrete Algorithms)* (konferenciaanyag). 1992.
- [129] Ramakrishnan Srikant–Rakesh Agrawal: Mining generalized association rules. *Proceedings of the 21st International Conference on Very Large Databases (VLDB), Zurich, Switzerland*, 1995.

- [130] Ramakrishnan Srikant–Rakesh Agrawal: Mining sequential patterns: Generalizations and performance improvements. In Peter M. G. Apers–Mokrane Bouzeghoub–Georges Gardarin (szerk.): *Proceedings of 5th International Conference Extending Database Technology, EDBT* (konferenciaanyag), 1057. köt. 1996. 25-29, Springer-Verlag, 3–17. p. ISBN 3-540-61057-X. URL <http://citeseer.nj.nec.com/article/srikant96mining.html>.
- [131] T-Online.: Általános szerződési feltételek, 2006.  
URL [http://www.t-online.hu/dokumentumok/toh\\_aszf\\_060331.pdf](http://www.t-online.hu/dokumentumok/toh_aszf_060331.pdf).
- [132] Lyn C. Thomas: A survey of credit and behavioural scoring; forecasting financial risk of lending to consumers. *International Journal of Forecasting* 16, 2000.
- [133] Lyn C. Thomas: A survey of credit and behavioural scoring; forecasting financial risk of lending to consumers. *International Journal of Forecasting*, 16. évf. (2000).
- [134] Shiby Thomas–Sreenath Bodagala–Khaled Alsabti–Sanjay Ranka: An efficient algorithm for the incremental updation of association rules in large databases. In *Knowledge Discovery and Data Mining* (konferenciaanyag). 1997, 263–266. p.  
URL <http://citeseer.nj.nec.com/thomas97efficient.html>.
- [135] Shiby Thomas–Sunita Sarawagi: Mining generalized association rules and sequential patterns using SQL queries. In *Knowledge Discovery and Data Mining* (konferenciaanyag). 1998, 344–348. p. URL <http://citeseer.nj.nec.com/thomas98mining.html>.
- [136] Hannu Toivonen: Sampling large databases for association rules. In *The VLDB Journal* (konferenciaanyag). 1996, 134–145. p.  
URL <http://citeseer.nj.nec.com/toivonen96sampling.html>.
- [137] Jeffrey K. Uhlmann: Satisfying general proximity/similarity queries with metric trees. *Inf. Process. Lett.*, 40. évf. (1991) 4. sz.
- [138] J. R. Ullmann: An algorithm for subgraph isomorphism. *J. ACM*, 23. évf. (1976) 1. sz. ISSN 0004-5411.
- [139] John von Neumann.: First draft of a report on the EDVAC. Contract No. W-670-ORD-4926 Between the United States Army Ordnance Department and the University of Pennsylvania, 1945. június.  
URL <http://qss.stanford.edu/~{}godfrey/vonNeumann/vnedvac.pdf>.
- [140] Jianyong Wang–Jiawei Han–Jian Pei: Closet+: Searching for the best strategies for mining frequent closed itemsets. In *In Proceedings of the Ninth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'03)* (konferenciaanyag). Washington, DC, USA, 2003. URL <http://citeseer.nj.nec.com/wang03closet.html>.
- [141] E. D. Wiener–J. O. Pedersen–A. S. Weigend: A neural network approach to topic spotting. In *Proc. of the SDAIR-95, 4th Annual Symposium on Document Analysis and Information Retrieval* (konferenciaanyag). Las Vegas, NV, 1995, 317–332. p.

- [142] Ian H. Witten – Eibe Frank: *Data Mining: Practical Machine Learning Tools and Techniques*. Morgan Kaufmann Series in Data Management Sys sorozat. Second. kiad. 2005. June, Morgan Kaufmann. ISBN 0120884070.  
URL <http://www.amazon.fr/exec/obidos/ASIN/0120884070/citeulike04-21>.
- [143] Ying Yang – Geoffrey I. Webb: Proportional k-interval discretization for naive-bayes classifiers. In *EMCL '01: Proceedings of the 12th European Conference on Machine Learning* (konferenciaanyag). London, UK, 2001, Springer-Verlag, 564–575. p. ISBN 3-540-42536-5.
- [144] S. B. Yao: Tree structures construction using key densities. In *Proceedings of the 1975 annual conference* (konferenciaanyag). 1975, ACM Press, 337–342. p.
- [145] Mohammed J. Zaki: Efficiently mining frequent trees in a forest. Jelentés, Troy, NY, 12180, 2001. July, Computer Science Department, Rensselaer Polytechnic Institute.
- [146] Mohammed J. Zaki: Efficiently mining frequent trees in a forest. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining* (konferenciaanyag). 2002, ACM Press, 71–80. p. ISBN 1-58113-567-X.
- [147] Mohammed J. Zaki – Karam Gouda: Fast vertical mining using difsets. In *Proceedings of the ninth ACM SIGKDD international conference on Knowledge discovery and data mining* (konferenciaanyag). 2003, ACM Press, 326–335. p. ISBN 1-58113-737-0.
- [148] Mohammed Javeed Zaki: Sequence mining in categorical domains: Incorporating constraints. In *CIKM* (konferenciaanyag). 2000, 422–429. p.  
URL <http://citeseer.nj.nec.com/zaki00sequence.html>.
- [149] Mohammed Javeed Zaki – Ching-Jui Hsiao: Charm: An efficient algorithm for closed itemset mining. In *Proceedings of 2nd SIAM International Conference on Data Mining* (konferenciaanyag). Arlington, VA, USA, 2002.
- [150] Mohammed Javeed Zaki – Mitsunori Ogihara: Theoretical foundations of association rules. In *Proceedings of third SIGMOD'98 Workshop on Research Issues in Data Mining and Knowledge Discovery (DMKD'98)* (konferenciaanyag). Seattle, Washington, 1998.  
URL <http://citeseer.nj.nec.com/zaki98theoretical.html>.
- [151] Mohammed Javeed Zaki – Srinivasan Parthasarathy – Mitsunori Ogihara – Wei Li: New algorithms for fast discovery of association rules. In David Heckerman – Heikki Mannila – Daryl Pregibon – Ramasamy Uthurusamy – Menlo Park (szerk.): *Proceedings of the third International Conference on Knowledge Discovery and Data Mining* (konferenciaanyag). 1997. 12-15, AAAI Press, 283–296. p. ISBN 1-57735-027-8.  
URL <http://http://citeseer.nj.nec.com/30063.html>.