

08:15

D.-1-105

Információkezelés 4.

1.EA

2009.02.11

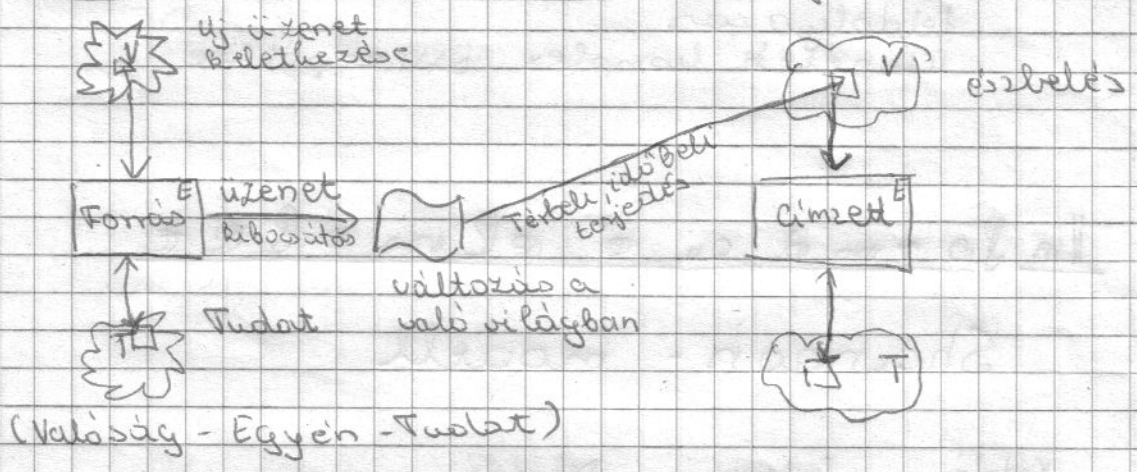
Győrfy : Információ és kódelmélet (Typotex)
Elementary Information Theory

Információelmélet alapjai

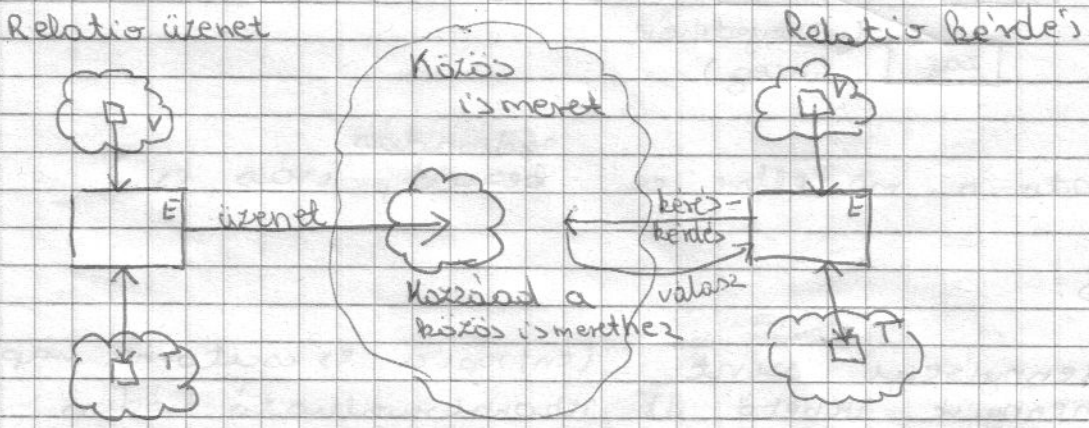
- Kommunikáció - információs rendszerek világa (véletlen jelenségképe épít) (jóvőre vonatkozik)
- Shannon-Jele entropia
- Kolmogorov entropia (számításelméletre épül) (múltre vonatkozik)

Kommunikáció lényege az emberi elmék kölcsönhatása (tudatok)

Elemi kommunikációs lépés. (információ átadás)



Információs-rendszer jellegű kommunikáció.



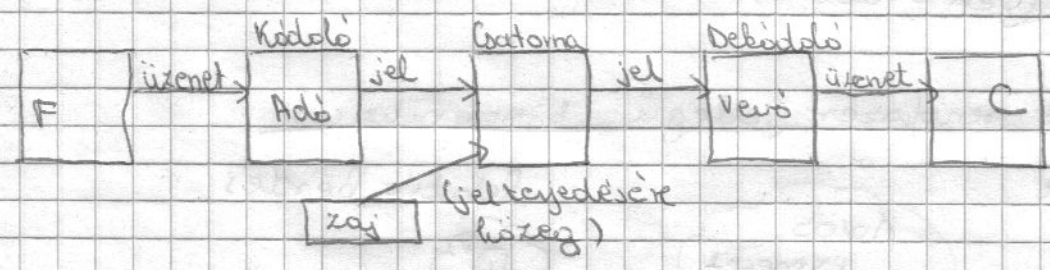
Történeti állomások:

- beszéd kialakulása
(hangképzés (megfelelő légmozgás), hallás, ^{tudat} agyrndszer,
(gestikuláció - testbeszéd)
- jeljegyzés - írás, rajz, kép
(írás megmarad)
- nyomtatás (szerszorosítás könnyen)
(komputer + szolgáltatás a híres- hírdésékre)
- jelenleg: előbb analóg →
 - rögép
 - telefon-adatátvitel
 - mozgóképek
 - számítógép
 - fénykép
 - rádiófrekvenciás műsorvezetés
 - hűtőszelvény
 - digitális kamera, műszerek
- háttérben, tárolás
számítás - mit lépés kezelni
feldolgozás
üzenetek komplex összerakása

hang, példa, érzékelés, zebra

Információ elmélet

Shannon - modell



WWW.FSZEK.HU

(Rohozható a modellre a rádióadás, beszéd, írás is...)

Elvárás:

- A) Mennyiségi szint** (entropia és csatorna kapacitás)
- mennyire vihető át rekonstruálható módon a kívánt üzenet a csatornában?
- B) Megértési szint** (jelentés, szemantika) 2/1

c) Hatékonyági szint (a kívánt hatás elérése)

Mennyiségi szint:

- Forrás jellemzője: mennyi üzenetet bocsát ki, időegység hál visszaüzeneti sebesség paraméter is hál
- Csatorna: időegység alatt átvihető különböző csatornajelek száma

Forrási egyezmény példái:

elemi üzenet: $\left\{ \begin{array}{l} \text{igen} \\ \text{nem} \end{array} \right.$

- ① pénzfeldobás: fej \rightarrow igen
írás \rightarrow nem

Átvitel: 0-1 sorozat segítségével

Feladat: $1024 \times$ ismételyük meg

1024 0-1 bit

2^{1024} különböző jel hál hogy átmenjen a csatornán

- ② ötösöm van a lottón - igen
nincs - nem

1024 het eredményét

- Ha nincs ötös - 0

- Ha 1 ötös van - $10 \frac{1024}{10 \text{ bit}}$

- Ha több ötös van - $11 \frac{1024}{\text{teljes}}$

$$P_0 = 0,999977$$

$$P_1 = 2,33 \cdot 10^{-5}$$

$$P_2 = \dots$$

MF. kódhossz várható értéke?

(Poisson határelosítás)

Kül ①-② - Mennyire bizonytalan a 2
esemény kioldásodása!

2009. 02. 18 , 2. előadás

ablinux.inf.elte.hu / infoke24

- segédanyagok

Mennyiségek

1. Csatorna: Kapacitás

a) bináris szimmetrikus csatornával való ekivalencia alapján
- egységnyi idő alatt átlagosan ugyanannyi különböző jelzést vihető át, mint egy C bit/sec sebességű bináris csatornában

2^C különböző bit vihető át

2. Forrás: véges sok különböző szimbólumból való választás (n -a szimbólumok száma)
választás - véletlen körülmények között történik

A választást valószínűség eloszlással jellemezzük

N hosszú üzenetet választ

Elvárás: az N hosszú üzenetek halmazán kell megadni

n^N elemű eloszlás
Shannon-entropiáját rendeljük hozzá

- n^N elemű halmazon, elemeit sorszámozzuk $i=1, \dots, n^N$
választási valószínűség: P_i

Shannon-entropia:
$$H_N = - \sum_{i=1}^{n^N} P_i \log_2 P_i$$

Forrás jellemzésére H_N átlagát keressük szimbólumonként:

Def:
$$H := \lim_{N \rightarrow \infty} \frac{1}{N} H_N$$

Mérvszám:
[bit / szimbólum]

Megvalósítható: a forrás sebessége - sec.-ként hány szimbólumot választ

$$V = \left[\frac{\text{szimbólum}}{\text{sec}} \right]$$

C : Bit/sec
 H : Bit/szimbólum
 V : szimbólum/sec

Mennyiségi feladat megoldhatósága: veszteségmentesen átvihető minden üzenet

Zajmentes csatorna alaptétele (Shannon)

⊕ Nem lehet veszteségmentesen működtetni a kommunikációs rendszert, ha $V > \frac{C}{H}$;

Tetszőleges ($\epsilon > 0$)-ra \exists kódolás, hogy $V < \frac{C}{H} - \epsilon$ sebességgel veszteségmentesen működtethető a rendszer.

Def: Csatorna kapacitás:

Diszkrét, zajmentes csatorna

Átvihető jelek: s_1, s_2, \dots, s_m
 Átviteli idők: t_1, t_2, \dots, t_n

Szim. bin. csatorna: $m = 2$
 $t_1 = t_2 = t$

Ha T időig működtetjük, akkor hány különböző jelsorozatot hildketek át?

T/t jel vihető át

$$N(T) = 2^{T/t} \quad \text{--> azes lehetőségek}$$

$$C = 1/t \quad \left. \begin{array}{l} N(T) = 2^{C \cdot T} \\ C = \frac{1}{T} \cdot \log_2 N(T) \end{array} \right\}$$

Két csatorna ekvivalens, ha az $N(T)$ mennyiségek "lényegében" megegyeznek

Def: Legyen a csatornán T idő alatt átvihető különböző jelsorozatok száma $N(T)$, ekkor a csatorna kapacitása:

$$C = \lim_{T \rightarrow \infty} \frac{\log_2 N(T)}{T}$$

$$2^{T \cdot (C - \delta)} < N(T) < 2^{T \cdot (C + \delta)} \quad \underline{\underline{C}}: \text{tetszőlegesen kicsi } \delta > 0 \quad \text{②/2}$$

S_1, \dots, S_m | ha $t_1 = t_2 = \dots = t_m = t$, akkor szimmetrikus
 t_1, \dots, t_m

C megoldása?

Általános eset:

$$N(T) = N(T-t_1) + N(T-t_2) + \dots + N(T-t_m)$$

(differencia egyenlet)
 Megoldása:

$$1 - X^{-t_1} - X^{-t_2} - \dots - X^{-t_m} = 0$$

X_0 : a legnagyobb pozitív gyök

$$X_0^T = (X_0^{-t_1} + X_0^{-t_2} + \dots + X_0^{-t_m}) \cdot X_0^T =$$

$$= X_0^{T-t_1} + X_0^{T-t_2} + \dots + X_0^{T-t_m}$$

[X_0 hatóánya úgy viselkedik mint $N(T)$]

$$N(T) = k \cdot X_0^T$$

$$C = \lim_{T \rightarrow \infty} \frac{\log_2 N(T)}{T} = \boxed{\log_2 X_0 = C}$$

// T lineáris növekedése exponenciálisan
 növeli az $N(T)$ -t

Forrás jellemzése:

Shannon-entropia levezetése

- minden lehetséges üzenetválasztásra működőn
- nem a szimbólumok fontosak, hanem a válaszítás bizonytalansága (milyen eloszlás szerint?)

Bizonytalanság mérésére: eloszlást jellemez

Valószínűségi eloszlás: (P_1, \dots, P_n) , ahol $\forall i: P_i > 0$ és
 $\sum_{i=1}^n P_i = 1$

Keressünk a (véges, diszkrét) eloszlásokon egy $H(P_1, \dots, P_n)$ fnt, melyre Shannon-elméletei:

(0) H az eloszlásokon van értelmezve, független a $\{P_1, \dots, P_n\}$ permutációtól

(1) H legyen folytonos minden változójában
(akkor is, ha $P_n \rightarrow 0$ $H(P_1, \dots, P_n) \rightarrow H(P_1, \dots, P_{n-1})$)

(2) Az $A(n) = H(\frac{1}{n}, \dots, \frac{1}{n})$ egyenletes eloszlás bizonytalansága

$$A(n) < A(m), \text{ ha } n < m$$

// Monotonitási elvárás a szimmetrikus
szituációra

(3) Elágazási (lépcsőzetési szabály)

$$H(P_1, P_2, \dots, \underbrace{P_{n-1}, P_n}_{Q_j}) = H(P_1, P_2, \dots, P_{n-2}, Q_j) + Q_j \cdot H\left(\frac{P_{n-1}}{Q_j}, \frac{P_n}{Q_j}\right)$$

$Q_j = P_{n-1} + P_n$

$$H(P_1, \dots, P_n) = H(P_1, \dots, P_{n-2}, Q_j) + Q_j \cdot H\left(\frac{P_{n-1}}{Q_j}, \frac{P_n}{Q_j}\right)$$

// Általános elágazási szabály: (indukcióval)

legyen: $Q_j = \sum_{i=1}^{n_j} P_{j,i} \quad j=1, \dots, m$

$$0 < P_{j,i} < 1, \quad \sum_{j=1}^m Q_j = 1$$

$$H(\underbrace{P_{j1}, \dots, P_{jn1}}_{Q_1}, \underbrace{P_{j1}, \dots, P_{jn2}}_{Q_2}, \dots, \underbrace{P_{j1}, \dots, P_{jn,m}}_{Q_m}) =$$

$$= H(Q_1, \dots, Q_m) + \sum_{j=1}^m Q_j \cdot H\left(\frac{P_{j1}}{Q_j}, \dots, \frac{P_{jn_j}}{Q_j}\right)$$

↑ elágazásig tartó bizonytalanság

↑ elágazás atomi bizonytalanság

Kaptunk egy ftt-egyenletet H -ra



$A(n) - (3)$ feltételeket kielégítő k , csak

$$k \cdot \sum_{i=1}^n P_i \cdot \log_2 P_i \text{ alakú lehet}$$

k megválasztása? amit 1 bittel tudunk mérni

$$1 = H\left(\frac{1}{2}, \frac{1}{2}\right)$$

$$k(-1) = 1, \quad \boxed{k = -1}$$

$$\boxed{H(P_1, \dots, P_n) = -\sum_{i=1}^n P_i \cdot \log_2 P_i}, \text{ ami a Shannon-entropia formula}$$

Biz: Racionális eloszlás: $+1 \leftrightarrow A$
 kifejezés segítségével

$$P_i = \frac{q_i}{m} \quad i=1, \dots, n$$

$$\sum_{i=1}^n P_i = 1 \Rightarrow \sum_{i=1}^n q_i = m$$

Az m elemű egyenletes eloszlásból indulunk - csoportok képzése

$$H\left(\underbrace{\frac{q_1}{m}, \dots, \frac{q_m}{m}}_{\text{elég. rász.}}\right) = \underbrace{A(m)}_{\text{teljes}} - \sum_{i=1}^n \frac{q_i}{m} \cdot \underbrace{A(q_i)}_{\text{elég. rász. után}} =$$

$$= -\sum_{i=1}^n \frac{q_i}{m} \cdot [A(q_i) - A(m)]$$

Ha $A(n) = k' \cdot \log_2 n$ alapú (be kell biz.)

\Rightarrow bővíthetjük a Schan. és lesz a racionális eset és a valósz. az már a folytonosságból kijön

Egyenletes eloszlás entropiája: $\boxed{A(n) = k' \cdot \log_2 n}$

$A(n)$ meghatározása:

$$A(n \cdot m) = H\left(\underbrace{\frac{1}{nm}, \dots, \frac{1}{nm}}_{\text{elég. r.}}\right) = A(n) + n \cdot \frac{1}{n} \cdot A(m) \quad \text{elég. után}$$

$$A(n \cdot m) = A(n) + A(m)$$

$$A(n^m) = m \cdot A(n)$$

Rögzített s egész számot: $s^m \leq t^n \leq s^{m+1}$
 (t paraméter, egész szám)

$$s^m \leq t^n \leq s^{m+1} \quad \left| \begin{array}{l} \text{alkalmazuk: } \log_2 x \\ A(x) \end{array} \right. \left. \begin{array}{l} \text{monoton} \\ \Rightarrow \leq \text{marad} \end{array} \right.$$

$$\log_2 s^m \leq \log_2 t^n < \log_2 s^{m+1} \quad \left| \quad A(s^m) \leq A(t^n) < A(s^{m+1}) \right.$$

$$m \cdot \log_2 s \leq n \cdot \log_2 t < (m+1) \cdot \log_2 s \quad \left| \quad m \cdot A(s) \leq n \cdot A(t) < (m+1) \cdot A(s) \right.$$

$$\frac{m}{n} \cdot \log_2 s \leq \log_2 t < \frac{m+1}{n} \cdot \log_2 s \quad \left| \quad \frac{m}{n} A(s) \leq A(t) < \frac{m}{n} A(s) + \frac{1}{n} A(s) \right.$$

$\frac{m}{n} \cdot \log_2 s + \frac{1}{n} \log_2 s$
 $\downarrow n \rightarrow \infty$
 0

$$\frac{m}{n} \leq \frac{\log_2 t}{\log_2 s} < \frac{m}{n} + \frac{1}{n} \quad \left| \quad \frac{m}{n} \leq \frac{A(t)}{A(s)} < \frac{m}{n} + \frac{1}{n} \right.$$

$$n \rightarrow \infty \quad \therefore \quad \frac{m}{n} \rightarrow \frac{\log_2 t}{\log_2 s} \quad \left| \quad \frac{m}{n} \rightarrow \frac{A(t)}{A(s)} \right.$$

$$\Rightarrow \quad \frac{A(t)}{A(s)} = \frac{\log_2(t)}{\log_2(s)}$$

$$A(t) = \underbrace{\frac{A(s)}{\log_2 s}}_R \cdot \log_2 t$$

$$H\left(\frac{1}{n}, \dots, \frac{1}{n}\right) = \log_2 n$$

3. előadás, 2009. 02. 25.

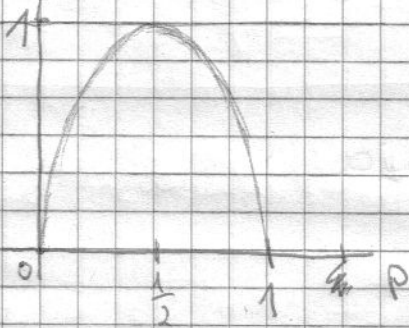
$$H(P_1, \dots, P_n) = - \sum_{i=1}^n P_i \cdot \log_2 P_i$$

Tulajdonságai:

① $H(\frac{1}{n}, \dots, \frac{1}{n}) =: A(n) = \log_2 n$

② Kételemű eloszlás:

$$H(P, 1-P) = -P \cdot \log_2 P - (1-P) \cdot \log_2 (1-P)$$



$$H(P, 1-P) \leq 1$$

Ha p nő, akkor H is nő

③ $H(P_1, \dots, P_n) \leq \log_2 n$ (egyenletesség teljes növekedés)

$$H(P_1, \dots, P_{n-1}, P_n) = H(P_1, \dots, P_{n-2}, P) + p \cdot H\left(\frac{P_{n-1}}{p}, \frac{P_n}{p}\right)$$

$$H(P_1, \dots, P_{n-2}, q_1, q_2) = H(P_1, \dots, P_{n-2}, P) + p \cdot H\left(\frac{q_1}{p}, \frac{q_2}{p}\right)$$

$$(P_{n-1} + P_n = q_1 + q_2 = p)$$

⇒ Kiegészítés növeli az entropiát

(Forrás amit kibocsát - val. szin. változó)

Valószínűségi változó entropiája

ξ , értékei: x_1, \dots, x_n (diszkrét, véges)

eloszlása: $P(\xi = x_i) \quad i=1, \dots, n$

$$\sum_{i=1}^n P(\xi = x_i) = 1$$

$$H(\xi) := - \sum_{i=1}^n P(\xi = x_i) \cdot \log_2 P(\xi = x_i)$$

①/3 Forrás jellemezése: N db v.sz. változóval

Két dimenziós eloszlás entropiája : ξ, η

Együttes eloszlás: $P(\xi = x_i, \eta = y_j)$ $i=1 \dots n$
 $j=1 \dots m$

$$H(\xi, \eta) = - \sum_{i=1}^n \sum_{j=1}^m P(\xi = x_i, \eta = y_j) \cdot \log_2 P(\xi = x_i, \eta = y_j)$$

ξ eloszlása?

$$P(\xi = x_i) = \sum_{j=1}^m P(\xi = x_i, \eta = y_j)$$

$$P(\eta = y_j) = \sum_{i=1}^n P(\xi = x_i, \eta = y_j)$$

Független eloszlások entropiája:

$$H(\xi, \eta), H(\xi), H(\eta)$$

ξ, η függetlenek, akkor $P(\xi = x_i, \eta = y_j) = P(\xi = x_i) \cdot P(\eta = y_j)$

$$H(\xi, \eta) = - \sum_{i=1}^n \sum_{j=1}^m P(\xi = x_i) \cdot P(\eta = y_j) \cdot [\log_2 P(\xi = x_i) + \log_2 P(\eta = y_j)] =$$

$$= - \sum_{i=1}^n \sum_{j=1}^m P(\xi = x_i) \cdot P(\eta = y_j) \cdot \log_2 P(\xi = x_i) +$$

$$+ - \sum_{i=1}^n \sum_{j=1}^m P(\xi = x_i) \cdot P(\eta = y_j) \cdot \log_2 P(\eta = y_j) =$$

$$= - \sum_{i=1}^n P(\xi = x_i) \cdot \log_2 P(\xi = x_i) \cdot \underbrace{\sum_{j=1}^m P(\eta = y_j)}_{=1}$$

$$+ \sum_{i=1}^n P(\xi = x_i) \cdot \sum_{j=1}^m (\log_2 P(\eta = y_j)) \cdot P(\eta = y_j)$$

$$= H(\xi) + H(\eta)$$

$$\boxed{H(\xi, \eta) = H(\xi) + H(\eta)}$$

(Vérhetően ez a legnagyobb bizonytalanság)

Tétel: $H(\xi, \eta) \leq H(\xi) + H(\eta)$ és
 $=$ a.c.s.a. ξ, η függetlenek

Biz:

Tegyük fel, h. megfigyeltük $\eta = y_j$ értéket

$$P(\xi | \eta = y_j) \quad j = 1, \dots, n$$

$$P(\xi = x_i | \eta = y_j) = \frac{P(\xi = x_i, \eta = y_j)}{P(\eta = y_j)}$$

$$P(A|B) = \frac{P(AB)}{P(B)}$$

$$H(\xi | \eta = y_j) = - \sum_{i=1}^n P(\xi = x_i | \eta = y_j) \cdot \log_2 P(\xi = x_i | \eta = y_j)$$

\neq

Feltételes entrópia:

Mennyi bizonytalanság marad ξ -re nézve, ha η -t megfigyeltém?

$H(\xi, \eta)$ -ha egyiket sem figyeltém meg belőle ebből eltűnik η megfigyeléseinek bizonytalansága:

$$H(\xi, \eta) - H(\eta) =: H(\xi | \eta)$$

ξ feltételes entrópiája η szerint

Chagorin szabályt felhasználva:

$H(\eta)$ - elágazásig tartó

$H(\xi | \eta)$ - elágazás utáni bizonytalanság

$$H(\xi | \eta) = - \sum_{j=1}^n P(\eta = y_j) \cdot \sum_{i=1}^n P(\xi = x_i | \eta = y_j) \cdot \log_2 P(\xi = x_i | \eta = y_j)$$

$$= - \sum_{j=1}^n P(\eta = y_j) \cdot \sum_{i=1}^n \frac{P(\xi = x_i, \eta = y_j)}{P(\eta = y_j)} \cdot \log_2 \frac{P(\xi = x_i, \eta = y_j)}{P(\eta = y_j)}$$

$$= \sum_{j=1}^n \sum_{i=1}^n P(\xi = x_i, \eta = y_j) \cdot \log_2 \frac{P(\eta = y_j)}{P(\xi = x_i, \eta = y_j)}$$

$$H(S, K) \leq H(S) + H(K)$$

$$H(S, K) - H(K) \leq H(S)$$

Tétel (másik alakja)

$$H(S|K) \leq H(S) \quad \text{a. o. a. = , ha függetlenek}$$

\forall f h megfigyelhetjük K -t, vagy K egy $f(K)$ függvényét.

Lemma

$$H(S|K) \leq H(S|f(K)) \quad (\text{kevesebb ismeret az entropiát kevésbé csökkenti})$$

és egyenlőség feltétele:

$$\forall z \text{-re, amire } f(y) = z$$

$$P(S=x|K=y) = P(S=x|f(K)=z)$$

$$(S, f(K))$$

$$P(S=x, f(K)=z) = \sum_{f(y)=z} P(S=x, K=y)$$

$$P(f(K)=z) = \sum_{f(y)=z} P(K=y)$$

$$H(S|f(K)) = \sum_x \sum_z \left[\sum_{f(y)=z} P(S=x, K=y) \right] \cdot \log_2 \frac{\sum_{f(y)=z} P(K=y)}{\sum_{f(y)=z} P(S=x, K=y)}$$

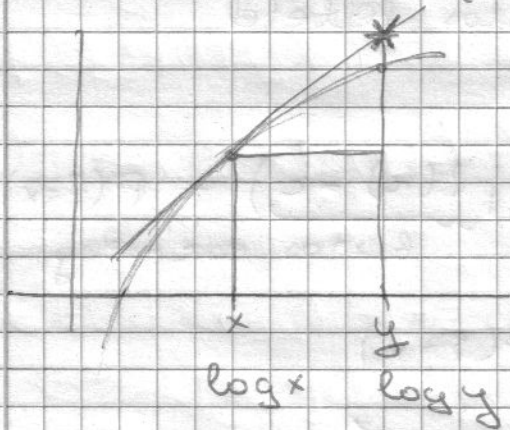
$$H(S|K) = \sum_x \underbrace{\sum_z \sum_{f(y)=z} P(S=x, K=y)}_{= \sum_y} \cdot \log_2 \frac{P(K=y)}{P(S=x, K=y)}$$

csak a $\sum_x \sum_z$ -n belüli összegek kell használni
 ezek típusa: $a: i=1, \dots, n$ $a = \sum a_i$
 $b: j=1, \dots, n$ $b = \sum b_j$

$$\sum_{i=1}^n a_i \cdot \log \frac{b_i}{a_i} \quad (H(S|K))$$

$$\leq a \cdot \log \frac{b}{a} \quad (H(S|P(K)))$$

= feltétele: $\frac{b_i}{a_i} = \frac{b}{a} \quad i=1 \dots n$



* közelről becüli $\log y$ -t

$$\log y \leq \log x + c_x (y-x)$$

$$= a \cdot \frac{b}{a} - a \cdot \frac{x}{a}$$

[előnyeg, ha konvex $x = \frac{b}{a}$
für $(\log z)$]

$$\log \frac{b_i}{a_i} \leq \log \frac{b}{a} + c \left(\frac{b_i}{a_i} - \frac{b}{a} \right)$$

$$= a \cdot \frac{b_i}{a_i} - a \cdot \frac{b}{a} \quad \forall i\text{-re}$$

$$* a_i \Rightarrow \sum_{i=1}^n a_i \log \frac{b_i}{a_i} \leq \sum_{i=1}^n a_i \log \frac{b}{a} + c \sum_{i=1}^n \left(\frac{b_i}{a_i} - \frac{b}{a} \right)$$

$$\leq \log \frac{b}{a} \cdot \sum_{i=1}^n a_i + \dots$$

$$\leq a \cdot \log \frac{b}{a} + \dots$$

$$z\text{-re: } \sum_{P(Y)=z} P(S=x, K=y) \log_2 \frac{P(K=y)}{P(S=x, K=y)} \leq$$

$$\leq \left(\sum_{P(Y)=z} P(S=x, K=y) \right) \cdot \log_2 \frac{\sum_{P(K)=y} P(S=x, K=y)}{\sum_{P(Y)=z} P(S=x, K=y)}$$

$$e_i = \frac{P(S=x, K=y)}{P(K=y)} = \frac{\sum_{P(Y)=z} P(S=x, K=y)}{\sum_{P(Y)=z} P(K=y)} \quad , \text{ha } P(Y)=z$$

$$P(S=x|K=y) = P(S=x|P(K)=z)$$

lemma

lemmából, hogyan bizonyítsuk bc, hogy:

$$H(S|K) \leq H(S)$$

válasszuk az $f(y) \equiv c$ fnt.

$$\Rightarrow H(S|K) \leq H(S|f(K)) \stackrel{\text{biztos feltétel}}{=} H(S)$$

= feltétele: $f(y) = c$

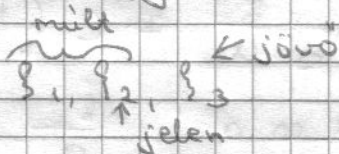
$$P(\xi=x|K=y) = P(\xi=x|f(K)=c) = P(\xi=x)$$

biztos esemény

$$\frac{P(\xi=x, K=y)}{P(K=y)} = P(\xi=x)$$

↓ függetlenek

lemma következménye:



■ tétel

$$H(\xi_3|\xi_1, \xi_2) \leq H(\xi_3|\xi_2)$$

$$= a \text{ cs } a, \text{ ha } P(\xi_3=x|\xi_1=y, \xi_2=z) = P(\xi_3=x|\xi_2=z)$$

"a legnagyobb bizonytalanság a jövőre

akkor van ha csak a jelen ismerné és a múltat nem"

(Markovlánc tulajdonság: "a jelen is-meretében a múlt és a jövő független")

"A jövő a múltól csak a jelenen keresztül függ"

2009.03.04., 4. előadás

Múlt órán: $\sum a_i \log \frac{b_i}{a_i} \leq a \log \frac{b}{a}$ N

ahol $a = \sum a_i$

$b = \sum b_i$

Következmények:

① (p_1, \dots, p_n) (q_1, \dots, q_n) isz. eloszlások
a: p_i b: q_i

$$\sum_{i=1}^n p_i \cdot \log_2 \frac{p_i}{q_i} \leq 0$$

(v. $\log_2 \frac{1}{1} = 0$)

$$-\sum_{i=1}^n p_i \log_2 p_i \leq -\sum_{i=1}^n p_i \log_2 q_i$$

② $b_i = 1 \quad i=1, \dots, n$

$a_i = (p_1, \dots, p_n)$

$$\sum_{i=1}^n p_i \cdot \log_2 \frac{1}{p_i} \leq \log_2 \frac{n}{1} = \log_2 n$$

(egyenletes entrópia maximális)

③ Jensen-egyenlőtlenség

$h(x)$ konvex sz., ξ - valsz. változó
várható érték

$$E(h(\xi)) \leq h(E(\xi))$$

(esetintésként h a logaritmus, E az entrópia)

Forrás jellemzése:

Tipikus közelítések: a lehető legbizonytalanabbra készülőnk, azaz felülről kell közelíteni a forrás entrópiáját

Válaztható szimbólumok: x_1, \dots, x_n

①/4 0-dol rendű közelítés - nincs megfigyelés

- Egyenletes eloszlást tételezzünk fel, ξ

ξ eloszlása: $P(\xi = x_i) = \frac{1}{n} \quad \forall i=1, \dots, n$

- függetlenek a választások, N hosszú választásnál (ξ_1, \dots, ξ_N)

$$P(\xi_1 = x_{i_1}, \dots, \xi_N = x_{i_N}) = P(\xi_1 = x_{i_1}) \cdot \dots \cdot P(\xi_N = x_{i_N}) = \left(\frac{1}{n}\right)^N$$

$$H = \lim_{N \rightarrow \infty} \frac{1}{N} \cdot H(\xi_1, \dots, \xi_N)$$

egy szimbólumra jutó átlagos entrópia

$$H(\xi_1, \dots, \xi_N) = N \cdot \log_2 n \quad (\text{egyenletes})$$

$$H = \lim_{N \rightarrow \infty} \frac{1}{N} N \log_2 n$$

$$H = \log_2 n$$

1-rendű közelítés:

M hosszú üzenetet megfigyelünk

x_i előfordulása: $m_i \quad \forall i=1, \dots, n$

Eloszlás közelítése: $P_i = P(\xi = x_i) = \frac{m_i}{M}, \quad i=1, \dots, n$

$$P(\xi_1 = x_{i_1}, \dots, \xi_N = x_{i_N}) = \prod_{s=1}^N P(\xi_s = x_{i_s}) = \prod_{i=1}^n P_i^{m_i}$$

(ahol a lehető legbizonytalanabb, ha függetlenek)

$$= \prod_{i=1}^n P_i^{m_i}$$

m_i - x_i szimbólum hányszor fordul elő a produktumban, azaz P_i hányszor fordul elő

$$H(\xi_1, \dots, \xi_N) = \sum_{i=1}^n H(\xi_i) = N \cdot H(\xi)$$

$$H = - \sum_{i=1}^n P(\xi = x_i) \cdot \log_2 P(\xi = x_i)$$

2-rendű közelítés:

M megfigyelés, m_i - x_i előfordulásainak száma

m_{ij} - az $x_i x_j$ egymás utáni előfordulásainak a száma

$P(\xi = x_i) = \frac{m_i}{M}$ - önmagában való választás

$P(\xi_t = x_i, \xi_{t+1} = x_j) = \frac{m_{ij}}{M}$ - egymás utáni párokra való vetületi (perem) eloszlás

2/4

$$P_{ij} := P(S_{t+1} = x_j | S_t = x_i) = \frac{P(S_t = x_i, S_{t+1} = x_j)}{P(S_t = x_i)} = \frac{\frac{m_{ij}}{n}}{\frac{m_i}{n}} = \frac{m_{ij}}{m_i}$$

$$P_{j|i} = \frac{m_{ij}}{m_i}$$

$$P(S_1 = x_{i_1}, \dots, S_N = x_{i_N}) = P(S_N = x_{i_N} | S_1 = x_{i_1}, \dots, S_{N-1} = x_{i_{N-1}}) \cdot P(S_1 = x_{i_1}, \dots, S_{N-1} = x_{i_{N-1}})$$

// Markovlan tulajd: $H(S_3 | S_1, S_2) \leq H(S_3 | S_2)$

Markov-lánc: (ehelyett van a legenyebb bizonyít.)

$$P(S_N = x_{i_N} | S_1 = x_{i_1}, \dots, S_{N-1} = x_{i_{N-1}}) = P(S_N = x_{i_N} | S_{N-1} = x_{i_{N-1}})$$

Teljesüljön $\forall x_{i_N}, x_{i_{N-1}}$ értékek N -től függetlenül.

$$= P_{i_{N-1} i_N}$$

homogén, egy lépéses Markovlánc:

$$P(S_{t+1} = x_i | S_1 = x_{i_1}, \dots, S_t = x_{i_t}, S_t = x_i) = P(S_{t+1} = x_i | S_t = x_i) = P_{ij}$$

P_{ij} = átmeneti valószínűségek

(átmenet valószínűség mátrix, \sum egy sor = 1)

$$P(S_1 = x_{i_1}, \dots, S_N = x_{i_N}) = P(S_1 = x_{i_1}, \dots, S_{N-1} = x_{i_{N-1}}) \cdot P(S_N = x_{i_N} | \dots) \\ = \dots = P(S_1 = x_{i_1}) \cdot P(S_2 = x_{i_2} | S_1 = x_{i_1}) \cdot \dots \cdot P(S_N = x_{i_N} | S_{N-1} = x_{i_{N-1}})$$

$M_{ij} = P(S_t = x_j | S_{t+1} = x_i)$ alakú tényező k száma

$$P(S_1 = x_{i_1}, \dots, S_N = x_{i_N}) = P(S_1 = x_{i_1}) \cdot \prod_{i,j} P_{ij}^{M_{ij}}$$

$$H(S_1, \dots, S_N) = H(S_1, \dots, S_{N-1}) + H(S_N | S_1, \dots, S_{N-1}) =$$

ahor max, ha csak az utolsótól függ

$$= H(S_1, \dots, S_{N-1}) + H(S_N | S_{N-1}) = \dots =$$

$$= H(S_1) + H(S_2 | S_1) + \dots + H(S_N | S_{N-1})$$

$$H(\mathcal{S}_t | \mathcal{S}_{t-1}) = \sum_{i=1}^n P(\mathcal{S}_{t-1} = x_i) \cdot H(\mathcal{S}_t | \mathcal{S}_{t-1} = x_i) =$$

$$H(\mathcal{S}_t | \mathcal{S}_{t-1} = x_i) = - \sum_{j=1}^n P(\mathcal{S}_t = j | \mathcal{S}_{t-1} = x_i) \cdot \log_2 P_{ji|i}$$

$$= - \sum_{j=1}^n P_{ji|i} \cdot \log_2 P_{ji|i} = H(\mathcal{S} | i)$$

$$H(\mathcal{S}_t | \mathcal{S}_{t-1}) = \sum_{i=1}^n P(\mathcal{S}_{t-1} = x_i) \cdot H(\mathcal{S} | i)$$

$$P(\mathcal{S}_t = x_i)$$

$$t \rightarrow \infty$$

Markov-lánc ... odikus tulajdonság

$$\lim_{t \rightarrow \infty} P(\mathcal{S}_t = x_i) = \pi_i$$

(Feltétel: minden állapot visszatérő legyen)
ne legyen periodikus

Markov folyamat határeloszlása - stationárius eloszl.

$$P(\mathcal{S}_t = x_i) = P(\mathcal{S}_{t-1} = x_j) \cdot P_{ij}$$

$$\lim_{t \rightarrow \infty} \pi_i = \sum_{j=1}^n \pi_j \cdot P_{ij} \quad \forall i \text{-re}$$

(ez egy lin. egy. rendszer)

(π_1, \dots, π_n) az átmenet valószínűségi mátrix
szajátvektora a $\lambda=1$ sajátértékkel.

legyen: $P(\mathcal{S}_t = x_i) = \pi_i$ - ettől kezdve $P(\mathcal{S}_t = x_i) = \pi_i$
marad, az időtől nem függ
(azaz station.)

$$H(\mathcal{S}_1, \dots, \mathcal{S}_n) = H(\mathcal{S}_1) + \sum_{j=1}^n H(\mathcal{S} | j) \cdot \left(\sum_{t=2}^n P(\mathcal{S}_{t-1} = x_j) \right)$$

$$\lim_{n \rightarrow \infty} \frac{1}{n} H(\mathcal{S}_1, \dots, \mathcal{S}_n) = \sum_{j=1}^n H(\mathcal{S} | j) \cdot \left(\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{t=2}^n P(\mathcal{S}_{t-1} = x_j) \right)$$

$$= \sum_{j=1}^n \pi_j \cdot H(\mathcal{S} | j)$$

$\downarrow \pi_j$

4/4

$$H(s_1, \dots, s_n) = \sum_{j=1}^n \pi_j \cdot H(s | j)$$

3-ahérendű:

m i k - $X_i X_j X_k$ egymás után

$$P(\xi_t = X_k | \xi_{t-1} = X_j, \xi_{t-2} = X_i) - \text{rögzítjük } \frac{m_{ijk}}{m_{ij}}$$

Kétlépéses Markov lánc:

$$P(\xi_t = X_k | \xi_1 = X_i, \dots, \xi_{t-3} = X_i, \xi_{t-2} = X_j, \xi_{t-1} = X_j)$$

$$= P(\xi_t = X_k | \xi_{t-2} = X_i, \xi_{t-1} = X_j)$$

(betűszimbólumból állóink a szóak megfigyelésére ...)

(Ez a levezetés nem kell vizsgára)

Infókeres 4., 5. előadás, 2009. 03. 11.

"nagy számok törvénye" típusú

Kosszi sorozatok jellemzése:

Forma: független, azonos eloszlású

// legrosszabb eset
// legbizonytalanabb

X_1, \dots, X_n - választható jelek

$$P(\xi = x_i) = P_i$$

$$P(\xi_1 = x_{i_1}, \dots, \xi_n = x_{i_n}) = \prod_{i=1}^n P_i^{M_i} \quad / = M_i \text{-es } x_i \text{ előfordulási száma}$$

$$P(\vec{\xi} = \vec{x})$$

(relatív gyakoriság tart a vsz. -hez) \leftarrow NSZT

$$M_i = N \cdot P_i, \quad i=1 \dots n$$

$$P(\vec{\xi} = \vec{x}) = \prod_{i=1}^n P_i^{N \cdot P_i} = 2^{-\sum N \cdot P_i \cdot \log_2 P_i} = 2^{-N \cdot H}$$

Tétel: Tetszőleges $\varepsilon > 0, \delta > 0$ - hoz létezik N_0 , hogy

$N > N_0$ esetén az N kosszi sorozatok két halmazba sorolhatók:

(1) lényegtelen halmaz, annak a vszge, hogy ebből választunk $< \varepsilon$

(2) Tipikus halmaz (ebbe $1-\varepsilon$ -nél nagyobb vszgeel esik a választás)

\vec{x} tipikus

$$\left| \frac{\log_2 P(\vec{\xi} = \vec{x})}{N} - H \right| < \delta$$

összetöleg H

Nagyon kosszi jel sorozat esetén a tipikus halmazban majdnem egyenletesen terjed a választás

Szerepe a tételnek:

Tipikus halmaz: $|M_n|$ - adjunk becslést az elemszámra

Becsljük a $P(\vec{\xi} = \vec{x})$ értéket, ahol \vec{x} tipikus:

$$2^{-N(H+\delta)} < P(\vec{\xi} = \vec{x}) < 2^{-N(H-\delta)}$$

$$|M_n| < 2^{N(H+\delta)}$$

Biz: (tétel)

ξ

$$\prod_{i=1}^n p_i \cdot n_i$$

$n_i \quad i=1, \dots, n$

n_i -t előállítani: $S_j^{(i)} \quad j=1, \dots, N \quad S_j^{(i)} = \begin{cases} 0, & \text{ha } \xi_j \neq X_i \\ 1, & \text{ha } \xi_j = X_i \end{cases}$

$$n_i = \sum_{j=1}^N S_j^{(i)}, \quad P(S_j^{(i)} = 1) = p_i$$

$$M(n_i) = N \cdot M(S_j^{(i)}) = N \cdot p_i \quad \text{M-várható érték}$$

$$D^2(n_i) = N \cdot D^2(S_j^{(i)}) = N \cdot (p_i - p_i^2) = N \cdot p_i \cdot (1 - p_i) \quad \text{D-szórás } [E(X^2) - (E(X))^2]$$

Chebisev-egyenlőtlenség:

$$P\left(\left(\xi - M(\xi)\right)^2 > \gamma^2\right) \leq \frac{D^2(\xi)}{\gamma^2}$$

Alkalm

$$P\left(\left(n_i - N p_i\right)^2 > (\xi' N)^2\right) \leq \frac{N \cdot p_i \cdot (1 - p_i)}{(\xi' N)^2} < \frac{\text{Konstans}}{N} < \xi'$$

Megválasztjuk N -et olyan nagyra, hogy ξ' -től függően ez teljesüljön $\forall i=1, \dots, n$ -re.

Erre az N -re

$$P(\text{legalább egy } i\text{-re } |n_i - N p_i| > \xi' N) \leq n \cdot \xi'$$

$$\xi - \text{leígyegtelen halmaz: } \xi' = \frac{\xi}{n}$$

A komplementer esemény tipikus halmaza:

$$\forall i\text{-re } |n_i - N p_i| < \xi' N$$

$$N(p_i - \xi') < n_i < N(p_i + \xi')$$

$$\prod_{i=1}^n p_i^{N(p_i + \xi')} < \prod_{i=1}^n p_i^{n_i} < \prod_{i=1}^n p_i^{N(p_i - \xi')} \quad \parallel p_i < 1$$

$$N \cdot \left(\sum_{i=1}^n \log_2 p_i (p_i + \xi')\right) < \log_2 P(\vec{\xi} = \vec{x}) < N \cdot \left(\sum_{i=1}^n (p_i - \xi') \cdot \log_2 p_i\right)$$

$$H + \xi' \cdot \sum_{i=1}^n \log_2 p_i < \frac{\log_2 P^{-1}(\vec{\xi} = \vec{x})}{N} < H - \xi' \cdot \sum_{i=1}^n \log_2 p_i$$

// H behoz egy miniszert és a $\log_2 P^{-1}$ fordít az előjel

S

Tipikus halmaz (tulajdonsága)

A legvalószínűbb választások megadásával:

$$P_1 \geq P_2 \dots \geq P_n$$

legszűrűbb választás N db X_i , $M_1 = N$
és nincs benne a tipikus halmazba.

$$\vec{x}_1, \vec{x}_2, \dots$$

$$P(\vec{X} = \vec{x}_i) \leq P(\vec{X} = \vec{x}_{i-1})$$

lényeges sorozat

$$\sum_{i=1}^L P(\vec{X} = \vec{x}_i) \geq 1 - \lambda$$

$$2^{N(H-\delta)} < L(\lambda) < 2^{N(H+\delta)}$$

$$\sum_{i=1}^{L-1} P(\vec{X} = \vec{x}_i) < 1 - \lambda$$

↑ tipikus halmaz
elemzámának becslése
igaz rá

$L(\lambda)$ és M_N nem ugyanaz, de kevés elem esik a metszetükön kívül.

Kódolás és csatornkapacitás viszonya:

Forrás: S_1, \dots, S_N kibocsátott

ezt a csatornában $M_1, \dots, M_M(S_1, \dots, S_N)$ formában vesszük át

legegyszerűbb kódolás: $f(S_1, \dots, S_N) = (M_1, \dots, M_M(S_1, \dots, S_N))$

$$H(f(S)) \leq H(S) \quad \text{és " = " } \Leftrightarrow \text{ ha } f \text{ invertálható}$$

(a közvetlen megfigyeléshez képest a közvetett megfigyelés csak vesztes információt)

⇓

$$\hookrightarrow H(S_1, \dots, S_N) \geq H(M_1, \dots, M_M(S_1, \dots, S_N))$$

Nincs veszteség: invertálható a kódolás, azaz $\exists f^{-1}$

$$f^{-1}(M_1, \dots, M_M) = (S_1, \dots, S_N)$$

$$\Rightarrow H(M_1, \dots, M_M) \geq H(S_1, \dots, S_N)$$

③ / 5

Teljes, ha nincs veszteség: $H(M_1, \dots, M_M) = H(S_1, \dots, S_N)$

Ha VDA-val kidolgoz, akkor a kódoló automa-
tának még van egy plusz állapota is van: $+d$

A csatorna kimenetén ezt nem észlelem, kisebb
az entrópia...

T ideig működik a csatorna, T hosszú jelsoxozat
Mekkora ekkor a csatorna kimenetén a maximális
entrópia? C bit/sec

Különböző T hosszú jelsoxozatok száma: $N(T)$

$$C = \lim_{T \rightarrow \infty} \frac{\log_2 N(T)}{T}$$

Maximális: egyenletes eloszlás esetén: $\log_2 N(T) = A$, ahol
 $T(C-S) < N(T) < T \cdot (C+S)$ $\Rightarrow T \cdot (C-S) < N < T \cdot (C+S)$

(T ideig működő forrás entrópiája nem halad-
hatja meg a T ideig működő csatorna ^{max} entrópiáját)

Tajmentes csatorna alapfeltételének bizonyítása

Forrás: H bit/szimbólum
Csatorna: C bit/sec

Tétel:

a) Nem lehet C/H -nél nagyobb sebességgel működtet-
ni a forrást úgy, hogy V üzenet a csatorna
kimenetéből visszaállítható legyen
(Nem lehet veszteségmentes)

b) Tetszőleges $S > 0$ -hoz létezik kódolás, hogy
 $C/H - S$ sebességgel veszteségmentesen lehet
üzenetet átadni.

Biz:

a) T idő, V sebesség szimb/sec
Mennyi a forrás entrópia T idő alatt: $H \cdot T \cdot V$ -nél kisebb

$$V \cdot T \cdot (H - S) < H(T) < V \cdot T \cdot (H + S)$$

Csatorna kimenet max. entrópiája: $(T(C-S), T \cdot (C+S))$ között

$$V \cdot T \cdot (H - S) > T \cdot (C + S) \quad \text{- ez biztos veszteséges}$$

$$V \cdot T \cdot (H - S) \leq T \cdot (C + S)$$

$$V \leq \frac{C + S}{H - S}, \quad S, s \text{ tetszőlegesen kicsi}$$

Tehát V tetszőlegesen közelebe kerülhetünk

2009.03.17., 6. előadás, Infóker 4.

$S > 0$ $\frac{C}{H} - S$ sebességgel beszűkítésmentes

legyen az üzenetátvitel

(görög betű: tetszőleges kicsi lehet, és ha elég nagy az ~~n~~ ^{üzenethosszáig} N , akkor az egyenlőtlenség teljesül)

N hosszú üzenet, kérdés: milyen hosszú csatornajelet tartozik hozzá

(Általában igaz ergodikus forrásra)

Független, azonos eloszlású forrás:

X_1, \dots, X_n

N hosszú üzenetek:

(i) Tipikus üzenetek, $1 - \varepsilon$ -nél nagyobb valószínűséggel elemszáma $\leq 2^{N(H+S)}$

(ii) Nem tipikus üzenetek, ε -nél kevesebb valósz-séggel elemszáma: $\leq n^N = 2^{N \log_2 n}$

(i) halmaz kötélszáma: T_1

$$2^{T_1(C-\delta)} \geq 2^{N(H+S)}$$

$$T_1 \geq N \cdot \frac{H+S}{C-\delta} = N \cdot \frac{H}{C} + \delta_1$$

(ii) Marad felhasználhatóan T_1 hosszú csatornajelet vegyünk egy ívet, után T_2 hosszú jellel kódoljuk a nem tipikus üzeneteket

$$T_2 > N \frac{\log_2 n}{C-\delta} = N \cdot \frac{\log_2 n}{C} + \delta_2$$

A szimbólumok eloszlása: (P_1, \dots, P_n)

$X_1 - P_1$ uszgel választom

\vdots
 $X_n - P_n$

Ködlössz várható értéke: $\sum_{i=1}^n P_i \cdot l(w_i)$

Tétel: $H(P_1, \dots, P_n) \leq \sum_{i=1}^n P_i \cdot l(w_i)$ minden prefix kódra

(Biz1 indirekt + de qkior nagyobb sebességet érhetnek el mint az alaptétel mond)

Biz2:

a_i, b_i, a, b

$$\sum_{i=1}^n a_i \log \frac{b_i}{a_i} \leq a \log \frac{b}{a}$$

$a_i \rightarrow (P_1, \dots, P_n)$ $(q_1, \dots, q_n) \leftarrow b_i$

$$\sum_{i=1}^n P_i \log_2 \frac{q_i}{P_i} \leq 1 \cdot \log_2 1 = 1 \cdot 0 = 0$$

$$H(P_1, \dots, P_n) \leq - \sum_{i=1}^n P_i \log_2 q_i$$

$$\sum_{i=1}^n 2^{-l(w_i)} = 1, \quad q_i := 2^{-l(w_i)}$$

$$H(P_1, \dots, P_n) \leq - \sum_{i=1}^n P_i \log_2 2^{-l(w_i)}$$

$$H(P_1, \dots, P_n) \leq - \sum_{i=1}^n P_i (-l(w_i))$$

$$H(P_1, \dots, P_n) \leq \sum_{i=1}^n P_i \cdot l(w_i)$$

Nevezetes kódolások:

1) Shannon-Fano kód:

$$\sum_{i=1}^n P_i \cdot l(w_i) < H + 1 \quad (1 \text{ bit veszteséggel teljes megközelítéssel az entrópiát})$$

Kód elkészítése: $P_1 \geq P_2 \geq \dots \geq P_n$

$$Q_1 := 0 \quad (i=1)$$
$$Q_i := \sum_{j=1}^{i-1} P_j \quad (i=2, \dots, n)$$

Kód: binárisan kódoljuk a a_i számokat,

$0, a_1, \dots, a_n$ l_i pontossággal megadjuk

úgy válasszuk l_i -t, hogy $2^{-l_i} \leq p_i \leq 2^{-l_i+1}$

($l_i - 1$ 0-as bittel kezdődik)

$$l_i \geq -\log_2 p_i \geq l_i - 1 \quad | \cdot p_i \text{ és } \sum$$

$$-\sum_{i=1}^n p_i \log_2 p_i \geq \sum_{i=1}^n p_i \cdot (l_i - 1)$$

$$H \geq \left(\sum_{i=1}^n p_i l_i \right) - 1$$

$$\sum_{i=1}^n p_i l_i < H + 1$$

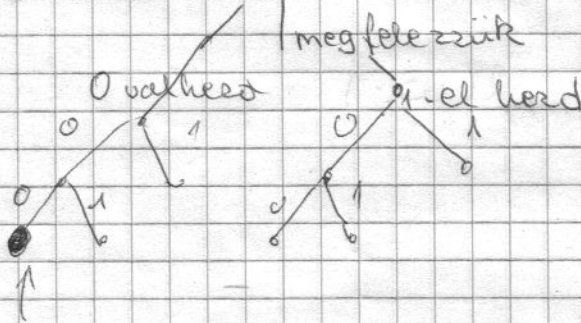
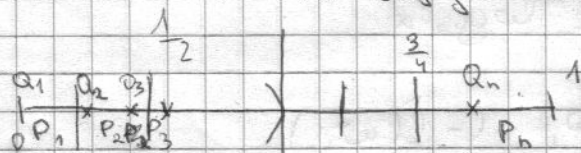
prefix kód?

- hosszabbak nőnek

- "érték" nő

w_i, w_{i+1} biztos teljesül a prefixmentesség?

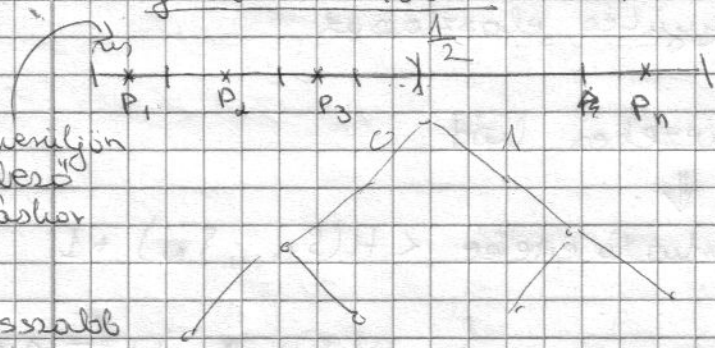
Igen, w_i hosszabbak az utolsóhoz valamit és így változik, nem lesz prefix.



leáll: az intervallumban egyetlen pont marad

(Mátrix: meghatározni kell az elosztást)

2) Gilbert - Moore kód : P_1, \dots, P_n nincs rendezve



intervallum feleső-
pontjára
felezéses kódol

ide kerüljön
a feleső
leálláshoz

leáll - ha egyedül
marad

1-el hosszabb
a kód

(1-el többször kell
felezni : $\frac{P_i}{2}$ legyen
benne)

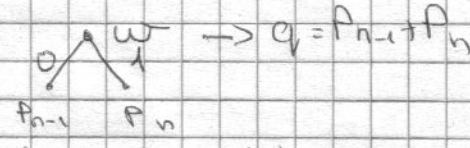
$$\sum_{i=1}^n P_i l_i < H + 2$$

3) Huffman - kód

Rekurzió: P_1, \dots, P_{n-1}, P_n

a két legkisebb : P_{n-1}, P_n

kódja w_0 P_{n-1}
 w_1 P_n



Aztán vesszük a (P_1, \dots, P_{n-2}, q) Huffman kódját
 w

Optimális: Bármely más prefix kódra : w'_i

$$\sum_{i=1}^n P_i l(w_i) \geq \underbrace{\sum_{i=1}^n P_i l(w'_i)}_{\text{Huffman-kód}}$$

$$\sum_{i=1}^n P_i l(w_i) < H + 1$$

Ha gyákoriság van : x_i n_i - szer fordul elő

$$\sum_{i=1}^n n_i = N$$

$$\min \sum_{i=1}^n n_i \cdot l(w_i)$$

$$\Leftrightarrow \min \sum_{i=1}^n \left(\frac{n_i}{N} \right) \cdot l(w_i)$$

↑
classis

Közelítsük $\frac{1}{N}$ hányadost:

Vegyük a (ξ_1, \dots, ξ_N) együttes eloszlását

$H(\xi_1, \dots, \xi_N) \rightarrow$ Nevezetes kód



Kódl hossz várható értéke $< H(\xi_1, \dots, \xi_N) + 1$

Egy szimbólumra jutó kódl hossz: $\frac{1}{N} \cdot H(\xi_1, \dots, \xi_N) + \frac{1}{N} \sim H + \frac{1}{N}$

(tetszőlegesen közel tudok lenni a $\frac{1}{N}$ sebességgel, ha nagyon sok hosszú kódsorozatot kódolok egyben)

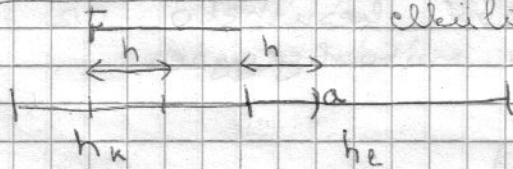
"Nagyon nagy kódl, valamelyen tipikus halmazon való viselkedéshez hasonlóan működnek."

4. előadás, 2009.03.25., Informatika 4.

Tómórités:

Lempel-Ziv kódok

- LZ77 kód: (hivatkozok az előzőekben már elemezte kódokra)



már volt est akaruk most hirtelen

- minden szimbólumhoz ^{kezd} kód, azaz ha n szimbólum van - $\log_2 n!$ hosszú tudjuk kódolni

Átkódolt kód: (b, h, c) ahol c az "a" kódja

$$L(b) \sim \log_2 h_k$$

$$L(h) \sim \log_2 h_e$$

időtől független "teljes" időbeni és térbeli átlazhatóság felismerhető

Stacionárius, ergodikus forrás esetén - ezzel elérjük az optimális kódolást

Szótár kódok - szótárépítő kódolás

LZ78 kódolás:

- kezdetben # szimbólumnak van egy kódja

- az új szöveg részen megnezi, milyen az előzőekkel leghosszabb egyezés van amivel egyezik, az már szótárban van és legyen a szótárkódja!

- Átkódolt: $\langle i, c \rangle$ - ahol c a folytató "a" kódja

- ebből készül új szótárbejegyzés

Stac + ergo. forrás esetén - ennél is elérhető az optimalitás, azaz a max sebesség.

①/4

(de egy idő után a késleltetés nagy lehet)

LZW-kód (T. Welch)

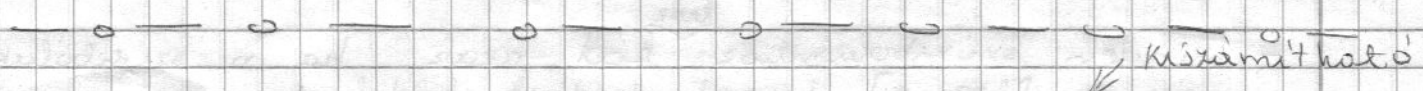
c legyen "a" kódja

az előző $\langle i, c \rangle$ párból csak i -t küldöm, a következő közbendő szö: "a"-val kezdődik

Majd a következő üzenetből készíthető az LZW kód szerint $\langle i, c \rangle$ szótárbejegyzés.

(Pl GIF is ilyen használ)

Korlátozzuk a szótármeletet!



Kolmogorov-Bonyolultság
(-u- -entropia)

Algoritmikus
információelmélet

(1960-as
évek közepéig)

(1964) - Műlt tömör leírása

(Sh: a kódhosszt minimalizálni)

(Kol: egy nagyon hosszú kódot, hogyan tudok tömörebbet előállítani, és hogy állítsam ebből vissza)

- Egy nagy adat lehető legrövidebb tömöreteje, amiből visszaállítható

(Hogyan lehet a kódban es lehet tömöríteni, hogy egy számítógép ebből visszaadja a szöveget nyomtatva?
0-ból visszaállítja
Nem elég a kód, kell a dekódolás is)

Tömörítetőség:

x tömöríteni,	kód P ,	$I(P) = x$
x alakja	$I(P)$	x kódja
		"hossz"
$x \in \mathbb{N}$	$I(P) = P$	x
$x = 2^k$	$I_1(P) = 2^P$	k
$x = 2^{2^k}$	$I_2(P) = 2^{2^P}$	k
$x = 2^{2^{2^k}}$	$I_3(P) = 2^{2^{2^P}}$	k
		$\log_2 \dots \log_2 x$

2/4

Egyek a tömörítő ω -eket használhatom.

x - tömörítése

i - melyik ω -t használom

p - híd (tömörítetten a legkövetkező)

(i, p)

$$X = 1025 = 1024 + 1 = 2^{10} + 1$$

Lehetne: $f_{\omega}(p) = 2^p + 1$

f - kiszámítható, és ezek megszámlálható sokan vannak, tehát az f -ek halmazára rekurzíván felsorolhatjuk.
Ilyen ω -ek: építkezések:
1. építkezés - Turing gép

2. Parciális rekurzív ω -ek:

$$\omega: \mathbb{N}^k \rightarrow \mathbb{N} \quad \mathbb{N} := \mathbb{N}_0$$

Alap ω -ek:

- zéró ω : $Z: \mathbb{N} \rightarrow \mathbb{N}, Z(x) = 0$

- növegető: $S: \mathbb{N} \rightarrow \mathbb{N}, S(x) = x + 1$

- vetítő (projekció): $P_{n,i}: \mathbb{N}^n \rightarrow \mathbb{N}, P(x_1, \dots, x_i, \dots, x_n) = x_i$

Függvény képzések:

- Kompozíció: $f: \mathbb{N}^n \rightarrow \mathbb{N}$ képződik
 $h: \mathbb{N}^k \rightarrow \mathbb{N}$
 $g_1, \dots, g_k: \mathbb{N}^n \rightarrow \mathbb{N}$

$$h(g_1(x_1, \dots, x_n), \dots, g_k(x_1, \dots, x_n)) =: f(x_1, \dots, x_n)$$

- Primitív rekurzió: $f: \mathbb{N}^{n+1} \rightarrow \mathbb{N}$
 $g: \mathbb{N}^n \rightarrow \mathbb{N}$
 $h: \mathbb{N}^{n+2} \rightarrow \mathbb{N}$

$$f(x_1, \dots, x_n, 0) =: g(x_1, \dots, x_n)$$

$$f(x_1, \dots, x_n, y+1) = h(x_1, \dots, x_n, y, f(x_1, \dots, x_n, y))$$

Primitív rekurzív ω -ek osztály a:

- zárt a kompozícióra, primitív rekurzióra (z.s.p.)

- generáló sorozat

(Lehet olyan ω -t ami nem primitív rekurzív, ezért "büvös" meg!)

- A generáló sorozathoz könnyű grammatikát szerkeszteni (véges írással)

- sorszámhatós a generáló sorozatok

- univerzális fr:

$u(k, x)$: - \forall generáló sorozathoz van k , amire $u(k, x) = f(x)$

- az $u(k, x)$ minden k -ra primitív rekurzív

- $u(k, x)$ kiszámítható

- $u(k, k) + 1 = h(x)$ fr-t tekintsük

-- kiszámítható

-- $h(x)$ milyen fr?

-- $h(x)$ nem primitív rekurzív

-- ha az lenne, h lenne a sorszáma

$$h(n) \leq u(n, n) + 1$$

$$= u(n, n) \quad (\text{hiszen } n \text{ a } h \text{ sorszáma})$$

-- ellentmondás $\rightarrow h(x)$ nem primitív rekurzív

- Pr. rek. fr. hátránya: lassan növekszenek

- Ackermann fr-ek

$$A : \mathbb{N}^2 \rightarrow \mathbb{N}, \quad A(0, y) = y + 1$$

$$A(x+1, 0) = A(x, 1)$$

$$A(x+1, y+1) = A(x, A(x+1, y))$$

$$a_m(x) = A(m, x)$$

a_4 : toronyhatvány, aminek a magassága X .

a_5 : -||- , -||- , $a_4^*(x)$

- Egy pr. rek. fr. az \forall valamelyik Ackermann fr-nél lassabban növekszik.

- $a(m) := a_m(m)$ - \forall pr. rek. fr.-nél gyorsabban növekszik

- Még egy lépés szükséges:

- minimalizáció λ -képzés:

$$f: \mathbb{N}^n \rightarrow \mathbb{N} \text{ képződik}$$

$$g: \mathbb{N}^{n+1} \rightarrow \mathbb{N}$$

$$f(\vec{x}) = \text{Mg}(|g(\vec{x}, y) = 0|)$$

az $f(x_1, \dots, x_n) = k$, ha:

1) $g(x_1, \dots, x_n, k) = 0$

2) $y < k : g(x_1, \dots, x_n, y) \neq 0$

3) $g(x_1, \dots, x_n, y)$ értelmezve van $\forall y \leq k$

(Teljesen értelmezhető 0-tól k -ig is)

// Pl: $g(x, y) = x + y + 1$

$f(x)$ - sehol nincs értelmezve

(Parciális λ -t kapunk)

- Minimalizáció hozzártelevel a parciális rekurzió λ -t lehet kapjuk.

- Ez az osztály azonos a Turing-biszámítható ~~halmaz~~ λ -ek osztályával (Church-tézis)

Szigorú változat: reguláris minimalizáció

- $\forall x_1, \dots, x_n$ -re van $h: g(x_1, \dots, x_n, h) = 0$

- g totális λ .

\Rightarrow Rekurzió λ -ek osztályát kapjuk.

Az unioverzális λ - nem veszt ki a parciális rekurzió λ -ek osztályából. $g(x)$

$$U(h, x) = g(x), \text{ ha } g \text{ sorozata}$$

\uparrow az ott, ahol értelmezve vannak

$$h(x) = U(x, x) + 1$$

Sorozat: n

$$U(n, x) = h(x)$$

$$h(n) \begin{cases} = U(n, n) + 1 \\ = U(n, n) \end{cases}$$

ha $h(n)$ értelmezve lenne \Rightarrow ellentmondás

$\Rightarrow h(n)$ nincs értelmezve

(általós λ -ek elrontása nem értékelhető)

-Ugyanez nem működik a reguláris λ -ekre

\Rightarrow nincs univerzális felsoroló λ -ünk

(regularitás eldönthetetlen algoritmikusán)
 \Rightarrow tehát nem lehet megkonstruálni a generáló sortól

// értékelhetőség azon működik, hogy egy minimalizáció után nincs feltétel olyan helyen, ahol a λ lehet nincs értelmezve.

2009.04.01, # 8. előadás, Indukció 4.

(XML-ről Ulman-Widow könyv alapján)

Kleene-jele normál alak: $\exists U, V$

$U: \mathbb{N} \rightarrow \mathbb{N}$, $V: \mathbb{N}^3 \rightarrow \mathbb{N}$ prim. rekurzió
függvények

úgy, hogy bármely $f: \mathbb{N} \rightarrow \mathbb{N}$ parciális rekurzió
függvényhez létezik n ,

$$f(x) = U(\mu t (V(n, x, t) \neq 0))) \quad \forall x \text{-re}$$

\mathbb{R}

$A \subset \mathbb{R}$

$x \in A$ véges időben megmondható

$$h_A(x) := \begin{cases} 0 & \text{ha } x \in A \\ 1 & \text{ha } x \notin A \end{cases} \quad (\text{ha } h \text{ totális f.})$$

A -rekurzió (elődítható), ha h_A rekurzió

$A \subseteq \mathbb{R}$ rekurzióval felsorolható, ha létezik $f: \mathbb{N} \rightarrow \mathbb{N}$
rekurzió f., hogy $x \in A \Leftrightarrow \exists n: f(n) = x$

ha $x \in A$ - véges lépésben bizonyítható,
ha $x \notin A$ - tudjuk-e igazolni

A rekurzió $\Leftrightarrow A$ és \bar{A} is rekurzióval felsorolható

Rice-tétel: (U-univerzális f. szerinti sorrendezés)
"A" tetszőleges részhalmaza a parciális rekurzió-függvényeknek

$A = \{n \mid f_n \in A\}$, ekkor A a.c.s.a. eldönthető
ha $A = \emptyset$ v. $A = \text{teljes}$

(ilyen tulajdonság pl., h. a Turing gépek
probléma megoldása, vagy hogy egy
parciális rek. f. rekurzió-e)

1/8

Unio Turing gép - i-edik Turing gépet
először p szimulálja, ha i -t először neki
paraméternek

Kolmogorov - Bonyolultság: (Algoritmikus entropia)

(Kolmogorov, Chaitin, Solomonoff) ~ 1964

deképezést készítsunk:

$$\Omega = \{0, 1\}, \quad \mathbb{N} = \{0, 1, \dots\}$$

Λ	-üres szó	0
0	} 1 bit hosszú	1
1		2
00		3
01		4
10		5
11		6
000		7
001	- - - - -	8
010		
011		
100		
101		
110		
111		

(100 - fordított sorrend)

Kölcsönösen egyértelmű

$$x \begin{cases} \in \Omega \\ \in \mathbb{N} \end{cases}$$

$$l(x) = x \text{ hossza } (\approx \log_2 x)$$

Konkaténáció: xy

Rendezett párok kódolása: $(x, y) \rightarrow z$

①

	0	1	2	3	4
0	00	01	10	11	
1		10	11		
2			10	11	
3				10	11

(x, y) hossza vagy lehet, bár x pl. lehet

② Önkorlátos kódolás: \bar{x} (prefix kód az \mathbb{N} fölött)

$$x = d_1 d_2 \dots d_n$$

$$\bar{x} = \underbrace{11\dots 1}_n 0 d_1 \dots d_n$$

$$l(\bar{x}) = 2l(x) + 1$$

$$x' = \bar{n} x$$

$$l(x') = 2 \log_2 n + l(x) + 1$$

$$2 \log_2(l(x)) + l(x) + 1$$

Celhetne folytatni \bar{n} helyett n' , stb. ...

② / 8

(x, y)

$z_0(x, y) = xy = z_0 \xrightarrow{\text{vizsfejtés}}$

$\Pi_{0,1}(z_0) = x, \Pi_{0,2}(z_0) = y$

$z_1(x, y) = x' y = z_1$

$\Pi_{1,1}(z_1) = x, \Pi_{1,2}(z_1) = y$

X - mennyire tömöríthető?

$f: \mathbb{N} \rightarrow \mathbb{N}$

Def: Az $x \in \mathbb{N}$ -re bonyolultsága az $f: \mathbb{N} \rightarrow \mathbb{N}$ függvény szerint

$C_f(x) = \sum \min \{l(p) \mid f(p) = x\}, \infty$, ha nincs (c-bonyolultság)

Alaptétel: (Invariencia tétel) (K-Ch-S)

létezik olyan optimális $f_0: \mathbb{N} \rightarrow \mathbb{N}$ par. rek, hogy $\forall f: \mathbb{N} \rightarrow \mathbb{N}$ -hez létezik k, c konstans, hogy $C_{f_0}(x) \leq C_f(x) + k, \forall x \in \mathbb{N}$ -re.

Következménye:

f_0, g_0 is optimális

$|C_{f_0}(x) - C_{g_0}(x)| \leq k$ $\forall x \in \mathbb{N}$

Optimális f_0 -ek szerinti bonyolultság konstans erejéig egyértelmű

\Rightarrow Rögzítsük az f_0 optimális f_0 -t.

Def: Az $x \in \mathbb{N}$ Kolmogorov bonyolultsága (entrópiája)

$C(x) = C_{f_0}(x)$

(tömöríthetőség első mérése)

Alaptétel bizonyítása:

Legyen $x, f: \mathbb{N} \rightarrow \mathbb{N}$ tetszőleges

$C_f(x) = k$, létezik $f(p) = x : l(p) = k$

Univerzális $U(y, n)$ f_0 -t, ahol n : sorok

f -hez $\exists n : U(y, n) = f(y)$
z-t választom úgy, hogy $z := n \cdot p$ -ből

$f_0(z) = U(\Pi_{1,2}(z), \Pi_{1,1}(z)) = U(p, n) = f(p) = x$

$$\Rightarrow C_{f_0}(x) \leq C(z) = \underbrace{C(n^*)}_{k_f} + C(p) = C(p) + k_f = C_f(x) + k_f$$

~~Egyenlős leírásának hossza: $C(n^*)$~~
~~Ködleírás: $C(n)$~~ } \Rightarrow

~~Tömítésnél \Rightarrow minimalizáljuk a leírás és a visszaírás hosszát~~

Alaptétel dohs dg ok $C(x)$ sz:

1. Mit jelent az, hogy $C(x) = k$?

\Rightarrow létezik egy p kód, $f_0(p) = x$ és $l(p) = k$

1a. $C(x) \leq k$? $f_0(p) = x$ és $l(p) \leq k$

(2^{k+1} ilyen kód lehet)

2. Alaptétel: felülről becsülhető azaz $f: \mathbb{N} \rightarrow \mathbb{N}$

$$C_{f_0}(x) \leq C_f(x) + k_f$$

3. Van-e minden x -nek bonyolultabbja?

$C(x)$ totális függvény \mathbb{N} -es felhasználással

$f(p) = p$ tot. használva

$$C(x) \leq C_f(x) + k_f = l(x) + k_f$$

(A Kolm. bony: max konstanssal nagyobb a hosszal)

"Semmiem bonyolultabb mint annyien"

4 $C(x)$ nem kiszámítható

$C_f(x)$ - ha f nincs értelmezve a p helyen, $l(p) = k$

\Rightarrow nem tudom megmondani, h p semmiem
nem a legjobb

Minimalizáció - nem tudom megmondani,

hogy nem kisebb-e k_2 -nél valószínű,

Inklúzió $C(x)$ kiszámítható:

x_0 : az első x , amire $C(x) \geq 1 = 2^0$

x_1 : az első x után, ami $C(x) \geq 2$

④/8

x_k az első x_{k-1} után, amire $C(x) \geq 2^k$

Legyen:

$$f(p) = x_p$$

$$C_f(x_p) \leq l(p) \quad \left. \begin{array}{l} \text{Inv. tul.} \\ C(x_p) \geq 2^p \end{array} \right\} C(x_p) \leq l(p) + k_f$$

$$C(x_p) \geq 2^p \quad \left. \begin{array}{l} \text{Inv. tul.} \\ C(x_p) \leq l(p) + k_f \end{array} \right\} C(x_p) \leq l(p) + k_f$$

$$2^p \leq l(p) + k_f \quad | \text{ellentmondás}$$

\Rightarrow nem kiszámítható

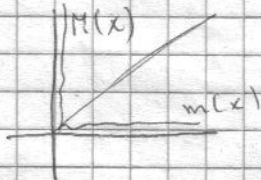
5. $\lim_{x \rightarrow \infty} C(x) = +\infty$ // (1)-es alapján

6. $m(x) := \min_{x \leq y} C(y)$ ($C(x)$ alsó burkolója)

$$\lim_{x \rightarrow \infty} m(x) = +\infty$$

$m(x)$: \uparrow monoton kiszámítható f.-nél lassabban

(Indirekt biz. szel)



$M - m(x)$ \rightarrow a túlközvetítés az $x=y$ -n

Ezek között van a bonyolultság

7. $(x+h)$ -t ábrázolni

$$C(x+h) \leq C(x) + 2 \cdot C(h) + k_f$$

(húszit változtatok az x -en húszit változik a bonyolultság)

$g(x,h) = y$ fut használók

$$C(g(x,h)) \leq C(x) + 2 \cdot C(h) + k_g$$

$\underbrace{k_g}_{\text{változtatás kioldás értéke}}$

8. $\log_2 x$ a $C(x)$ alsó burkolója



alsó burkoló



alsó burkoló

2009. 04. 22, Infóház 4.,

10. előadás

- Ullmann könyv utolsó két fejezete XML

XSLT - Extensible Stylesheet Language for Transformations

// Dokumentumok átszerkesztését biztosítja

XSLT-stíluslap

Stíluslapokat adunk meg, ezek is XML dokumentumok formátumúak

xsl: jelölő - nevére hivatkozik

```
<?xml ... >
```

```
<xsl:stylesheet xmlns:xsl="http://www.w3.org/..."
```

```
  : (stylesheet hivatkozása)
```

```
</xsl:stylesheet>
```

Stíluslap - sablonokból áll (template)

sablon jelölője: xsl:template

speciális attribútumai:

- match (milyen elemekre kell majd alkalmazni)

```
<xsl:template match="XPath kifejezés" >
```

```
  : (sablonon belül, amit először alkalmazni)
```

(ez lehet az a template alkalmazás helye)

```
</xsl:template>
```

XPath itt lehet - abszolút /

- relatív: egy elem feloldozásán belül használjuk

Templaten belül egy speciális jelölőjezési lehetőség, az elem gyerekeire kérni feloldozást:

```
<xsl:apply-templates select="kifejezés"/>
```

- olyan template keres, aminek a match felté-

tele illeszkedik a gyerekekre
- relatív útvonallal lehet megadni a matchet

10.

Ha illeszhető a select feltételkifejezésre,
akkor hajtja végre a template-t.

- Rekurzív hívást is lehet végezni...

XML-ből értékek elérése:

`<xsl:value-of select="kifejezés" />`
kifejezés - XPath, eredménye szöveg-nek tekintendő

Iteráció:

`<xsl:for-each select="kifejezés">`
⋮ (- amit végre kell hajtani)
`</xsl:for-each>`

Feltételes parancsok:

`<xsl:if test="logikai kif.">`
[törzs - mit akarunk csinálni] `<then>`
`</xsl:if>`
`<xsl:if test="negált log. kif.">` `<else>`
`</xsl:if>`

// Van `choose` is (otherwise ággal)

output (html, szöveg, XML)

↑ ez alá a node alá kerül az eredmény

- Bejárás az XML dűcsin is a legjobban
illeszkedő template-t választja.

(Vissza a bonyolultság elméletéhez)

Példa:

$$\{x \mid l(x) = n \wedge C(x) \leq n - \Delta\}$$

(legalább Δ -val jobban tömöríthetőek, mint amilyen hosszúak)

Mekkora lehet ennek az elemszáma?

A lehetséges kódok, amit f_0 használhat hossz szerint számozzuk (le: = (+)) = (-)

$$\begin{aligned} & 0 - 1 \\ + & 1 - 2 \\ + & 2 - 2^2 \\ + & \vdots \\ + & 2^{n-\Delta} \end{aligned}$$

$C(x) \leq n - \Delta$, azt jelenti tömöríthető $n - \Delta$ hosszú kóddal

$$\underline{2^{n-\Delta+1} - 1}$$

max elemszám
 Hogyan aránylik 2^n -hez: $2^{-\Delta+1}$

"A dolgok többsége olyan bonyolult amilyen"

" $\Delta = 2$: a halmas ^{max} fele tömöríthető $n-2$ hosszú kóddal"

Változik az x , hogyan változik a bony.

$$\begin{aligned} y &= x+h \\ C(x) & \\ C(y) & \text{ -?} \end{aligned}$$

$$\begin{aligned} g(x, h) &= y \\ C(y) &\leq C(x) + C(h) + 2 \cdot \log_2 C(h) + k_g \end{aligned}$$

u, rendezett párt kell kóddolni

Biz: Invariancia tételre alapul, kell adnunk g -ra egy ügyes kóddolást

$$l(x) = k, \quad C(h) = m$$

$$f_0(p) = x, \quad f_0(q) = h$$

$$l(p) = k, \quad l(q) = m$$

p és q ismeretében x ki kell számolnunk

$y = t!$ De nekünk nem a p -ról, hanem egy h -ról kell... $x = (x) \mid (x)$

→ standard invariáns h -ról

$$r = \{q\}^p$$

$$f(r) = g(\underbrace{f_0(\pi_{12}(r))}_p, \underbrace{f_0(\pi_{11}(r))}_q) = y$$

$$C_{\downarrow}(y) = l(r) = l(p) + l(q) + 2 \cdot l(l(q)) = k + m + 2 \cdot \log_2 m = C(x) + C(h) + 2 \log_2 C(h)$$

"Kicsit változik x , kicsit változik a bonyolultság"

Feltételes Kolmogorov- bonyolultság:

(relatív)

Mennyit segít y kiszámításához az x ismerete?

Ha ismerjük egy $f: \mathbb{N}^2 \rightarrow \mathbb{N}$ lvt: $f(x, p) = y$, akkor elég a p ismerete.

Def: Az $f: \mathbb{N}^2 \rightarrow \mathbb{N}$ felhasználásával az y feltételes bonyolultsága x ismeretében:

$$C(y|x) = \min_{f(x,p)=y} l(p), \text{ illetve } \infty, \text{ ha nincs ilyen } p$$

Kérdés: van-e ebben az esetben is rögzíthető f ?

Alaptétel-invariencia tétel: létezik optimális $f_0: \mathbb{N}^2 \rightarrow \mathbb{N}$ parciális rekurzió f_0 , hogy bármely $f: \mathbb{N}^2 \rightarrow \mathbb{N}$ lvt-hoz megadható egy k_f konstans, hogy

$$C_{f_0}(y|x) \leq C_f(y|x) + k_f \quad \forall x, y \in \mathbb{N}$$

(k_f nem függ x -től, csak f -től)

④ / 10

Biz: Kétváltozós ^{f_0 -ek} univerzális parc. reb. f_0 -t
bell. hasonlónak: $U(n, x, p)$

(ca. a konstrukció mint a lényesebbes esetben)

Invariancia-tétel: $f_0^{(2)}, g_0^{(2)}$ optimális:

$$|C_{f_0}(y|x) - C_{g_0}(y|x)| \leq R_{f_0, g_0}$$

- Konstans erejű megfigyelések.

Köszönetül $f_0^{(2)}: \mathbb{N}^2 \rightarrow \mathbb{N}$ optimális f_0 -t.

Def: $C(y|x) = C_{f_0}(y|x)$

2009. 04. 29. A, 11. előadás, InfoCéz 4.

Feltételes Kolmogorov bonyolultság:

$$C_{f_0}^{(2)}(y|x) = C(y|x)$$

Alaptulajdonságai:

- invariancia-tétel alapján: $f: \mathbb{N}^2 \rightarrow \mathbb{N}$

$$C(y|x) \leq C_f(y|x) + k_f$$

\uparrow x, y -től független konstans

- $C(y|x) = k$ jelentése: $\exists p: f_0(p, x) = y$ és $l(p) = k$

- $C(y) \leq C(y|x) + C(x) + 2 \cdot \log_2 C(x) + k$

(hasonlóan mint: $g = g(h)$ inv. tétellel $C(y)$ -ra)

Mire használható? - Halmas demének bonyolultsága

- Halmas szerinti feltételes bonyolultság

$$\mathcal{A} = \{a_1, \dots, a_N\}$$

$$C(a|\mathcal{A}) = ? \quad a \in \mathcal{A}$$

$$\langle \mathcal{A} \rangle = a_1 a_2 \dots a_{N-1} a_N^*$$

$\leftarrow \mathcal{A}$ kódolása önkéntes kódokkal

$$C(a|\mathcal{A}) = C(a|\langle \mathcal{A} \rangle)$$

Adott \mathcal{A} , $|\mathcal{A}| = N$ Mennyire tömören ábrázolható?

$\{x \mid x \in \mathcal{A}, C(x|\mathcal{A}) < \log_2 N - \Delta\}$ -? Milyen elemszám

kényül adott bonyolultsági elem lehet?

Lehetséges p kódok: $f_0(p, \langle \mathcal{A} \rangle) = x \in \mathcal{A}$

$$0 \quad 1 \quad 2 \quad \dots \quad \log_2 N - \Delta - 1 \quad \leftarrow \text{kódkörök}$$
$$1 + 2 + 2^2 + \dots + 2^{\log_2 N - \Delta - 1} = 2^{\log_2 N - \Delta} - 1$$

$$\frac{|\{x \mid x \in \mathcal{A}, C(x|\mathcal{A}) < \log_2 N - \Delta\}|}{|\{x \mid x \in \mathcal{A}\}|} = \frac{2^{\log_2 N - \Delta}}{2^{\log_2 N}} \leq 2^{-\Delta}$$

A halmas is méretében sem lehet tömöröbben ábrázolható az elemek az egyenletes kódolásnál.
a legtöbb elemet

Keresem a lehető legkisebb halmaszt, ami befeleli a jelenségem és egyenletesen kódolom.

Példa: Személyek adatainak ^{aira} relációs AB: 10.2026

R (csnév, unév)

csnév: 20 byte = 160 bit

unév: 20 byte = 160 bit ($\times 16$) = ($\times 16$) 320

A: az összes lehetséges relációelőfordulás

- lehetséges különálló sorok száma: $2^{160} \cdot 2^{160} = 2^{320}$

- lehetséges előfordulások száma: 2^{320} (részhalmoz) $\times 2^{320}$

$d = (n) | \mathcal{A} | = 2^{320}$

$\{I \mid I \in \mathcal{A}\}$, $C(I | \mathcal{A}) = \log_2 2^{320} = 2^{320}$

Sorok számira van egy korlát a való életben:
< 10 milliárd

10 milliárd 320 bit elég kódolnia a sorokat

Példa: csnév: elő 256 leggyakoribb

unév: -11-

Magyar statisztika szerint a populáció 80%-nak ez a csnév illeti 80%-nak ez az utónév.

4. frekvencia:

1. Minolbet név gyakori: 26% - 2 byte

2. Csnév gyakori: 16% 22 byte

3. Unév gyakori: 16% 22 byte

4. egyik sem: 4% 4 byte

(Jól kell homogén részhalmozokra szétválasztani)

Paraméteres halmazsereg:

(pl. adatbázis (reláció) semmel rögzítet egy paramétert és az előfordulás a halmazsereg)

$a \rightarrow B_a$ véges halmaz

$h(a) = B_a$ - h egy totális fun (rekurzív)

} meg tudom mondani mi van benne és mi nincs

legyen $A \subseteq N \times N$ rekurzívan felsorolható relációhalmaz

$B_a = \{x \mid (x, a) \in A\}$

megszorítás: minden a -ra $|B_a| = m_a$ véges

m_a -t nem tudom megmondani, mert nem

tudom megmondani mi nincs benne

Tétel: létezik k_a , hogy $C(x|B_a) = C(x|a) \leq \log_2 m_a + k_a$

$\forall x \in B_a$ (és $\forall a$)

A halmaz elemeinek a tömörítetősége mindig eléri a \log_2 (elemszám) -ot.

(Alulról + felülről is becsülni tudjuk)

Biz:

B_a halmaz elemei:

A-hoz van felsoroló függvény

$f(a, p)$ konst. nyaljuka:

Futtatjuk A felsorolását, keressük az első $x_i \in B_a$ -t

$\Rightarrow f(a, 0) = x_1$

Tovább: $x_2 \in B_a$ -t megtaláltuk: $f(a, 1) = x_2$

:

$x_{m_a} \in B_a$ elem: $f(a, m_a) = x_{m_a}$

és után az $f(a, p)$ nincs értelmezve

$x \in B_a$: $C_f(x|a)$ - kód egy $(0, \dots, m_a)$ közötti érték:

hossza $\log_2 m_a$ -val felülől becsülhető!

Erdőes halmassereg:

$$D_k = \{x \mid C(x) \leq k\}$$

k paraméterben ez egy halmassereg...

Meg tudjuk-e konstruálni a felsőbőítőt?

f_0 optimális f_0 : hely és lépésszám bejárása

$f_0 : 0$

$f_0 : 1$

x^n	0	1	2	3
1	1	1	2	3
2	1	2	3	
3	2	3		
	3			

ha nem ért véget

veszem a leív. állít.

ha $x \in D_k$: ha valamelyik k hosszú értékre a pont $f_0(p) = x$ eredményt.

$D_k^+ = \{x \mid C(x) = k\}$ nem kiszámítható halmaz
(wi. akkor a $C(x)$ is kiszámítható lenne, de általában nem lehet bebizonyítani, h. a halmaz bony. ga valaminek k)

$D_k^- = \{x \mid C(x) < k\}$ ez sem kiszámítható

Nem tudom megmondani nincs-e jobb mód nála...

Kolmogorov-entropia — Schannon-entropia viszonya

Nagyon hosszú jelenségekre

Shannon-entropia: N hosszú üzenetek halmaza

(stationarius, ergodikus forrás)

két halmazba sorolható:

-lényeges, tipikus halmaz: ebbe $1-\epsilon$ vss. gel esik

$$(1-\epsilon) 2^{N(H-\delta)} < m_N < 2^{N(H+\delta)}$$

\bar{x} tipikus: statisztika jellemzi B_N halmazba

esik (ez lesz a par. halmassereg)

$$C(\bar{x} | B_N) = C(\bar{x} | N) \leq \log_2 m_N + k_{\text{forrás}}$$

N hosszú üzenetek kódját összevesszük N -nel:

Shannon-entropia H

$$H \sim \frac{K_{\text{olm}}}{N \cdot (H - \delta)} < \frac{C(\bar{x} | B_N)}{N} < \frac{N \cdot (H + \delta)}{N} \sim H$$

Ha nagyon nagy jelenségeket vesszük, akkor a tip. halmazzal Kolm. Bonyolultsága visszaverődik a Sch. entropiát.

Hügyötte: egyenletessel jellemezhetőek

-o-

Összes közül melyek a tipikusak és melyek nem?

-o-

0-1 sorozatok véletlenszerűségét vizsgáljuk.

Lehet-e pénzfeldobás eredménye?

1000 pénzfeldobás: 0010100... 1... 101 -véletlenszerű
00010100...
00000100...

egyeseket ~~na~~ egyesével σ -ra cseréltem:

00... 0 000 -nem véletlen

Van-e határvonal a véletlen és nem véletlen között?

Nem...

Szabályszerűség

Ami igazán véletlen az nem tömöríthető...
(Közösséget bonyolultsága a hosszal nő!)

Mennyire tipikus?

Véletlenség defektusa jellemezhető: Δ

$$\exists x \mid x \in A, C(x | A) < \log_2 |A| - \Delta$$

~~XXXXXXXXXX~~

↑ véletlen defektusa legalább Δ

Példa: Fekete - fehér szinezések

(homogén ill. "máskor" színűség) $G(\mathbb{R}^n) = C(\mathbb{R}^n)$

fele fekete - fele fehér

Egymásra helyezve a fölötti érdekes jelenség...

→ színgörbék ill. színes "transzformáció" -
sík...

2009. 05. 06., Infóképzés 4. 12. előadás

(Az információ nem növekedési törvénye!)

Adatbázis - x (x az adatbázis tartalma)

Kérdés: y

Válasz: $a = f(y, x)$

Mennyi az a információ tartalma?

$$C(a|x) \leq C(y) + h$$

ui:

$$C(a|x) = C_f(a|x) + K_f \leq C(y) + K_2$$

Shannon-entropia esetén:

$$- \sum_{i=1}^n p_i \log_2 p_i$$

p_i valószínűség - $-\log_2 p_i$ (mennyiséget rendeljük)

x - $C(x)$ (kolmogorov esetben ez rendeljük)
(kódl hossz)

$$p_i = 2^{-l_i}$$

Egész számokon adjunk meg egy v.s.z. eloszlást:

$$P(i) = 2^{-2^i}$$

$$\sum_{i=0}^{\infty} P(i) = \infty \quad (\text{tehát nem jó eloszlás})$$

$$\text{Ha } x\text{-et rögzíttem: } \sum_{f(p)=x} 2^{-l(p)} = \infty$$

összes lehetőség es kódok összege az összeg

Miért? Mert amit az f felhasznál, nem lehet prefixmentes kódrendszer.

f_0 - univerzális T-gép segítségével

$$U(i, p) = T_i(p)$$

Pényfeldobással állítsuk elő a T-gép bemenetét,
és legyen ez a w_2 .

Nem jó, mert ...

Nem egyértelmű, hogy mire vezetett a pénzfeldobás

Prefix Kolmogorov-bonyolultság:

- Csak prefixmentes f_0 -eket használunk:

$f: \mathbb{N} \rightarrow \mathbb{N}$
 $\Omega \rightarrow \Omega$ } ha az f értelmezve van p és q -n
akkor egyik sem prefixe a másiknak

(Kraft-egyenlőtlenség: $\{p_i, \dots\}$ rendszer prefixmentes,
akkor $\sum_{i=0}^{\infty} 2^{-l(p_i)} \leq 1$) (prefixmentes kódolás)

Ha f prefixmentes, a $C_f(x)$ értelmezhető

Létezik optimális $g_0: \Omega \rightarrow \Omega$ prefixmentes f_0 , amire

$\forall f$ prefixmentes f_0 -hez létezik k_f konstans:

$$C_{g_0}(x) \leq C_f(x) + k_f$$

(vii. létezik univerzális felsoroló f_0)

Vetzőleges $h: \Omega \rightarrow \Omega$ parc. rek. f_0 ,
 $\exists h_p$ prefixmentes f_0

Idő tér konstansnak tekintve feltatjuk a hirtelkeles
(átlósan) és ha találunk ^{értelmezést} akkor az tartozik
a kódhoz, ha a követő, amit találunk
prefixe az előzőnek, akkor hi hagyjuk.

T_p - prefixmentes T-gép - véletlen sorozatok
 k bit hosszú helyen, 2^{-k} valószínűséggel áll meg

$\sum 2^{-k} \leq 1$ a Kraft-egyenletlenség miatt

$K(x) = C_{g_0}(x)$ prefixmentes Kolmogorov-bonyolultság!

$\pi(x) = 2^{-K(x)}$ (valószínűség)

$$\sum_{i=0}^{\infty} 2^{-K(i)} < 1$$

$$\pi^*(x) := \sum_{g_0(p)=x} 2^{-l(p)} < 1$$

$x = g_0$ szerinti összes lehetséges kódja

g_0 -t képzeljük el T-gépnek, akkor annak a vsz-ge, hogy x -et ad éppen $\pi(x)$ mennyi

Azaz: a megállás valószínűsége x eredménnyel

(Azok fordulnak elő nagyobb vsz-ge, amit egyszerűbben tudunk értelmezni)

$\pi(x)$, $\pi^*(x)$ - az egész számok fölött értelmezett eloszlások

— 0 —

A legnagyobb bizonytalanság: egyenletes eloszlás. Tudunk-e az egész számok halmán egyenletes eloszlást megadni? Nem!

Mi a lehető leg-egyenletesebb eloszlás?

(Ami a lehető leglassabban nő)

3) Felig kiszámítható (semi) eloszlások között
 12. a lehető legegyenletesebb eloszlás a $\pi(x)$ ($\pi^*(x)$)

Félig hisz: nem tudom megmondani, de tetőlegesen jól közelítem a pontos értéket.

Ha $p(x)$ félig hisz. eloszlás, akkor $\exists c$ konstans, hogy $p(x) < c \cdot \pi(x)$
(Nem tud a $\pi(x)$ -nél gyorsabb lenni)

- Magoráló eloszlásnak nevezik ezt a tulajdonságot

(Algoritmikus tanulás háttérében ez van)

Probléma: nem tudjuk kiszámolni a $\pi(x)$ -et, csak közelítő értékekkel dolgozhatunk.

$$C(x) = k \Rightarrow \exists p : f_0(p) = x \text{ és } l(p) = k$$

$$K(x) \quad g(p') = f_0(p) = x \\ l(p') = 2 \cdot \log_2 k + k$$

$$\Rightarrow K(x) \leq C(k) + 2 \cdot \log_2 C(k) + c$$

-o-

Feltételes prefix-bonyolultság:

$f(x, p)$ p -ben prefixmentes

$$C_f(y|x)$$

$$g_0^{(2)}: N^2 \rightarrow N \text{ optimális}$$

$$C_{g_0}(y|x) \leq C_f(y|x) + k_f$$

$$\Rightarrow \underline{K(y|x)}$$

Van értelme: $K(x|A)$, $x \in A$ - nak

Shannon-entropia esetében: $H(S, \pi) = H(S|\pi) + H(\pi)$

vez. Kolmogorov-bonyolultságra:

$$(x, y) \text{ rendezett párja: } (x, y) = x'y$$

$$C(x, y) = C(x|y)$$

$$C(x, y) \stackrel{?}{\leq} \underbrace{C(y|x)} + C(x)$$

til nagy lehet

Es az additív tulajdonság nem teljesül, de k-ra

$$K(x, y) \leq K(x) + K(y|x) + 2 \log_2 K(x)$$

Bmeg: Ha tudnánk az x bonyolultságát, akkor konstans

oregig: $K(x, y) \stackrel{+}{=} K(x) + K(\#|x, K(x))$

↑ additív

van x -nek tanúja:

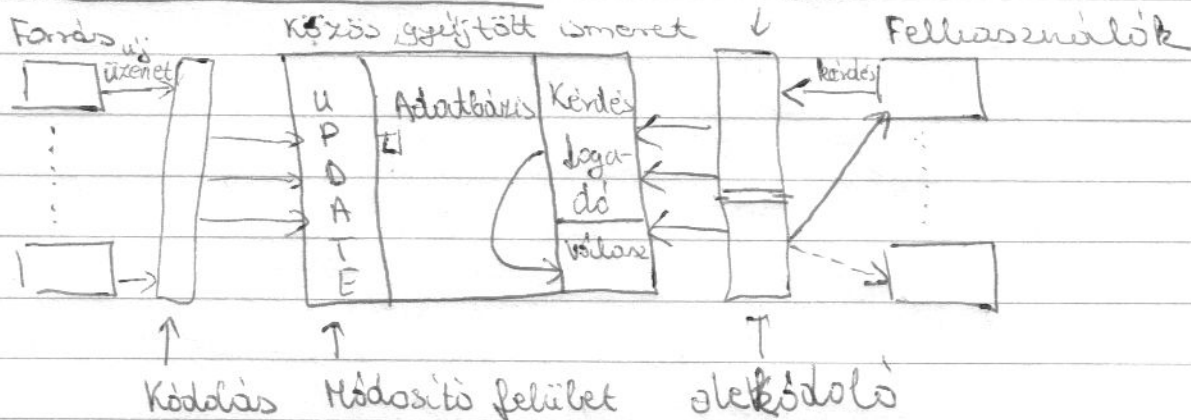
$$x^* : g_0(x^*) = x$$

$$l(x^*) = K(x)$$

Elvi modell az adatbázis kezelő rendszeréhez

Diagram: (Sch. modellhez hasonlóan)

Információ-rendszer



Megszorítások:

- mi lehet az adatbázis tartalma?

(valamilyen módon formalizálni kell)

- milyen új üzenetek lehetségesek, és hogyan épülnek be a régi ismeretekbe?

- milyen kérdések szolgalmozhatnak meg?
válaszok?

5/12. (tömörítettség, hatékony lekérdezés - ellentmond a módosíthatóság -nak)

13. előadás, 2009. 05. 13., Infókéz4.

ABKR formális modellje:

1. Logikai - fogalmi modell
Valóság modellezése - kódrendszer

Céja:

közötti kapcsolat

Kódrendszer: kiszámítási, bonyolultságelmélettel
elemzéssel

Valóság mod. : Absztrakt adat típusok
adatmodellek

előfordulásait \rightarrow kódrendszer

Kódrendszer: bináris szavak rendszere

① Mi lehet az adatbázis lehetőségei és állapota
a kód. alapján?
 \downarrow véges bináris szavak

$$\Sigma \subseteq \Omega$$

Σ - rekurzív v. rekurzívan
felismerhető

Sémák: leírása adat, tehát:

$$S \subseteq \Omega$$

$s \in S$, Σ_s : séma szerint lehetséges adatok

$$\langle s, y \rangle = x = s' y$$

(s séma belüli y előford.

$s_1 \dots s_n$ sémák

$y_1 \dots y_n$ állapotok

$\langle s_i, y_i \rangle$ felsorolás kód

② Hogyan lehet módosítani az AB tartalmát?

Y - módosító adatok halmaza

a) $M: Y \times \Sigma \rightarrow \Sigma$ - módosító fű, kiszámítható fű.

b) M : módosító fűek halmaza

$m \in M$, y -célszerű módosító adat

($y \in Y_m$)

$$m(y, x) = x'$$

- Zártsági elvárás: $M(Y \times \Sigma) \subseteq \Sigma$

(szemantika adatközlésben elől el) generikus dyan legyen a módosítás, h az eredmény értelmes legyen

- ekvivalencia osztályok bevezetése $\Sigma - n$

-üres módosítás: Λ , $x \equiv x'$ ha

$$M(\Lambda, x) = M(\Lambda, x')$$

Sémánkénti módosítások: $s \in S$, Σ_s , M_s - módosító fű halmaza

Adat - mennyiségek:

{ - hossz (al. mértékegység) - redundancia
- $K(x)$, $C(x)$

→ ezek alapján elemezhető a módosítások komplexitása

Séma megválasztása mikor jó? (kol. bony. alapján)

$x \in \Sigma_s$, $K(x|s) \sim \log_2 N$

↑ véges

Σ_s -en egyszerű kódot használunk

$$|\Sigma_s| = N$$

Σ_s -en homogén leaslat lesz

Fixikai modell

- beábrázolódik egy bináris kód címesletű, véges kapacitású rekeszekbe

Ha van egy fizikai kód: (c_1, l_1, y_1)
 (c_2, l_2, y_2)
 \vdots
 (c_k, l_k, y_k)

↓ hossz ↓ csatlak

Ψ - fizikai kód

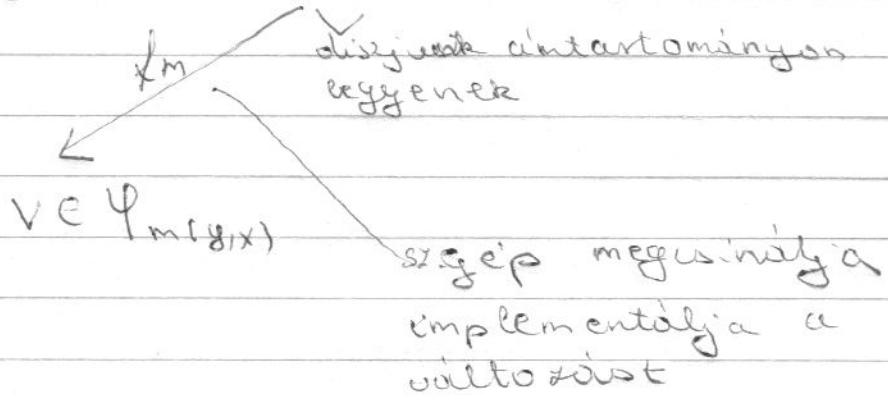
$\Psi, z \in \Phi$

$\Psi: \Sigma \rightarrow \Phi, x \in \Sigma, \Psi_x \in \Phi$

Kódoltások rendszere: $y \in \Upsilon, \Psi_y$

$w \in \Psi_y, z \in \Psi_x$
 ↓ összeküszes

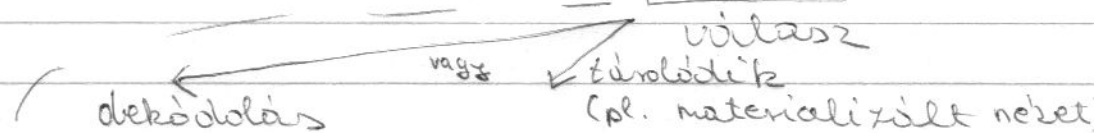
$$\Psi_{x,y} = \Psi_y \circ \Psi_x = wz$$



Eljárásainknak is lesz egy kódolt formája, ugyanabban a címtartományban mozog.

Lekeódolás: $q - k \in \Psi_q$

$$\Psi_q \Psi_x \Psi_{ABKR} \rightarrow \Psi_{A(q,x)} (\Psi_x \Psi_{ABKR})$$



$\Sigma_q \rightarrow$ felhasználói üzenet