



Információs rendszerek elméleti alapjai

Információelmélet



Az entrópia tulajdonságai

Jensen egyenlőtlenség

$h(x)$ konvex függvény,

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

$h(x)$ szigorúan konvex, akkor az egyenlő

feltétele: $P(\xi = E(\xi)) = 1$ – $\xi = 1$ val. gel konstans

konvexitás

$$x, y \in [a, b], 0 < \lambda < 1$$

$$h(\lambda x + (1 - \lambda)y) \leq \lambda h(x) + (1 - \lambda)h(y)$$

szigorú konvexitás

$$E(\xi) \text{ pontban, } \lambda x + (1 - \lambda)y = E(\xi), \forall x, y$$

h logaritmus, E az entrópia

Forrás jellemzése



A forrás közelítései:

Megfigyeljük a forrást M hosszú üzenetet véve. Ebből közelítjük a forrás valószínűségi eloszlását, statisztikákból. (Stacionáris forrásokot vizsgálunk. Azaz időtől ne függjön az eloszlás

Az üzenet (diszkrét) lehetséges szimbólumai, x_1, x_2, \dots, x_n .



Forrás jellemzése

0-d rendű közelítés: csak n értékét használjuk fel.

A legbizonytalanabb az egyenletes eloszlás;

Egy szimbólum választás: ξ_j valószínűségi változó

$$P(\xi_i = x_j) = \frac{1}{n}, \quad j = 1, \dots, n$$

N szimbólum egymás után (n^N lehetőség)

$$P(\xi_1 = x_{i_1}, \xi_2 = x_{i_2}, \dots, \xi_n = x_{i_N}) = n^{-N}, \quad j = 1, \dots, n$$

Független, egyenletes, azonos eloszlású változó.



Forrás jellemzése

Entrópia : $H(\xi_i) = \log_2 n$

$H(\xi_1, \xi_2, \dots, \xi_N) = NH(\xi) = N \log_2 n, \quad j = 1$

$$H(\xi) = \log_2 n = \lim_{N \rightarrow +\infty} \left(\frac{1}{N} H(\xi_1, \xi_2, \dots, \xi_N) \right)$$

1 szimbólumra jutó
entrópia/bizonytalanság



Forrás jellemzése

Első-rendű közelítés

M hosszú üzenetet megfigyelünk

x_i előfordulásainak száma: $m_i \forall i=1, \dots, n$

Rögzítjük ξ_i eloszlásának közelítését,

$$P(\xi_i=x_j)=m_j/M=p_j \forall j=1, \dots, n$$

A vetületi eloszlások rögzítése mellett a legbizonytalanabb együttes eloszlás a független eloszlás

Forrás jellemzése

Első-rendű közelítés

Független, azonos eloszlású változók.

$$P(\xi_1 = x_{i_1}, \xi_2 = x_{i_2}, \dots, \xi_n = x_{i_N}) = \prod_{j=1}^N P(\xi_j = x_{i_j}) =$$

$$\prod_{i=1}^n p_i^{\mu_i}, \quad \mu_i = \text{az } x_i \text{ száma } \{x_{i_1}, x_{i_2}, \dots, x_{i_N}\} \text{ - ben}$$

$$H(\xi_i) = - \sum_{j=1}^m p_j \log_2 p_j \dots \text{Entrópia}$$

Forrás jellemzése



Együttes entrópia (függetlenség miatt additív)

Egy szimbólumra jutó entrópia:

$$H(\xi_1, \dots, \xi_N) = -N \sum_{j=1}^n p_j \log_2 p_j$$

$$H(\xi) = \lim_{N \rightarrow \infty} \frac{1}{N} H(\xi_1, \dots, \xi_N) = -\sum_{j=1}^n p_j \log_2 p_j$$

Forrás jellemzése

Másod-rendű közelítés:

Megfigyeljük az egymás utáni párok gyakoriságát,

m_{ij} : az $x_i x_j$ előfordulásainak száma

m_i : az x_i előfordulásainak száma

A $p_{i|j} = m_{ij} / m_j$ hányados megfelel feltételes valószínűségnek:

$$P(\xi_t = x_j | \xi_{t-1} = x_i) = \frac{P(\xi_t = x_j, \xi_{t-1} = x_i)}{P(\xi_{t-1} = x_i)}$$

Ezzel rögzítjük az átmeneti valószínűségeket. (átmenet valószínűség mátrix, sorok összege 1.

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

A 2' Tétel szerint a legbizonytalanabb eloszlás, ami ennek eleget tesz az egy lépéses homogén Markov lánc

Forrás jellemzése

$$\forall t - re \quad P(\xi_t = x_i | \xi_1 = 1, \dots, \xi_{t-1} = x_j) = P(\xi_t = x_i | \xi_{t-1} = x_j) = p_{i|j}$$

Az együttes eloszlás :

$$P(\xi_1 = x_{i_1}, \dots, \xi_N = x_{i_N}) =$$

$$P(\xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) \cdot P(\xi_N = x_{i_N} | \xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) =$$

$$P(\xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) P(\xi_N = x_{i_N} | \xi_{N-1} = x_{i_{N-1}}) =$$

$$P(\xi_1 = x_{i_1}) P(\xi_2 = x_{i_2} | \xi_1 = x_{i_1}) \dots P(\xi_N = x_{i_N} | \xi_{N-1} = x_{i_{N-1}}) = P(\xi_1 = x_{i_1}) \prod_{i,j} p_{i|j}^{\mu_{i,j}}$$

ahol $\mu_{i,j}$ az x_i, x_j átmenetek száma, $P(\xi_t = x_i | \xi_{t-1} = x_j)$

Forrás jellemzése

Markov lánc
entrópiája
határérték lesz
Szükséges
fogalom az
ergodikus (határ)
eloszlás.
(Ergodikus
Markov-lánc:
minden
állapotból
minden állapot
elérhető (véges
lépéssel, pozitív
valószínűséggel)
és nem
periodikus,
Minden érték
visszatérő legyen.

$$\underbrace{H(\xi_1, \dots, \xi_N)}_{\text{Mi a}} = \underbrace{H(\xi_N | \xi_1, \dots, \xi_{N-1})}_{\text{viszony?}} + \underbrace{H(\xi_1, \dots, \xi_{N-1})}_{\text{viszony?}} =$$

$$H(\xi, \eta) = H(\xi) + H(\eta) \wedge H(\xi) = H(\xi | \eta)$$

$$H(\xi_N | \xi_{N-1}) + H(\xi_1, \dots, \xi_{N-1}) =$$

$$H(\xi_1) + H(\xi_2 | \xi_1) + \dots + H(\xi_N | \xi_{N-1})$$

$$H(\xi_t | \xi_{t-1}) = \sum_{i=1}^n P(\xi_{t-1} = x_i) H(\xi_t | \xi_{t-1} = x_i)$$

$$H(\xi_t | \xi_{t-1} = x_i) = \sum_{j=1}^n P(\xi_t = x_j | \xi_{t-1} = x_i) \cdot \log_2 P(\xi_t = x_j | \xi_{t-1} = x_i) =$$

$$-\sum_{j=1}^n p_{j|i} \log_2 p_{j|i} = H(\xi | x_i)$$

$$H(\xi_t | \xi_{t-1}) = \sum_{i=1}^n P(\xi_{t-1} = x_i) \cdot H(\xi | x_i)$$

Forrás jellemzése

A Markov lánc tulajdonság miatt :

$$H(\xi_N | \xi_1, \dots, \xi_{N-1}) = \sum_{i_1, i_2, \dots, i_{N-1}} P(\xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}).$$

$$H(\xi_N | \xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) = \sum_{i_1, i_2, \dots, i_{N-1}} P(\xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}).$$

$$\left(- \sum_{i=1}^n P(\xi_N = x_i | \xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) \cdot \log_2 P(\xi_N = x_i | \xi_1 = x_{i_1}, \dots, \xi_{N-1} = x_{i_{N-1}}) \right)$$

Markov lánc és ξ_{N-1} eloszlása az együttes eloszlásból:

$$\sum_{j=1}^n P(\xi_{N-1} = x_j) \left(- \sum_{i=1}^n p_{i|j} \log_2 p_{i|j} \right) =$$

$$\sum_{j=1}^n P(\xi_{N-1} = x_j) H(\xi | x_j)$$

Forrás jellemzése

Visszahelyettesítve $H(\xi_1, \dots, \xi_N)$ -be,

és alkalmasan visszafelé, $N = 2$ -ig, kapjuk :

$$H(\xi_1, \dots, \xi_{N-1}, \xi_N) = \sum_{t=2}^N \sum_{j=1}^n P(\xi_{t-1} = x_j) H(\xi | x_j) =$$

$$\sum_{j=1}^n H(\xi | x_j) \left(\sum_{t=1}^N P(\xi_{t-1} = x_j) \right)$$

Kihasználva, hogy ergodikus Markov lánc

$$\exists \lim_{t \rightarrow \infty} P(\xi_t = x_i) = \pi_i \text{ határeloszlás, } \forall i = 1, \dots, n$$

Forrás jellemzése

π_i az átmenet valsz. mátrix sajátvektora ($\lambda = 1$ sajátértékkel):

$$P(\xi_t = x_i) = \sum_{j=1}^n P(\xi_{t-1} = x_j) p_{i|j}$$

Mindkét oldal határértékét véve :

$$\pi_i = \sum_{j=1}^n \pi_j p_{i|j} \quad \forall i$$

Ez egy lineáris egyenletrendszer

π_1, \dots, π_n , az átmenet valsz. mátrix sajátvektorai,

$P(\xi_t = x_i) = \pi_i \dots$ ettől kezdve π_i nem függ az időtől
azaz stacionáris.

Forrás jellemzése

visszatérve az entrópiára

$$H(\xi_1, \dots, \xi_N) = H(\xi_1) + \sum_{j=1}^n H(\xi | x_j) \left(\sum_{t=1}^N P(\xi_{t-1} = x_j) \right)$$

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^n P(\xi_{t-1} = x_j) = \pi_j$$

ezért

$$H(\xi) = \lim_{N \rightarrow \infty} \frac{1}{N} H(\xi_1, \dots, \xi_N) = \sum_{j=1}^n \pi_j H(\xi | x_j)$$

Forrás jellemzése



Lehet 3-ad, 4-ed stb rendű közelítés, ami két, három stb. lépéses Markov lánc.

Nyelvek esetén szokás/lehetőség, szavakat tekinteni a forrás választásának. Lehet a szavak felett 0-ad, 1-ső, 2-od, stb. közelítés.

(2-od rendű alapján szimulálva értelmes mondattöredékek nagy gyakorisággal keletkeznek.)

Matematikai kitérő – A valószínűségszámítás alapfogalmairól

