

Elemente de teoria informației¹⁹

29.

(Exercițiu cu caracter teoretic:
proprietăți dezirabile ale entropiei)

■ □ • CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2.2

Prin definiție, entropia (în sens Shannon) a unei variabile aleatoare discrete X ale cărei valori sunt luate cu probabilitățile p_1, p_2, \dots, p_n este $H(X) = -\sum_i p_i \log p_i$. Însă legătura dintre această definiție formală și obiectivul avut în vedere — și anume, acela de a exprima gradul de *incertitudine* cu care se produc valorile unei astfel de variabile aleatoare — nu este foarte intuitivă.

Scopul acestui exercițiu este de a arăta că orice funcție $\psi_n(p_1, \dots, p_n)$ care satisface trei proprietăți dezirabile pentru entropie este în mod necesar de forma $-K \sum_i p_i \log p_i$ unde K este o constantă reală pozitivă. Iată care sunt aceste *proprietăți*:

A1. Funcția $\psi_n(p_1, \dots, p_n)$ este continuă în fiecare din argumentele ei și simetrică.

Din punct de vedere formal, în acest caz simetria se traduce prin egalitatea $\psi_n(p_1, \dots, p_i, \dots, p_j, \dots, p_n) = \psi_n(p_1, \dots, p_j, \dots, p_i, \dots, p_n)$ pentru orice $i \neq j$. Informal spus, dacă două dintre valorile care sunt luate de variabila aleatoare X (și anume x_i și x_j) își schimbă între ele probabilitățile (p_i și respectiv p_j), valoarea entropiei lui X nu se schimbă.

A2. Funcția $\psi_n(1/n, \dots, 1/n)$ este monoton crescătoare în raport cu n .

Altfel spus, dacă toate evenimentele sunt echiprobabile, atunci entropia crește odată cu numărul de evenimente posibile.

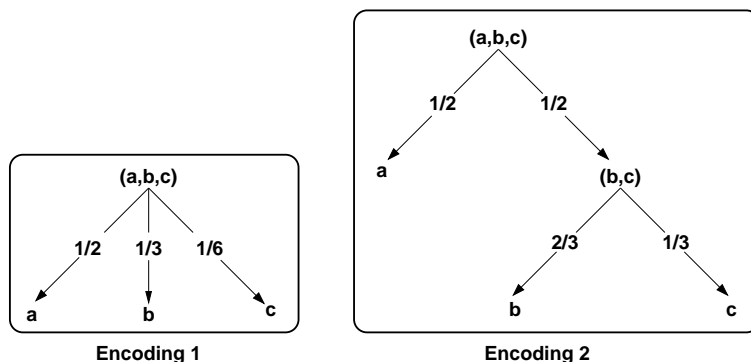
A3. Dacă faptul de a alege între mai multe evenimente posibile poate fi realizat prin mai multe alegeri succesive, atunci $\psi_n(p_1, \dots, p_n)$ trebuie să se poată scrie/calcula ca o sumă ponderată.

De *exemplu*, dacă evenimentele (a, b, c) se produc respectiv cu probabilitățile $(1/2, 1/3, 1/6)$, atunci acest fapt poate fi echivalat cu

- a alege mai întâi cu probabilitate de $1/2$ între a și (b, c) ,
 - urmat de a alege între b și c cu probabilitățile $2/3$ și $1/3$ respectiv.
- (A se vedea imaginile de mai jos, Encoding 1 și Encoding 2.)

Din punct de vedere formal, proprietatea A3 impune ca, pe acest exemplu, $\psi_3(1/2, 1/3, 1/6)$ să fie egal cu $\psi_2(1/2, 1/2) + 1/2 \cdot \psi_2(2/3, 1/3)$.

¹⁹Observație importantă: În toate problemele care urmează, referitor la entropie / teoria informației se va considera în mod implicit că notația 'log' desemnează logaritmul în baza 2. De asemenea, prin convenție, se va considera $p \cdot \log p = 0$ pentru $p = 0$.



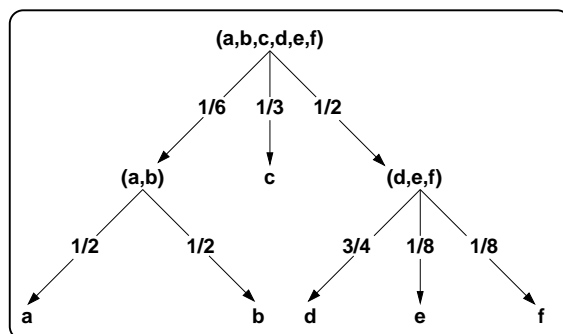
Așadar, în acest exercițiu vi se cere să arătați că dacă o funcție de n variabile $\psi_n(p_1, \dots, p_n)$ satisface proprietățile A1, A2 și A3 de mai sus, atunci există $K \in \mathbb{R}^+$ astfel încât $\psi_n(p_1, \dots, p_n) = -K \sum_i p_i \log p_i$ unde K depinde de ψ .

Indicație:

Veți face rezolvarea acestei probleme în mod gradual, parcurgând următoarele puncte (dintre care primele două puncte au rolul de a vă acomoda cu noțiunile din enunț):

a. Arătați că $H(1/2, 1/3, 1/6) = H(1/2, 1/2) + \frac{1}{2}H(2/3, 1/3)$. Altfel spus, verificați faptul că definiția clasică a entropiei, $H(X) = \sum_i p_i \log 1/p_i$, satisface proprietatea A3 pe *exemplul* care a fost dat mai sus.

b. Calculați entropia în cazul distribuției „codificării” din figura de mai jos, folosind din nou proprietatea A3.



Următoarele întrebări tratează cazul particular $A(n) \stackrel{\text{not.}}{=} \psi(1/n, 1/n, \dots, 1/n)$.

c. Arătați că

$$A(s^m) = m A(s) \text{ pentru orice } s, m \in \mathbb{N}^*. \quad (1)$$

Mai departe, pentru $s, m \in \mathbb{N}^*$ fixați, vom considera $t, n \in \mathbb{N}^*$ astfel încât $s^m \leq t^n \leq s^{m+1}$.

d. Verificați că, prin logaritmare a acestei duble inegalități și apoi prin rearanjare, obținem $\frac{m}{n} \leq \frac{\log t}{\log s} \leq \frac{m}{n} + \frac{1}{n}$ pentru $s \neq 1$, și deci

$$\left| \frac{m}{n} - \frac{\log t}{\log s} \right| \leq \frac{1}{n}. \quad (2)$$

e. Explicați de ce $A(s^m) \leq A(t^n) \leq A(s^{m+1})$.

f. Combinând ultima inegalitate de mai sus cu inegalitatea (1), avem $A(s^m) \leq A(t^n) \leq A(s^{m+1}) \Rightarrow mA(s) \leq nA(t) \leq (m+1)A(s)$. Verificați că

$$\left| \frac{m}{n} - \frac{A(t)}{A(s)} \right| \leq \frac{1}{n} \text{ pentru } s \neq 1. \quad (3)$$

g. Combinând inegalitățile (2) și (3), arătați că

$$\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n} \text{ pentru } s \neq 1 \quad (4)$$

și, în consecință

$$A(t) = K \log t \text{ cu } K > 0 \text{ (din cauza proprietății A2)}. \quad (5)$$

Observație:

Rezultatul de mai sus ($A(t) = K \log t \Leftrightarrow \psi_t(1/t, \dots, 1/t) = Kt \frac{1}{t} \log t$) se generalizează ușor la cazul $\psi(p_1, \dots, p_n)$ cu $p_i \in \mathbb{Q}$ pentru $i = 1, \dots, n$ (cazul $p_i \notin \mathbb{Q}$ nu este tratat aici):

Considerăm o mulțime de N evenimente echiprobabile. Fie $\mathcal{P} = (S_1, S_2, \dots, S_k)$ o partiționare a acestei mulțimi de evenimente. Notăm $p_i = |S_i|/N$.

Propunem următoarea *codificare*: Vom alege mai întâi S_i , una din submulțimile din partiția \mathcal{P} , în funcție de probabilitățile p_1, \dots, p_k . Extragem apoi unul din elementele mulțimii S_i , cu probabilitate uniformă.

Conform egalității (5), avem $A(N) = K \log N$. Folosind proprietatea A3 și *codificarea* în doi pași propusă mai sus, rezultă că

$$A(N) = \psi_k(p_1, \dots, p_k) + \sum_i p_i A(|S_i|).$$

Așadar,

$$K \log N = \psi_k(p_1, \dots, p_k) + K \sum_i p_i \log |S_i|.$$

Prin urmare,

$$\begin{aligned} \psi_k(p_1, \dots, p_k) &= K \left[\log N - \sum_i p_i \log |S_i| \right] \\ &= K \left[\log N \sum_i p_i - \sum_i p_i \log |S_i| \right] = -K \sum_i p_i \log \frac{|S_i|}{N} \\ &= -K \sum_i p_i \log p_i \end{aligned}$$

Răspuns:

a. Facem calculele:

$$\begin{aligned} H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) &= \frac{1}{2} \log 2 + \frac{1}{3} \log 3 + \frac{1}{6} \log 6 = \left(\frac{1}{2} + \frac{1}{6}\right) \log 2 + \left(\frac{1}{3} + \frac{1}{6}\right) \log 3 \\ &= \frac{2}{3} + \frac{1}{2} \log 3 \end{aligned}$$

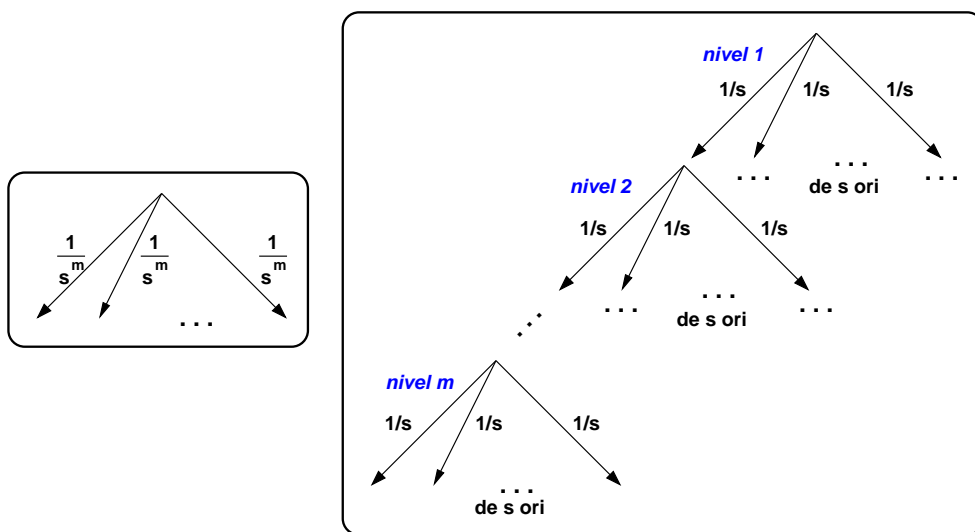
$$\begin{aligned}
H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2}H\left(\frac{2}{3}, \frac{1}{3}\right) &= 1 + \frac{1}{2}\left(\frac{2}{3}\log\frac{3}{2} + \frac{1}{3}\log 3\right) = 1 + \frac{1}{2}\left(\log 3 - \frac{2}{3}\right) \\
&= \frac{2}{3} + \frac{1}{2}\log 3
\end{aligned}$$

și rezultă că $H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) = H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2}H\left(\frac{2}{3}, \frac{1}{3}\right)$.

b. Folosind proprietatea A3, entropia „codificării” din figura dată în enunț este:

$$\begin{aligned}
H\left(\frac{1}{6}, \frac{1}{3}, \frac{1}{2}\right) &+ \frac{1}{6}H\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{1}{2}H\left(\frac{3}{4}, \frac{1}{8}, \frac{1}{8}\right) \\
&= \frac{1}{6}\log 6 + \frac{1}{3}\log 3 + \frac{1}{2}\log 2 + \frac{1}{6} + \frac{1}{2}\left(\frac{3}{4}\log\frac{4}{3} + \frac{2}{8}\log 8\right) \\
&= \frac{1}{6}\log 2 + \frac{1}{6}\log 3 + \frac{1}{3}\log 3 + \frac{1}{2} + \frac{1}{6} + \frac{3}{8}(2 - \log 3) + \frac{3}{8} \\
&= \frac{1}{6} + \frac{1}{2} + \frac{1}{6} + \frac{3}{4} + \frac{3}{8} + \left(\frac{1}{6} + \frac{1}{3} - \frac{3}{8}\right)\log 3 \\
&= \frac{47}{24} + \frac{1}{8}\log 3 = 1.958 + 0.125\log 3 = 2.156
\end{aligned}$$

c. Pentru calculul lui $A(s^m)$ se poate folosi atât o „codificare” imediată cât și una (des)compusă, ca în figura următoare:



Aplicând proprietatea A3 pe „codificarea” din figura de mai sus, partea dreaptă, avem:

$$\begin{aligned}
A(s^m) &= A(s) + s \cdot \frac{1}{s}A(s) + s^2 \cdot \frac{1}{s^2}A(s) + \dots + s^{m-1} \cdot \frac{1}{s^{m-1}}A(s) \\
&= \underbrace{A(s) + A(s) + A(s) + \dots + A(s)}_{\text{de } m \text{ ori}} = mA(s)
\end{aligned}$$

d. Aplicând funcția log fiecărui termen al inegalității $s^m \leq t^n \leq s^{m+1}$ obținem $m \log s \leq n \log t \leq (m+1) \log s$. Apoi, pentru $s \neq 1$, împărțind prin $n \log s$, rezultă:

$$\frac{m}{n} \leq \frac{\log t}{\log s} \leq \frac{m}{n} + \frac{1}{n} \Rightarrow 0 \leq \frac{\log t}{\log s} - \frac{m}{n} \leq \frac{1}{n} \Rightarrow \left| \frac{\log t}{\log s} - \frac{m}{n} \right| \leq \frac{1}{n}$$

e. Datorită proprietății A2 din enunț, inegalitatea $s^m \leq t^n \leq s^{m+1}$ implică

$$\psi_{s^m} \left(\frac{1}{s^m}, \dots, \frac{1}{s^m} \right) \leq \psi_{t^n} \left(\frac{1}{t^n}, \dots, \frac{1}{t^n} \right) \leq \psi_{s^{m+1}} \left(\frac{1}{s^{m+1}}, \dots, \frac{1}{s^{m+1}} \right)$$

ceea ce reprezintă exact $A(s^m) \leq A(t^n) \leq A(s^{m+1})$.

f. Datorită proprietății (1), dubla inegalitate $A(s^m) \leq A(t^n) \leq A(s^{m+1})$ devine $m A(s) \leq n A(t) \leq (m+1) A(s)$. Împărțind această inegalitate prin $n A(s)$, despre care se poate spune că este nenul pentru orice $s \neq 1$, rezultă:

$$\frac{m}{n} \leq \frac{A(t)}{A(s)} \leq \frac{m}{n} + \frac{1}{n} \Rightarrow 0 \leq \frac{A(t)}{A(s)} - \frac{m}{n} \leq \frac{1}{n} \Rightarrow \left| \frac{A(t)}{A(s)} - \frac{m}{n} \right| \leq \frac{1}{n}$$

g. Inegalitățile duble de mai jos rescriu convenabil proprietățile (2) și (3):

$$-\frac{1}{n} \leq \frac{m}{n} - \frac{\log t}{\log s} \leq \frac{1}{n} \quad \text{și} \quad -\frac{1}{n} \leq \frac{A(t)}{A(s)} - \frac{m}{n} \leq \frac{1}{n}$$

Însumându-le membru cu membru, rezultă

$$-\frac{2}{n} \leq \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \leq \frac{2}{n} \Rightarrow \left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n}$$

Din inegalitatea dublă $s^m \leq t^n \leq s^{m+1}$ rezultă că odată cu m crește și n (acesta din urmă depinzând de valorile lui s, t și m).²⁰ Așadar, atunci când m tinde la infinit vom avea și $n \rightarrow \infty$. Dacă trecem la limită inegalitatea $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n}$ pentru $n \rightarrow \infty$, rezultă $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \rightarrow 0$, de unde avem $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| = 0$ și deci $\frac{A(t)}{A(s)} = \frac{\log t}{\log s}$. Rezultă că $A(t) = \frac{A(s)}{\log s} \log t = K \log t$. Evident, constanta $K = \frac{A(s)}{\log s}$ nu depinde de t . Variind valorile lui t , rezultă că $A(t) = K \log t$ pentru orice $t \in \mathbb{N}^*$.

O *observație* finală: din inegalitatea $s^m \leq t^n \leq s^{m+1}$, dacă $s \neq 1$ va rezulta că de fapt și $t \neq 1$. Atunci când pentru A se folosește formula clasică a entropiei, egalitatea $A(t) = K \log t$ se verifică și în cazul $t = 1$, fiindcă $A(1) = \sum_p p \log p = 1 \log 1 = 0$, iar $K \log 1 = 0$.

30.

(Entropie, entropie corelată,
entropie condițională, câștig de informație:
definiții și proprietăți imediate)

□ L. Ciortuz, pornind de la
CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2

Fie X și Y variabile aleatoare discrete. Dăm pe scurt următoarele definiții:

²⁰Mai precis, este util să observăm o succesiune de forma $s^m \leq t^{n_1} \leq s^{m+1} \leq t^{n_2} \leq s^{m+2} \leq \dots$.

• **Entropia variabilei X :**

$$H(X) \stackrel{\text{def.}}{=} -\sum_i P(X = x_i) \log P(X = x_i) \stackrel{\text{not.}}{=} E_X[-\log P(X)].$$

Prin *convenție*, dacă $p(x) = 0$ atunci vom considera $p(x) \log p(x) = 0$.

• **Entropia condițională specifică a variabilei Y în raport cu valoarea x_k a variabilei X :**

$$H(Y | X = x_k) \stackrel{\text{def.}}{=} -\sum_j P(Y = y_j | X = x_k) \log P(Y = y_j | X = x_k) \\ \stackrel{\text{not.}}{=} E_{Y|X=x_k}[-\log P(Y | X = x_k)].$$

• **Entropia condițională medie a variabilei Y în raport cu variabila X :**

$$H(Y | X) \stackrel{\text{def.}}{=} \sum_k P(X = x_k) H(Y | X = x_k) \stackrel{\text{not.}}{=} E_X[H(Y | X)].$$

• **Entropia corelată a variabilelor X și Y :**

$$H(X, Y) \stackrel{\text{def.}}{=} -\sum_{i,j} P(X = x_i, Y = y_j) \log P(X = x_i, Y = y_j) \\ \stackrel{\text{not.}}{=} E_{X,Y}[-\log P(X, Y)].$$

• **Informația mutuală a variabilelor X și Y , numită de asemenea *câștigul de informație* al variabilei X în raport cu variabila Y (sau invers):**

$$MI(X, Y) \stackrel{\text{not.}}{=} IG(X, Y) \stackrel{\text{def.}}{=} H(X) - H(X | Y) = H(Y) - H(Y | X)$$

(Observație: ultima egalitate de mai sus are loc datorită rezultatului de la punctul c de mai jos.)

Arătați că:

a. $H(X) \geq 0$. În particular, $H(X) = 0$ dacă și numai dacă variabila X este constantă.

b. $H(Y | X) = -\sum_i \sum_j P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i)$.

c. $H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$.

Mai general: $H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})$ (regula de înlănțuire).

Răspuns:

a. Este ușor să arătăm că $H(X) = -\sum_i P(X = x_i) \log P(X = x_i) \geq 0$.

Știm că $\log x \leq 0$ pentru $\forall x \leq 1$ și $\log x \geq 0$ pentru $\forall x \geq 1$. De asemenea știm că $P(X = x_i) \in [0, 1]$ (fiind o probabilitate). Așadar,

$$H(X) = -\sum_i P(X = x_i) \log P(X = x_i) = \sum_i \underbrace{P(X = x_i)}_{\geq 0} \underbrace{\log \frac{1}{P(X = x_i)}}_{\geq 0} \geq 0$$

Pentru a arăta că $H(X) = 0$ dacă și numai dacă X este constantă vom demonstra că ambele implicații au loc:

„ \Rightarrow “ Presupunem că $H(X) = 0$, adică $\sum_i P(X = x_i) \log \frac{1}{P(X = x_i)} = 0$. Datorită faptului că fiecare termen din această sumă este mai mare sau egal cu 0, rezultă că $H(X) = 0$ doar dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $\log \frac{1}{P(X = x_i)} = 0$, adică dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $P(X = x_i) = 1$. Cum însă $\sum_i P(X = x_i) = 1$ rezultă că există o singură valoare x_1 pentru X astfel

încât $P(X = x_1) = 1$, iar $P(X = x) = 0$ pentru orice $x \neq x_1$. Altfel spus, variabila aleatoare discretă X este constantă.²¹

„ \Leftarrow ” Presupunem că variabila X este constantă, ceea ce înseamnă că X ia o singură valoare x_1 , cu probabilitatea $P(X = x_1) = 1$. Prin urmare, $H(X) = -1 \cdot \log 1 = 0$.

b. Pentru a demonstra egalitatea cerută vom porni de la definiția lui $H(Y | X)$ și apoi vom efectua câteva transformări elementare:

$$\begin{aligned}
 H(Y | X) &= \sum_i P(X = x_i) H(Y | X = x_i) \\
 &= \sum_i P(X = x_i) \left[- \sum_j P(Y = y_j | X = x_i) \log P(Y = y_j | X = x_i) \right] \\
 &= - \sum_i \sum_j \underbrace{P(X = x_i) P(Y = y_j | X = x_i)}_{=P(X=x_i, Y=y_j)} \log P(Y = y_j | X = x_i) \\
 &= - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i)
 \end{aligned}$$

c. În primul rând, trebuie să demonstrăm că

$$H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$$

Din definiția entropiei corelate știm că: $H(X, Y) = - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(X = x_i, Y = y_j)$. Vom aplica mai întâi regula de multiplicare, $P(X, Y) = P(X) \cdot P(Y | X)$, după care vom transforma logaritmul produsului în sumă de logaritmi. Pentru claritatea demonstrației vom nota prescurtat $p(x_i) = P(X = x_i)$, $p(x_i, y_j) = P(X = x_i, Y = y_j)$, etc.

$$\begin{aligned}
 H(X, Y) &= - \sum_i \sum_j p(x_i, y_j) \log p(x_i, y_j) \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log[p(x_i) \cdot p(y_j | x_i)] \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) [\log p(x_i) + \log p(y_j | x_i)] \\
 &= - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log p(x_i) - \sum_i \sum_j p(x_i) \cdot p(y_j | x_i) \log p(y_j | x_i) \\
 &= - \sum_i p(x_i) \log p(x_i) \cdot \underbrace{\sum_j p(y_j | x_i)}_{=1} - \sum_i p(x_i) \sum_j p(y_j | x_i) \log p(y_j | x_i) \\
 &= H(X) + \sum_i p(x_i) H(Y | X = x_i) = H(X) + H(Y | X)
 \end{aligned}$$

Pentru a demonstra egalitatea $H(X, Y) = H(Y) + H(X | Y)$, se procedează analog, înlocuind $p(x_i, y_j)$ nu cu $p(x_i) \cdot p(y_j | x_i)$, ci cu $p(y_i) \cdot p(x_j | y_i)$.

²¹Mai corect spus, X este constantă pe tot domeniul de definiție, eventual cu excepția unei mulțimi de probabilitate 0.

Pentru cazul general

$$H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})$$

vom folosi regula de înlănțuire de la variabile aleatoare

$$P(X_1, \dots, X_n) = P(X_1) \cdot P(X_2 | X_1) \cdot P(X_3 | X_1, X_2) \cdot \dots \cdot P(X_n | X_1, \dots, X_{n-1}),$$

precum și scrierea entropiei sub formă de medie, $H(X) = E \left[\log \frac{1}{P(X)} \right]$:

$$\begin{aligned} H(X_1, \dots, X_n) &= E \left[\log \frac{1}{p(x_1, \dots, x_n)} \right] \\ &= - E_{p(x_1, \dots, x_n)} [\log p(x_1, \dots, x_n)] \\ &= - E_{p(x_1, \dots, x_n)} [\log p(x_1) + \log p(x_2 | x_1) + \dots + \log p(x_n | x_1, \dots, x_{n-1})] \\ &= - E_{p(x_1)} [\log p(x_1)] - E_{p(x_1, x_2)} [\log p(x_2 | x_1)] - \dots \\ &\quad - E_{p(x_1, \dots, x_n)} [\log p(x_n | x_1, \dots, x_{n-1})] \\ &= H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1}) \end{aligned}$$

31. (Entropie, entropie condițională, câștig de informație: exemplificare [clasică])

■ □ ● ○ CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 2

Problema aceasta se referă la aruncarea a două zaruri perfecte, cu 6 fețe.

a. Calculează distribuția probabilistă a sumei numerelor de pe cele două fețe care au fost obținute/„observate“ în urma aruncării zarurilor.

În continuare, suma aceasta va fi asimilată cu o variabilă aleatoare, notată cu S .

b. Cantitatea de *informație* obținută (sau: *surpriza* pe care o resimțim) la „observarea“ producerii valorii x a unei variabile aleatoare X oarecare este prin *definiție*

$$\text{Information}(P(X = x)) = \text{Surprise}(P(X = x)) = \log_2 \frac{1}{P(X = x)} = -\log_2 P(X = x).$$

Această cantitate este exprimată (numeric) în *biți de informație*.

Cât de surprins vei fi atunci când vei „observa“ $S = 2$, respectiv $S = 11$, $S = 5$ și $S = 7$? (Vei exprima de fiecare dată rezultatul în biți. Puteți folosi $\log_2 3 = 1.584962501$.)

c. Calculează entropia variabilei S .

d. Să presupunem acum că vei arunca aceste două zaruri pe rând, iar la aruncarea primului zar se obține numărul 4. Cât este entropia lui S în urma acestei „observații“? S-a pierdut, ori s-a câștigat informație în acest proces? Calculează cât de multă informație (exprimată în biți) s-a pierdut ori s-a câștigat.

Răspuns:

a. Redăm distribuția lui S (ușor de calculat) în următorul tabel:

| | | | | | | | | | | | |
|--------|------|------|------|------|------|------|------|------|------|------|------|
| S | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| $P(S)$ | 1/36 | 2/36 | 3/36 | 4/36 | 5/36 | 6/36 | 5/36 | 4/36 | 3/36 | 2/36 | 1/36 |

b. Conform definiției date, vom avea:

$$\begin{aligned} \text{Information}(S = 2) &= -\log_2(1/36) = \log_2 36 = 2 \log_2 6 = 2(1 + \log_2 3) \\ &= 5.169925001 \text{ biți} \end{aligned}$$

$$\text{Information}(S = 11) = -\log_2 2/36 = \log_2 18 = 1 + 2 \log_2 3 = 4.169925001 \text{ biți}$$

$$\text{Information}(S = 5) = -\log_2 4/36 = \log_2 9 = 2 \log_2 3 = 3.169925001 \text{ biți}$$

$$\text{Information}(S = 7) = -\log_2 6/36 = \log_2 6 = 1 + \log_2 3 = 2.584962501 \text{ biți}$$

c. Conform definiției pentru entropie (vezi problema 30), $H(S)$ este media ponderată (cu ajutorul probabilităților) a „surprizelor” / cantităților de informație produse la „observarea” tuturor valorilor variabilei S . Făcând calculele, vom obține:

$$\begin{aligned} H(S) &= -\sum_{i=1}^n p_i \log p_i \\ &= -\left(2 \cdot \frac{1}{36} \log \frac{1}{36} + 2 \cdot \frac{2}{36} \log \frac{2}{36} + 2 \cdot \frac{3}{36} \log \frac{3}{36} + 2 \cdot \frac{4}{36} \log \frac{4}{36} + \right. \\ &\quad \left. 2 \cdot \frac{5}{36} \log \frac{5}{36} + \frac{6}{36} \log \frac{6}{36}\right) \\ &= \frac{1}{36}(2 \log_2 36 + 4 \log_2 18 + 6 \log_2 12 + 8 \log_2 9 + 10 \log_2 \frac{36}{5} + 6 \log_2 6) \\ &= \frac{1}{36}(2 \log_2 6^2 + 4 \log_2 6 \cdot 3 + 6 \log_2 6 \cdot 2 + 8 \log_2 3^2 + 10 \log_2 \frac{6^2}{5} + 6 \log_2 6) \\ &= \frac{1}{36}(40 \log_2 6 + 20 \log_2 3 + 6 - 10 \log_2 5) \\ &= \frac{1}{36}(60 \log_2 3 + 46 - 10 \log_2 5) = 3.274401919 \text{ biți.} \end{aligned}$$

d. Distribuția variabilei S condiționată de observarea feței 4 la prima aruncare este:

| | | | | | | | | | | | |
|------------|---|---|---|-----|-----|-----|-----|-----|-----|----|----|
| S | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 |
| $P(S ...)$ | 0 | 0 | 0 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 1/6 | 0 | 0 |

În consecință, folosind definiția entropiei condiționale specifice (vezi de asemenea problema 30), vom avea:

$$H(S|First-die-shows-4) = -6 \cdot \frac{1}{6} \log_2 \frac{1}{6} = \log_2 6 = 2.58 \text{ biți,}$$

ceea ce înseamnă că se obține următorul *câștig* de informație:

$$IG(S; First-die-shows-4) = H(S) - H(S|First-die-shows-4) = 3.27 - 2.58 = 0.69 \text{ biți.}$$

Altfel spus, atunci când ni se comunică faptul că la aruncarea celor două zaruri primul dintre ele produce fața 4, această informație va reduce ulterior entropia variabilei S (sau, am putea spune, „surpriza” medie provocată de valorile ei) cu 0.69 biți.

32.

(Probabilități marginale,
entropii, entropii condiționale medii)

■ □ ● CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 3

Un doctor trebuie să pună un diagnostic unui pacient care are simptome de răceală (C , de la engl. cold). Factorul principal pe care doctorul îl ia în considerare pentru a elabora diagnosticul este timpul, adică starea vremii de afară (T). Variabila aleatoare C ia două valori, *yes* și *no*, iar variabila aleatoare T ia 3 valori: *sunny* (însorit), *rainy* (ploios) și *snowy* (foarte rece, să zicem). Distribuția corelată a celor două variabile este dată în tabelul următor:

| | $T = \textit{sunny}$ | $T = \textit{rainy}$ | $T = \textit{snowy}$ |
|--------------------|----------------------|----------------------|----------------------|
| $C = \textit{no}$ | 0.30 | 0.20 | 0.10 |
| $C = \textit{yes}$ | 0.05 | 0.15 | 0.20 |

a. Calculați probabilitățile marginale $P(C)$ și $P(T)$.

Sugestie: Folosiți formula $P(X = x) = \sum_y P(X = x; Y = y)$. De exemplu,

$$P(C = \textit{no}) = P(C = \textit{no}, T = \textit{sunny}) + P(C = \textit{no}, T = \textit{rainy}) + P(C = \textit{no}, T = \textit{snowy}).$$

b. Calculați entropiile $H(C)$ și $H(T)$.

c. Calculați entropiile condiționale medii $H(C|T)$ și $H(T|C)$.

Răspuns:

a. Folosind formula dată, vom obține: $P_C = (0.6, 0.4)$ și $P_T = (0.35, 0.35, 0.30)$.

b. Aplicând definiția pentru entropie (vezi problema 30), rezultă:

$$\begin{aligned}
 H(C) &= 0.6 \log \frac{5}{3} + 0.4 \log \frac{5}{2} = \log 5 - 0.6 \log 3 - 0.4 = 0.971 \text{ biți} \\
 H(T) &= 2 \cdot 0.35 \log \frac{20}{7} + 0.3 \log \frac{10}{3} \\
 &= 0.7(2 + \log 5 - \log 7) + 0.3(1 + \log 5 - \log 3) \\
 &= 1.7 + \log 5 - 0.7 \log 7 - 0.3 \log 3 = 1.581 \text{ biți}.
 \end{aligned}$$

c. Aplicând definiția pentru entropie condițională medie (vezi de asemenea

problema 30), vom avea:

$$\begin{aligned}
 H(C|T) &= 0.35 \cdot H\left(\frac{0.30}{0.30+0.05}, \frac{0.05}{0.30+0.05}\right) + 0.35 \cdot H\left(\frac{0.20}{0.20+0.15}, \frac{0.15}{0.20+0.15}\right) + \\
 &\quad 0.30 \cdot H\left(\frac{0.10}{0.10+0.20}, \frac{0.20}{0.20+0.10}\right) \\
 &= \frac{7}{20} \cdot H\left(\frac{6}{7}, \frac{1}{7}\right) + \frac{7}{20} \cdot H\left(\frac{4}{7}, \frac{3}{7}\right) + \frac{3}{10} \cdot H\left(\frac{1}{3}, \frac{2}{3}\right) \\
 &= \frac{7}{20} \cdot \left(\frac{6}{7} \log \frac{7}{6} + \frac{1}{7} \log 7\right) + \frac{7}{20} \cdot \left(\frac{4}{7} \log \frac{7}{4} + \frac{3}{7} \log \frac{7}{3}\right) + \frac{3}{10} \cdot \left(\frac{1}{3} \log 3 + \frac{2}{3} \log \frac{3}{2}\right) \\
 &= \frac{7}{20} \cdot \left(\log 7 - \frac{6}{7} - \frac{6}{7} \log 3\right) + \frac{7}{20} \cdot \left(\log 7 - \frac{8}{7} - \frac{3}{7} \log 3\right) + \frac{3}{10} \cdot \left(\log 3 - \frac{2}{3}\right) \\
 &= \frac{7}{10} \log 7 - \left(\frac{3}{10} + \frac{4}{10} + \frac{2}{10}\right) - \log 3 \left(\frac{6}{20} + \frac{3}{20} - \frac{3}{10}\right) \\
 &= \frac{7}{10} \log 7 - \frac{3}{20} \log 3 - \frac{9}{10} = 0.82715 \text{ biți.}
 \end{aligned}$$

Similar,

$$\begin{aligned}
 H(T|C) &= 0.60 \cdot H\left(\frac{0.30}{0.30+0.20+0.10}, \frac{0.20}{0.30+0.20+0.10}, \frac{0.10}{0.30+0.20+0.10}\right) + \\
 &\quad 0.40 \cdot H\left(\frac{0.05}{0.05+0.15+0.20}, \frac{0.15}{0.05+0.15+0.20}, \frac{0.20}{0.05+0.15+0.20}\right) \\
 &= \frac{3}{5} \cdot H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) + \frac{2}{5} \cdot H\left(\frac{1}{8}, \frac{3}{8}, \frac{1}{2}\right) \\
 &= \frac{3}{5} \left(\frac{1}{2} + \frac{1}{3} \log 3 + \frac{1}{6} (1 + \log 3)\right) + \frac{2}{5} \left(\frac{1}{8} \cdot 3 + \frac{3}{8} (3 - \log 3) + \frac{1}{2}\right) \\
 &= \frac{3}{5} \left(\frac{2}{3} + \frac{1}{2} \log 3\right) + \frac{2}{5} \left(2 - \frac{3}{8} \log 3\right) \\
 &= \frac{6}{5} + \frac{3}{20} \log 3 = 1.43774 \text{ biți.}
 \end{aligned}$$

33. (Calcularea entropiei unei variabile aleatoare continue:
cazul distribuției exponențiale)

■ □ CMU, 2011 spring, R. Rosenfeld, HW2, pr. 2.c

Pentru o distribuție de probabilitate continuă P , entropia se definește astfel:

$$H(P) = \int_{-\infty}^{+\infty} P(x) \log_2 \frac{1}{P(x)} dx$$

Calculați entropia *distribuției* continue *exponențiale* de parametru $\lambda > 0$. Vă reamintim că definiția acestei distribuții este următoarea:

$$P(x) = \begin{cases} \lambda e^{-\lambda x}, & \text{dacă } x \geq 0; \\ 0, & \text{dacă } x < 0. \end{cases}$$

Indicație: Dacă $P(x) = 0$, veți presupune că $-P(x) \log P(x) = 0$.

Răspuns:

Dat fiind faptul că funcția P se anulează pe intervalul $(-\infty, 0)$, este natural ca mai întâi să „rupem” intervalul de integrare pentru $\int_{-\infty}^{\infty} P(x) \log_2 \frac{1}{P(x)} dx$ în două: $(-\infty, 0]$ și $[0, \infty)$. Așadar,

$$\begin{aligned} H(P) &= \int_{-\infty}^0 P(x) \log_2 \frac{1}{P(x)} dx + \int_0^{\infty} P(x) \log_2 \frac{1}{P(x)} dx \\ &\stackrel{\text{def. } P}{=} \int_{-\infty}^0 0 \log_2 0 dx + \int_0^{\infty} \lambda e^{-\lambda x} \log_2 \frac{1}{\lambda e^{-\lambda x}} dx \end{aligned}$$

Prima dintre aceste două ultime integrale este 0, conform *indicației* din enunț. Pentru a putea calcula cea de-a doua integrală (în expresia căreia apare numărul e), vom schimba baza logaritmului, și anume vom trece din baza 2 în baza e (baza logaritmului natural, \ln).²²

Prin urmare,

$$\begin{aligned} H(P) &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \ln \frac{1}{\lambda e^{-\lambda x}} dx = \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \left(\ln \frac{1}{\lambda} + \ln \frac{1}{e^{-\lambda x}} \right) dx \\ &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda + \ln e^{\lambda x}) dx \\ &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda + \lambda x) dx \\ &= \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} (-\ln \lambda) dx + \frac{1}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} \lambda x dx \\ &= \frac{-\ln \lambda}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} dx + \frac{\lambda}{\ln 2} \int_0^{\infty} \lambda e^{-\lambda x} x dx \\ &= \frac{\ln \lambda}{\ln 2} \int_0^{\infty} (e^{-\lambda x})' dx - \frac{\lambda}{\ln 2} \int_0^{\infty} (e^{-\lambda x})' x dx \end{aligned}$$

Prima integrală se rezolvă foarte ușor:

$$\int_0^{\infty} (e^{-\lambda x})' dx = e^{-\lambda x} \Big|_0^{\infty} = e^{-\infty} - e^0 = 0 - 1 = -1$$

Pentru a rezolva cea de-a doua integrală se poate folosi *formula de integrare prin părți*:

$$\int_0^{\infty} (e^{-\lambda x})' x dx = e^{-\lambda x} x \Big|_0^{\infty} - \int_0^{\infty} e^{-\lambda x} x' dx = e^{-\lambda x} x \Big|_0^{\infty} - \int_0^{\infty} e^{-\lambda x} dx$$

Integrala definită $e^{-\lambda x} x \Big|_0^{\infty}$ nu se poate calcula direct (din cauza conflictului $0 \cdot \infty$ care se produce atunci când lui x i se atribuie valoarea-limită ∞), ci se calculează folosind *regula lui l'Hôpital*:

$$\lim_{x \rightarrow \infty} x e^{-\lambda x} = \lim_{x \rightarrow \infty} \frac{x}{e^{\lambda x}} = \lim_{x \rightarrow \infty} \frac{x'}{(e^{\lambda x})'} = \lim_{x \rightarrow \infty} \frac{1}{\lambda e^{\lambda x}} = \frac{1}{\lambda} \lim_{x \rightarrow \infty} e^{-\lambda x} = e^{-\infty} = 0,$$

²² Pentru aceasta, vom folosi formula $\log_a b = \frac{\log_c b}{\log_c a}$, valabilă pentru orice $a > 0$, $b > 0$ și $c > 0$, cu $a \neq 1$ și $c \neq 1$. În calculele care urmează vom folosi și alte formule de la logaritmi.

deci

$$e^{-\lambda x} x \Big|_0^\infty = 0 - 0 = 0.$$

Integrala $\int_0^\infty e^{-\lambda x} dx$ se calculează ușor:

$$\int_0^\infty e^{-\lambda x} dx = -\frac{1}{\lambda} \int_0^\infty (e^{-\lambda x})' dx = -\frac{1}{\lambda} e^{-\lambda x} \Big|_0^\infty = -\frac{1}{\lambda}(0 - 1) = \frac{1}{\lambda}$$

Prin urmare,

$$\int_0^\infty (e^{-\lambda x})' x dx = 0 - \frac{1}{\lambda} = -\frac{1}{\lambda},$$

ceea ce conduce la rezultatul final:

$$H(P) = \frac{\ln \lambda}{\ln 2}(-1) - \frac{\lambda}{\ln 2} \left(-\frac{1}{\lambda}\right) = -\frac{\ln \lambda}{\ln 2} + \frac{1}{\ln 2} = \frac{1 - \ln \lambda}{\ln 2}.$$

34.

(Entropia relativă: definiție și proprietăți elementare;
exprimarea câștigului de informație
cu ajutorul entropiei relative)

■ □ ○ prelucrare de Liviu Ciortuz după

CMU, 2007 fall, Carlos Guestrin, HW1, pr. 1.2

Entropia relativă sau divergența Kullback-Leibler (KL) a unei distribuții p în raport cu o altă distribuție q — ambele distribuții fiind discrete — se definește astfel:

$$KL(p||q) \stackrel{\text{def.}}{=} - \sum_{x \in X} p(x) \log \frac{q(x)}{p(x)}$$

Din perspectiva teoriei informației, divergența KL specifică numărul de *biți adiționali* care sunt necesari în medie pentru a transmite valorile variabilei X atunci când presupunem că aceste valori sunt distribuite conform distribuției („model“) q , dar în realitate ele urmează o altă distribuție, p .²³

a. Demonstrați că $KL(p||q) \geq 0$ și apoi arătați că egalitatea are loc dacă și numai dacă $p = q$.²⁴

Indicație:

Pentru a demonstra punctul acesta puteți folosi *inegalitatea lui Jensen*:

Dacă $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ este o funcție convexă, atunci pentru orice $t \in [0, 1]$ și orice $x_1, x_2 \in \mathbb{R}$ urmează $\varphi(tx_1 + (1-t)x_2) \leq t\varphi(x_1) + (1-t)\varphi(x_2)$. Dacă φ este funcție strict convexă, atunci egalitatea are loc doar dacă $x_1 = x_2$.

Mai general, pentru orice $a_i \geq 0$, $i = 1, \dots, n$ cu $\sum_i a_i \neq 0$ și orice $x_i \in \mathbb{R}$, $i =$

²³ *Atenție:* Divergența KL nu este o măsură de *distanță* între două distribuții probabiliste, fiindcă în general ea nu este simetrică ($KL(p||q) \neq KL(q||p)$). Pentru „simetrizare“, se consideră $JSD(p||q) = \frac{1}{2}(KL(p||q) + KL(q||p))$, care se numește *divergența Jensen-Shannon*.

Măsura

$$\begin{aligned} d(X, Y) &\stackrel{\text{def}}{=} H(X, Y) - IG(X; Y) = H(X) + H(Y) - 2IG(X; Y) \\ &= H(X | Y) + H(Y | X), \end{aligned}$$

cunoscută ca *variația informației*, este o măsură de distanță (metrică), adică este nengativă, simetrică, implică egalitatea indiscernabililor și satisface inegalitatea triunghiului.

²⁴ Mai general, $KL(p||q)$ este cu atât mai mică cu cât „asemănarea“ dintre distribuțiile p și q este mai mare.

$1, \dots, n$, avem $\varphi\left(\frac{\sum_i a_i x_i}{\sum_j a_j}\right) \leq \frac{\sum_i a_i \varphi(x_i)}{\sum_j a_j}$. Dacă φ este strict convexă, atunci egalitatea are loc doar dacă $x_1 = \dots = x_n$. Evident, rezultate similare pot fi formulate și pentru funcții concave.

b. Câștigul de informație poate fi definit ca fiind entropia relativă dintre distribuția corelată observată a lui X și Y pe de o parte, și produsul distribuțiilor marginale p_X și p_Y pe de altă parte:

$$IG(X, Y) \stackrel{\text{def.}}{=} KL(p_{X,Y} \parallel (p_X p_Y)) = - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right)$$

Arătați că această nouă definiție a câștigului de informație este echivalentă cu definiția dată anterior (vezi problema 30). Cu alte cuvinte, arătați că

$$KL(p_{X,Y} \parallel (p_X p_Y)) = H[X] - H[X | Y]$$

Observație: Din noua definiție introdusă mai sus pentru câștigul de informație, rezultă imediat că

$$\begin{aligned} IG(X, Y) &= \sum_y p(y) \sum_x p(x | y) \log \frac{p(x | y)}{p(x)} = \sum_y p(y) KL(p_{X|Y} \parallel p_X) \\ &= E_Y[KL(p_{X|Y} \parallel p_X)] \end{aligned}$$

ceea ce înseamnă că $IG(X, Y)$ poate fi văzută ca o medie (în raport cu distribuția lui Y) a divergenței KL dintre distribuția condițională a lui X în raport cu Y pe de o parte, și distribuția lui X pe de altă parte.

c. Arătați că $IG(X, Y) \geq 0$ pentru orice variabile aleatoare discrete X și Y . În particular, $IG(X, Y) = 0$ dacă și numai dacă X și Y sunt independente.

Răspuns:

a. Vom dovedi inegalitatea $KL(p \parallel q) \geq 0$ folosind inegalitatea lui Jensen, în expresia căreia vom înlocui φ cu funcția convexă $-\log_2$, pe a_i cu $p(x_i)$ și pe x_i cu $\frac{q(x_i)}{p(x_i)}$. (Pentru conveniență, în cele ce urmează vor renunța la indicele variabilei x .) Vom avea:

$$\begin{aligned} KL(p \parallel q) &\stackrel{\text{def.}}{=} - \sum_x p(x) \log \frac{p(x)}{q(x)} \\ &\stackrel{\text{Jensen}}{\geq} - \log \left(\sum_x p(x) \frac{q(x)}{p(x)} \right) = - \log \left(\underbrace{\sum_x q(x)}_1 \right) = - \log 1 = 0 \end{aligned}$$

Vom demonstra acum că $KL(p \parallel q) \geq 0 \Leftrightarrow p = q$.

Egalitatea $p(x) = q(x)$ implică $\frac{q(x)}{p(x)} = 1$, deci $\log \frac{q(x)}{p(x)} = 0$ pentru orice x , de unde rezultă imediat $KL(p \parallel q) = 0$.

Pentru a demonstra implicația inversă, se ține cont că în inegalitatea lui Jensen are loc egalitatea doar în cazul în care $x_i = x_j$ pentru orice i și j . În cazul de

față, această condiție se traduce prin faptul că raportul $\frac{q(x)}{p(x)}$ este același pentru orice valoare a lui x . Ținând cont că $\sum_x p(x) = 1$ și $\sum_x p(x) \frac{q(x)}{p(x)} = \sum_x q(x) = 1$, rezultă că $\frac{q(x)}{p(x)} = 1$ sau, altfel spus, $p(x) = q(x)$ pentru orice x , ceea ce înseamnă că distribuțiile p și q sunt identice.

b. Vom folosi regula de multiplicare, și anume $p(x, y) = p(x | y)p(y)$:

$$\begin{aligned}
 KL(p_{XY} || (p_X p_Y)) &\stackrel{\text{def.}}{=} KL - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right) \\
 &= - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x | y)p(y)} \right) = - \sum_x \sum_y p(x, y) [\log p(x) - \log p(x | y)] \\
 &= - \sum_x \sum_y p(x, y) \log p(x) - \left(- \sum_x \sum_y p(x, y) \log p(x | y) \right) \\
 &= - \sum_x \log p(x) \underbrace{\sum_y p(x, y)}_{=p(x)} - H[X | Y] \\
 &= H[X] - H[X | Y] = IG(X, Y)
 \end{aligned}$$

c. Atât inegalitatea că $IG(X, Y) \geq 0$ cât și echivalența $IG(X, Y) = 0 \Leftrightarrow X$ și Y sunt independente constituie consecințe imediate ale punctelor b și a discutate mai sus.

35.

(Informația mutuală, aplicație:
selecția de trăsături)

□ CMU, 2009 spring, Ziv Bar-Joseph, HW5, pr. 6

În tabelul alăturat se dă un set de opt observații/instanțe, reprezentate ca tupluri de valori ale variabilelor aleatoare binare de „intrare“ X_1, X_2, X_3, X_4, X_5 și ale variabilei aleatoare binare de „ieșire“ Y .

Am dori să reducem spațiul de trăsături $\{X_1, X_2, X_3, X_4, X_5\}$ folosind o metodă de selecție de tip *filtru*.

| X_1 | X_2 | X_3 | X_4 | X_5 | Y |
|-------|-------|-------|-------|-------|-----|
| 0 | 1 | 1 | 0 | 1 | 0 |
| 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | 1 | 1 | 0 | 1 |
| 0 | 1 | 1 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 1 | 1 |
| 1 | 0 | 0 | 1 | 0 | 1 |
| 1 | 1 | 1 | 0 | 1 | 1 |

a. Calculați informația mutuală $MI(X_i, Y)$ pentru fiecare i .

b. Ținând cont de rezultatul de la punctul precedent, alegeți cel mai mic subset de trăsături în așa fel încât cel mai bun clasificator antrenat pe acest spațiu (reduc) de trăsături să fie cel puțin la fel de bun ca și cel mai bun clasificator antrenat pe întreg spațiul de trăsături. Justificați alegerea pe care ați făcut-o.

Răspuns:

a. Pentru calculul informației mutuale putem folosi formula din problema 34:

$$MI(X, Y) = \sum_x \sum_y p_{X,Y}(x, y) \log \left(\frac{p_{X,Y}(x, y)}{p_X(x)p_Y(y)} \right)$$

Probabilitățile marginale, estimate în sensul verosimilității maxime (MLE), sunt:

| | P_{X_1} | P_{X_2} | P_{X_3} | P_{X_4} | P_{X_5} | P_Y |
|---|-----------|-----------|-----------|-----------|-----------|-------|
| 0 | 0.5 | 3/8 | 0.5 | 0.5 | 0.5 | 0.25 |
| 1 | 0.5 | 5/8 | 0.5 | 0.5 | 0.5 | 0.75 |

iar probabilitățile corelate sunt:

| X_i | Y | $P_{X_1,Y}$ | $P_{X_2,Y}$ | $P_{X_3,Y}$ | $P_{X_4,Y}$ | $P_{X_5,Y}$ |
|-------|-----|-------------|-------------|-------------|-------------|-------------|
| 0 | 0 | 1/8 | 1/8 | 1/8 | 1/4 | 0 |
| 0 | 1 | 3/8 | 1/4 | 3/8 | 1/4 | 1/2 |
| 1 | 0 | 1/8 | 1/8 | 1/8 | 0 | 1/4 |
| 1 | 1 | 3/8 | 1/2 | 3/8 | 1/2 | 1/4 |

Se poate observa că X_1 și Y sunt independente, deci $MI(X_1, Y) = 0$, conform problemei 34.c. Similar, $MI(X_3, Y) = 0$. În rest, efectuând calculele obținem $MI(X_2, Y) = 0.01571$, $MI(X_4, Y) = 0.3113$ și $MI(X_5, Y) = 0.3113$.

b. La selecția de trăsături vom alege acele trăsături X_i care au informație mutuală nenulă în raport cu Y . Acestea sunt X_2, X_4 și X_5 . Celelalte două trăsături, X_1 și X_3 sunt independente în raport cu Y .

Totuși, inspectând datele, observăm că dacă vom selecta doar trăsăturile X_2, X_4 și X_5 vom avea două instanțe (vezi prima și ultima linie din tabel) care au aceleași trăsături ($X_2 = 1, X_4 = 0, X_5 = 1$) dar au etichete/ieșiri diferite: $Y = 0$, respectiv $Y = 1$. Așadar, vom adăuga la setul de trăsături selectate anterior și variabila X_1 , care va permite dezambiguizarea în cazul acestor două instanțe.

Observație: Deși $MI(X_1, Y) = 0$, nu rezultă că $MI(X_1, X_j) = 0$ pentru $j \in \{2, 4, 5\}$. Aceasta explică de ce adăugarea variabilei X_1 la setul $\{X_2, X_4, X_5\}$ poate aduce un câștig de informație.

36. (Entropia corelată: forma particulară a relației de „înlănțuire” în cazul variabilelor aleatoare independente)

□ prelucrare de Liviu Ciortuz după
CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 7.b

Demonstrați că dacă X și Y sunt variabile aleatoare independente discrete, atunci $H(X, Y) = H(X) + H(Y)$.

Este adevărată și reciprocă acestei afirmații? Adică, atunci când are loc egalitatea $H(X, Y) = H(X) + H(Y)$ rezultă că variabilele X și Y sunt independente?

Răspuns:

Conform problemei 30.c, *formula de înlănțuire* a entropiilor pentru cazul general (adică, indiferent dacă X și Y sunt sau nu independente) este:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y). \quad (8)$$

Conform definiției câștigului de informație (vezi problema 30),

$$IG(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (9)$$

De asemenea, conform problemei 34.c,

$$IG(X, Y) = 0 \Leftrightarrow X \text{ și } Y \text{ sunt independente.} \quad (10)$$

Din relațiile (9) și (10) rezultă că

$$H(Y) = H(Y|X) \Leftrightarrow X \text{ și } Y \text{ sunt independente.} \quad (11)$$

Așadar, dacă X și Y sunt independente, coroborând relațiile (11) și (8) vom avea $H(X, Y) = H(Y) + H(X)$.

Invers, dacă $H(X, Y) = H(X) + H(Y)$, din relația (8) rezultă că $H(Y) = H(Y|X)$, ceea ce implică faptul că X și Y sunt independente, conform relației (11).

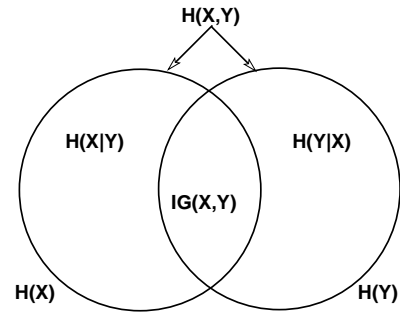
Observație: Din egalitățile (8) și (9), rezultă

$$H(X, Y) = H(X) + H(Y) - IG(X, Y).$$

Conform proprietății de ne-negativitate a câștigului de informație ($IG(X, Y) \geq 0$; vezi problema 34.c), rezultă

$$H(X, Y) \leq H(X) + H(Y).$$

Această ultimă relație, precum și relațiile (8) și (9) sunt ilustrate în figura alăturată.



37.

(Cross-entropie: definiție, o proprietate (ne-negativitatea) și un exemplu simplu de calculare [a valorii cross-entropiei])

□ • ◦ CMU, 2011 spring, Roni Rosenfeld, HW2, pr. 3.c

Cross-entropia a două distribuții p și q , desemnată prin $CH(p, q)$, reprezintă numărul mediu de biți necesari pentru a codifica un eveniment dintr-o mulțime oarecare de posibilități, atunci când schema de codificare folosită se bazează pe o distribuție de probabilitate dată q , în loc să se bazeze pe distribuția „adevărată” p . În cazul în care distribuțiile p și q sunt discrete, această noțiune se definește formal astfel:

$$CH(p, q) = - \sum_x p(x) \log q(x).$$

În cazul distribuțiilor continue, definiția se obține/construiește prin analogie:

$$CH(p, q) = - \int_X p(x) \log q(x) dx.$$

a. Poate oare cross-entropia să ia valori negative? Faceți o demonstrație sau dați un contra-exemplu.

b. În multe experimente, pentru a stabili calitatea diferitelor ipoteze/modele, se procedează la evaluarea/compararea lor pe un set de date. Să presupunem că, urmărind să faci predicția *funcției de probabilitate* asociate unei anumite *variabile aleatoare* care are 7 valori posibile, ai obținut (printr-un procedeu oarecare) două *modele* diferite, iar *distribuțiile de probabilitate* prezise de către aceste două modele sunt respectiv:

$$q_1 = \left(\frac{1}{10}, \frac{1}{10}, \frac{1}{5}, \frac{3}{10}, \frac{1}{5}, \frac{1}{20}, \frac{1}{20} \right) \text{ și } q_2 = \left(\frac{1}{20}, \frac{1}{10}, \frac{3}{20}, \frac{7}{20}, \frac{1}{5}, \frac{1}{10}, \frac{1}{20} \right).$$

Să zicem că pentru evaluare folosești un set de date caracterizat de următoarea distribuție *empirică*:

$$p_{\text{empiric}} = \left(\frac{1}{20}, \frac{1}{10}, \frac{1}{5}, \frac{3}{10}, \frac{1}{5}, \frac{1}{10}, \frac{1}{20} \right).$$

Calculează cross-entropiile $CH(p_{\text{empiric}}, q_1)$ și $CH(p_{\text{empiric}}, q_2)$.

Care dintre aceste două modele va conduce la o cross-entropie mai mică? Putem oare garanta că acest model este într-adevăr [cel] mai bun? Explică/justifică răspunsul [pe care l-ai] dat.

Răspuns:

a. Nu, cross-entropia nu poate lua valori negative. Iată cum demonstrăm: Știm că pentru orice funcții de probabilitate p și q și pentru orice x (care aparține domeniului de valori al unei variabile aleatoare care are o astfel de distribuție de probabilitate), valorile $p(x)$ și $q(x)$ satisfac inegalitățile $0 \leq p(x) \leq 1$ and $0 \leq q(x) \leq 1$. Inegalitatea $q(x) \leq 1$ implică faptul că $\log q(x) \leq 0$. Din $0 \leq p(x)$ și $-\log q(x) \geq 0$, rezultă că $0 \leq -p(x) \log q(x)$. În consecință, suma tuturor acestor termeni va fi de asemenea mai mare sau egală cu 0, deci cross-entropia nu poate fi niciodată negativă.

Observație importantă: Spre deosebire de entropie (vezi problema 68), cross-entropia nu este mărginită. Ea poate crește la infinit; vezi cazul când pentru o anumită valoare x sunt adevărate simultan relațiile $p(x) \neq 0$ și $q(x) = 0$.²⁵

²⁵ Mai precis, $\lim_{q(x) \rightarrow +0} (-p(x) \cdot \log_2 q(x)) = -p(x)(-\infty) = +\infty$.

b. Facem calculele, folosind formula cross-entropiei:

$$\begin{aligned}
 CH(p_{\text{empiric}}, q_1) &= \\
 &- \left(\frac{1}{20} \log_2 \frac{1}{10} + \frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{5} \log_2 \frac{1}{5} + \frac{3}{10} \log_2 \frac{3}{10} + \frac{1}{5} \log_2 \frac{1}{5} + \frac{1}{10} \log_2 \frac{1}{20} \right. \\
 &\quad \left. + \frac{1}{20} \log_2 \frac{1}{20} \right) = \frac{3}{20} \log_2 10 + \frac{2}{5} \log_2 5 + \frac{3}{10} \log_2 \frac{10}{3} + \frac{3}{20} \log_2 20 = \\
 &\frac{3}{20} \log_2 2 \cdot 5 + \frac{2}{5} \log_2 5 + \frac{3}{10} \log_2 \frac{2 \cdot 5}{3} + \frac{3}{20} \log_2 2^2 \cdot 5 = \\
 &\left(\frac{3}{20} + \frac{3}{10} + 2 \cdot \frac{3}{20} \right) + \left(\frac{3}{20} + \frac{2}{5} + \frac{3}{10} + \frac{3}{20} \right) \log_2 5 - \frac{3}{10} \log_2 3 = \\
 &\frac{3}{4} + \log_2 5 - \frac{3}{10} \log_2 3 = 2.596439345 \text{ biți}
 \end{aligned}$$

$$\begin{aligned}
 CH(p_{\text{empiric}}, q_2) &= \\
 &- \left(\frac{1}{20} \log_2 \frac{1}{20} + \frac{1}{10} \log_2 \frac{1}{10} + \frac{1}{5} \log_2 \frac{3}{20} + \frac{3}{10} \log_2 \frac{7}{20} + \frac{1}{5} \log_2 \frac{1}{5} + \frac{1}{10} \log_2 \frac{1}{10} \right. \\
 &\quad \left. + \frac{1}{20} \log_2 \frac{1}{20} \right) = \\
 &\frac{1}{10} \log_2 20 + \frac{1}{5} \log_2 10 + \frac{1}{5} \log_2 \frac{20}{3} + \frac{3}{10} \log_2 \frac{20}{7} + \frac{1}{5} \log_2 5 = \\
 &\frac{1}{10} \log_2 2^2 \cdot 5 + \frac{1}{5} \log_2 2 \cdot 5 + \frac{1}{5} \log_2 \frac{2^2 \cdot 5}{3} + \frac{3}{10} \log_2 \frac{2^2 \cdot 5}{7} + \frac{1}{5} \log_2 5 = \\
 &\left(2 \cdot \frac{1}{10} + \frac{1}{5} + 2 \cdot \frac{1}{5} + 2 \cdot \frac{3}{10} \right) + \left(\frac{1}{10} + 3 \cdot \frac{1}{5} + \frac{3}{10} \right) \log_2 5 - \frac{1}{5} \log_2 3 - \frac{3}{10} \log_2 7 = \\
 &\frac{7}{5} + \log_2 5 - \frac{1}{5} \log_2 3 - \frac{3}{10} \log_2 7 = 2.562729118 \text{ biți}.
 \end{aligned}$$

Așadar, distribuția p_{empiric} are o cross-entropie mai mică în [raport cu] modelul q_2 . Este deci rezonabil să afirmăm că alegerea modelului q_2 este mai bună.

Totuși, nu putem garanta că acest model este întotdeauna cel mai bun, fiindcă aici lucrăm cu o distribuție „empirică“, iar distribuția „adevărată“ nu neapărat se reflectă în mod complet/perfect în această distribuție empirică.

De obicei, *bias-ul de eșantionare* (engl., sampling bias), precum și *insuficiența datelor de antrenament* vor contribui la lărgirea „spațiului“ care diferențiază distribuția adevărată de distribuția empirică. Prin urmare, în practică, atunci când concepem un [astfel de] experiment de evaluare a mai multor distribuții probabiliste, trebuie să avem permanent în minte faptul acesta și, dacă este posibil, să folosim tehnici care reduc/minimizează aceste riscuri.

38. (Inegalitatea lui Gibbs: un caz particular;
comparație între valorile entropiei și ale cross-entropiei)

□ *Liviu Ciortuz, 2012, după www.en.wikipedia.org*

Fie $P = \{p_1, \dots, p_n\}$ o distribuție de probabilitate discretă.

a. Arătați că pentru orice distribuție de probabilitate $Q = \{q_1, \dots, q_n\}$ are loc inegalitatea:

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

Altfel spus, $H(P) \leq CH(P, Q)$, unde $H(P)$ este entropia distribuției P , iar $CH(P, Q)$ este *cross-entropia* lui P în raport cu Q .

b. Arătați că în formula de mai sus egalitatea are loc dacă și numai dacă $p_i = q_i$ pentru $i = 1, \dots, n$.

Observație: În formula din enunț, în locul bazei 2 pentru logaritm poate fi folosită orice bază supraunitară.

Indicație: Dacă în inegalitatea dată se trece termenul din partea stângă în partea dreaptă, obținem $0 \leq \sum_{i=1}^n p_i \log_2 p_i - \sum_i p_i \log_2 q_i \Leftrightarrow 0 \leq -\sum_{i=1}^n p_i \log_2 \frac{q_i}{p_i}$.

Puteți face legătura dintre expresia din partea dreaptă a acestei ultime inegalități și definiția *entropiei relative* (numită de asemenea *divergența Kullback-Leibler*, vedeți problema 34) și apoi să folosiți proprietățile entropiei relative.

Răspuns:

Expresia $-\sum_{i=1}^n p_i \log_2 \frac{q_i}{p_i}$ la care s-a ajuns în *Indicație* este exact divergența Kullback-Leibler dintre distribuțiile P și Q . Formal, scriem acest lucru astfel: $KL(P||Q) = -\sum_{i=1}^n p_i \log_2 \frac{q_i}{p_i} = CH(P, Q) - H(P)$.

a. La problema 34.a am demonstrat inegalitatea $KL(P||Q) \geq 0$, care are loc pentru orice distribuții probabiliste discrete P și Q . Aceasta este exact proprietatea de care avem nevoie pentru a justifica inegalitatea dată în enunț la acest punct ($H(P) \leq CH(P, Q)$).

b. Tot la problema 34.a s-a demonstrat că $KL(P||Q)$ are valoarea 0 dacă și numai dacă distribuțiile P și Q sunt identice. În contextul nostru, această proprietate se transpune imediat sub forma $H(P) = CH(P, Q) \Leftrightarrow p_i = q_i$ pentru $i = 1, \dots, n$.

Elemente de teoria informației

65. (Probabilități marginale, entropii, entropii condiționale medii)

• ◦ CMU, 2011 spring, Roni Rosenfeld, HW2, pr. 1.d

Echipa de fotbal american The Steelers (din Pittsburgh) va juca în cupa Superbowl XLV contra echipei The Green Bay Packers. Pregătindu-și meciul, ei (fotbaliștii echipei The Steelers) se gândesc să-și definească strategia de joc în funcție de doi factori majori:

- dacă jucătorul Ben Roethlisberger va fi (sau nu) accidentat la vremea meciului (Injured = yes / no), și
- cum anume va fi vremea (Weather = foggy / rainy / clear sky).

Iată distribuția corelată a acestor două tipuri de evenimente:

| | Weather = foggy | rainy | clear sky | P(Injured) |
|---------------|-----------------|-------|-----------|------------|
| Injured = no | 0.1 | 0.25 | 0.35 | |
| Injured = yes | 0.05 | 0.1 | 0.15 | |
| P(Weather) | | | | |

- a. Pornind de la distribuția corelată dată, completați ultima linie și ultima coloană din tabelul de mai sus cu valorile corespunzătoare distribuțiilor marginale $P(\text{Weather})$ și $P(\text{Injured})$.
- b. Calculați entropiile $H(\text{Weather})$ și $H(\text{Injured})$.
- c. Calculați entropiile condiționale medii $H(\text{Injured} \mid \text{Weather})$ și $H(\text{Weather} \mid \text{Injured})$.

Veți putea folosi următoarele aproximări: $\log_2 3 = 1.585$, $\log_2 5 = 2.322$, $\log_2 7 = 2.807$, $\log_2 11 = 3.459$ și $\log_2 13 = 3.700$.

66. (O margine superioară pentru valoarea entropiei unei variabile aleatoare discrete)

■ □ * CMU, 2003 fall, T. Mitchell, A. Moore, HW1, pr. 1.1

Comentariu: La problema 30 s-a demonstrat că entropia oricărei variabile aleatoare discrete este ne-negativă ($H(X) \geq 0$).²⁷ La acest exercițiu veți demonstra — tot pentru cazul discret — că există și o margine superioară pentru $H(X)$.

Așadar, fie X o variabilă aleatoare discretă care ia n valori și urmează distribuția probabilistă P . Conform definiției, entropia lui X este

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

Arătați că $H(X) \leq \log_2 n$.

Sugestie: Puteți folosi inegalitatea $\ln x \leq x - 1$ care are loc pentru orice $x > 0$.

²⁷Extensia acestei proprietăți la cazul variabilelor aleatoare continue este imediată.

67. (Entropia corelată: forma particulară a relației de „înlănțuire” în cazul variabilelor aleatoare independente)

□ • ○ * prelucrare de Liviu Ciortuz după CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 7.b

La problema 36 s-a demonstrat că dacă X și Y sunt variabile aleatoare independente, atunci are loc egalitatea $H(X, Y) = H(X) + H(Y)$ și reciproc. Demonstrația a folosit proprietăți/rezultate obținute în problemele anterioare (30 și 34).

Implicația directă (X și Y independente $\Rightarrow H(X, Y) = H(X) + H(Y)$) se poate obține însă și în mod direct, pornind de la definiția independenței variabilelor aleatoare. Vă cerem să faceți astfel demonstrația acestei implicații. Veți trata mai întâi cazul variabilelor aleatoare discrete și apoi cazul variabilelor aleatoare continue.

68. (Entropie corelată și condiționată: formula de „înlănțuire” condițională)

□ • ○ * CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 4.b

Demonstrați că proprietatea

$$H(X, Y|A) = H(Y|A) + H(X|Y, A)$$

este adevărată pentru oricare 3 variabile aleatoare discrete X , Y și A .

Explicați în mod intuitiv, într-o singură frază, care este semnificația proprietății de mai sus.

69. (O proprietate a câștigului de informație: ne-negativitatea)

■ □ * Liviu Ciortuz

Definiția *câștigului de informație* (sau: a *informației mutuale*) al unei variabile aleatoare X în raport cu o altă variabilă aleatoare Y este $IG(X, Y) = H(X) - H(X | Y) = H(Y) - H(Y | X)$.²⁸ La problema 34 s-a demonstrat — pentru cazul în care X și Y sunt discrete — că $IG(X, Y) = KL(P_{X,Y} || P_X P_Y)$, unde KL desemnează *entropia relativă* (sau: *divergența Kullback-Leibler*), P_X și P_Y sunt distribuțiile variabilelor X și, respectiv, Y , iar $P_{X,Y}$ este distribuția corelată a acestor variabile. Tot la problema 34 s-a arătat că divergența KL este întotdeauna ne-negativă. În consecință, $IG(X, Y) \geq 0$ pentru orice X și Y . La acest exercițiu vă cerem să demonstrați inegalitatea $IG(X, Y) \geq 0$ în manieră directă, adică fără a apela la divergența Kullback-Leibler.

Sugestie: Puteți folosi următoarea formă a inegalității lui Jensen:

$$\sum_{i=1}^n \alpha_i \log x_i \leq \log \left(\sum_{i=1}^n \alpha_i x_i \right)$$

unde baza logaritmului se consideră supraunitară, $\alpha_i \geq 0$ pentru $i = 1, \dots, n$ și $\sum_{i=1}^n \alpha_i = 1$.

²⁸ Vezi problema 30.

70.

(Cross-entropia — o aplicație:
selecția modelelor probabiliste)

□ • ◦ CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 10

Avem un zar măsluit (engl., unfair die). Probabilitățile [reale] de apariție pentru fiecare dintre fețele de la 1 la 6 sunt date de distribuția

$$P_{true} = (0.08, 0.55, 0.15, 0.12, 0.05, 0.05).$$

Fără să cunoască acest fapt, două persoane, identificate cu A și respectiv B , ne sugerează următoarele modele de probabilitate pentru zarul măsluit:

$$P_A = (0.07, 0.14, 0.24, 0.24, 0.05, 0.26)$$

$$P_B = (0.25, 0.13, 0.21, 0.03, 0.11, 0.27)$$

a. Elaborați câteva idei relativ la cum am putea măsura/determina care dintre aceste două modele este mai bun.

b. Calculați cross-entropiile $CH(P_{true}, P_A)$, $CH(P_{true}, P_B)$ și $CH(P_A, P_B)$. Presupunând că alegem ca măsură/mijloc de evaluare a modelelor cross-entropia, care dintre cele două modele (P_A și P_B) credeți că este mai bun?

71.

(Proprietăți ale entropiei: Adevărat sau Fals?)

□ * CMU, 2011 spring, Roni Rosenfeld, HW2, pr. 2.a.1
 CMU, 2008 fall, Eric Xing, final exam, pr. 1.4
 CMU, 2012 spring, Roni Rosenfeld, HW2, pr. 7

Stabiliți dacă următoarele propoziții sunt adevărate sau false.

a. Entropia nu este negativă.

b. $H(X, Y) \geq H(X) + H(Y)$ pentru orice două variabile aleatoare X și Y .

c. Dacă X și Y sunt variabile aleatoare independente, atunci $H(X, Y) = H(X) + H(Y)$.