

Învățare automată

— Licență, anul III, 2021-2022 —
examen „suplimentar“ I

Nume student:

Grupa:

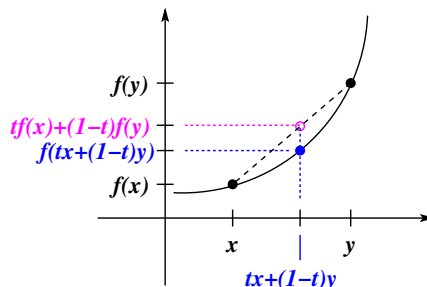
1. (Inegalitatea lui Jensen și câteva consecințe ale ei)

Liviu Ciortuz, 2019

Dacă $f : \mathbb{R} \rightarrow \mathbb{R}$ este o **funcție convexă**, atunci, conform **definiției**,¹ pentru orice $t \in [0, 1]$ și orice $x_1, x_2 \in \mathbb{R}$ urmează

$$f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2). \quad (1)$$

Dacă f este funcție strict convexă, atunci egalitatea are loc doar dacă $x_1 = x_2$.



a. Folosind definiția de mai sus, demonstrați **inegalitatea lui Jensen**:²

Pentru orice $a_i \geq 0$, $i = 1, \dots, n$ cu $\sum_i a_i = 1$ și orice $x_i \in \mathbb{R}$, $i = 1, \dots, n$, dacă f este funcție convexă,³ atunci

$$f\left(\sum_i a_i x_i\right) \leq \sum_i a_i f(x_i). \quad (2)$$

Mai general, pentru orice $a'_i \geq 0$, cu $i = 1, \dots, n$ și $\sum_i a'_i \neq 0$ avem

$$f\left(\frac{\sum_i a'_i x_i}{\sum_j a'_j}\right) \leq \frac{\sum_i a'_i f(x_i)}{\sum_j a'_j}. \quad (3)$$

Observații:

1. Dacă f este strict convexă, atunci în relațiile de mai sus egalitatea are loc doar dacă $x_1 = \dots = x_n$.
2. Evident, rezultate similare cu cele de mai sus pot fi formulate și pentru funcții concave, înlocuind în relațiile (2) și (3) semnul \leq cu \geq .

b. Demonstrați **inegalitatea mediilor** folosind inegalitatea lui Jensen:⁴

$$\frac{x_1 + x_2 + \dots + x_n}{n} \geq \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} \quad \text{pentru orice } x_i \geq 0, i = 1, \dots, n.$$

c. În contextul teoriei probabilităților, inegalitatea lui Jensen este exprimată astfel: dacă X este o variabilă aleatoare și f este o funcție convexă, atunci $f(E[X]) \leq E[f(X)]$. (Similar, dacă f este funcție concavă, atunci $f(E[X]) \geq E[f(X)]$.)

Demonstrați această inegalitate în cazul în care X este variabilă aleatoare discretă cu un număr finit de valori (adică, $|Val(X)| < \infty$).

²Johan Jensen, inginer și matematician danez (1859–1925).

³Dacă numerele a_i , cu $i = 1, \dots, n$, satisfac cele două proprietăți indicate, spunem că suma $\sum_i a_i x_i$ este o **combinație convexă** a numerelor x_1, \dots, x_n . Atunci când se folosește doar prima condiție ($a_i \geq 0$, $i = 1, \dots, n$), suma respectivă se numește **combinație canonică**. În fine, dacă se renunță la ambele condiții, suma se numește **combinație afină**. (Cf. https://en.wikipedia.org/wiki/Linear_combination.)

⁴Primul membru al **inegalității mediilor** este media aritmetică, iar cel de-al doilea este media geometrică.

2.

(Variabile aleatoare: proprietăți de bază pentru medii, varianță, covarianță)

Fie variabila aleatoare $X : \Omega \rightarrow \mathbb{R}$, cu funcția de probabilitate P .

Dacă X este variabilă aleatoare *discretă*, atunci prin definiție $P(x) \stackrel{not.}{=} P(X = x) \stackrel{not.}{=} P(\{\omega \mid X(\omega) = x\}) \geq 0$ pentru orice $x \in \mathbb{R}$, și $\sum_{x_i \in \text{Val}(X)} P(x_i) = 1$, unde $\text{Val}(X)$ este mulțimea valorilor variabilei aleatoare X .

Dacă X este variabilă aleatoare *continuă*, având funcția densitate de probabilitate p , atunci prin definiție $p(X = x) \geq 0$ pentru orice $x \in \mathbb{R}$, și $\int_{-\infty}^{\infty} p(X = x) dx = 1$ (sau, scris mai simplu: $\int_{-\infty}^{+\infty} p(x) dx = 1$).

a. Dacă X este variabilă aleatoare discretă, *media* sa se definește ca fiind numărul real $E[X] = \sum_{x_i \in \text{Val}(X)} x_i \cdot P(X = x_i)$. Dacă X este variabilă aleatoare continuă, media sa este $E[X] = \int_{-\infty}^{\infty} x \cdot p(X = x) dx$.

Arătați că pentru orice două variabile aleatoare W și Z de același tip (adică fie ambele discrete fie ambele continue), având același domeniu de definiție (Ω), avem

$$E[W + Z] = E[W] + E[Z].$$

De asemenea, demonstrați că pentru orice constantă $a \in \mathbb{R}$, are loc egalitatea

$$E[aX] = aE[X]. \quad (4)$$

Notăți că aX este o variabilă aleatoare definită pe același domeniu (Ω) ca și variabila X , cu proprietatea că $(aX)(\omega) \stackrel{def.}{=} aX(\omega)$ pentru orice $\omega \in \Omega$.

Observație: Cele două egalități de mai sus se pot combina sub o formă mai generală: pentru orice variabile aleatoare (fie toate discrete fie toate continue) X_1, \dots, X_n și pentru orice constante $a_1, \dots, a_n \in \mathbb{R}$, cu $n \geq 1$, are loc egalitatea

$$E[a_1X_1 + \dots + a_nX_n] = a_1E[X_1] + \dots + a_nE[X_n].$$

Această egalitate este cunoscută sub numele de *proprietatea de liniaritate a mediei*.

b. Fie X o variabilă aleatoare. Notăm $\bar{X} = E[X]$. *Varianța* lui X se definește ca fiind $\text{Var}(X) = E[(X - \bar{X})^2]$. Arătați că:

$$\text{Var}(X) = E[X^2] - (E[X])^2.$$

Observație importantă: Veți vedea că această proprietate este adeseori folosită (în locul definiției varianței) în diverse demonstrații care vor urma.

De asemenea, demonstrați că pentru orice constantă $a \in \mathbb{R}$, are loc egalitatea

$$\text{Var}(aX) = a^2 \text{Var}(X). \quad (5)$$

Prin urmare, în cazul varianței nu avem o proprietate de liniaritate similară cu cea din cazul mediei.

Indicație: La acest punct nu este necesar să faceți demonstrațiile separat pentru cele două cazuri, discret și respectiv continuu.

c. *Covarianța* a două variabile aleatoare X și Y care au același domeniu de definiție (Ω) se definește astfel: $\text{Cov}(X, Y) = E[(X - E[X])(Y - E[Y])]$, unde $E[X]$ este media lui X . (Este imediat faptul că noțiunea de covarianță generalizează noțiunea de varianță.)

Demonstrați egalitatea:

$$\text{Cov}(X, Y) = E[XY] - E[X] \cdot E[Y]. \quad (6)$$

Consecință imediată (din relațiile (4) și (6)):

$$\text{Cov}(aX, bY) = ab \text{Cov}(X, Y) \quad \forall a, b \in \mathbb{R}. \quad (7)$$

Observații:

1. Este imediat faptul că proprietatea (7) generalizează proprietatea (5).
2. Spre deosebire de varianță, care poate lua doar valori mai mari sau egale cu 0 (ceea ce decurge imediat din definiția de la punctul b), covarianța poate lua și valori negative. Mai mult, se poate demonstra că atunci când $\text{Var}(X) \neq 0$ și $\text{Var}(Y) \neq 0$, avem următoarele margini (una inferioară și cealaltă superioară) pentru $\text{Cov}(X, Y)$:

$$-\sqrt{\text{Var}(X) \cdot \text{Var}(Y)} \leq \text{Cov}(X, Y) \leq +\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}.$$

3. (Distribuția gaussiană unidimensională: verificarea condițiilor de definiție pentru p.d.f., calculul mediei și al varianței)

Considerăm o variabilă aleatoare X care urmează distribuția normală (gaussiană):

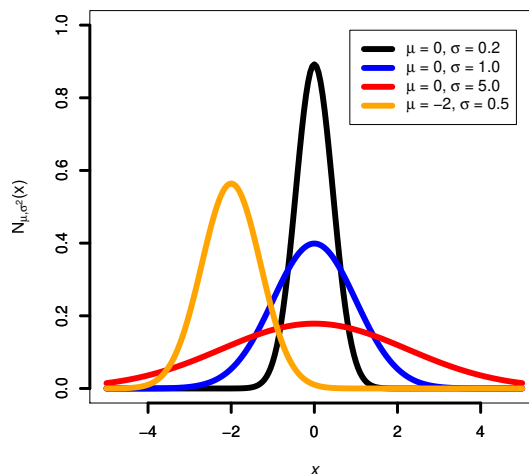
Distribuția gaussiană: p.d.f.

$$\mathcal{N}(X = x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x - \mu)^2}{2\sigma^2}}$$

unde σ poate fi orice număr real pozitiv, iar μ orice număr real.

Arătați că:

- a. \mathcal{N} este într-adevăr o funcție de densitate de probabilitate (p.d.f.), adică $\int_{-\infty}^{+\infty} \mathcal{N}(x) dx = 1$.
- b. $E[X] = \mu$.
- c. $\text{Var}[X] = \sigma^2$.



Sugestie (pentru punctul a): Pentru cazul distribuției gaussiene *standard* — la care veți face „reducere”, de la cazul *general*, printr-o schimbare liniară de variabilă —, proprietatea de demonstrat devine

$$\int_{-\infty}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 1 \text{ sau, echivalent } \int_{-\infty}^{+\infty} e^{-\frac{x^2}{2}} dx = \sqrt{2\pi}.$$

Pentru demonstrarea ultimei egalități vă recomandăm să arătați că

$$\left(\int_{-\infty}^{+\infty} e^{-\frac{x^2}{2}} dx \right)^2 = 2\pi \text{ sau, echivalent } \left(\int_{-\infty}^{+\infty} e^{-\frac{x^2}{2}} dx \right) \cdot \left(\int_{-\infty}^{+\infty} e^{-\frac{y^2}{2}} dy \right) = 2\pi,$$

deci

$$\int_{x=-\infty}^{+\infty} \int_{y=-\infty}^{+\infty} e^{-\frac{x^2+y^2}{2}} dx dy = 2\pi.$$

Ultima egalitate poate fi demonstrată prin trecerea din sistemul de coordonate cartezian în sistemul de coordonate polare. În mod concret, $(x, y) \mapsto (r, \theta)$, unde $r \in [0, +\infty)$ și $\theta \in [0, 2\pi)$, cu $x = r \cos \theta$ și $y = r \sin \theta$, ceea ce constituie o corespondență bijectivă. Vă reamintim că *regula de schimbare de variabilă* pentru cazul vectorial, formulată (aici) pentru cazul probabilist, este următoarea:⁵

Presupunem că $V = [V_1 \dots V_n]^\top \in \mathbb{R}^n$ este un vector de variabile aleatoare având funcția de densitate de probabilitate comună $f_V : \mathbb{R}^n \rightarrow \mathbb{R}$. Dacă definim un alt vector de variabile aleatoare, Z , obținut prin compunerea $Z = H \circ V$, unde $H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ este o funcție bijectivă și derivabilă [pe componente] în raport cu fiecare dintre argumentele sale, atunci Z va avea funcția de densitate de probabilitate comună $f_Z : \mathbb{R}^n \rightarrow \mathbb{R}$, unde

$$f_Z(z) = f_V(v) \cdot \left| \det \begin{pmatrix} \frac{\partial v_1}{\partial z_1} & \cdots & \frac{\partial v_1}{\partial z_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial v_n}{\partial z_1} & \cdots & \frac{\partial v_n}{\partial z_n} \end{pmatrix} \right|.$$

Matricea al cărei determinant este calculat în expresia de mai sus se numește matrice *jacobiană*, iar determinantul respectiv se numește *determinant jacobian*, de la numele matematicianului german Carl Gustav Jacob Jacobi (1804 – 1851).

4. (Distribuția exponențială și distribuția Gamma: verificarea condițiilor de definiție pentru p.d.f., calculul mediilor și al varianțelor)

a. Distribuția *exponențială* este o distribuție continuă, care are funcția densitate de probabilitate

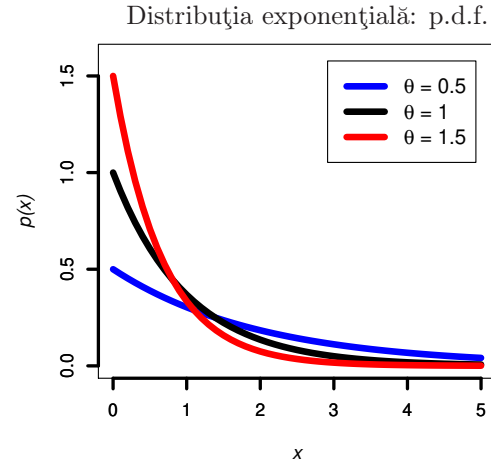
$$p(x | \theta) = \begin{cases} \theta e^{-\theta x} & \text{pentru } x \geq 0 \\ 0 & \text{pentru } x < 0. \end{cases}$$

unde $\theta > 0$ este un parametru real. Arătați mai întâi că funcția $p(\cdot)$ este într-adevăr funcție densitate de probabilitate (engl., probability density function, p.d.f.) și apoi că media și respectiv varianța distribuției exponențiale sunt $\frac{1}{\theta}$ și respectiv $\frac{1}{\theta^2}$.

b. Funcția Γ a lui Euler (1707-1783) este o generalizare în \mathbb{R}^+ a definiției numerelor factoriale din \mathbb{N} (și anume, $\Gamma(r) = (r-1)!$ pentru orice $r \in \mathbb{N}^*$).

Formula de definiție a acestei funcții este următoarea:

$$\Gamma(r) \stackrel{\text{def.}}{=} \int_0^{+\infty} t^{r-1} e^{-t} dt \text{ pentru orice } r > 0.$$



⁵Cf. *The Multivariate Gaussian Distribution*, Chuong Do, Stanford University, 2008.

Demonstrați următoarele proprietăți ale funcției Γ :

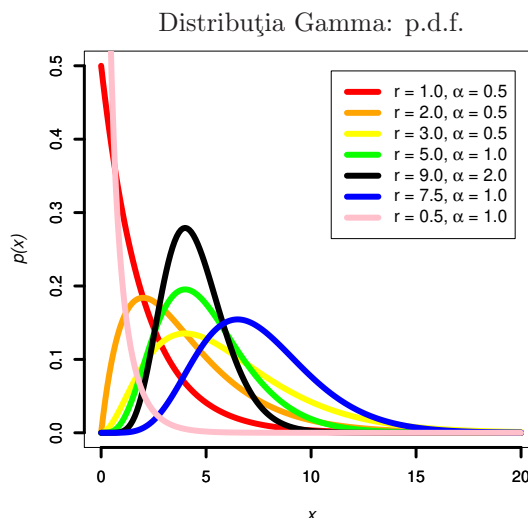
i. $\Gamma(r+1) = r\Gamma(r)$ pentru orice $r > 0$.

ii. $\Gamma(1) = 1$ și $\Gamma(1/2) = \sqrt{\pi}$.⁶

c. Distribuția *Gamma* este o distribuție continuă, de parametri $r > 0$ (care dă forma distribuției, engl., shape) și $\alpha > 0$ (numit *rata*, engl., rate), cu funcția densitate de probabilitate definită pe \mathbb{R}^+ prin expresia următoare:

$$p(x) \stackrel{\text{not.}}{=} \text{Gamma}(x|r, \alpha) \\ \stackrel{\text{def.}}{=} \frac{\alpha^r}{\Gamma(r)} x^{r-1} e^{-\alpha x} \text{ pentru } x \geq 0,$$

unde constanta de normalizare este $\frac{\alpha^r}{\Gamma(r)}$.



Demonstrați mai întâi că funcția $p(\cdot)$ este într-adevăr p.d.f. și apoi că media distribuției Gamma este $\frac{r}{\alpha}$, iar varianța ei este $\frac{r}{\alpha^2}$.

Sugestie: Puteți folosi următoarea proprietate, care este demonstrată la problema 3:

$$\int_{-\infty}^{\infty} v^2 e^{-\frac{v^2}{2}} dv = 2 \int_0^{\infty} v^2 e^{-\frac{v^2}{2}} dv = \sqrt{2\pi}. \quad (8)$$

Observație: Se poate vedea imediat că $\text{Gamma}(x|1, \alpha) = \alpha e^{-\alpha x}$ pentru orice $x \geq 0$, ceea ce corespunde [definiției funcției de densitate a] distribuției exponențiale (vedeți punctul a). Așadar, se poate spune că distribuția exponențială este membru al familiei de distribuții Gamma.

5.

(Distribuții de tip Bernoulli;
calculul verosimilității datelor;
estimarea parametrilor în sensul MLE)

CMU, 2005 fall, T. Mitchell, A. Moore, midterm, pr. 1.3

Avem două monede. Probabilitatea de apariție a stemei este θ în cazul primei monede și 2θ în cazul celei de-a doua monede.

Presupunem că aruncăm aceste două monede de mai multe ori, în mod independent una de cealaltă, și obținem rezultatele din tabelul alăturat.

Moneda	Rezultat
1	stemă
2	ban
2	ban
2	ban
2	stemă

a. Care este log-verosimilitatea acestor date în funcție de θ ?

b. Cât este estimarea / valoarea de verosimilitate maximă (engl., maximum likelihood, ML) a lui θ ?

⁶Din relația de recurență i. și din relația $\Gamma(1) = 1$, rezultă $\Gamma(r+1) = r \cdot \Gamma(r) = r \cdot (r-1) \cdot \Gamma(r-1) = \dots = r \cdot (r-1) \cdot \dots \cdot 2 \cdot 1 = r!$. Așadar, într-adevăr funcția Γ a lui Euler generalizează noțiunea de *produs factorial*.

6. (Entropie, entropie comună, entropie condițională, câștig de informație: definiții și proprietăți imediate)

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

- 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 comună a variabilelor X și Y :

$$H(X, Y) \stackrel{\text{def.}}{=} -\sum_i \sum_j P(X = x_i, Y = y_j) \log P(X = x_i, Y = y_j)$$

$$\stackrel{\text{not.}}{=} E_{X,Y}[-\log P(X, Y)].$$

- *Câștigul de informație* al variabilei X în raport cu variabila Y (sau invers), numit de asemenea *informația mutuală* a variabilelor X și Y :

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

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

Arătați că:

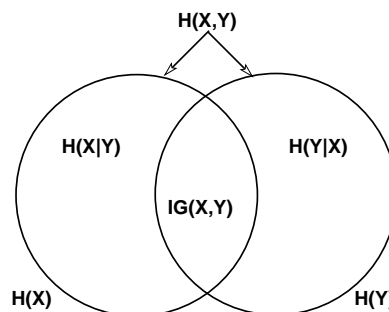
- a. $H(X) \geq 0$. În particular, $H(X) = 0$ dacă și numai dacă variabila X este constantă.

- 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 înlanțuire).

Observație: Relația precedentă, precum și relația de definiție pentru câștigul de informație sunt ilustrate în figura alăturată.



7. (Entropie, entropie condițională specifică, câștig de informație: exemplificare)

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

- a. Calculează distribuția probabilistică 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*

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

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.

8. (Câștigul de informație / informația mutuală, o aplicație: selecția de trăsături)

În tabelul următor se dă un set de opt observații / instanțe, reprezentate ca tuple 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*.

a. Calculați câștigul de informație / 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.

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

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

Comentariu: La problema 6.a am demonstrat că entropia oricărei variabile aleatoare discrete este nenegativă ($H(X) \geq 0$).⁷ 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

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

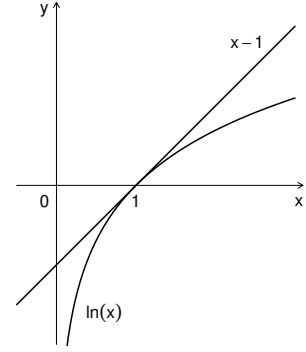
⁷Extensia acestei proprietăți la cazul variabilelor aleatoare continue este imediată.

Arătați că $H(X) \leq \log_2 n$.

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

Sugestie (2): Puteți folosi inegalitatea lui Jensen (vedeți problema 1).

Sugestie (3): Puteți folosi metoda multiplicatorilor lui Lagrange.



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

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

Conform problemei 6.c, *formula de înlănțuire* a entropiilor pentru cazul general (adică, indiferent dacă X și Y sunt sau nu independente) este:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y). \quad (9)$$

Demonstrați că dacă X și Y sunt variabile aleatoare independente discrete, atunci $H(X, Y) = H(X) + H(Y)$, și reciproc: atunci când are loc egalitatea $H(X, Y) = H(X) + H(Y)$ rezultă că variabilele X și Y sunt independente.

11. (Nenegativitatea câștigului de informație; o condiție necesară și suficientă pentru anularea lui)

Definiția câștigului de informație (sau: a *informației mutuale*) al unei variabile aleatoare X în raport cu o altă variabilă aleatoare Y este $IG(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$.⁸ La problema 15 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 comună a acestor variabile. Tot la problema 15 s-a arătat că divergența KL este întotdeauna nenegativă. În consecință, $IG(X, Y) \geq 0$ pentru orice X și Y .

a. Aici vă cerem să demonstrați inegalitatea $IG(X, Y) \geq 0$ în manieră directă, plecând de la prima definiție dată mai sus, fără a mai apela la divergența Kullback-Leibler.

b. Arătați tot într-o manieră directă că $IG(X, Y) = 0$ dacă și numai dacă X și Y sunt independente. (Într-o manieră indirectă, acest rezultat a fost demonstrat la problema 15.c.)

Sugestie: Puteți folosi următoarea formă [particulară] a inegalității lui Jensen:⁹

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

⁸Vedeți problema 6.

⁹Vedeți problema 1.

unde baza logaritmului se consideră supraunitară, $a_i \geq 0$ pentru $i = 1, \dots, n$ și $\sum_{i=1}^n a_i = 1$.¹⁰ Egalitatea are loc dacă și numai dacă $x_1 = x_2 = \dots = x_n$.

12. (Entropia distribuției continue uniforme și a distribuției gaussiene uni-dimensionale)

- a. Calculați entropia distribuției continue uniforme definite pe intervalul $[a, b]$.
b. Calculați entropia distribuției gaussiene unidimensionale de parametri $\mu \in \mathbb{R}$ și $\sigma^2 \in \mathbb{R}_+$, pentru care funcția de densitate de probabilitate (p.d.f.) este:

$$\mathcal{N}(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \text{ pentru } x \in \mathbb{R}.$$

13. (Calcularea entropiei unor variabile aleatoare continue: cazul distribuției exponențiale și al distribuției Gamma)

Pentru o variabilă aleatoare X care urmează o distribuție continuă cu funcția densitate de probabilitate (p.d.f.) p , entropia se definește astfel:

$$H(X) = \int_{-\infty}^{+\infty} p(x) \log_2 \frac{1}{p(x)} dx \stackrel{p(x) \neq 0}{=} - \int_{-\infty}^{+\infty} p(x) \log_2 p(x) dx.$$

Indicație: Dacă $p(x) = 0$, veți presupune că $-p(x) \log_2 p(x) = 0$.

- a. Calculați entropia distribuției continue exponențiale de parametru $\lambda > 0$. Vă reamintim că definiția p.d.f.-ului acestei distribuții este următoarea:¹¹

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

- b. Calculați entropia distribuției Gamma, pentru care funcția de densitate de probabilitate este:¹²

$$p(x) \stackrel{\text{not.}}{=} \text{Gamma}(x|r, \alpha) \stackrel{\text{def.}}{=} \frac{\alpha^r}{\Gamma(r)} x^{r-1} e^{-\alpha x} \text{ pentru } x \geq 0,$$

cu $r > 0$ (forma), $\alpha > 0$ (rata) și $\Gamma(r) \stackrel{\text{def.}}{=} \int_0^{+\infty} x^{r-1} e^{-x} dx$ pentru orice $r > 0$ (funcția lui Euler). (Vedeți ex. 4.b.)

Indicație:

Este posibil să fie nevoie să folosiți o proprietate specială, numită *regula de derivare sub semnul de integrală* (engl., differentiation under the integral sign), care constituie obiectul pentru următoarea

Teoremă: Dacă $f(r, t)$ este o funcție cu valori reale, continuă și derivabilă în raport

¹⁰Avantajul la această problemă, comparativ cu problema 15.a, este că aici se lucrează cu o singură distribuție (p), nu cu două distribuții (p și q). Totuși, demonstrația de aici va fi mai laborioasă.

¹¹La ex. 4.a puteți vedea graficul acestei funcții de densitate pentru câteva valori ale parametrului (λ).

¹²La ex. 4.c puteți vedea graficul acestei funcții de densitate pentru câteva valori ale parametrilor r și α .

cu r pe intervalul (a, b) , iar derivata parțială $\frac{\partial}{\partial r} f(r, t)$ este de asemenea continuă pe intervalul (a, b) , atunci¹³

$$\int_a^b \frac{\partial}{\partial r} f(r, t) dt = \frac{\partial}{\partial r} \int_a^b f(r, t) dt.$$

Menționăm că regula aceasta, datorată lui Gottfried Leibnitz (1646-1716), se generalizează — sub o formă semnificativ mai elaborată — la cazul când limitele de integrare a și b depind de r , adică sunt de forma $a(r)$ și respectiv $b(r)$, iar aceste două funcții sunt continue și au derivatele continue.

14. (Redescoperirea definiției entropiei pornind de la un set de *proprietăți dezirabile* ale ei)

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 $H(X) = -\sum_i p_i \log p_i$. Înșă 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.

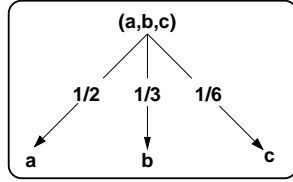
Scopul acestui exercițiu este de a arăta că orice funcție $\psi_n(p_1, \dots, p_n)$ care satisface trei proprietăți dezirabile pentru entropie este în mod necesar de forma $-K \sum_i p_i \log p_i$ unde K este o constantă reală pozitivă. Iată care sunt aceste *proprietăți*:¹⁴

- 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 ca o sumă ponderată a entropiilor calculate la fiecare stadiu / alegere.
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.)

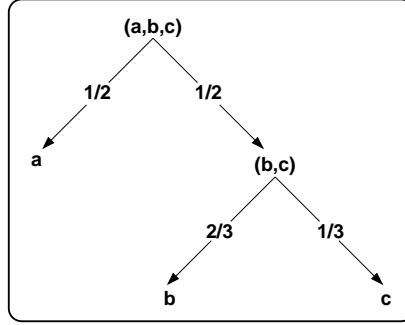
Din punct de vedere formal, proprietatea A3 impune ca, pe acest exemplu, $\psi_3(1/2, 1/3, 1/6)$ să fie egal cu $\psi_2(1/2, 1/2) + 1/2 \cdot \psi_2(2/3, 1/3)$.

¹³Pentru demonstrație, vedeți „Proof of basic form“ pe site-ul https://en.wikipedia.org/wiki/Leibniz_integral_rule#Alternative_Proof_of_General_Form_with_Variable_Limits_using_the_Chain_Rule (accesat la data de 15.12.2020).

¹⁴LC: Deși nu se specifică în enunțul original al problemei, este necesar / natural să considerăm și proprietatea următoare: [A0.] $\psi_n(p_1, \dots, p_n) \geq 0$ pentru orice $n \in \mathbb{N}^*$ și orice $p_1, \dots, p_n \in [0, 1]$ astfel încât $\sum_i p_i = 1$, pentru că ψ_n este văzută că măsură a *dezordinii*; de asemenea, $\psi_1(1) = 0$ pentru că în acest caz particular nu există niciun fel de dezordine.



Encoding 1



Encoding 2

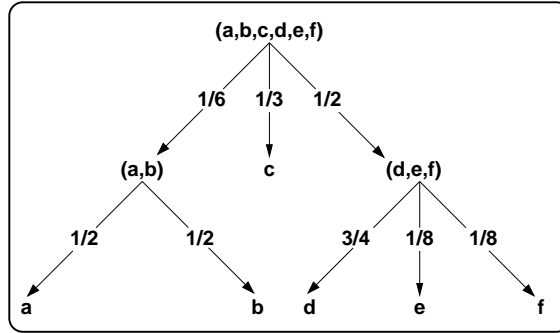
Așadar, în acest exercițiu vi se cere să arătați că dacă o funcție de n variabile $\psi_n(p_1, \dots, p_n)$ satisface proprietățile A1, A2 și A3 de mai sus, atunci există $K \in \mathbb{R}^+$ astfel încât $\psi_n(p_1, \dots, p_n) = -K \sum_i p_i \log p_i$ unde $K \in \mathbb{R}_+$ este o constantă.¹⁵

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 $A(n) \stackrel{not.}{=} \psi(1/n, 1/n, \dots, 1/n)$.

c. Arătați că

$$A(s^m) = m A(s) \text{ pentru orice } s, m \in \mathbb{N}^*. \quad (10)$$

La punctele d – g de mai jos, pentru orice pereche de numere $s, m \in \mathbb{N}^*$ fixate, cu $s > 1$, vom considera $t, n \in \mathbb{N}^*$ astfel încât

$$s^m \leq t^n \leq s^{m+1}. \quad (11)$$

¹⁵În demonstrație, vom vedea, va rezulta $K = \frac{1}{\log s} \psi_s\left(\frac{1}{s}, \dots, \frac{1}{s}\right)$ pentru un număr oarecare $s \in \mathbb{N}^* \setminus \{1\}$, fixat.

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

$$\left| \frac{m}{n} - \frac{\log t}{\log s} \right| \leq \frac{1}{n}. \quad (12)$$

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

f. Combinând ultima inegalitate de mai sus cu egalitatea (10), avem $A(s^m) \leq A(t^n) \leq A(s^{m+1}) \Rightarrow m A(s) \leq n A(t) \leq (m+1) A(s)$. Verificați că

$$\left| \frac{m}{n} - \frac{A(t)}{A(s)} \right| \leq \frac{1}{n} \text{ pentru } s \neq 1. \quad (13)$$

g. Combinând inegalitățile (12) și (13), arătați că

$$\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n} \text{ pentru } s \neq 1 \quad (14)$$

și, în consecință

$$A(t) = K \log t \text{ cu } K > 0 \text{ (din cauza proprietății A2)}. \quad (15)$$

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 $\psi(p_1, \dots, p_k)$ cu $p_i \in \mathbb{Q}$ pentru $i = 1, \dots, k$.

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

Conform egalității (15), avem $A(N) = K \log N$. Folosind proprietatea A3 și *codificarea* în doi pași propusă mai sus, rezultă că

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

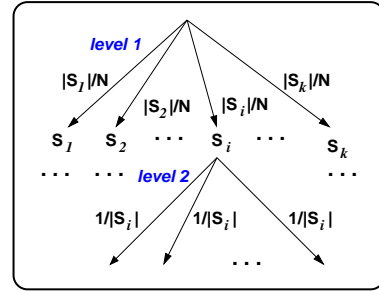
Așadar,

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

Prin urmare,

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

Putem trata în sfârșit și cazul $p_i \in \mathbb{R}$, cu $\sum_i p_i = 1$.¹⁶ Considerăm pentru fiecare p_i , $i = 1, \dots, k-1$ un șir de aproximări succesive $\{q_i^{(n)}\}_{n \geq 1} \subset \mathbb{Q}$ astfel încât $\lim_{n \rightarrow \infty} q_i^{(n)} = p_i$.



¹⁶Cf. Ștefan Bălăucă, student, FII, 2020f.

(De exemplu, putem defini $q_i^{(n)} = \text{numărul format luând în considerare doar din primele } n \text{ zecimale ale lui } p_i$.)

Definim de asemenea şirul $\{q_k^{(n)}\}_{n \geq 1} \subset \mathbb{Q}$ astfel: $q_k^{(n)} = 1 - \sum_{i=1}^{k-1} q_i^{(n)}$. Trecând la limită în această relaţie de definiţie, obţinem:

$$\lim_{n \rightarrow \infty} q_k^{(n)} = 1 - \sum_{i=1}^{k-1} \lim_{n \rightarrow \infty} q_i^{(n)} = 1 - \sum_{i=1}^{k-1} p_i = p_k.$$

Aşadar, avem $\lim_{n \rightarrow \infty} q_i^{(n)} = p_i$, $\forall i = 1, \dots, k$. Deoarece $q_i^{(n)} \in \mathbb{Q}$, pentru $i = 1, \dots, k$, cu $k \geq 1$, este valabilă relaţia

$$\psi(q_1^{(n)}, \dots, q_k^{(n)}) = -K \sum_i q_i^{(n)} \ln q_i^{(n)}, \quad \forall n \geq 1.$$

Trecând la limită în această relaţie, vom avea:

$$\lim_{n \rightarrow \infty} \psi(q_1^{(n)}, \dots, q_k^{(n)}) = -K \sum_i \lim_{n \rightarrow \infty} q_i^{(n)} \ln q_i^{(n)} = -K \sum_i p_i \ln p_i. \quad (16)$$

Datorită proprietăţii A1, funcţia ψ este continuă în fiecare dintre argumentele ei, astfel că

$$\lim_{n \rightarrow \infty} \psi(q_1^{(n)}, \dots, q_k^{(n)}) = \psi(\lim_{n \rightarrow \infty} q_1^{(n)}, \dots, \lim_{n \rightarrow \infty} q_k^{(n)}) = \psi(p_1, \dots, p_k). \quad (17)$$

Din relaţiile (16) şi (17) putem concluziona că

$$\psi(p_1, \dots, p_k) = -K \sum_i p_i \ln p_i, \quad \text{pentru } \forall p_i \in [0, 1], \text{ a.î. } \sum_i p_i = 1.$$

15. (Entropia relativă: definiţie şi proprietăţi elementare; exprimarea câştigului de informaţie cu ajutorul entropiei relative)

Entropia relativă sau divergenţa Kullback-Leibler (KL) a unei distribuţii p în raport cu o altă distribuţie q — ambele distribuţii fiind discrete — se defineşte astfel:

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

Din perspectiva teoriei informaţiei, divergenţa KL specifică numărul de *biţi adiţionali* care sunt necesari în medie pentru a transmite valorile variabilei X atunci când presupunem că aceste valori sunt distribuite conform distribuţiei („model“) q , dar în realitate ele urmează o altă distribuţie, p .¹⁷

¹⁷ *Atenţie:* Divergenţa KL nu este o măsură de *distanţă* între două distribuţii probabiliste, fiindcă în general ea nu este simetrică ($KL(p||q) \neq KL(q||p)$) şi nici nu satisface inegalitatea triunghiului. Pentru „simetrizare“, se consideră $M(p, q) = \frac{1}{2}(p+q)$, apoi se defineşte funcţia $JSD(p||q) = \frac{1}{2}KL(p||M) + \frac{1}{2}KL(q||M)$, care se numeşte *divergenţa Jensen-Shannon*. În sfârşit, se poate arăta că $\sqrt{JSD(p||q)}$ defineşte o măsură de distanţă (metrică), adică este nenegativă, simetrică, implică identitatea indiscernabililor şi satisface inegalitatea triunghiului; ea este numită *distanţa Jensen-Shannon*.

Variaţia informaţiei, definită prin

$$VI(X, Y) \stackrel{\text{def.}}{=} H(X, Y) - IG(X, Y) = H(X) + H(Y) - 2IG(X, Y) = H(X | Y) + H(Y | X),$$

este de asemenea o măsură de distanţă.

a. Demonstrați inegalitatea $KL(p||q) \geq 0$ și apoi arătați că egalitatea are loc dacă și numai dacă $p = q$.¹⁸

Indicație:

Pentru a demonstra punctul acesta puteți folosi *inegalitatea lui Jensen*.¹⁹

Dacă $f : \mathbb{R} \rightarrow \mathbb{R}$ este o *funcție convexă*, atunci pentru orice $a_i \geq 0$, $i = 1, \dots, n$ cu $\sum_i a_i = 1$ și orice $x_i \in \mathbb{R}$, $i = 1, \dots, n$, avem $f(\sum_i a_i x_i) \leq \sum_i a_i f(x_i)$. Dacă f este strict convexă, atunci egalitatea are loc doar dacă $x_1 = \dots = x_n$. Pentru funcții concave, semnul inegalității este \geq .

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

$$\begin{aligned} IG(X, Y) &\stackrel{def.}{=} KL(p_{X,Y} || (p_X p_Y)) = - \sum_x \sum_y p_{X,Y}(x, y) \log \left(\frac{p_X(x)p_Y(y)}{p_{X,Y}(x, y)} \right) \\ &\stackrel{not.}{=} - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right) \end{aligned}$$

Arătați că această nouă definiție a câștigului de informație este echivalentă cu definiția dată anterior (vedeți problema 6). Cu alte cuvinte, demonstrați egalitatea

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

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

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

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. O consecință imediată a punctelor a și b este faptul că $IG(X, Y) \geq 0$ (deci $H(X) \geq H(X|Y)$ și $H(Y) \geq H(Y|X)$) pentru orice variabile aleatoare discrete X și Y . Folosind din nou rezultatele de la punctele a și b, arătați că $IG(X, Y) = 0$ dacă și numai dacă X și Y sunt independente.

16. (Cross-entropie: definiție, o proprietate (nenegativitatea) și un exemplu simplu de calculare a valorii cross-entropiei)

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:

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

¹⁸Mai general, $KL(p||q)$ este cu atât mai mică cu cât „asemănarea” dintre distribuțiile p și q este mai mare.

¹⁹Vedeți problema 1.

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

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

Observație: Ținând cont de definiția entropiei relative (cunoscută și sub numele de divergența Kullback-Leibler), vedeți pr. 15, putem scrie:

$$KL(p||q) = CH(p, q) - H(p).$$

Cross-entropia — ca și entropia relativă; vedeți problema 15 —, spre deosebire de entropia comună, nu este simetrică în raport cu cele două distribuții / argumente: în general, $CH(p, q) \neq CH(q, p)$.

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

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.

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

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:

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

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.

Indicații:

1. Dacă în inegalitatea dată se trece termenul din partea stângă în partea dreaptă, obținem $0 \leq \sum_{i=1}^n p_i \log_2 p_i - \sum_{i=1}^n p_i \log_2 q_i \Leftrightarrow 0 \leq -\sum_{i=1}^n p_i \log_2 \frac{q_i}{p_i}$. Puteți face legătura dintre expresia din partea dreaptă a acestei ultime inegalități și definiția *entropiei relative* (numită de asemenea *divergența Kullback-Leibler*; vedeți problema 15) și apoi să folosiți proprietățile entropiei relative.
2. Pentru a demonstra într-o manieră mai directă inegalitatea lui Gibbs, puteți folosi inegalitatea lui Jensen (vedeți problema 1).

18. (O „regulă de înlănțuire“ pentru entropia relativă / divergența KL)

Vă readucem aminte — vedeți problema 15 — că entropia relativă sau divergența Kullback-Leibler (KL) dintre două distribuții discrete $P(X)$ și $Q(X)$ este definită astfel:²⁰

$$KL(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)}.$$

Pentru conveniență, vom presupune că $P(x) > 0$ și $Q(x) > 0$ pentru orice x . Uneori în cele ce urmează vom scrie $KL(P||Q)$ sub forma $KL(P(X)||Q(X))$.

Entropia relativă (sau, divergența KL) dintre două distribuții condiționale $P(X|Y)$ și $Q(X|Y)$ se definește în felul următor:

$$KL(P(X|Y)||Q(X|Y)) = \sum_y P(y) \left(\sum_x P(x|y) \log \frac{P(x|y)}{Q(x|y)} \right).$$

Aceasta poate fi văzută ca fiind media [probabilistă a] divergenței KL dintre distribuțiile condiționale corespunzătoare pentru x (adică, dintre $P(X|Y=y)$ și $Q(X|Y=y)$), media fiind calculată în raport cu valorile y ale variabilei aleatoare Y .

Demonstrați că are loc următoarea proprietate (de tip „regulă de înlănțuire“) pentru divergența KL:

$$KL(P(X,Y)||Q(X,Y)) = KL(P(X)||Q(X)) + KL(P(Y|X)||Q(Y|X)). \quad (18)$$

19. (Echivalența dintre minimizarea entropiei relative (divergența KL) și maximizarea funcției de log-verosimilitate)

Presupunem că ni se dă o problemă de estimare a parametrilor unei distribuții probabiliste și că dispunem de setul de date de antrenament $\{x_i; i = 1, \dots, m\}$. Considerăm distribuția empirică $\hat{P}(x) = \frac{1}{m}$, adică distribuția uniformă peste setul de date de antrenament. Așadar, a face eșantionare (engl., sampling) cu distribuția empirică înseamnă a alege în mod aleatoriu un exemplu din setul de antrenament.

Presupunem că avem o anumită familie de distribuții P_θ , indexată după parametrul θ . (Dacă doriți, puteți considera $P_\theta(x)$ ca notație alternativă pentru $P(x; \theta)$.)

Demonstrați că a găsi valoarea de verosimilitate maximă (engl., maximum likelihood estimate) pentru parametrul θ este echivalent cu a găsi cea distribuție P_θ pentru

²⁰ Atunci când P și Q sunt p.d.f.-uri (funcții de densitate de probabilitate) pentru variabile aleatoare continue, dacă în formula de mai sus înlocuim simbolul de sumare cu integrala, proprietatea enunțată în acest exercițiu rămâne valabilă. Totuși, din motive de simplitate, aici vom lucra doar cu funcții masă de probabilitate / distribuții discrete.

care entropia relativă (divergența KL) față de distribuția empirică \hat{P} este minimă. Așadar, vă cerem să demonstrați relația următoare:

$$\arg \min_{\theta} KL(\hat{P}||P_{\theta}) = \arg \max_{\theta} \sum_{i=1}^m \ln P_{\theta}(x_i).$$

Comentariu: Ne vom referi acum la relația dintre proprietatea enunțată în acest exercițiu pe de o parte și estimarea parametrilor făcută de algoritmul de clasificare binară Bayes Naiv cu atribute de tip Bernoulli de cealaltă parte (vedeți capitolul de *Clasificare bayesiană*). În modelul de clasificare Bayes Naiv se presupune că P_{θ} este de forma următoare: $P_{\theta}(x, y) = p(y) \prod_{i=1}^n p(x_i|y)$. Conform regulei de înlănțuire pentru divergența KL (vedeți relația (18) de la ex. 18), avem:

$$KL(\hat{P}||P_{\theta}) = KL(\hat{P}(y)||p(y)) + \sum_{i=1}^n KL(\hat{P}(x_i|y)||p(x_i|y)).$$

Această relație arată că problema găsirii maximumului verosimilității / minimumului divergenței KL (pentru a estima valorile parametrilor) se descompune în $2n + 1$ probleme de optimizare independente: una pentru distribuția a priori a clasei, $p(y)$, și câte una pentru fiecare distribuție condițională $p(x_i|y)$, corespunzătoare fiecărei trăsături x_i în raport cu fiecare din cele două valori posibile pentru eticheta y . Concret, găsirea estimării de verosimilitate maximă pentru fiecare dintre aceste probleme în mod individual rezultă în maximizarea verosimilității distribuției comune. (O observație similară se poate formula și în legătură cu estimarea parametrilor rețelelor bayesiene (engl., Bayesian networks.)

20. (Proprietăți ale entropiei: Adevărat sau Fals?)

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

21. (Cât de multe date de antrenament necesită algoritmul Bayes Naiv vs. algoritmul Bayes Optimal?
[LC: complexitatea la eșantionare])

Unul dintre motivele pentru care folosim clasificatorul Bayes Naiv este faptul că el necesită mult mai puține date de antrenament (în vederea estimării parametrilor) decât clasificatorul Bayes Optimal.

Acest exercițiu te va ajuta să înțelegi cât de importantă este această diferență dintre cei doi algoritmi.

Presupunem că o *observație* / *instanță* este o valoare generată în mod aleatoriu de către vectorul de variabile aleatoare $\bar{X} = (X_1, \dots, X_{d-1}, X_d)$, unde fiecare X_i este o variabilă aleatoare urmând distribuția probabilistică Bernoulli de parametru $p = 0.5$. Considerăm X_1, \dots, X_{d-1} variabilele de intrare, iar $X_d = Y$ variabila de ieșire.

Pentru a estima în sensul verosimilității maxime (MLE) parametrii clasificatorului Bayes Optimal, avem nevoie să *observăm* / *întâlnim* fiecare valoare a lui \bar{X} de un număr rezonabil de ori. Similar, pentru a antrena clasificatorul Bayes Naiv avem

nevoie să întâlnim fiecare valoare a fiecărei variabile X_i ($i = \overline{1, d}$) de un număr rezonabil de ori.

Ne întrebăm cât de multe observații sunt necesare (a fi generate) pentru ca fiecare valoare a variabilei comune \bar{X} în cazul algoritmului Bayes Optimal, și respectiv fiecare valoare a variabilelor X_i ($i = \overline{1, d}$) în cazul Bayes Naiv să fie întâlnită cel puțin o dată. (În practică este nevoie de mult mai multe observații, dar în acest exercițiu ne limităm la câte o singură observație pentru fiecare valoare în parte.)

Indicație: La rezolvarea punctelor de mai jos vă sugerăm să folosiți următoarele două inegalități:

- pentru orice evenimente E_1, \dots, E_n , avem $P(E_1 \cup \dots \cup E_n) \leq \sum_{i=1}^n P(E_i)$.²¹
- $(1 - \frac{1}{e})^k \leq \frac{1}{e}$ pentru orice $k \geq 1$, unde $e \approx 2.71828$ este baza logaritmului natural.

a. Începem cu algoritmul Bayes Naiv. Fie $i \in \{1, \dots, d\}$ fixat. Arătați că dacă s-au făcut N observații (având forma $\bar{x}_j = (x_1^j, \dots, x_{d-1}^j, x_d^j)$ cu $j = 1, \dots, N$), atunci probabilitatea să nu fi întâlnit ambele valori ale variabilei X_i este $\frac{1}{2^{N-1}}$. (Observați că această fracție reprezintă un număr foarte mic atunci când N este suficient de mare.)

b. Fie $\varepsilon > 0$ fixat. Folosind prima inegalitate din *indicația* de mai sus, arătați că dacă au fost făcute câte $N_{NB} = 1 + \log_2 \frac{d}{\varepsilon}$ observații de aceeași formă ca mai sus, atunci probabilitatea să nu fi întâlnit ambele valori pentru fiecare dintre variabilele X_i ($i = \overline{1, d}$) este mai mică sau egală cu ε .

c. Acum trecem la algoritmul Bayes Optimal. Fie \bar{x} o instanță (fixată) a variabilei comune \bar{X} . Folosind a doua inegalitate din *indicația* de mai sus, arătați că dacă s-au făcut N observații (fiecare observație implicând simultan toate variabilele X_i cu $i = \overline{1, d}$), atunci probabilitatea ca să nu se fi întâlnit niciodată \bar{x} este mai mică sau egală cu $e^{-\frac{N}{2^d}}$.

d. Arătați că dacă au fost făcute cel puțin $N_{JB} = 2^d \ln \frac{2^d}{\varepsilon}$ observații, atunci probabilitatea ca să nu se fi întâlnit toate instanțele variabilei comune \bar{X} este mai mică sau egală cu ε .

e. Dacă se fixează $\varepsilon = 0.1$, calculați valorile N_{NB} și N_{JB} pentru $d = 2$, $d = 5$ și $d = 10$. Ce concluzie puteți trage? (Altfel spus, cum interpretați rezultatele?)

Observații:

1. Știm că pentru clasificatorul Bayes Naiv, este necesar să estimăm din date probabilitățile $P(Y = y)$ — adică, $P(X_d = x_d)$ — și $P(X_i = x_i | Y = y)$, pentru $i = 1, \dots, d-1$. Putem considera că este suficient să întâlnim cu o probabilitate de cel puțin $1 - \varepsilon$ toate valorile y , precum și perechile de forma (x_i, y) , cu $i = 1, \dots, d-1$. În mod implicit, problema noastră simplifică și mai mult cerințele, considerând că este suficient să găsim cu probabilitate de cel puțin $1 - \varepsilon$ toate valorile x_i , pentru $i = 1, \dots, d$ (subînțelegând că atunci vor apărea în date, cu probabilități semnificative, atât valorile y cât și fiecare dintre combinațiile (x_i, y) , cu $i = 1, \dots, d-1$).

2. Pentru clasificatorul Bayes Optimal, în mod similar, este necesar să estimăm probabilitățile $P(Y = y)$ și $P(X_1 = x_1, \dots, X_{d-1} = x_{d-1} | Y = y)$, ceea ce, evident, va permite calcularea probabilităților de forma $P(x_1, \dots, x_{d-1}, y)$. Problema noastră

²¹ Aceasta se numește *proprietatea de subaditivitate* a probabilităților.

consideră în mod implicit că este suficient să găsim cu probabilitate de cel puțin $1 - \varepsilon$ toate combinațiile (x_1, \dots, x_{d-1}, y) .

22. (Algoritmul Bayes Naiv: raportul cu regresia logistică și natura separatorului decizional; cazul când variabilele de intrare sunt de tip boolean)

a. [Bayes Naiv și Regresia Logistică: relația dintre regulile de decizie]

Fie Y o variabilă aleatoare Bernoulli, iar $X = (X_1, \dots, X_d)$ un vector de variabile booleene. Demonstrați că distribuția condițională $P(Y|X)$ are forma funcției logistice de argument $z = -(w_0 + w_1X_1 + \dots + w_dX_d)$, cu parametri $w_0, w_1, \dots, w_d \in \mathbb{R}$,²² adică

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^d w_i X_i)}$$

și, prin urmare

$$P(Y = 0|X) = \frac{\exp(w_0 + \sum_{i=1}^d w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^d w_i X_i)}.$$

Vă reamintim că *funcția logistică* (sau *sigmoidală*) este definită prin expresia $\sigma(z) = 1/(1 + e^{-z})$ pentru orice $z \in \mathbb{R}$.

Comentariu:

Regresia logistică (și, mai general, *clasificatorii „discriminativi“*) învață [în mod] direct parametri distribuției $P(Y|X)$,²³ pe când algoritmul Bayes Naiv (și, mai general, *clasificatorii „generativi“*) învață [parametrii pentru] distribuțiile $P(X|Y)$ și $P(Y)$, cu ajutorul cărora va calcula apoi $P(Y|X)$ și cea mai probabilă valoare pentru Y (atunci când X are o valoare fixată / dată). Vom spune că regresia logistică este corespondentul „discriminativ“ al clasificatorului „generativ“ Bayes Naiv.

Indicații:

1. Vom introduce o *notație* simplă, care ne va fi de folos în continuare. Întrucât variabilele X_i sunt booleene, odată fixată o valoare y_k pentru variabila Y , vom avea nevoie de un singur parametru pentru a defini distribuția condițională $P(X_i|Y = y_k)$, pentru fiecare $i = 1, \dots, d$. Așadar, vom desemna cu θ_{i1} probabilitatea $P(X_i = 1|Y = 1)$ și, prin urmare $P(X_i = 0|Y = 1) = 1 - \theta_{i1}$. În mod similar, vom desemna cu θ_{i0} probabilitatea $P(X_i = 1|Y = 0)$.

2. Remarcați că odată ce am introdus notațiile de mai sus, vom putea scrie $P(X_i|Y = 1)$ după cum urmează:

$$P(X_i|Y = 1) = \theta_{i1}^{X_i} (1 - \theta_{i1})^{(1-X_i)}, \quad (19)$$

bineînțeles, cu excepția cazurilor când $\theta_{i1} = 0$ și $X_i = 0$, respectiv $\theta_{i1} = 1$ și $X_i = 1$. Observați că atunci când X_i are valoarea 1, cel de-al doilea factor din partea dreaptă a egalității (19) este 1, pentru că exponentul lui este zero. Deci $P(X_i|Y = 1) = \theta_{i1}^{X_i} = \theta_{i1}$ pentru $X_i = 1$. În mod similar, atunci când $X_i = 0$ primul factor este egal cu 1, pentru că exponentul lui este zero. Deci $P(X_i|Y = 1) = (1 - \theta_{i1})^{1-X_i} = 1 - \theta_{i1}$ pentru $X_i = 0$.

²²LC: Prin urmare, *separatorul decizional* (sau, *granița de decizie*) pentru algoritmul Bayes Naiv este — într-o astfel de situație — liniar (în funcție de argumentele X_1, \dots, X_d). Ecuația separatorului decizional va fi $w_0 + w_1X_1 + \dots + w_dX_d = 0$.

²³LC: Parametrii distribuției $P(Y|X)$ sunt în acest caz $w_i \in \mathbb{R}$, cu $i = 0, 1, \dots, d$, iar învățarea lor se face prin maximizarea funcției de verosimilitate $\mathcal{L}(w) \stackrel{\text{not.}}{=} P(D|w)$, unde D este setul de date de antrenament. La rândul ei, maximizarea aceasta se realizează prin aplicarea unei metode de optimizare, de exemplu metoda gradientului ascendent sau metoda lui Newton.

b. [Relaxarea presupuziției de independență condițională]

Pentru a putea exprima interacțiunile dintre trăsături, modelul regresiei logistice poate fi extins cu niște termeni suplimentari. De exemplu, putem adăuga un termen care să exprime dependența dintre trăsăturile X_1 și X_2 :

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + w_{1,2}X_1X_2 + \sum_{i=1}^d w_iX_i)}.$$

În mod similar, presupuziția de independență condițională asumată de către algoritmul Bayes Naiv poate fi relaxată astfel încât trăsăturile X_1 și X_2 să nu mai trebuiască să satisfacă independența condițională. Așadar, vom putea scrie:

$$P(Y|X) = \frac{P(Y) P(X_1, X_2|Y) \prod_{i=3}^d P(X_i|Y)}{P(X)}.$$

Demonstrați că în acest caz distribuția $P(Y|X)$ are aceeași formă ca și modelul de regresie logistică augmentat cu un termen suplimentar, care exprimă dependența dintre X_1 și X_2 (și, în acest fel, modelul extins al regresiei logistice rămâne corespundentul discriminativ al clasificatorului nostru generativ).

Indicații:

3. De data aceasta o altă notație simplă ne va ajuta. Vom avea nevoie de mai mulți parametri decât la punctul a pentru a defini distribuția comună $P(X_1, X_2|Y)$. Așa că vom nota $\beta_{ijk} = P(X_1 = i, X_2 = j|Y = k)$, pentru fiecare combinație posibilă de valori pentru indicii i, j și k .

4. Această nouă notație poate fi folosită acum pentru a exprima probabilitatea $P(X_1, X_2|Y = k)$ după cum urmează:

$$P(X_1, X_2|Y = k) = (\beta_{11k})^{X_1X_2} (\beta_{10k})^{X_1(1-X_2)} (\beta_{01k})^{(1-X_1)X_2} (\beta_{00k})^{(1-X_1)(1-X_2)} \quad (20)$$

pentru $k \in \{0, 1\}$, cu excepția următoarelor cazuri: *i.* $\beta_{11k} = 0$ și $X_1X_2 = 0$, *ii.* $\beta_{10k} = 0$ și $X_1(1-X_2) = 0$, *iii.* $\beta_{01k} = 0$ și $(1-X_1)X_2 = 0$ și *iv.* $\beta_{00k} = 0$ și $(1-X_1)(1-X_2) = 0$.

23.

(Asupra folosirii algoritmului k -NN în spații (\mathbb{R}^p) de dimensiune (p) mare: un avertisment: „blestemul marilor dimensiuni“)

Considerăm punctele x_1, x_2, \dots, x_n distribuite în mod independent și uniform într-o sferă (notată cu B) care are raza egală cu unitatea²⁴ și centrul în O , originea spațiului \mathbb{R}^p . Așadar, $B = \{x : \|x\|^2 \leq 1\} \subset \mathbb{R}^p$, unde $\|x\| = \sqrt{x \cdot x}$, iar operatorul \cdot desemnează produsul scalar din \mathbb{R}^p .

În această problemă veți studia „mărimea“ vecinătății de tip 1-NN pentru originea O și cum anume variază ea în raport cu dimensiunea p . În acest fel, veți putea vedea care sunt dezavantajele folosirii algoritmului k -NN într-un spațiu de dimensiune mare.

Din punct de vedere formal, „mărimea“ menționată mai sus va fi identificată cu d^* , distanța de la O la cel mai apropiat vecin din mulțimea $\{x_1, x_2, \dots, x_n\}$:

$$d^* \stackrel{not.}{=} \min_{1 \leq i \leq n} \|x_i\|.$$

Observație: Din moment ce eșantionul $\{x_1, x_2, \dots, x_n\}$ este generat în mod aleatoriu, distanța d^* poate fi văzută ca fiind [produsă de către] o variabilă aleatoare.

²⁴Termenul folosit în limba engleză pentru o astfel de sferă este *unit ball*.

a. În cazul particular $p = 1$, calculați expresia *funcției de distribuție cumulativă*²⁵ a lui d^* (văzută ca variabilă aleatoare), și anume $P(d^* \leq t)$ pentru $t \in [0, 1]$.

b. Determinați expresia *funcției de distribuție cumulativă* (c.d.f.) a lui d^* în cazul general, adică pentru $p \in \{1, 2, 3, \dots\}$.

Sugestie: Puteți folosi următoarea formulă pentru volumul unei sfere de rază r din \mathbb{R}^p :

$$V_p(r) = \frac{(r\sqrt{\pi})^p}{\Gamma\left(\frac{p}{2} + 1\right)},$$

unde Γ reprezintă funcția Gamma a lui Euler, care are proprietățile:

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(1) = 1, \quad \text{iar } \Gamma(x+1) = x\Gamma(x) \text{ pentru } x > 0.^{26}$$

c. Care este *mediana* variabilei aleatoare d^* (adică, valoarea lui t pentru care $P(d^* \leq t) = 1/2$)? Va trebui ca răspunsul să fie formulat în funcție de n și p (dimensiunea eșantionului și, respectiv, dimensiunea spațiului din care se face extragerea instanțelor, \mathbb{R}^p).

Pentru $n = 100$, alcătuiți un grafic cu valorile [funcției] mediane pentru $p = 1, 2, 3, \dots, 100$. Valorile lui p vor fi plasate pe axa Ox , iar valorile mediane pe axa Oy . Ce observați?

d. Folosind funcția de distribuție cumulativă (c.d.f.) de la punctul b, determinați cât de mare ar trebui să fie n (mărimea eșantionului) astfel încât

$$P(d^* \leq 0.5) \geq 0.9,$$

adică, cu probabilitate de cel puțin $9/10$, distanța d^* de la originea O la cel mai apropiat vecin să fie mai mică decât $1/2$ (adică, jumătate din distanța de la O la marginea sferei). Va trebui să formulați răspunsul ca expresie a unei funcții în raport cu variabila p .

Reprezentați grafic valorile acestei funcții pentru $p = 1, 2, \dots, 20$, plasând valorile lui p pe axa Ox și valorile funcției pe axa Oy . Ce observați?

Sugestie: Pentru $\ln(1-x)$, puteți face apel la dezvoltarea sa sub formă de *serie Taylor*:

$$\ln(1-x) = -\sum_{i=1}^{\infty} \frac{x^i}{i} \text{ pentru } -1 \leq x < 1.$$

e. În urma rezolvării punctelor de mai sus, ce puteți spune despre dezavantajele algoritmului k -NN în raport cu [diferitele valori posibile pentru] n și p ?

²⁵Engl., cumulative distribution function, c.d.f.

²⁶Se verifică ușor că pentru $p = 3$ se obține volumul sferei: $V_3(r) = \frac{(r\sqrt{\pi})^3}{\frac{3}{4}\sqrt{\pi}} = \frac{4\pi r^3}{3}$.