

Învățare automată

— Licență, anul III, 2021-2022 —

examen „suplimentar“ II (exerciții teoretice)

Nume student:

Grupa:

1. (Distribuția gaussiană unidimensională: estimarea mediei în sensul MLE și respectiv MAP, atunci când varianța este cunoscută)

În această problemă ne propunem să calculăm *estimatorul de verosimilitate maximă* (engl., maximum likelihood estimator, MLE) și *estimatorul de probabilitate maximă a posteriori* (engl., maximum a posteriori probability (MAP) estimator) pentru media unei distribuții gaussiene unidimensionale. Concret, presupunem că avem n instanțe, x_1, \dots, x_n generate în mod independent de către o distribuție normală cu varianța *cunoscută*, σ^2 , și media *necunoscută*, μ .

Calculați estimatorul MLE pentru media μ . Elaborati calculele în mod detaliat.

2. (Distribuția gaussiană, cazul $\mu = 0$: estimarea varianței în sensul MLE)

Fie X o variabilă aleatoare de distribuție normală (gaussiană) cu media 0 și varianța σ^2 , adică $X \sim \mathcal{N}(0, \sigma^2)$.

Găsiți estimarea de verosimilitate maximă (engl., maximum likelihood estimate) pentru parametrul σ^2 , adică σ_{MLE}^2 .

3. (Distribuția gaussiană unidimensională: estimarea varianței în sensul MLE [în cazul când nu se impun restricții asupra mediei μ])

Fie $x_1, \dots, x_n \in \mathbb{R}$ instanțe generate în mod independent și identic distribuite conform gaussienei $N(x|\mu, \sigma^2)$.

Comentarii:

1. Vă reamintim că funcția densitate de probabilitate (p.d.f.) a distribuției gaussiene unidimensionale este definită de expresia:

$$p(x) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right).$$

2. La problema 1.a am calculat estimatorul MLE al mediei μ pentru această distribuție. (Am notat acest estimator cu μ_{MLE} .)

Calculați estimatorul MLE al varianței σ^2 (notat cu σ_{MLE}^2), presupunând că nu se impune nicio restricție asupra lui μ .¹

¹Este util de știut de asemenea că la problema 2 se cere ca (tot pentru distribuția gaussiană unidimensională) să se estimeze varianța σ^2 în sens MLE, presupunând însă că media μ este 0.

4.

(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 $H(X) = -\sum_i p_i \log p_i$. 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 (sau, axiome) 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*:²

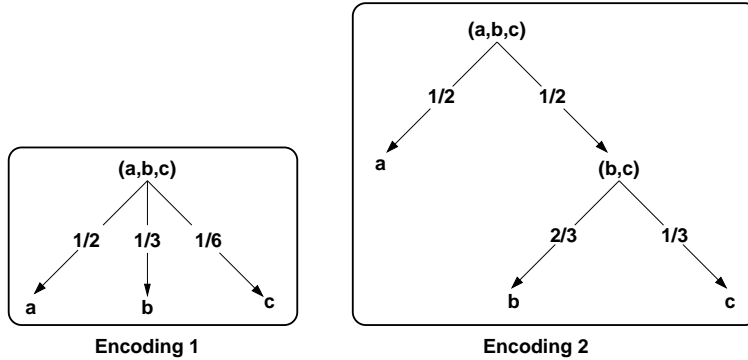
A1. Funcția $\psi_n(p_1, \dots, p_n)$ este continuă în fiecare din argumentele ei și *simetrică*.
Din punct de vedere formal, în acest caz simetria se traduce prin egalitatea $\psi_n(p_1, \dots, p_i, \dots, p_j, \dots, p_n) = \psi_n(p_1, \dots, p_j, \dots, p_i, \dots, p_n)$ pentru orice $i \neq j$. Informal 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 strict *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 $\psi_2(1/2, 1/2) + 1/2 \cdot \psi_2(2/3, 1/3)$.



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.³

²LC: Deși nu se specifică în enunțul original al problemei, este necesar / natural să considerăm și proprietatea următoare: [A0.] $\psi_1(1) = 0$ pentru că ψ_n este văzută ca măsură a *dezordinii* / *incertitudinii*, iar în cazul particular $n = 1$ nu există niciun fel de dezordine / incertitudine.

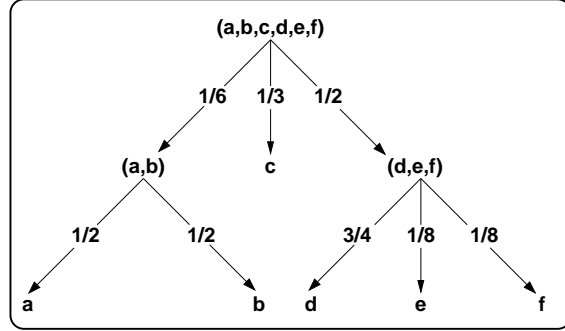
³În demonstrație, vom vedea, va rezulta $K = \frac{1}{\log s} \psi_s\left(\frac{1}{s}, \dots, \frac{1}{s}\right)$ pentru un număr oarecare $s \in \mathbb{N}^* \setminus \{1\}$, fixat.

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 alăturată, folosind din nou proprietatea A3.



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

c. Arătați că

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

La punctele d – g de mai jos, pentru orice număr $t \in \mathbb{N}^*$ (fixat), vom considera — pe lângă $n \in \mathbb{N}^*$, care a fost de fapt introdus atunci când am spus că aici ne ocupăm de $A(n)$ — numerele $s, m \in \mathbb{N}$, cu $s > 1$ astfel încât⁴

$$s^m \leq t^n \leq s^{m+1}. \quad (2)$$

d. Verificați că, prin logaritmare⁵ dublei inegalități (2) 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

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

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 (1), 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

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

g. Combinând inegalitățile (3) și (4), arătați că

$$\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n} \text{ pentru } s \neq 1 \quad (5)$$

și, în consecință

$$A(t) = K \log t \text{ cu } K > 0 \text{ (din cauza proprietății A2)}. \quad (6)$$

⁴LC: În idee, am putea să fixăm $s = 2$ și apoi să alegem $m \in \mathbb{N}$ (în funcție de t și n) astfel încât relația (2) să fie satisfăcută. Însă raționamentul următor nu depinde (în esență) de valoarea aleasă (și fixată) pentru s .

⁵Puteți folosi ca bază a logaritmului orice număr supra-unitar, arbitrar ales, dar fixat.

h. Arătați că 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 $\psi_k(p_1, \dots, p_k)$ cu $p_i \in \mathbb{Q}^+$ pentru $i = 1, \dots, k$ și $\sum_{i=1}^k p_i = 1$:

$$\psi_k(p_1, \dots, p_k) = -K \sum_i p_i \log p_i.$$

i. În sfârșit, tratați și cazul $p_i \in \mathbb{R}^+$, $i = 1, \dots, k$, cu $\sum_i p_i = 1$.

5. (Algoritmul Bayes [Naiv] gaussian: aplicare pe date din \mathbb{R})

Presupunem că dispunem de setul de date de antrenament din tabelul alăturat; singurul atribut de intrare (X) ia valori reale, iar atributul de ieșire (Y) este de tip Bernoulli, deci ia două valori, notate cu A și respectiv B .

X	Y
0	A
2	A
3	B
4	B
5	B
6	B
7	B

a. Pornind de la acest set de date, va trebui mai întâi să învățați *parametrii* clasificatorului Bayes gaussian, prin metoda estimării de verosimilitate maximă (MLE).⁶ Centralizați rezultatele, completând tabelul următor:

$\mu_A =$	$\sigma_A^2 =$	$P(Y = A) =$
$\mu_B =$	$\sigma_B^2 =$	$P(Y = B) =$

b. Notăm $\alpha = p(X = 2|Y = A)$ și $\beta = p(X = 2|Y = B)$.

- Cât este $p(X = 2, Y = A)$ în funcție de α ?
- Cât este $p(X = 2, Y = B)$ în funcție de β ?
- Cât este $p(X = 2)$ în funcție de α și β ?
- Cât este $p(Y = A|X = 2)$ în funcție de α și β ?

c. Cum va clasifica algoritmul Bayes [Naiv] gaussian punctul $X = 2$? Puteți exprima răspunsul fie în funcție de α și β , fie — mai bine! — calculând în prealabil valorile lui α și β în funcție de parametrii calculați / estimați la punctul precedent.

6. (Clasificatorul Bayes [Naiv] gaussian, cazul când se folosește un singur atribut de intrare: zone de decizie și granițe de decizie; analiza diferitelor cazuri specifice)

În acest exercițiu ne propunem să identificăm *zonele de decizie* și *granițele de decizie* determinate de clasificatorul Bayes [Naiv] gaussian (abreviat *GB*), atunci când folosim un singur atribut de intrare X , iar atributul de ieșire, pe care îl vom nota cu Y , este binar, deci poate lua două valori, desemnate în continuare cu A și B . Vom desemna prin $p_A = p$ *probabilitatea de selecție* pentru clasa A . Prin urmare,

⁶Vedeți secțiunea corespunzătoare din capitolul de *Fundamente*, în speță (pentru acest caz) problemele 1.a și 3.a.

probabilitatea de selecție pentru clasa B este $p_B = 1 - p$. Vom presupune că $p \in (0, 1)$ și $X|(Y = A) \sim \mathcal{N}(x|\mu_A, \sigma_A^2)$, iar $X|(Y = B) \sim \mathcal{N}(x|\mu_B, \sigma_B^2)$.

a. Arătați că regula de decizie a acestui clasificator Bayes [Naiv] gaussian

$$\hat{Y}_{GB}(X = x) = A \Leftrightarrow p \cdot \mathcal{N}(x|\mu_A, \sigma_A^2) \geq (1 - p) \cdot \mathcal{N}(x|\mu_B, \sigma_B^2) \quad (7)$$

devine echivalentă cu

$$(\sigma_A^2 - \sigma_B^2)(x - x_1)(x - x_2) \geq 0,$$

unde

$$x_1 = \frac{\sigma_A^2 \mu_B - \sigma_B^2 \mu_A - \sqrt{\Delta'}}{\sigma_A^2 - \sigma_B^2} \quad \text{și} \quad x_2 = \frac{\sigma_A^2 \mu_B - \sigma_B^2 \mu_A + \sqrt{\Delta'}}{\sigma_A^2 - \sigma_B^2},$$

cu

$$\Delta' \stackrel{\text{not.}}{=} \sigma_A^2 \sigma_B^2 \left[(\mu_A - \mu_B)^2 + (\sigma_A^2 - \sigma_B^2) \ln \left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B} \right)^2 \right],$$

în condițiile în care $\sigma_A^2 \neq \sigma_B^2$ și $\Delta' \geq 0$.

b. Arătați că atunci când $\sigma_A^2 = \sigma_B^2 \stackrel{\text{not.}}{=} \sigma^2$ și $\mu_A \neq \mu_B$, regula de decizie (7) este echivalentă cu

$$x \geq x_0 \text{ dacă } \mu_A > \mu_B$$

și, respectiv

$$x \leq x_0 \text{ dacă } \mu_A < \mu_B,$$

unde

$$x_0 = \frac{\mu_A + \mu_B}{2} + \frac{\sigma^2}{\mu_A - \mu_B} \cdot \ln \frac{1-p}{p}.$$

c. Arătați că este posibil ca în anumite cazuri să avem $\Delta' < 0$, adică inegalitatea (7) să fie ori adevărată pentru orice $x \in \mathbb{R}$, ori falsă pentru orice $x \in \mathbb{R}$. Altfel spus, arătați că pot exista combinații de valori pentru parametrii σ_A și σ_B (ambii strict pozitivi), μ_A și μ_B din \mathbb{R} , precum și $p \in (0, 1)$, astfel încât

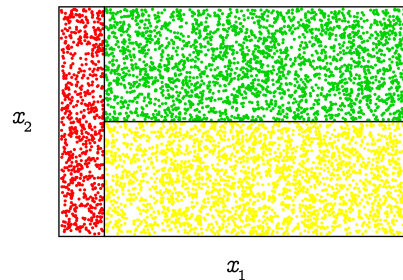
$$(\mu_A - \mu_B)^2 + (\sigma_A^2 - \sigma_B^2) \ln \left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B} \right)^2 < 0. \quad (8)$$

7.

(Arbori de decizie cu variabile continue:
ID3 ca algoritm “greedy”)

Fie următoarea problemă de clasificare ternară.

Considerăm că în figura alăturată regiunea dreptunghiulară este populată în mod dens cu puncte caracterizate de două atribute numerice (continue), x_1 și x_2 . Cele trei sub-dreptunghiuri (roșu, verde și galben) reprezintă trei clase de puncte, C_1, C_2 , și C_3 .



Dimensiunile $x_1 \times x_2$ ale dreptunghiurilor roșu, verde și galben sunt 1×6 , 7×3 și respectiv 7×3 . Dreptunghiul roșu este populat în mod uniform cu 6000 de puncte din clasa C_1 . Dreptunghiul verde este populat în mod uniform cu 42000 de puncte din clasa C_2 . Dreptunghiul galben este populat în mod uniform cu 42000 de puncte din clasa C_3 . Pentru simplitate, nu vom considera alte puncte decât acestea.

a. Care este numărul minim de noduri de test pe care trebuie să le aibă un arbore de decizie pentru a clasifica în mod corect acest set de date?

b. Câte noduri de test are arborele de decizie obținut în urma antrenării algoritmului ID3 pe acest set de date, folosind criteriul maximizării câștigului de informație?

Indicație: Pentru a determina entropiile condiționale minime, puteți folosi proprietatea A3 de la problema 4 de la capitolul de *Fundamente*.

c. Avem același număr de noduri în cele două cazuri de mai sus, sau nu? Care credeți că este explicația?

d. Un arbore de decizie poate să clasifice setul de date din figura de mai sus cu 100% acuratețe (presupunând că nu există zgomote la nivel de etichete). Ce condiții trebuie să satisfacă în general un set de date de acest gen astfel încât arborele de decizie rezultat în urma antrenării să fie cât mai compact și să producă o acuratețe de 100%?

Indicație: Fiecare nod intern al arborelui de decizie corespunde unui test bazat pe o singură trăsătură. Gândiți-vă ce fel de clase de funcții / granițe de separare corespund unui astfel de arbore de decizie.

8. (O clasă de concepte învățabile în sens empiric γ -slab cu ajutorul compașilor de decizie: seturile de instanțe din \mathbb{R} , care sunt etichetate în mod consistent)

Rezumat: În această problemă vom arăta că, în ipoteza că instanțele de antrenament cu care lucrăm au un singur atribut, x , care ia valori în \mathbb{R} , iar aceste instanțe (x_i , cu $i = 1, \dots, m$) sunt etichetate în mod „consistent“ (adică, necontradictoriu), există o constantă $\gamma > 0$ care reprezintă o garanție (engl., guarantee, or edge) pentru învățabilitate empirică „slabă“, folosind compași de decizie în conjuncție cu algoritmul AdaBoost.⁷

Din punct de vedere *analitic*, compașii de decizie determinați de praguri de separare (engl., thresholding-based decision stumps) pot fi definiți ca niște funcții-prag (sau, funcții-treaptă), care sunt indexate pe de o parte de un prag $s \in \mathbb{R}$ și pe de altă parte de un semn $+/-$, în felul următor:

$$\phi_{s,+}(x) = \begin{cases} 1 & \text{dacă } x \geq s \\ -1 & \text{dacă } x < s \end{cases} \quad \text{iar} \quad \phi_{s,-}(x) = \begin{cases} -1 & \text{dacă } x \geq s \\ 1 & \text{dacă } x < s. \end{cases}$$

Prin urmare, $\phi_{s,+}(x) = -\phi_{s,-}(x)$ pentru orice $x \in \mathbb{R}$.

Dat fiind un set de date de antrenament consistent etichetate $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, cu $x_i \in \mathbb{R}$ și $y_i \in \{-1, +1\}$ pentru $i = 1, \dots, m$, *obiectivul* nostru este să investigăm existența unei *garanții* pentru învățabilitate „slabă“ pe setul S . Vom demonstra că există o valoare $\gamma > 0$, astfel încât pentru *orice* distribuție probabilistă p definită pe S , putem găsi un prag $s \in \mathbb{R}$ cu proprietatea următoare:

$$\text{error}_p(\phi_{s,+}) \leq \frac{1}{2} - \gamma \quad \text{sau} \quad \text{error}_p(\phi_{s,-}) \leq \frac{1}{2} - \gamma,$$

⁷Noțiunile de învățabilitate empirică „slabă“ și garanție (γ) pentru învățabilitate empirică „slabă“ au fost definite la problema 23.e. Folosind notațiile de acolo, vă reamintim că *garanția* $\gamma \in (0, 1/2)$ are proprietatea că la orice iterație t a algoritmului AdaBoost există o ipoteză „slabă“ h_t astfel încât

$$\text{error}_{D_t}(h_t) \stackrel{\text{def.}}{=} \sum_{i=1}^m D_t(i) \cdot 1_{\{h_t(x_i) \neq y_i\}} \leq \frac{1}{2} - \gamma.$$

În formula aceasta, simbolul $1_{\{\dots\}}$ desemnează binecunoscuta funcție indicator. Concret, funcția $1_{\{h_t(x_i) \neq y_i\}}$ ia valoarea 1 pentru acei i pentru care $h_t(x_i) \neq y_i$, și 0 în caz contrar. Vă reamintim că p este o distribuție probabilistă discretă dacă $p_i \geq 0$ pentru $i = 1, \dots, m$ și $\sum_{i=1}^m p_i = 1$.

unde

$$error_p(\phi_{s,+}) \stackrel{\text{def.}}{=} \sum_{i=1}^m p_i \cdot 1_{\{y_i \neq \phi_{s,+}(x_i)\}} \quad \text{și} \quad error_p(\phi_{s,-}) \stackrel{\text{def.}}{=} \sum_{i=1}^m p_i \cdot 1_{\{y_i \neq \phi_{s,-}(x_i)\}}.$$

Facem *presupunerea* că toate instanțele de antrenament sunt distincte, deci nu există $x_i = x_j$ cu $i \neq j$. De asemenea, vom presupune — fără a reduce generalitatea; însă acest fapt ne va permite să simplificăm notațiile din problemă — că

$$x_1 > x_2 > \dots > x_m.$$

a. Arătați că, dat fiind S , pentru orice prag $s \in \mathbb{R}$ (fixat), dacă notăm cu $m_0(s) \in \{0, 1, \dots, m\}$ indicele care are proprietatea că $x_i \geq s$ pentru orice $i \leq m_0(s)$ și $x_i < s$ pentru orice $i > m_0(s)$,⁸ atunci au loc egalitățile

$$\begin{aligned} error_p(\phi_{s,+}) &= \frac{1}{2} - \frac{1}{2} \left(\sum_{i=1}^{m_0(s)} y_i p_i - \sum_{i=m_0(s)+1}^m y_i p_i \right) \\ error_p(\phi_{s,-}) &= \frac{1}{2} - \frac{1}{2} \left(\sum_{i=m_0(s)+1}^m y_i p_i - \sum_{i=1}^{m_0(s)} y_i p_i \right). \end{aligned}$$

Convenție: În expresii precum cele de mai sus, veți trata sumele indexate pe mulțimi vide ca fiind egale cu 0. Așadar, $\sum_{i=1}^0 a_i = 0$, oricare ar fi termenii a_i și, similar, $\sum_{i=m+1}^m a_i = 0$.

b. Arătați că, dat fiind setul de date de antrenament S , există o valoare $\gamma > 0$, care poate depinde de m , numărul de instanțe din S (dar nu și de distribuțiile probabiliste p care pot fi definite pe S), astfel încât pentru orice astfel de distribuție p , putem găsi un indice m_0 cu proprietatea următoare:

$$|f(m_0)| \geq 2\gamma, \text{ unde } f(m_0) \stackrel{\text{def.}}{=} \sum_{i=1}^{m_0} y_i p_i - \sum_{i=m_0+1}^m y_i p_i.$$

Sugestie: Analizați diferența $f(m_0 + 1) - f(m_0)$. Care este valoarea pe care ați găsit-o pentru γ ?

c. Ținând cont de răspunsurile pe care le-ați dat la punctele a și b, cât anume este *garanția* [pentru învățabilitate empirică „slabă”] pe care compașii de decizie o furnizează pentru orice set de date de antrenament $S = \{(x_1, y_1), \dots, (x_m, y_m)\}$, considerând că toate instanțele $x_i \in \mathbb{R}$, cu $i = 1, \dots, m$ sunt distincte?

d. Puteți indica o margine superioară (engl., upper bound) pentru numărul de compași de decizie necesari pentru a obține eroare la antrenare 0 pe un astfel de set de date de antrenament (S)?

9. (Algoritmul 1-NN [comparativ cu clasificatorul Bayes Optimal]: o margine superioară pentru *eroarea medie asimptotică* [la antrenare])

Un rezultat interesant obținut de Cover și Hart (1967) arată că, atunci când numărul datelor de antrenament tinde la infinit, iar datele de antrenament umplu spațiul în mod dens, rata medie a erorii produsă de către clasificatorul 1-NN este mărginită

⁸Știm că un astfel de indice există fiindcă $x_1 > x_2 > \dots > x_m$.

superior de dublul ratei medii a erorii pentru clasificatorul Bayes Comun (engl., Joint Bayes), care este numit adeseori și *Bayes Optimal* (engl., Optimal Bayes).

În acest exercițiu vi se va arăta, pas cu pas, cum se demonstrează rezultatul lui Cover și Hart în cazul particular al clasificării binare. Așadar, fie x_1, x_2, \dots instanțele de antrenament, iar y_1, y_2, \dots etichetele corespunzătoare, cu $y_i \in \{0, 1\}$. Putem considera instanțele x_i ca fiind puncte într-un spațiu euclidian d -dimensional.

Notăm $p_y(x) = P(X = x \mid Y = y)$ probabilitatea condiționată care reprezintă distribuția instanțelor din clasa y . Vom presupune că aceste probabilități condiționate sunt continue în raport cu variabila x și că $p_y(x) \in (0, 1)$ pentru orice x și orice y . Notăm cu θ probabilitatea ca un exemplu de antrenament selectat în mod aleatoriu să fie din clasa 1, așadar $\theta \stackrel{\text{not.}}{=} P(Y = 1)$. Din nou, presupunem că $\theta \in (0, 1)$.

a. Calculați probabilitatea ca o instanță oarecare x să aparțină clasei 1: $q(x) \stackrel{\text{not.}}{=} P(Y = 1 \mid X = x)$. Exprimați $q(x)$ în funcție de $p_0(x)$, $p_1(x)$ și θ .

b. Clasificatorul Bayes Optimal asignează unui punct dat x cea mai probabilă clasă, $\arg \max_y P(Y = y \mid X = x)$. (Aceasta implică faptul că algoritmul Bayes Optimal maximizează probabilitatea clasificării corecte a tuturor datelor.) Considerând o instanță oarecare x , calculați probabilitatea ca x să fie clasificat greșit folosind clasificatorul Bayes Optimal, în funcție de probabilitatea $q(x) \stackrel{\text{not.}}{=} P(Y = 1 \mid X = x)$ care tocmai a fost calculată la punctul precedent. Veți desemna această nouă probabilitate cu $Error_{Bayes}(x)$.

c. Acum considerăm clasificatorul 1-NN. Acesta îi asignează unei instanțe oarecare de test x eticheta celei mai apropiate instanțe de antrenament x' . Dată fiind o instanță de antrenament x (aleasă în mod arbitrar, dar fixată), calculați eroarea „așteptată” (engl., expected error) produsă de către clasificatorul 1-NN, adică probabilitatea ca instanța x să fie clasificată greșit. Notați această probabilitate cu $Error_{1-NN}(x)$ și exprimați-o sub forma unei funcții definite în raport cu probabilitățile $q(x)$ și $q(x')$.

d. În cazul asimptotic, numărul de exemple de antrenament al fiecărei clase tinde la infinit, iar datele de antrenament umplu spațiul în mod dens. Atunci $q(x') \rightarrow q(x)$, unde, ca și mai sus, x' este cel mai apropiat vecin al lui x .⁹ Făcând această substituție în rezultatul obținut la punctul anterior, deduceți expresia erorii asimptotice pentru clasificatorul 1-NN în punctul x , adică $\lim_{x' \rightarrow x} Error_{1-NN}(x)$, în funcție de probabilitatea $q(x)$.

e. Arătați că eroarea asimptotică obținută la punctul d este mai mică decât dublul erorii clasificatorului Bayes Optimal obținută la punctul b, adică:

$$\lim_{x' \rightarrow x} Error_{1-NN}(x) \leq 2Error_{Bayes}(x).$$

În final, din această inegalitate deduceți relația corespunzătoare între ratele medii ale erorilor:¹⁰

$$E[\lim_{n \rightarrow \infty} Error_{1-NN}] \leq 2E[Error_{Bayes}].$$

⁹Adică, $P(Y = 1 \mid X = x') \rightarrow P(Y = 1 \mid X = x)$. Aceasta se justifică ținând cont de continuitatea lui $p_y(x) \stackrel{\text{not.}}{=} P(X = x \mid Y = y)$ care a fost asumată în enunț și, de asemenea, de rezultatul obținut la punctul a.

¹⁰Cititorul atent va remarca faptul că în expresia de mai jos ($E[\lim_{n \rightarrow \infty} Error_{1-NN}]$) s-a înlocuit $\lim_{x \rightarrow x'}$ (folosită anterior) cu $\lim_{n \rightarrow \infty}$, pentru că se face trecerea la medii. Când $n \rightarrow \infty$, conform presupunziilor din enunț, rezultă $x \rightarrow x'$ pentru orice x .

10. (Compararea clasificatorilor 1-NN și Bayes Optimal:
o margine superioară mai bună
pentru *rata medie a erorii asimptotice* a lui 1-NN)

La problema 9 am demonstrat că *rata medie a erorii* clasificatorului 1-NN este mărginită asimptotic¹¹ de dublul ratei medii a erorii clasificatorului Bayes Optimal.

Arătați că — în aceleași condiții ca la problema 9 — se poate obține o margine chiar mai bună:

$$E\left[\lim_{n \rightarrow \infty} Error_{1-NN}\right] \leq 2E[Error_{Bayes}](1 - E[Error_{Bayes}]).$$

11. (Algoritmul EM:
chestiuni calitative / metodologice privind aplicarea
în cazul unei mixturi de două distribuții gaussiene unidimensionale
presupunând probabilitățile de selecție egale și fixate)

A fost odată, cu mult timp în urmă, un sat care era situat într-o regiune cu sute de lacuri. În acele lacuri trăiau două specii de pești, însă în fiecare lac, nu erau decât pești dintr-o singură specie. (Peștii din lacuri diferite puteau fi de specii diferite.) Peștii din aceste două specii arătau identic, miroseau identic atunci când erau gătiți și aveau exact același gust delicios, însă una dintre specii era otrăvitoare și oricine mânca pești din specia respectivă murea. Singura diferență observabilă la aceste specii de pești consta în efectul asupra nivelului pH-ului (aciditatea) apei din lacul în care trăiau. Lacurile cu pești neotrăvitori aveau pH-ul distribuit conform unei distribuții gaussiene de medie (μ_{sigur}) și varianță (σ_{sigur}^2) necunoscute, iar pH-ul din lacurile cu pești otrăvitori era distribuit conform unei alte distribuții gaussiene de medie (μ_{mortal}) și varianță (σ_{mortal}^2) de asemenea necunoscute. (Peștii otrăvitori cauzau o aciditate care era în general mai mare / ridicată decât în cazul celorlalți pești.)

Așa cum era firesc, pentru a ieși din încurcătură — adică pentru a determina caracterul otrăvitor ori neotrăvitor al peștilor din fiecare lac (sau dintr-un lac oarecare, din care încă nu s-au consumat pești) în funcție de nivelul pH-ului apei —, sătenii au apelat la învățarea automată. Cu toate acestea, au avut loc dezbateri aprinse cu privire la modalitatea corectă de aplicare a algoritmului Expectation-Maximization la problema lor.

Pentru fiecare dintre modalitățile prezentate mai jos, indicați dacă ea constituie o aplicare corectă a algoritmului EM și dacă va conduce la o estimare rezonabilă pentru parametrii μ și σ^2 ai fiecărei clase.

a. Se ghicesc valorile inițiale pentru parametrii μ și σ^2 corespunzătorii fiecărei specii. (1) Pentru fiecare lac, se determină — folosind teorema lui Bayes și distribuțiile gaussiene de parametri μ și σ^2 — care este cea mai probabilă specie de pești asociată lacului respectiv. (2) Se recalculează valorile pentru μ și σ^2 utilizând metoda estimării de verosimilitate maximă. Se repetă iterativ pașii (1) și (2) până se ajunge la convergență.

b. Pentru fiecare lac, se ghicește probabilitatea (inițială) că lacul respectiv ar fi populat cu pești neotrăvitori. (1) Folosind aceste probabilități, se estimează în sensul verosimilității maxime valorile μ și σ corespunzătoare fiecărei clase. (2) Utilizând aceste estimări pentru μ și σ , se recalculează probabilitățile de ‘siguranță’ pentru fiecare lac. Se repetă iterativ pașii (1) și (2) până se ajunge la convergență.

¹¹ Adică, atunci când $n \rightarrow \infty$, unde n este numărul de instanțe de antrenament.

c. Se calculează media și varianța nivelului de pH pentru toată mulțimea lacurilor. Se folosesc aceste valori pentru a inițializa parametrii μ și σ^2 corespunzători fiecărei specii de pești. (1) Folosind aceste valori pentru μ și σ^2 , se calculează pentru fiecare lac probabilitatea să conțină pești otrăvitori. (2) Se găsesc valorile de verosimilitate maximă pentru μ și σ^2 . Se repetă iterativ pașii (1) și (2) până se ajunge la convergență.

12.

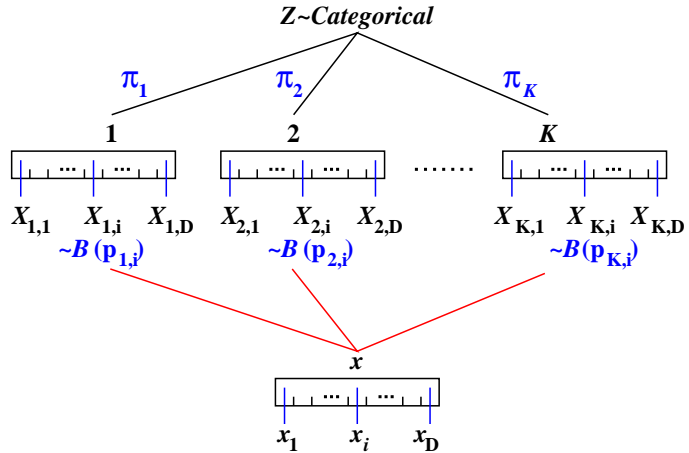
(Algoritmul EM pentru rezolvarea unei mixturi de vectori de variabile Bernoulli [independente]; aplicare la recunoșterea cifrelor scrise de mână)

În acest exercițiu, veți crea un algoritm de tip EM (expectation-maximization) care să clusterizeze imagini alb-negru. Input-urile $x^{(i)}$ pot fi văzute ca vectori de valori binare corespunzătoare culorilor alb și negru. Obiectivul este să clusterizăm aceste imagini în [mai multe] grupuri. Pentru a rezolva această problemă, veți folosi un model de tip mixtură de distribuții Bernoulli.

A. Mixtura de distribuții Bernoulli

a. Considerăm un vector de variabile aleatoare binare, $x \in \{0,1\}^D$. Presupunem că fiecare variabilă x_d urmează o distribuție *Bernoulli*(p_d), deci $P(x_d = 1) = p_d$. Fie $p \in (0,1)^D$ vectorul rezultat de parametri Bernoulli. Scrieți expresia probabilității $P(x|p)$, considerând că variabilele x_d sunt independente între ele.

b. Să presupunem acum că avem o mixtură de K vectori de distribuții Bernoulli: fiecare vector $x^{(i)}$ este generat folosind un vector de variabile Bernoulli independente, de parametru $p^{(k)}$ *not.* $(p_1^{(k)}, \dots, p_D^{(k)})$.



Presupunem că dispunem de o distribuție aleatoare π *not.* (π_1, \dots, π_K) , cu π_k indicând probabilitatea de selecție a setului de parametri Bernoulli $p^{(k)}$, pentru $k = 1, \dots, K$. Scrieți expresia probabilității $P(x^{(i)}|p, \pi)$, unde cu p notăm (de acum încolo) ansamblul $(p^{(1)}, \dots, p^{(K)})$.

c. Presupunem că avem input-urile $X = \{x^{(i)}\}_{i=1, \dots, n}$. Folosind rezultatele de la punctele anterioare, scrieți expresia log-verosimilității datelor X , $\ln P(X|\pi, p)$.

B. Pasul de estimare (engl., expectation step)

d. Acum vom introduce variabilele latente pentru algoritmul EM. Fie $z^{(i)} \in \{0,1\}^K$ un vector indicator, astfel încât $z_k^{(i)} = 1$ dacă vectorul $x^{(i)}$ a fost generat de vectorul de distribuții *Bernoulli*($p^{(k)}$), și 0 în caz contrar. Fie $Z = \{z^{(i)}\}_{i=1, \dots, n}$. Cât este probabilitatea $P(z^{(i)}|\pi)$? Dar $P(x^{(i)}|z^{(i)}, p, \pi)$?

Indicație: Folosiți artificul ridicării la putere (engl., exponentiation trick), ținând cont de faptul că $z_k^{(i)} \in \{0,1\}$.

e. Folosind cele două probabilități calculate la punctul d , scrieți expresia verosimilității datelor complete, $P(X, Z|\pi, p)$.

f. Fie $\mu(z_k) \stackrel{not.}{=} E[z_k^{(i)}|x^{(i)}, \pi, p]$. Demonstrați că

$$\mu(z_k^{(i)}) = \frac{\pi_k \prod_{d=1}^D (p_d^{(k)})^{x_d^{(i)}} (1 - p_d^{(k)})^{1-x_d^{(i)}}}{\sum_{j=1}^K \pi_j \prod_{d=1}^D (p_d^{(j)})^{x_d^{(i)}} (1 - p_d^{(j)})^{1-x_d^{(i)}}}.$$

Fie \bar{p} și $\bar{\pi}$ noile valori ale parametrilor pe care dorim să le obținem prin maximizare, iar p și π valorile lor de la iterația precedentă. Folosind aceste notații, arătați că funcția „auxiliară“, care este necesară pentru pasul M al algoritmului de estimare-maximizare, are expresia următoare:

$$\begin{aligned} & E[\ln P(X, Z|\bar{p}, \bar{\pi})|X, p, \pi] \\ &= \sum_{i=1}^N \sum_{k=1}^K \mu(z_k^{(i)}) \left[\ln \bar{\pi}_k + \sum_{d=1}^D \left(x_d^{(i)} \ln \bar{p}_d^{(k)} + (1 - x_d^{(i)}) \ln(1 - \bar{p}_d^{(k)}) \right) \right]. \end{aligned}$$

C. Pasul de maximizare

g. Acum trebuie să maximizăm expresia funcției „auxiliare“ care a fost dedusă la punctul f , în raport cu $\bar{\pi}$ și \bar{p} . Mai întâi, arătați că valoarea \bar{p} care maximizează funcția „auxiliară“ este

$$\bar{p}^{(k)} = \frac{\sum_{i=1}^N \mu(z_k^{(i)}) x^{(i)}}{N_k},$$

unde $N_k = \sum_{i=1}^N \mu(z_k^{(i)})$.

h. Demonstrați că valoarea $\bar{\pi}$ care maximizează funcția „auxiliară“ este:

$$\bar{\pi}_k = \frac{N_k}{\sum_{k'} N_{k'}}.$$

Sugestie: Dat fiind faptul că este vorba de rezolvarea unei probleme de optimizare cu restricții, puteți apela la metoda multiplicatorilor Lagrange.

13. (EM pentru combinarea unei clasificări de tip Bayes Naiv cu o mixtură de vectori de distribuții Bernoulli, folosind și un termen de penalizare / regularizare. Aplicație: recunoașterea cifrelor scrise de mână)

• ◦ U. Toronto, Radford Neal,
“Statistical Methods for Machine Learning and Data Mining” course,
2014 spring, HW2

Obiectivul dumneavoastră în această problemă va fi să clasificați cifre scrise de mână, cu ajutorul unor modele de mixturi [de distribuții probabile] care vor fi obținute prin maximizarea verosimilității datelor însoțită de un termen de penalizare, folosind algoritmul EM.

Datele pe care le veți folosi constau din 800 de imagini de antrenament și 1000 de imagini de test pentru cifre scrise de mână (din codurile poștale de pe plicuri din SUA). Aceste imagini au fost extrase din bine-cunoscutul set de date MNIST, prin selectare aleatorie dintr-un total de 60000 de imagini de antrenament furnizate, reducând rezoluția acestor imagini de la 28×28 de pixeli la 14×14 de pixeli, făcând media valorilor pixelilor pentru fiecare bloc de pixeli de

dimensiune 2×2 , iar apoi folosind niște praguri pentru aceste medii în așa fel încât să obținem valori binare. Pe pagina web a acestei culegeri¹² este furnizat un fișier de date de antrenament având 800 de linii, fiecare linie conținând 196 de valori pentru pixeli și anume, fie 0, fie 1. Tot acolo este furnizat și un alt fișier care conține etichetele corespunzătoare acestor 800 de cifre (de la 0 la 9). În mod similar, vă punem la dispoziție și un fișier cu 1000 de imagini de test, precum și un fișier conținând etichetele acestor 1000 de imagini de test. (Nu veți avea voie să vă uitați la etichetele acestor imagini de test decât la sfârșit de tot, și anume atunci când va trebui să analizați cât de bine se comportă algoritmul de învățare automată.)

Veți clasifica aceste imagini pentru cifre scrise de mână folosind un *model generativ*, din care, dată fiind imaginea unei astfel de cifre, veți deriva probabilitățile pentru cele 10 clase posibile. *Clasa asociată unei cifre de test* este cea care are probabilitatea maximă.

Modelul generativ pe care îl vom folosi estimează *probabilitățile claselor* ca [fiind] *frecvențele* lor din setul de date de antrenament (ceea ce va determina o distribuție aproximativ uniformă — însă nu exact uniformă — peste cele 10 cifre). De asemenea, acest model estimează *distribuțiile de probabilitate ale imaginilor din cadrul fiecărei clase* cu ajutorul unui *model de mixtură cu K componente*, fiecare componentă modelând valorile celor 196 de pixeli ca [fiind] variabile aleatoare *independente*. Va fi convenabil să combinăm toate aceste 10 modele de mixturi [de distribuții probabiliste] într-un singur *model de mixtură cu $10K$ componente*, care modelează atât valorile pixelilor cât și etichetele claselor. Totuși, *probabilitățile [condiționate ale] componentelor* din cadrul modelului fiecărei clase vor fi *fixate*, astfel încât K componente vor atribui probabilitatea condiționată 1 cifrei 0, alte K componente vor atribui probabilitatea condiționată 1 cifrei 1, alte K vor atribui probabilitatea condiționată 1 cifrei 2 ș.a.m.d.

Așadar, modelul probabilist pentru *distribuția comună* asupra etichetelor y_i , precum și asupra valorilor pixelilor $x_{i,1}, \dots, x_{i,196}$ pentru $i = 0, \dots, 9$ poate fi exprimat sub forma următoare:

$$P(y_i, x_i | \pi, \theta) = \sum_{k=1}^{10K} \left(\pi_k q_{k,y_i} \prod_{j=1}^{196} \theta_{k,j}^{x_{i,j}} (1 - \theta_{k,j})^{1-x_{i,j}} \right).$$

Instanțele etichetate (x_i, y_i) se presupune a fi independente.

Parametrii acestui model sunt *probabilitățile de mixare / selecție*, π_1, \dots, π_{10K} , precum și probabilitățile $\theta_{k,j}$ asociate pixelilor (mai precis, pentru ca aceștia să fie setați la valoarea 1), pentru fiecare componentă în parte, adică pentru $k = 1, \dots, 10K$ și $j = 1, \dots, 196$. Probabilitățile condiționate $q_{k,y}$ ale componentelor fiecărei clase sunt fixate astfel:

$$q_{k,y} = \begin{cases} 1 & \text{atunci când } k \in \{Ky + 1, \dots, Ky + K\} \\ 0 & \text{în caz contrar,} \end{cases}$$

pentru $k = 1, \dots, 10K$ și $y = 0, \dots, 9$.

La sfârșitul acestui exercițiu veți implementa în Python / R / Matlab o funcție care identifică valorile parametrilor care maximizează log-verosimilitatea datelor de antrenament plus un *termen de penalizare* (engl., penalty term).¹³

Algoritmul EM poate fi ușor adaptat pentru a găsi [în locul estimării verosimilității maxime] estimarea verosimilității maxime penalizate. În raport cu versiunea generală a algoritmului EM [care a fost prezentată la curs], pasul E se păstrează

¹²Vedeți <https://profs.info.uaic.ro/~ciortuz/ML.ex-book/implementation-exercises/> U Toronto. 2014s. RNeal.HW2.EM-for-BernoulliMM-using-the-NBayes-assumption.handwritten-char-reco.data+R-code+sol/

¹³Cu cât valorile acestui termen de penalizare vor fi mai mari, cu atât va fi mai bine.

neschimbat, însă pasul M va maximiza acum o funcție care reprezintă media log-verosimilității penalizate a datelor complete, $E_Q[\ln P(x, z|\theta) + G(\theta)]$, unde $G(\theta)$ este termenul de penalizare.

Scopul penalizării este acela de a evita ca estimările probabilităților asociate pixelilor să fie 0 sau aproape de 0, situații care pot cauza probleme la clasificarea instanțelor de test. (Este imediat că dacă o astfel de probabilitate este 0, rezultă că probabilitatea oricărei instanțe de test va fi 0 pentru clasa respectivă.) Termenul de penalizare care se sumează la log-verosimilitate trebuie să fie

$$G(\theta) = \alpha \sum_{k=1}^{10K} \sum_{j=1}^{196} [\log(\theta_{k,j}) + \log(1 - \theta_{k,j})],$$

unde constanta α controlează „magnitudinea” penalizării.¹⁴

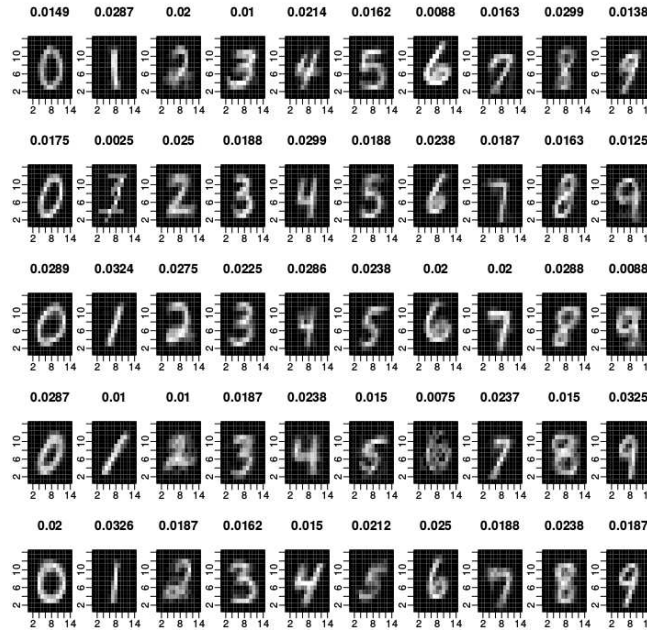
a. Demonstrați că într-adevăr [regula de actualizare de la] pasul E rămâne la fel ca în cazul versiunii generale a algoritmului EM.

b. La pasul M, reestimarea probabilităților de selecție π va rămâne de asemenea neschimbată [în raport cu cazul general],¹⁵ însă pentru reestimarea parametrilor θ va trebui să luați în considerare și termenul de penalizare. Demonstrați că formula pentru actualizarea parametrilor $\theta_{k,j}$ la pasul M este

$$\hat{\theta}_{k,j} = \frac{\alpha + \sum_{i=1}^n r_{i,k} x_{i,j}}{2\alpha + \sum_{i=1}^n r_{i,k}},$$

unde $r_{i,k}$ este probabilitatea ca instanța i să fi fost generată de către componenta k a mixturii, probabilitate care a fost estimată la pasul E.¹⁶

Comentariu: Pentru $K = 5$, putem ilustra sub forma unor imagini valorile parametrilor θ care au fost învățate pentru cele 50 de componente ale mixturii:



¹⁴Pentru acest exercițiu, α va putea fi fixat la valoarea 0.05, deși în aplicații reale el va trebui probabil să fie determinat folosind, de exemplu, metoda cross-validării.

¹⁵Va trebui să vă convingeți singuri de acest fapt.

¹⁶Ca să obțineți această formulă, ați putea să porniți de la demonstrația care a fost prezentată la problema 12, modificând-o în așa fel încât să includă și termenul de penalizare.

Este clar că în general cele cinci componente corespunzătoare fiecărei cifre au identificat diferite variante rezonabile de caligrafie a cifrelor, exceptând probabil câteva cazuri, cărora le corespund probabilități de selecție destul de mici (indicate sub formă numerică, deasupra imaginilor), așa cum este caracterul “1” de pe cel de-al doilea rând de imagini.