

0.1 Fundamente — Probleme rezolvate

0.1.1 Evenimente aleatoare și formula lui Bayes

1. (Proprietăți derivate din definiția funcției de probabilitate)

• CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.1

Fie două evenimente A și B .

a. Folosind doar proprietățile din definiția funcției de probabilitate, arătați că $P(A \cap \bar{B}) = P(A) - P(A \cap B)$.

b. Demonstrați *inegalitatea lui Bonferroni*: $P(A \cap B) \geq P(A) + P(B) - 1$.

c. Spunem că evenimentele A și B sunt *incompatibile* dacă $P(A \cap B) = 0$.

În ipoteza în care $P(A) = 1/3$ și $P(B) = 5/6$, este posibil ca evenimentele A și B să fie incompatibile? Justificați răspunsul.

Răspuns:

a. $A = A \cap \Omega = A \cap (B \cup \bar{B}) = (A \cap B) \cup (A \cap \bar{B})$.

Cum evenimentele $(A \cap B)$ și $(A \cap \bar{B})$, văzute ca mulțimi, sunt disjuncte, conform proprietății de „aditivitate numărabilă” din definiția funcției de probabilitate putem scrie $P(A) = P(A \cap B) + P(A \cap \bar{B})$, deci $P(A \cap \bar{B}) = P(A) - P(A \cap B)$.

b. Întrucât $A \cup B = (A \cap \bar{B}) \cup (A \cap B) \cup (\bar{A} \cap B)$, este imediat că $P(A \cup B) = P(A \cap \bar{B}) + P(A \cap B) + P(\bar{A} \cap B)$. Deci $P(A \cap B) = P(A \cup B) - P(A \cap \bar{B}) - P(\bar{A} \cap B) = [P(A \cup B) - P(A \cap \bar{B})] + [P(A \cup B) - P(\bar{A} \cap B)] - P(A \cup B) = P(A) + P(B) - P(A \cup B)$. Cum $P(A \cup B) \leq 1$ (fiind o probabilitate) putem scrie că: $P(A \cap B) = P(A) + P(B) - P(A \cup B) \geq P(A) + P(B) - 1$, deci $P(A \cap B) \geq P(A) + P(B) - 1$.

c. Pentru a studia posibilitatea ca evenimentele A și B să fie incompatibile vom aplica inegalitatea lui Bonferroni:

$$P(A) + P(B) = \frac{1}{3} + \frac{5}{6} = \frac{7}{6} \Rightarrow P(A \cap B) \geq \frac{1}{6} > 0$$

Deci $P(A \cap B)$ nu poate fi 0 și, în consecință, evenimentele A și B nu sunt incompatibile.

2. (Calcularea de probabilități elementare și probabilități condiționate)

• CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.4

În acest exercițiu vom arăta că — într-un anumit sens — probabilitatea unui eveniment se poate schimba (într-un anumit sens) dacă știm probabilitatea unui alt eveniment, legat de cel dintâi.

Se aruncă simultan două zaruri. Notăm cu S variabila aleatoare care desemnează suma numerelor rezultate din aruncarea celor două zaruri.

a. Calculați $P(S = 11)$.

b. Dacă știm că S este număr prim, cât devine probabilitatea de mai sus?

Răspuns:

a. Probabilitatea unui eveniment precum cel din enunț este dată de raportul dintre numărul de cazuri favorabile și numărul cazurilor posibile.

La aruncarea simultană a două zaruri, fiecare zar poate cădea pe una dintre cele 6 fețe, independent de celălalt zar. Prin urmare, există $6 \cdot 6 = 36$ posibilități (cazuri posibile). Pentru a se obține $S = 11$ cele două zaruri trebuie să aibă fie valorile $(5, 6)$, fie $(6, 5)$, deci există 2 posibilități (cazuri favorabile). Prin urmare,

$$P(S = 11) = \frac{2}{36} = \frac{1}{18}$$

b. Probabilitatea căutată este:

$$P(S = 11 \mid S = \text{prim}) = \frac{P(S = 11 \cap S = \text{prim})}{P(S = \text{prim})} = \frac{P(S = 11)}{P(S = \text{prim})}$$

Trebuie calculată probabilitatea ca S să fie număr prim. Pentru aceasta este necesar numărul de cazuri favorabile, adică numărul cazurilor pentru care $S \in \{2, 3, 5, 7, 11\}$. Există următoarele 15 posibilități:

- $S = 2$: $(1, 1)$
- $S = 3$: $(1, 2), (2, 1)$
- $S = 5$: $(1, 4), (4, 1), (2, 3), (3, 2)$
- $S = 7$: $(1, 6), (6, 1), (2, 5), (5, 2), (3, 4), (4, 3)$
- $S = 11$: $(5, 6), (6, 5)$

Deci $P(S = \text{prim}) = 15/36$. Prin urmare,

$$P(S = 11 \mid S = \text{prim}) = \frac{2/36}{15/36} = \frac{2}{15}$$

Întrucât $2/15 > 1/18$, putem spune că probabilitatea evenimentului $S = 11$ a crescut după ce am aflat un fapt suplimentar, și anume că suma rezultată la aruncarea celor două zaruri este un număr prim ($S = \text{prim}$).

Ca o *observație* cu caracter general: dacă $A \subseteq B$ (așa cum a fost cazul în această problemă), atunci rezultă imediat că $P(A) \leq P(A|B)$.

3.

(Spațiu de eșantionare – exemplificare;
calcul de probabilități condiționate)

• *CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 1.4*

O cutie conține trei carduri. Un card este roșu pe ambele părți, un alt card este verde pe ambele părți iar cel care a rămas este roșu pe o parte și verde

pe partea cealaltă. Selectăm în mod aleatoriu un card din această cutie; presupunem că nu-i vedem decât culoarea de pe fața superioară. Dacă această față este verde, care este probabilitatea ca și cealaltă față să fie verde?

Răspuns:

Pentru a rezolva această problemă este foarte util să stabilim mai întâi care sunt elementele ce compun *spațiul de eșantionare* (engl., sample space), Ω .¹⁴ Contrar intuiției primare, Ω nu va fi constituit din cele trei carduri, ci din fețele lor, fiindcă ceea ce observăm după o extragere este doar o față a unui card, nu ambele fețe ale cardului extras.

Din punct de vedere formal, vom folosi următoarea *notație* pentru carduri:

$$C_1 = (R1, R2), C_2 = (R3, V4), C_3 = (V5, V6),$$

unde $R1, R2, R3, V4, V5$ și $V6$ desemnează cele 6 fețe ale cardurilor. Așadar, $\Omega = \{R1, R2, R3, V4, V5, V6\}$.

Observație: Am fi putut nota fețele cardurilor folosind pur și simplu numerele $1, \dots, 6$, însă am preferat să însoțim fiecare dintre aceste numere cu litera R sau V care desemnează culoarea feței respective.

După ce am făcut această pregătire, probabilitatea cerută în enunțul problemei se calculează simplu, folosind regula clasică: $p = m/n$, unde m este numărul de cazuri favorabile, iar n este numărul de cazuri posibile.

Cazurile posibile sunt $V4, V5, V6$, iar cazurile favorabile sunt $V5$ și $V6$ deoarece doar pentru ele fața cealaltă a cardului este verde (este vorba de $V6$ și respectiv $V5$). Așadar, probabilitatea cerută este $\frac{2}{3}$.

4. (Evenimente aleatoare independente)

- CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.2.1
- CMU, 2009 spring, Ziv Bar-Joseph, HW1, pr. 1.1

Două evenimente A și B sunt independente statistic dacă $P(A \cap B) = P(A) \cdot P(B)$.

a. Arătați că dacă A și B sunt evenimente independente, atunci:

- A și \bar{B} sunt independente;
- \bar{A} și \bar{B} sunt independente.

b. Dacă evenimentul A este independent în raport cu el însuși, ce puteți spune despre $P(A)$?

Răspuns:

a. $P(A \cap \bar{B}) = P(A) - P(A \cap B) = P(A) - P(A) \cdot P(B) = P(A) \cdot (1 - P(B)) = P(A) \cdot P(\bar{B})$, deci A și \bar{B} sunt independente. (S-a folosit independența evenimentelor A și B .)

Pentru independența evenimentelor \bar{A} și \bar{B} se procedează analog:

$P(\bar{A} \cap \bar{B}) = P(\bar{B}) - P(A \cap \bar{B}) = P(\bar{B}) - P(A) \cdot P(\bar{B}) = P(\bar{B}) \cdot (1 - P(A)) = P(\bar{A}) \cdot P(\bar{B})$, deci

¹⁴A se vedea noțiunea prezentată la curs.

\bar{A} și \bar{B} sunt independente. (La cea de-a doua egalitate s-a folosit independența evenimentelor A și B demonstrată mai sus.)

b. Condiția de independență a evenimentului A față de el însuși se scrie astfel:

$$P(A \cap A) = P(A) \cdot P(A) \Rightarrow P(A) = [P(A)]^2 \Rightarrow P(A)[P(A) - 1] = 0$$

Ținând cont și de restricția $P(A) \in [0, 1]$, rezultă că $P(A) = 0$ sau $P(A) = 1$. (Atenție: nu rezultă neapărat că $A = \emptyset$ respectiv $A = \Omega$!)

5. (Evenimente aleatoare independente;
aplicarea proprietăților din definiția funcției de probabilitate)
• CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.2.3

Robert și Alina dau cu banul alternativ. Cel dintâi dintre ei care va obține stema (în engleză: head) câștigă jocul. Alina este prima care va da cu banul.

a. Dacă $P(\text{stemă}) = 1/2$, care este probabilitatea ca Alina să câștige jocul?

Indicație (1): Încercați să enumerați toate situațiile în care Alina poate câștiga.

Indicație (2): Pentru orice $a \in [0, 1]$, avem $\sum_{i=0}^{i=\infty} a^i = 1 + a + a^2 + \dots + a^n + \dots = \lim_{n \rightarrow +\infty} \frac{1 - a^{n+1}}{1 - a} = \frac{1}{1 - a}$.

b. Dacă $P(\text{stemă}) = p \in (0, 1]$, care este probabilitatea ca Alina să câștige jocul?

c. Ținând cont de expresia obținută la punctul b, dacă ar fi ca tu să joci acest joc, cum ai decide să intri în joc: primul ori al doilea (presupunând, bineînțeles, că ai avea posibilitatea să alegi)?

Răspuns:

a. Alina aruncă moneda în „pași” impari. Ea câștigă jocul dacă la pasul $2n + 1$ obține stema și până la pasul respectiv nimeni nu a obținut stema.

Notăm cu A evenimentul ca Alina să câștige jocul și cu A_i evenimentul ca Alina să câștige jocul la a i -a aruncare a banului. Întrucât evenimentele A_i , văzute ca mulțimi, sunt mutual disjuncte ($A_i \cap A_j = \emptyset$ pentru orice $i \neq j$) și $A = A_1 \cup A_3 \cup \dots$, conform proprietății de aditivitate numărabilă din definiția funcției de probabilitate rezultă că

$$P(A) = P(A_1) + P(A_3) + P(A_5) + \dots$$

Întrucât $p = \frac{1}{2}$, ținând cont [și] de faptul că toate aruncările sunt independente, probabilitățile corespunzătoare evenimentelor A_i se calculează astfel:

$$P(A_1) = P(\text{stemă}) = \frac{1}{2}$$

$$P(A_3) = (1 - P(\text{stemă})) \cdot (1 - P(\text{stemă})) \cdot P(\text{stemă}) = \left(\frac{1}{2}\right)^3$$

$$P(A_5) = \left(1 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) \cdot \left(1 - \frac{1}{2}\right) \cdot P(\text{stemă}) = \left(\frac{1}{2}\right)^5$$

...

$$P(A_{2i+1}) = (1 - P(\text{stemă}))^{2i} \cdot P(\text{stemă}) = \left(\frac{1}{2}\right)^{2i+1}$$

Prin urmare,

$$\begin{aligned} P(A) &= \sum_{i=0}^{\infty} P(A_{2i+1}) = \sum_{i=0}^{\infty} \left(\frac{1}{2}\right)^{2i+1} = \sum_{i=0}^{\infty} \left(\frac{1}{2}\right)^{2i} \cdot \frac{1}{2} = \frac{1}{2} \cdot \sum_{i=0}^{\infty} \left(\frac{1}{4}\right)^i = \frac{1}{2} \cdot \frac{1}{1-1/4} \\ &= \frac{1}{2} \cdot \frac{4}{3} = \frac{2}{3}. \end{aligned}$$

Așadar, pentru $p = 1/2$ probabilitatea ca Alina să câștige jocul este $2/3$ (deci de două ori mai mare decât probabilitatea ca Robert să câștige jocul), pur și simplu datorită faptului că ea este prima care dă cu banul. (Avantajul primului jucător!)

b. Dacă $P(\text{stemă}) = p \in (0, 1]$, se urmează același raționament ca mai sus, cu modificarea valorilor $P(A_i)$ astfel:

$$\begin{aligned} P(A_1) &= P(\text{stemă}) = p \\ P(A_3) &= (1 - P(\text{stemă})) \cdot (1 - P(\text{stemă})) \cdot P(\text{stemă}) = (1 - p)^2 \cdot p \\ P(A_5) &= (1 - p) \cdot (1 - p) \cdot (1 - p) \cdot (1 - p) \cdot P(\text{stemă}) = (1 - p)^4 \cdot p \\ &\dots \\ P(A_{2i+1}) &= (1 - P(\text{stemă}))^{2i} \cdot P(\text{stemă}) = (1 - p)^{2i} \cdot p \end{aligned}$$

Prin urmare, probabilitatea ca Alina să câștige devine:

$$P(A) = \sum_{i=0}^{\infty} P(A_{2i+1}) = \sum_{i=0}^{\infty} (1-p)^{2i} \cdot p = p \cdot \sum_{i=0}^{\infty} (1-p)^{2i} \stackrel{p \neq 0}{=} p \cdot \frac{1}{1 - (1-p)^2} = \frac{p}{2p - p^2} = \frac{1}{2 - p}$$

c. Jucătorul care aruncă primul banul câștigă jocul cu o probabilitate de $1/(2-p)$ (calculată mai sus). Cum $p \geq 0$, rezultă că $\frac{1}{2-p} \geq \frac{1}{2}$. Așadar, în cazul în care există posibilitatea de a alege, este de preferat să arunci primul.

6.

(Formula lui Bayes)

• CMU, 2001 fall, Andrew Moore, midterm exam, pr. 5

Se consideră două variabile aleatoare A și B despre care știm următoarele informații:

- a. $P(A | B) = 2/3$
- b. $P(A | B) = 2/3$ și $P(A | \bar{B}) = 1/3$
- c. $P(A | B) = 2/3$, $P(A | \bar{B}) = 1/3$ și $P(B) = 1/3$
- d. $P(A | B) = 2/3$, $P(A | \bar{B}) = 1/3$, $P(B) = 1/3$ și $P(A) = 4/9$.

În care din cele patru cazuri informațiile date sunt suficiente pentru a calcula $P(B | A)$? Există vreun caz în care apar informații superflue (i.e., informații care pot fi deduse din celelalte informații furnizate în cazul respectiv)?

Răspuns:

Conform teoremei lui Bayes, vom avea:

$$\begin{aligned} P(B | A) &= \frac{P(A | B) \cdot P(B)}{P(A)} = \frac{P(A | B) \cdot P(B)}{P(A | B) \cdot P(B) + P(A | \bar{B}) \cdot P(\bar{B})} \\ &= \frac{P(A | B) \cdot P(B)}{P(A | B) \cdot P(B) + P(A | \bar{B}) \cdot (1 - P(B))} \end{aligned}$$

Devine astfel evident că în cazurile c și d informațiile deținute sunt suficiente, pe când în celelalte două cazuri nu. În cazul d , informația $P(A) = 4/9$ este superfluă.

7.

(Formula lui Bayes)

• ◦ CMU, 2008 spring, T. Mitchell, W. Cohen, midterm, pr. 1.5

Presupunem că, răspunzând la o întrebare cu răspuns de genul adevărat / fals, un student fie cunoaște răspunsul, fie ghicește răspunsul. Probabilitatea de a cunoaște răspunsul este p , iar probabilitatea de a ghici răspunsul este $1 - p$.

Presupunem că probabilitatea de a răspunde corect la întrebare este

1 în cazul în care studentul cunoaște răspunsul

și 0.5 dacă studentul ghicește răspunsul.

Exprimați în funcție de p care este probabilitatea ca studentul examinat să cunoască răspunsul la întrebare, în ipoteza că el a răspuns corect (notație: $P(knew | correct)$).

Răspuns:

Evenimentele de interes în problema dată sunt:

$correct$ = studentul a răspuns corect la întrebare,

$knew$ = studentul cunoștea răspunsul corect

și complementarul acestuia din urmă:

$guess \stackrel{not.}{=} \neg knew$ = studentul ghicește răspunsul.

Aplicând formula lui Bayes obținem:

$$\begin{aligned} P(knew | correct) &= \frac{P(correct | knew) \cdot P(knew)}{P(correct | knew) \cdot P(knew) + P(correct | guess) \cdot P(guess)} \\ &= \frac{1 \cdot p}{1 \cdot p + 0.5 \cdot (1 - p)} = \frac{p}{0.5p + 0.5} = \frac{p}{0.5(p + 1)} = \frac{2p}{p + 1} \end{aligned}$$

8.

(Probabilități, chestiuni elementare:
Adevărat sau Fals?)

◻ • ◦ CMU, 2016 fall, N. Balcan, M. Gormley, HW1, pr. 6.1.1-4

În cele de mai jos vom nota cu Ω spațiul de eșantionare, iar cu \bar{A} complementul evenimentului A .

Marcați cu *adevărat* sau *fals* fiecare dintre afirmațiile următoare:

a. Pentru orice $A, B \subseteq \Omega$ astfel încât $P(A) > 0$ și $P(B) > 0$, are loc egalitatea $P(A|B)P(B) = P(B|A)P(A)$.

b. Pentru orice $A, B \subseteq \Omega$ astfel încât $P(B) > 0$, are loc egalitatea $P(A \cup B) = P(A) + P(B) - P(A|B)$.

c. Pentru orice $A, B, C \subseteq \Omega$ astfel încât $P(B \cup C) > 0$, urmează că $\frac{P(A \cup B \cup C)}{P(B \cup C)} \geq P(A|B \cup C)P(B \cup C)$

d. Pentru orice $A, B \subseteq \Omega$ astfel încât $P(A) > 0$ și $P(\bar{A}) > 0$, urmează că $P(B|A) + P(B|\bar{A}) = 1$.

Indicație:

Pentru fiecare afirmație adevărată, faceți demonstrația proprietății respective. Pentru fiecare afirmație falsă, dați fie un contraexemplu fie o justificare riguroasă.

Răspuns:

a. Adevărat. Demonstrația este imediată.

b. Fals.

Știm că pentru orice $A, B \subseteq \Omega$, are loc egalitatea $P(A \cup B) = P(A) + P(B) - P(A \cap B)$. Așadar, egalitatea din enunț ar fi adevărată dacă pentru orice $A, B \subseteq \Omega$ am avea $P(A \cap B) = P(A|B)$. Însă,

$$P(A \cap B) = P(A|B) \Leftrightarrow P(A \cap B) = \frac{P(A \cap B)}{P(B)} \Leftrightarrow P(B) = 1 \text{ sau } P(A \cap B) = 0.$$

Deci egalitatea nu este adevărată pentru orice A și B .

c. Adevărat.

Observăm că dacă se notează $D = B \cup C$, atunci inegalitatea din enunț devine mai ușor de manipulat: $\frac{P(A \cup D)}{P(D)} \geq P(A|D) \cdot P(D)$. Este imediat că ea este echivalentă cu $\frac{P(A \cup D)}{P(D)} \geq P(A \cap D)$, ceea ce implică $P(A \cup D) \geq P(A \cap D) \cdot P(D)$.

Ultima inegalitate este adevărată fiindcă pe de o parte $P(A \cup D) \geq P(A \cap D)$ și pe de altă parte $P(D) \in [0, 1]$.

d. Fals.

Știm că pentru orice $A, B \subseteq \Omega$ cu $P(A) > 0$, are loc egalitatea $P(B|A) + P(\bar{B}|A) = 1$ (o puteți demonstra imediat). Așadar, egalitatea din enunț ar fi adevărată dacă pentru orice $A, B \subseteq \Omega$ astfel încât $P(A) > 0$ și $P(\bar{A}) > 0$ am avea $P(B|\bar{A}) = P(\bar{B}|A)$.

Vom construi următorul contraexemplu: considerăm că la aruncarea unui zar perfect, evenimentul A este apariția unei fețe pare, iar evenimentul B este apariția feței 1. Urmează că

$$P(B|\bar{A}) = \frac{1}{3} \text{ și } P(\bar{B}|A) = 1.$$

Deci egalitatea nu este adevărată pentru orice evenimente A (cu $P(A) > 0$ și $P(\bar{A}) > 0$) și B .

0.1.2 Variabile aleatoare

9. (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 $P(x) \stackrel{\text{not.}}{=} P(X = x) \stackrel{\text{not.}}{=} P(\{\omega \mid X(\omega) = x\}) \geq 0$ 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. ■ CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW1, pr. 1.1

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

$$E[aX] = aE[X]. \quad (1)$$

Notați că aX este o variabilă aleatoare definită pe același domeniu (Ω) ca și variabila X , cu proprietatea că $(aX)(\omega) \stackrel{\text{def.}}{=} aX(\omega)$ 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. CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW1, pr. 1.3

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

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

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. *CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.3.1*

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

Demonstrați egalitatea:

$$Cov(X, Y) = E[XY] - E[X] \cdot E[Y]. \quad (3)$$

Consecință imediată (din relațiile (1) și (3)):

$$Cov(aX, bY) = ab Cov(X, Y) \quad \forall a, b \in \mathbb{R}. \quad (4)$$

Observații:

1. Este imediat faptul că proprietatea (4) generalizează proprietatea (2).

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, la problema 19 se va demonstra că atunci când $Var(X) \neq 0$ și $Var(Y) \neq 0$, avem următoarele margini (una inferioară și cealaltă superioară) pentru $Cov(X, Y)$:

$$-Var(X) \cdot Var(Y) \leq Cov(X, Y) \leq +Var(X) \cdot Var(Y).$$

Răspuns:

a. Pentru cazul discret vom folosi o formă echivalentă a formulei pentru media unei variabile aleatoare și anume $E[X] = \sum_{\omega \in \Omega} X(\omega) \cdot P(\omega)$. Prin urmare, putem scrie:

$$\begin{aligned} E[W + Z] &= \sum_{u \in Val(W+Z)} u \cdot P(W + Z = u) \\ &= \sum_{w \in Val(W), z \in Val(Z)} (w + z) \cdot P(W + Z = w + z) \\ &= \sum_{w \in Val(W)} \sum_{z \in Val(Z)} (w + z) \cdot P(\{\omega \in \Omega \mid (W + Z)(\omega) = w + z\}) \\ &= \sum_{w \in Val(W)} \sum_{z \in Val(Z)} w \cdot P(\{\omega \in \Omega \mid (W + Z)(\omega) = w + z\}) + \\ &\quad \sum_{w \in Val(W)} \sum_{z \in Val(Z)} z \cdot P(\{\omega \in \Omega \mid (W + Z)(\omega) = w + z\}) \\ &= \sum_{w \in Val(W)} w \sum_{z \in Val(Z)} P(\{\omega \in \Omega \mid W(\omega) = w, Z(\omega) = z\}) + \\ &\quad \sum_{z \in Val(Z)} z \sum_{w \in Val(W)} P(\{\omega \in \Omega \mid W(\omega) = w, Z(\omega) = z\}) \\ &= \sum_{w \in Val(W)} w P(\{\omega \in \Omega \mid W(\omega) = w\}) + \sum_{z \in Val(Z)} z P(\{\omega \in \Omega \mid Z(\omega) = z\}) \end{aligned}$$

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

Pentru cazul continuu:

$$\begin{aligned} E[W + Z] &= \int_w \int_z (w + z) p_{WZ}(w, z) dz dw \\ &= \int_w \int_z w p_{WZ}(w, z) dz dw + \int_w \int_z z p_{WZ}(w, z) dz dw \\ &= \int_w w \int_z p_{WZ}(w, z) dz dw + \int_z z \int_w p_{WZ}(w, z) dw dz \\ &= \int_w w p_W(w) dw + \int_z z p_Z(z) dz \\ &= E[W] + E[Z]. \end{aligned}$$

În cazul în care variabila aleatoare X este discretă, egalitatea $E[aX] = aE[X]$ se poate demonstra imediat, aplicând definiția mediei. Vom considera mai întâi subcazul $a \neq 0$:

$$\begin{aligned} E[aX] &= \sum_{u \in \text{Val}(aX)} u P(aX = u) = \sum_{x \in \text{Val}(X)} ax P(X = x), \quad \text{unde } x = \frac{1}{a}u \\ &= a \sum_{x \in \text{Val}(X)} x P(X = x) = aE[X]. \end{aligned}$$

Pentru subcazul $a = 0$, egalitatea $E[aX] = aE[X]$ este trivială.

Similar se face demonstrația egalității $E[aX] = aE[X]$ în cazul în care variabila aleatoare X este continuă.

b. Pentru a demonstra această proprietate vom ține cont de liniaritatea mediei (vedeți *Observația* de la punctul a):

$$\begin{aligned} \text{Var}(X) &= E[(X - \bar{X})^2] = E[X^2 - 2X\bar{X} + \bar{X}^2] \\ &= E[X^2] - E[2X\bar{X}] + E[\bar{X}^2] = E[X^2] - 2\bar{X}E[X] + \bar{X}^2 \\ &= E[X^2] - 2\bar{X}^2 + \bar{X}^2 = E[X^2] - \bar{X}^2 = E[X^2] - (E[X])^2. \end{aligned}$$

Egalitatea $\text{Var}(aX) = a^2 \text{Var}(X)$ rezultă imediat, aplicând definiția varianței și ținând cont de proprietatea $E[aX] = aE[X]$ pe care am demonstrat-o la punctul a.

$$\begin{aligned} \text{Var}(aX) &\stackrel{\text{def.}}{=} E[(aX - \underbrace{E[aX]}_{aE[X]})^2] = E[(a(X - E[X]))^2] \\ &= E[a^2(X - E[X])^2] = a^2 E[(X - E[X])^2] \stackrel{\text{def.}}{=} a^2 \text{Var}(X). \end{aligned}$$

Alternativ, putem folosi proprietatea pe care am demonstrat-o mai sus:

$$\begin{aligned} \text{Var}(aX) &= E[(aX)^2] - (E[aX])^2 = E[a^2 X^2] - (aE[X])^2 = a^2 E[X^2] - a^2 (E[X])^2 \\ &= a^2 (E[X^2] - (E[X])^2) = a^2 \text{Var}(X). \end{aligned}$$

c. Se folosește același gen de raționament ca la punctul precedent (prima parte):

$$\begin{aligned} \text{Cov}(X, Y) &= E[(X - E[X])(Y - E[Y])] \\ &= E[XY - XE[Y] - E[X]Y + E[X]E[Y]] \\ &= E[XY] - E[XE[Y]] - E[E[X]Y] + E[E[X]E[Y]] \\ &= E[XY] - E[X]E[Y] - E[X]E[Y] + E[X]E[Y] = E[XY] - E[X]E[Y]. \end{aligned}$$

10.

(Rezultat teoretic:

covarianța oricăror 2 variabile aleatoare independente este 0)

■ CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW1, pr. 1.2

În mod intuitiv, două variabile aleatoare X și Y sunt *independente* atunci când cunoașterea valorii uneia dintre ele (de exemplu X) nu furnizează niciun indiciu despre valoarea celeilalte variabile (Y , în acest caz).

Formal, dacă X și Y sunt variabile aleatoare discrete, independența lor revine la egalitatea $P(X = x, Y = y) = P(X = x) \cdot P(Y = y)$ pentru orice $x \in \text{Val}(X)$ și orice $y \in \text{Val}(Y)$.

Similar, dacă X și Y sunt variabile aleatoare continue, atunci $p(X = x, Y = y) = p(X = x) \cdot p(Y = y)$ pentru orice valori x și y posibile.

Arătați că dacă X și Y sunt variabile aleatoare independente de același tip (adică fie discret fie continuu), atunci

$$E[XY] = E[X] \cdot E[Y].$$

Echivalent: X, Y independente $\Rightarrow \text{Cov}(X, Y) = 0$. (A se vedea problema 9 punctul c.)

Observație: Reciproca implicației de mai sus nu este în general adevărată. A se vedea problema 11.

Răspuns:

Pentru cazul în care variabilele aleatoare independente X și Y sunt discrete, putem scrie:

$$\begin{aligned} E[XY] &= \sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} xy P(X = x, Y = y) \\ &\stackrel{\text{indep.}}{=} \sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} xy P(X = x) \cdot P(Y = y) \\ &= \sum_{x \in \text{Val}(X)} \left(x P(X = x) \sum_{y \in \text{Val}(Y)} y P(Y = y) \right) \\ &= \sum_{x \in \text{Val}(X)} x P(X = x) E[Y] = \left(\sum_{x \in \text{Val}(X)} x P(X = x) \right) E[Y] \\ &= E[X] \cdot E[Y] \end{aligned}$$

Pentru cazul continuu demonstrația este similară:

$$\begin{aligned} E[XY] &= \int_x \int_y xy p(X = x, Y = y) dy dx \\ &\stackrel{\text{indep.}}{=} \int_x \int_y xy p(X = x) \cdot p(Y = y) dy dx \\ &= \int_x x p(X = x) \int_y y p(Y = y) dy dx \\ &= \int_x x p(X = x) E[Y] dx = E[Y] \cdot \int_x x p(X = x) dx \\ &= E[Y] \cdot E[X] = E[X] \cdot E[Y] \end{aligned}$$

11.

(Covarianța nulă nu implică în mod neapărat independența variabilelor aleatoare)

CMU, 2009 spring, Tom Mitchell, HW2, pr. 1.3.2

CMU, 2009 fall, Geoff Gordon, HW1, pr. 3.1

a. Reciproca afirmației din problema 10 nu este în general adevărată, deci $Cov(X, Y) = 0 \not\Rightarrow X$ și Y sunt independente. Arătați aceasta folosind ca exemplu variabilele aleatoare ale căror distribuții sunt date în tabelul alăturat.

X	Y	$P(X, Y)$
0	0	1/3
1	0	0
2	0	1/3
0	1	0
1	1	1/3
2	1	0

b. Totuși, dacă X și Y sunt variabile aleatoare binare luând valori în mulțimea $\{0, 1\}$, iar $E[XY] = E[X] \cdot E[Y]$, atunci X și Y sunt independente. Justificați.

(Așadar, două variabile aleatoare binare cu valori în mulțimea $\{0, 1\}$ sunt independente atunci și numai atunci când au covarianța nulă.)

Răspuns:

a. Din datele furnizate în exemplul a rezultă următoarele probabilități:

$$\begin{aligned} P(X=0) &= 1/3 & P(Y=0) &= 2/3 & P(XY=0) &= 2/3 \\ P(X=1) &= 1/3 & P(Y=1) &= 1/3 & P(XY=1) &= 1/3 \\ P(X=2) &= 1/3 & & & P(XY=2) &= 0 \end{aligned}$$

Putem calcula mediile acestor variabile aleatoare folosind formula de definiție pentru variabile discrete $E[X] = \sum_x x \cdot P(X=x)$:

$$\begin{aligned} E[X] &= 0 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} = 1 \\ E[Y] &= 0 \cdot \frac{2}{3} + 1 \cdot \frac{1}{3} = \frac{1}{3} \\ E[XY] &= 0 \cdot \frac{2}{3} + 1 \cdot \frac{1}{3} + 2 \cdot 0 = \frac{1}{3}. \end{aligned}$$

Conform formulei care a fost demonstrată la problema 10, covarianța variabilelor X și Y este: $Cov(X, Y) = E[XY] - E[X] \cdot E[Y] = \frac{1}{3} - 1 \cdot \frac{1}{3} = 0$. Cu toate acestea, variabilele X și Y nu sunt independente. Într-adevăr, pentru $X=0$ și $Y=0$ se observă că $P(X=0, Y=0) = 1/3$ dar $P(X=0)P(Y=0) = 1/3 \cdot 2/3 = 2/9 \neq \frac{1}{3}$.

Prin acest contraexemplu am arătat că implicația „ $Cov(X, Y) = 0 \Rightarrow X$ și Y sunt independente” nu este în general adevărată.

b. În continuare vom demonstra că în cazul în care X și Y iau valori în mulțimea $\{0, 1\}$ implicația de mai sus este adevărată.

Dacă X și Y sunt variabile aleatoare binare, atunci:

$$\begin{aligned} E[X] &= 0 \cdot P(X=0) + 1 \cdot P(X=1) = P(X=1) \\ E[Y] &= P(Y=1) \\ E[XY] &= P(X=1, Y=1) \end{aligned}$$

Covarianța nulă înseamnă — conform problemelor 9 și 10 — că $E[XY] = E[X] \cdot E[Y]$, adică:

$$P(X=1, Y=1) = P(X=1)P(Y=1)$$

Însă, a demonstra independența variabilelor aleatoare X și Y revine la a arăta că $P(X, Y) = P(X)P(Y)$ pentru toate combinațiile posibile de valori ale variabilelor. Unul din cazuri ($X = 1, Y = 1$) este deja demonstrat, deci mai există încă 3 cazuri. Pentru acestea vom utiliza formulele: $P(A \cap B) = P(A) - P(A \cap \bar{B})$ și $P(\bar{A}) = 1 - P(A)$.

Cazul $X = 1, Y = 0$:

$$\begin{aligned} P(X = 1, Y = 0) &= P(X = 1) - P(X = 1, Y = 1) \\ &= P(X = 1) - P(X = 1)P(Y = 1) \\ &= P(X = 1)(1 - P(Y = 1)) \\ &= P(X = 1)P(Y = 0) \end{aligned}$$

Cazul $X = 0, Y = 1$ se tratează similar cu cazul anterior:

$$\begin{aligned} P(X = 0, Y = 1) &= P(Y = 1) - P(X = 1, Y = 1) \\ &= P(Y = 1) - P(X = 1)P(Y = 1) \\ &= P(Y = 1)(1 - P(X = 1)) \\ &= P(Y = 1)P(X = 0) \end{aligned}$$

Cazul $X = 0, Y = 0$:

$$\begin{aligned} P(X = 0, Y = 0) &= P(X = 0) - P(X = 0, Y = 1) \\ &= P(X = 0) - P(X = 0)P(Y = 1) \\ &= P(X = 0)(1 - P(Y = 1)) \\ &= P(X = 0)P(Y = 0) \end{aligned}$$

Prin urmare, egalitatea $P(X, Y) = P(X)P(Y)$ este adevărată pentru toate cazurile, deci variabilele X și Y sunt independente.

12. (Variabile aleatoare: calcul de medii)

Fie X o variabilă aleatoare pentru care $E(X) = \mu$ și $Var(X) = \sigma^2$.

a. *CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 2.4*

Cât este $E[X(X - 1)]$ în funcție de μ și σ ?

b. *CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, midterm, pr. 1.c*

Fie $c \in \mathbb{R}$. Care din următoarele variante sunt adevărate?

- | | |
|--------------------------------------------|---------------------------------------------|
| A. $E[(X - c)^2] = (\mu - c)^2 + \sigma^2$ | D. $E[(X - c)^2] = (\mu - c)^2 + 2\sigma^2$ |
| B. $E[(X - c)^2] = (\mu - c)^2$ | E. $E[(X - c)^2] = \mu^2 + c^2 + 2\sigma^2$ |
| C. $E[(X - c)^2] = (\mu - c)^2 - \sigma^2$ | F. $E[(X - c)^2] = \mu^2 + c^2 - 2\sigma^2$ |

Răspuns:

a. De la problema 9.a știm că media sumei a două variabile aleatoare este suma mediilor variabilelor respective. De asemenea, este imediat demonstrabil că

$E[c \cdot X] = c \cdot E[X]$, unde X, Y sunt variabile aleatoare iar c este o constantă. Ținând cont că $Var(X) = E[X^2] - (E[X])^2$ (vedeți problema 9.b), putem scrie:

$$\begin{aligned} E[X(X-1)] &= E[X^2 - X] = E[X^2] - E[X] \\ &= E[X^2] - (E[X])^2 + (E[X])^2 - E[X] \\ &= Var(X) + (E[X])^2 - E[X] \\ &= \sigma^2 + \mu^2 - \mu \end{aligned}$$

b. Pentru a găsi varianta adevărată vom calcula $E[(X-c)^2]$:

$$\begin{aligned} E[(X-c)^2] &= E[X^2 - 2cX + c^2] = E[X^2] - 2cE[X] + c^2 \\ &= E[X^2] - (E[X])^2 + (E[X])^2 - 2cE[X] + c^2 \\ &= \sigma^2 + \mu^2 - 2c\mu + c^2 \\ &= \sigma^2 + (\mu - c)^2 \end{aligned}$$

Deci varianta A este adevărată. (Toate celelalte variante sunt, în general, false.)

13.

(Variabile aleatoare discrete:
distribuții de probabilitate comune, marginale, condiționate;
regula de înlănțuire)

CMU, 2002 fall, Andrew Moore, final exam, pr. 4.a

Considerăm un set de date definit cu ajutorul a 3 variabile aleatoare cu valori booleene X, Y și Z . Care dintre seturile de informații de mai jos sunt suficiente pentru a specifica distribuția comună $P(x, y, z)$?

A.	B.	C.	D.
$P(\neg X \mid Z)$	$P(\neg X \mid \neg Z)$	$P(X \mid Z)$	$P(X \mid Z)$
$P(\neg X \mid \neg Z)$	$P(X \mid \neg Z)$	$P(X \mid \neg Z)$	$P(X \mid \neg Z)$
$P(\neg Y \mid X, Z)$	$P(Y \mid X, Z)$	$P(Y \mid X, Z)$	$P(Y \mid X, Z)$
$P(\neg Y \mid X, \neg Z)$	$P(Y \mid X, \neg Z)$	$P(Y \mid X, \neg Z)$	$P(Y \mid X, \neg Z)$
$P(\neg Y \mid \neg X, Z)$	$P(Y \mid \neg X, Z)$	$P(Y \mid \neg X, Z)$	$P(\neg Y \mid \neg X, \neg Z)$
$P(\neg Y \mid \neg X, \neg Z)$	$P(Y \mid \neg X, \neg Z)$	$P(\neg Y \mid \neg X, \neg Z)$	$P(Y \mid \neg X, \neg Z)$
$P(Z)$	$P(Z)$	$P(\neg Z)$	$P(Z)$

Răspuns:

Pentru a calcula distribuția comună a mai multor variabile aleatoare se poate aplica regula de înlănțuire.¹⁵ În cazul nostru, putem scrie:

$$P(X, Y, Z) = P(Z) \cdot P(X \mid Z) \cdot P(Y \mid X, Z)$$

Deoarece variabilele aleatoare X, Y și Z au valori booleene, pentru a specifica distribuția comună $P(X, Y, Z)$ este nevoie să se calculeze valoarea acestora în fiecare din cele 8 cazuri posibile:

¹⁵Pentru variabile aleatoare, regula de înlănțuire

$$P(A_1, A_2, \dots, A_n) = P(A_1) \cdot P(A_2 \mid A_1) \cdot P(A_n \mid A_1, A_2, \dots, A_{n-1})$$

se demonstrează imediat pornind de la regula de înlănțuire pentru (probabilități de) evenimente aleatoare. A se vedea problema 14.

$$\begin{array}{cccc} P(X, Y, Z) & P(X, Y, \neg Z) & P(X, \neg Y, Z) & P(X, \neg Y, \neg Z) \\ P(\neg X, Y, Z) & P(\neg X, Y, \neg Z) & P(\neg X, \neg Y, Z) & P(\neg X, \neg Y, \neg Z) \end{array}$$

Cunoaștem de asemenea relații de calcul de forma:

$$\begin{aligned} P(\neg X) &= 1 - P(X) \\ P(\neg X | Y) &= 1 - P(X | Y) \end{aligned}$$

Așadar, pentru a aplica regula de înlănțuire de mai sus este nevoie să cunoaștem

$$\begin{array}{ll} P(Z) \text{ sau } P(\neg Z); & P(Y | X, Z) \text{ sau } P(\neg Y | X, Z); \\ P(X | Z) \text{ sau } P(\neg X | Z); & P(Y | \neg X, Z) \text{ sau } P(\neg Y | \neg X, Z); \\ P(X | \neg Z) \text{ sau } P(\neg X | \neg Z); & P(Y | X, \neg Z) \text{ sau } P(\neg Y | X, \neg Z); \\ & P(Y | \neg X, \neg Z) \text{ sau } P(\neg Y | \neg X, \neg Z). \end{array}$$

Cu aceste precizări, putem specifica pentru fiecare dintre seturile de informații din enunț dacă sunt suficiente pentru a calcula distribuția comună $P(X, Y, Z)$.

Cazul A. Da. Se observă că putem calcula $P(X | Z)$ și $P(X | \neg Z)$ din primele două probabilități din enunț. De asemenea, utilizând următoarele 4 probabilități se pot deduce: $P(Y | X, Z)$, $P(Y | X, \neg Z)$, $P(Y | \neg X, Z)$ și $P(Y | \neg X, \neg Z)$. Iar din $P(Z)$ se obține $P(\neg Z)$. Prin urmare, există toate informațiile necesare distribuției comune $P(X, Y, Z)$.

Cazul B. Nu, informațiile din enunț nu sunt suficiente. Nu putem deduce valoarea pentru $P(X | Z)$.

Cazul C. Da, informațiile din enunț sunt suficiente. Din $P(X | Z)$ se obține $P(\neg X | Z)$, iar din $P(X | \neg Z)$ se obține $P(\neg X | \neg Z)$. Din următoarele 4 probabilități se obțin celelalte 4 necesare pentru $P(Y | X, Z)$, și anume: $P(\neg Y | X, Z)$, $P(\neg Y | X, \neg Z)$, $P(\neg Y | \neg X, Z)$ și respectiv $P(Y | \neg X, \neg Z)$.

Cazul D. Nu, informațiile din enunț nu sunt suficiente. Nu putem deduce valoarea pentru $P(Y | \neg X, Z)$.

14. (Variabile aleatoare discrete:
regula de multiplicare, varianta condițională)
CMU, 2009 fall, Geoff Gordon, HW2, pr. 2.1

Arătați că pentru orice valori x, y și z ale variabilelor aleatoare X, Y și Z respectiv, avem:

$$P(X = x, Y = y | Z = z) = P(X = x | Y = y, Z = z) \cdot P(Y = y | Z = z).$$

În notație simplificată: $P(X, Y | Z) = P(X | Y, Z) \cdot P(Y | Z)$.

Indicație: Folosiți regula de înlănțuire pentru evenimente aleatoare.

Răspuns:

Folosind definiția probabilității condiționate și regula de înlănțuire (cu termenii ordonați în mod convenabil), egalitatea cerută se obține astfel:

$$\begin{aligned}
P(X, Y | Z) &\stackrel{\text{def.}}{=} \frac{P(X, Y, Z)}{P(Z)} \stackrel{\text{not.}}{=} \frac{P(Z \cap Y \cap X)}{P(Z)} \\
&= \frac{P(Z)P(Y | Z)P(X | Y, Z)}{P(Z)} = P(Y | Z) \cdot P(X | Y, Z)
\end{aligned}$$

Observație: O altă metodă de rezolvare constă în a aplica pentru fiecare membru al egalității din enunț definiția probabilității condiționate, după simplificări obținându-se pentru ambii membri aceeași valoare:

$$\begin{aligned}
P(X, Y | Z) &= \frac{P(X, Y, Z)}{P(Z)} \\
P(X | Y, Z) \cdot P(Y | Z) &= \frac{P(X, Y, Z)}{P(Y, Z)} \cdot \frac{P(Y, Z)}{P(Z)} = \frac{P(X, Y, Z)}{P(Z)}
\end{aligned}$$

15. (Variabile aleatoare discrete: independență condițională)
CMU, 2005 fall, T. Mitchell, A. Moore, midterm, pr. 4

Fie variabilele aleatoare discrete A , B și C având distribuția comună conform tabelului de mai jos.

- a. Este variabila A independentă condițional de B în raport cu variabila C ?
- b. Dacă ați răspuns afirmativ la întrebarea a , faceți o schimbare în primele două linii ale tabelului de mai sus pentru a obține o distribuție pentru care răspunsul la aceeași întrebare să devină negativ. Invers, dacă ați răspuns negativ la întrebarea a , faceți o schimbare în primele două linii ale tabelului încât răspunsul să devină afirmativ.

A	B	C	$P(A, B, C)$
0	0	0	1/8
0	0	1	1/8
0	1	0	1/8
0	1	1	1/8
1	0	0	1/8
1	0	1	1/8
1	1	0	1/8
1	1	1	1/8

Răspuns:

- a. Faptul că variabila A este independentă condițional de B în raport cu variabila C se mai notează prin $A \perp B | C$ și poate fi demonstrat prin una din următoarele două relații:

$$\begin{aligned}
P(A = a, B = b | C = c) &= P(A = a | C = c) \cdot P(B = b | C = c) \text{ sau} \\
P(A = a | B = b, C = c) &= P(A = a | C = c), \text{ dacă } P(B = b, C = c) \neq 0.
\end{aligned}$$

pentru orice $a \in \text{Val}(A)$, $b \in \text{Val}(B)$, $c \in \text{Val}(C)$. Deși în general se folosește prima relație, pentru acest exercițiu este mai ușor să utilizăm cea de-a doua relație. Conform tabelului dat în enunț, rezultă imediat că $P(B = b, C = c) \neq 0$ pentru orice $b, c \in \{0, 1\}$. Vom demonstra că pentru orice $a, b, c \in \{0, 1\}$ este adevărat că $P(A = a | B = b, C = c) = P(A = a | C = c)$. Cele două probabilități condiționate vor fi calculate folosind datele din tabel:

$$\text{Cazul } (0, 0, 0): P(A = 0 | B = 0, C = 0) = \frac{1 \cdot \frac{1}{8}}{2 \cdot \frac{1}{8}} = \frac{1}{2} \text{ și } P(A = 0 | C = 0) = \frac{2 \cdot \frac{1}{8}}{4 \cdot \frac{1}{8}} = \frac{1}{2}.$$

Cazul $(0, 0, 1)$: $P(A = 0 \mid B = 0, C = 1) = \frac{1 \cdot \frac{1}{8}}{2 \cdot \frac{1}{8}} = \frac{1}{2}$ și $P(A = 0 \mid C = 1) = \frac{2 \cdot \frac{1}{8}}{4 \cdot \frac{1}{8}} = \frac{1}{2}$.

Se observă că pentru toate celelalte cazuri se obțin aceleași valori, deci

$$P(A = a \mid B = b, C = c) = P(A = a \mid C = c), \forall a, b, c \in \{0, 1\}.$$

Așadar, variabila A este independentă condițional de B în raport cu variabila C .

b. Schimbarea trebuie făcută în așa fel încât să se păstreze relația $\sum P(A, B, C) = 1$. O variantă posibilă este:

A	B	C	$P(A, B, C)$
0	0	0	$1/4$
0	0	1	0
...

Pentru aceste noi valori se observă că: $P(A = 0 \mid B = 0, C = 0) = \frac{1 \cdot 1/4}{1 \cdot 1/4 + 1 \cdot 1/8} = \frac{2}{3}$ și $P(A = 0 \mid C = 0) = \frac{1 \cdot 1/4 + 1 \cdot 1/8}{1 \cdot 1/4 + 3 \cdot 1/8} = \frac{3}{5}$. Este suficient un singur caz în care probabilitățile respective nu sunt egale, prin urmare variabila A nu este independentă condițional de variabila B în raport cu a treia variabilă, C .

16.

(Variabile aleatoare discrete:
distribuții comune, distribuții marginale, distribuții condiționale;
independență, independență condițională)

• CMU, 2016 fall, N. Balcan, M. Gormley, HW2, pr. 1.4

Fie trei variabile aleatoare X, Y și Z care iau valori în mulțimea $\{0, 1\}$. În tabelul de mai jos este dată distribuția probabilistă comună a acestor trei variabile, $P(X, Y, Z)$.

	$Z = 0$		$Z = 1$	
	$X = 0$	$X = 1$	$X = 0$	$X = 1$
$Y = 0$	$1/24$	$1/12$	$1/12$	$5/24$
$Y = 1$	$1/12$	p	q	$7/24$

a. Considerând că X și Y sunt independente, găsiți valorile lui p și q .

Indicație: Este util (deși nu obligatoriu) să calculați mai întâi $P(X, Y)$, distribuția comună a variabilelor X și Y , completând tabelul alăturat, după care veți calcula și distribuțiile (marginale) pentru X și Y , de preferință ca o linie și respectiv o coloană suplimentară la acest tabel.

	$X = 0$	$X = 1$
$Y = 0$		
$Y = 1$		

b. Considerând valorile lui p și q determinate la punctul a, sunt X și Y independente condițional în raport cu Z ? De ce?

Răspuns:

a. Conform definiției, variabilele aleatoare X și Y sunt independente dacă și numai dacă

$$P(X = x, Y = y) = P(X = x) \cdot P(Y = y), \forall x \in \text{Val}(X), \forall y \in \text{Val}(Y). \quad (5)$$

Observație: Evident, pentru a determina p și q , ar fi de dorit să scriem un sistem de două ecuații cu aceste două necunoscute. Cele două ecuații pot fi obținute prin particularizarea relației (5) pentru două perechi distincte de valori $(x, y) \in \text{Val}(X) \times \text{Val}(Y)$. Alternativ, va fi suficient să reținem doar o astfel de ecuație, pentru că în locul celeilalte putem folosi una din proprietățile din definiția funcției de probabilitate, și anume

$$\sum_{x,y,z} P(X = x, Y = y, Z = z) = 1. \quad (6)$$

Aceasta revine, practic la a însuma elementele din tabelul dat în enunț pentru distribuția $P(X, Y, Z)$.¹⁶

Pentru a calcula probabilitățile marginale $P(X = x, Y = y)$, vom proceda conform definiției: $P(X = x, Y = y) = \sum_{z \in \text{Val}(Z)} P(X = x, Y = y, Z = z)$. Și anume:

$$\begin{aligned} P(X = 0, Y = 0) &= P(X = 0, Y = 0, Z = 0) + P(X = 0, Y = 0, Z = 1) = \frac{1}{24} + \frac{1}{12} = \frac{1}{8} \\ P(X = 0, Y = 1) &= P(X = 0, Y = 1, Z = 0) + P(X = 0, Y = 1, Z = 1) = \frac{1}{12} + q \\ P(X = 1, Y = 0) &= P(X = 1, Y = 0, Z = 0) + P(X = 1, Y = 0, Z = 1) = \frac{1}{12} + \frac{5}{24} = \frac{7}{24} \\ P(X = 1, Y = 1) &= P(X = 1, Y = 1, Z = 0) + P(X = 1, Y = 1, Z = 1) = \frac{7}{24} + p. \end{aligned}$$

În mod similar,

$$\begin{aligned} P(X = 0) &= \sum_{y \in \text{Val}(Y)} P(X = x, Y = y) = \frac{1}{8} + \frac{1}{12} + q = \frac{5}{24} + q \\ P(X = 1) &= \frac{7}{24} + p + \frac{7}{24} = \frac{7}{12} + p \end{aligned}$$

și

$$\begin{aligned} P(Y = 0) &= \sum_{x \in \text{Val}(X)} P(X = x, Y = y) = \frac{1}{8} + \frac{7}{24} = \frac{5}{12} \\ P(Y = 1) &= 1 - P(Y = 0) = \frac{7}{12}. \end{aligned}$$

Punând acum toate aceste rezultate împreună, obținem:

	$X = 0$	$X = 1$	
$Y = 0$	$1/8$	$7/24$	$5/12$
$Y = 1$	$1/12 + q$	$7/24 + p$	$7/12$
	$5/24 + q$	$7/12 + p$	

¹⁶Remarcăm că, dacă vom urma *Indicația* din enunț, după ce vom scrie probabilitățile marginale $P(X, Y)$, $P(X)$ și $P(Y)$, va fi ușor să scriem încă alte trei proprietăți similare cu (6). Deci, o vom putea alege atunci pe cea mai „directă”/simplă dintre ele.

Acum este simplu de văzut că $5/24 + q + p + 14/24 = 1 \Leftrightarrow p + q = 5/24$. (Aceasta constituie prima noastră ecuație în p și q .) În consecință, am putea chiar să eliminăm q din tabelul de mai sus, scriind $P(X = 0, Y = 1) = 1/12 + q = 1/12 + 5/24 - p = 7/24 - p$ și $P(X = 0) = 5/24 + q = 5/24 + 5/24 - p = 5/12 - p$.

Aplicând definiția independenței variabilelor X și Y pentru $x = 0$ și $y = 0$, vom avea:

$$\begin{aligned} P(X = 0, Y = 0) &= P(X = 0) \cdot P(Y = 0) \Leftrightarrow \frac{1}{8} = \left(\frac{5}{12} - p\right) \cdot \frac{5}{12} \Leftrightarrow \frac{5}{12} - p = \frac{3}{2} \cdot \frac{1}{5} \\ &\Leftrightarrow p = \frac{5}{12} - \frac{3}{10} = \frac{25 - 18}{60} = \frac{7}{60}, \end{aligned}$$

de unde rezultă imediat

$$q = \frac{5}{24} - p = \frac{5}{24} - \frac{7}{60} = \frac{25 - 14}{120} = \frac{11}{120}.$$

De asemenea, vom avea $P(X = 0) = 5/12 - p = 5/12 - 7/60 = 3/10$ și $P(X = 1) = 7/10$.

În final, va trebui să verificăm și celelalte trei egalități din definiția independenței lui X și Y , pentru că, în general, este posibil ca sistemul de ecuații corespunzător relației (5) să fie supra-rescricționat (deci, incompatibil):¹⁷

$$P(X = 0, Y = 1) = P(X = 0) \cdot P(Y = 1) :$$

$$\frac{7}{24} - p = \frac{3}{10} \cdot \frac{7}{12} \Leftrightarrow \frac{7}{24} - \frac{7}{60} = \frac{7}{40} \Leftrightarrow 7 \cdot \left(\frac{1}{24} - \frac{1}{60}\right) = \frac{7}{40} \Leftrightarrow 7 \cdot \frac{5 - 2}{120} = \frac{7}{40} \quad (A)$$

$$P(X = 1, Y = 0) = P(X = 1) \cdot P(Y = 0) :$$

$$\frac{7}{24} = \frac{7}{10} \cdot \frac{5}{12} \quad (A)$$

$$P(X = 1, Y = 1) = P(X = 1) \cdot P(Y = 1) :$$

$$\frac{7}{24} + p = \frac{7}{10} \cdot \frac{7}{12} \Leftrightarrow \frac{7}{24} + \frac{7}{60} = 7^2 \cdot \frac{1}{10} \cdot \frac{1}{12} \Leftrightarrow \frac{1}{24} + \frac{1}{60} = \frac{7}{120} \Leftrightarrow \frac{5 + 2}{120} = \frac{7}{120} \quad (A)$$

Așadar, date fiind cele două valori ale lui p și q , variabilele X și Y sunt independente.

b. Conform definiției, variabilele aleatoare X și Y sunt independente condițional în raport cu variabila Z dacă și numai dacă

$$\begin{aligned} P(X = x, Y = y | Z = z) &= P(X = x | Z = z) \cdot P(Y = y | Z = z), \\ &\forall x \in \text{Val}(X), \forall y \in \text{Val}(Y), \forall z \in \text{Val}(Z). \end{aligned} \quad (7)$$

Dacă X și Y nu sunt independente condițional în raport cu Z , atunci $\exists x \in \text{Val}(X), \exists y \in \text{Val}(Y), \exists z \in \text{Val}(Z)$ astfel încât $P(X = x, Y = y | Z = z) \neq P(X = x | Z = z) \cdot P(Y = y | Z = z)$.

Pentru a putea exprima probabilitățile condiționate care apar în formulele de mai sus, vom calcula în prealabil distribuția marginală a lui Z pornind de la tabelul din enunț:

$$\begin{aligned} P(Z = 0) &= \sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} P(X = x, Y = y, Z = 0) \\ &= 5/24 + p = 5/24 + 7/60 = 39/120 = 13/40 \\ P(Z = 1) &= 1 - P(Z = 0) = 27/40. \end{aligned}$$

¹⁷Ar fi suficient chiar să verificăm doar două dintre cele trei egalități, fiindcă ultima decurge automat din celelalte, ținând cont de proprietățile $\sum_{x,y} P(X = x, Y = y) = 1$, $\sum_x P(X = x) = 1$ și $\sum_y P(Y = y) = 1$.

Vom verifica relația (7) pentru $x = y = z = 0$:

$$\begin{aligned}
 P(X = 0, Y = 0 | Z = 0) &= P(X = 0 | Z = 0) \cdot P(Y = 0 | Z = 0) \\
 \Leftrightarrow \frac{P(X = 0, Y = 0, Z = 0)}{P(Z = 0)} &= \frac{P(X = 0, Z = 0)}{P(Z = 0)} \cdot \frac{P(Y = 0, Z = 0)}{P(Z = 0)} \\
 \Leftrightarrow P(X = 0, Y = 0, Z = 0) \cdot P(Z = 0) &= P(X = 0, Z = 0) \cdot P(Y = 0, Z = 0) \\
 \Leftrightarrow \frac{1}{24} \cdot \frac{13}{40} &= \left(\frac{1}{24} + \frac{1}{12} \right) \cdot \left(\frac{1}{24} + \frac{1}{12} \right) \\
 \Leftrightarrow \frac{13}{24 \cdot 40} = \frac{1}{8} \cdot \frac{1}{8} &\Leftrightarrow \frac{13}{3 \cdot 5} = 1 \quad (F)
 \end{aligned}$$

Prin urmare, variabilele X și Y nu sunt independente condițional în raport cu variabila Z .

Observație: Această problemă pune în evidență următoarea *proprietate*, care merită a fi reținută: independența a două variabile aleatoare nu implică (în general) independența lor condițională în raport cu o a treia variabilă aleatoare.

17.

(Variabile aleatoare continue:
funcția densitate de probabilitate)

CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 1.5

Fie variabila aleatoare continuă X a cărei funcție densitate de probabilitate (în limba engleză: “probability density functions”; scris, sub formă abreviată, p.d.f.) este:

$$p(x) = \begin{cases} cx^2 & \text{pentru } 1 \leq x \leq 2 \\ 0 & \text{în caz contrar.} \end{cases}$$

a. Ținând cont de proprietățile funcției densitate de probabilitate (a se vedea notițele de la curs), cât trebuie să fie valoarea constantei c ?

b. Desenați graficul funcției de mai sus.

c. Calculați $P(X > 3/2)$.

Răspuns:

a. Faptul că $p(x)$ este funcție densitate de probabilitate pentru variabila aleatoare continuă X înseamnă că $p(x) \geq 0, \forall x$ și că $\int_{-\infty}^{+\infty} p(x)dx = 1$. Pentru a afla constanta c vom calcula valoarea integralei:

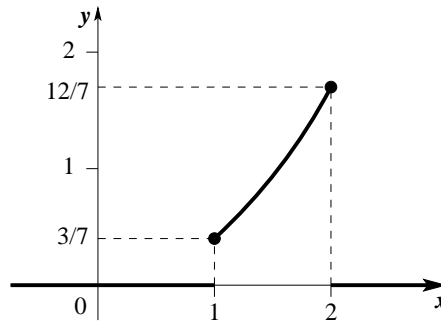
$$\begin{aligned}
 \int_{-\infty}^{+\infty} p(x)dx &= \underbrace{\int_{-\infty}^1 p(x)dx}_{=0} + \int_1^2 p(x)dx + \underbrace{\int_2^{+\infty} p(x)dx}_{=0} \\
 &= \int_1^2 cx^2dx = c \cdot \int_1^2 x^2dx = c \cdot \left. \frac{x^3}{3} \right|_1^2 = c \left(\frac{2^3}{3} - \frac{1^3}{3} \right) = c \cdot \frac{7}{3}
 \end{aligned}$$

Prin urmare, $c \cdot \frac{7}{3} = 1 \Rightarrow c = \frac{3}{7}$.

b. Trebuie să reprezentăm grafic funcția:

$$p(x) = \begin{cases} \frac{3}{7}x^2 & \text{pentru } 1 \leq x \leq 2 \\ 0 & \text{în caz contrar.} \end{cases}$$

În intervalul $[1, 2]$, această funcție este un fragment din parabola corespunzătoare funcției de gradul al doilea: $\frac{3}{7}x^2$, parabolă care are vârful în punctul $(0, 0)$. Putem calcula $p(1) = 3/7 \approx 0.42$ și $p(2) = 12/7 \approx 1.71$.



c. Valoarea probabilității cerute se poate calcula astfel:

$$\begin{aligned} P(X > 3/2) &\stackrel{\text{not.}}{=} P(\{\omega \mid X(\omega) > 3/2\}) \stackrel{\text{def.}}{=} \int_{X(\omega)=x>3/2} p(x)dx = \int_{3/2}^{+\infty} p(x)dx \\ &= \int_{3/2}^2 p(x)dx + \underbrace{\int_2^{+\infty} p(x)dx}_{=0} = \int_{3/2}^2 \frac{3}{7}x^2dx = \frac{3}{7} \cdot \int_{3/2}^2 x^2dx \\ &= \frac{3}{7} \cdot \left. \frac{x^3}{3} \right|_{3/2}^2 = \frac{1}{7} \cdot x^3 \Big|_{3/2}^2 = \frac{1}{7} \cdot \left(8 - \frac{27}{8} \right) = \frac{1}{7} \cdot \frac{64 - 27}{8} = \frac{37}{56}. \end{aligned}$$

O altă variantă de rezolvare este bazată pe folosirea *funcției cumulative de distribuție*, care, după cum știm, se definește prin relația $F(x) \stackrel{\text{def.}}{=} P(X \leq x)$, pentru orice $x \in \mathbb{R}$:

$$\begin{aligned} P(X > 3/2) &= 1 - P(X \leq 3/2) = 1 - F\left(\frac{3}{2}\right) = 1 - \int_{-\infty}^{3/2} p(x)dx \\ &= 1 - \left(\underbrace{\int_{-\infty}^1 p(x)dx}_0 + \int_1^{3/2} p(x)dx \right) = 1 - \left(0 + \int_1^{3/2} p(x)dx \right) \\ &= 1 - \int_1^{3/2} p(x)dx = 1 - \int_1^{3/2} \frac{3}{7}x^2dx = 1 - \frac{1}{7} \cdot x^3 \Big|_1^{3/2} \\ &= 1 - \frac{1}{7} \cdot \left(\frac{27}{8} - 1 \right) = 1 - \frac{1}{7} \cdot \frac{19}{8} \\ &= 1 - \frac{19}{56} = \frac{37}{56}. \end{aligned}$$

18.

(Variabile aleatoare continue:
funcția densitate de probabilitate)

CMU, 2008 spring, Eric Xing, HW1, pr. 1.1.b

Fie funcția

$$p(x) = \begin{cases} cx^{-d} & \text{pentru } x > 1 \\ 0 & \text{în caz contrar.} \end{cases}$$

Care sunt valorile posibile pentru c și d în așa fel ca p să poată reprezenta o funcție densitate de probabilitate?

Răspuns:

O funcție p poate reprezenta o funcție densitate de probabilitate dacă $p(x) \geq 0$ pentru $\forall x$, și $\int_{-\infty}^{\infty} p(x)dx = 1$. Vom folosi aceste două condiții pentru a calcula valorile posibile pentru c și d .

Prima condiție $p(x) \geq 0$ implică faptul că $c \geq 0$, asupra lui d neimpunând nicio restricție. Mai mult, ținând cont de forma lui p și de cea de-a doua condiție, vom avea chiar $c > 0$, fiindcă $c = 0$ ar implica $\int_{-\infty}^{\infty} p(x)dx = 0 \neq 1$.

Pentru a aplica cea de-a doua condiție, calculăm integrala:

$$\int_{-\infty}^{\infty} p(x)dx = \int_1^{\infty} cx^{-d}dx = c \int_1^{\infty} x^{-d}dx.$$

Vom avea de tratat două cazuri:

$$\text{Cazul 1: } d = 1 \Rightarrow c \int_1^{\infty} x^{-d}dx = c \int_1^{\infty} \frac{1}{x}dx = c \cdot \ln x \Big|_1^{\infty} = \infty$$

$$\text{Cazul 2: } d \neq 1 \Rightarrow c \int_1^{\infty} x^{-d}dx = c \cdot \frac{x^{-d+1}}{-d+1} \Big|_1^{\infty} = \frac{c}{1-d} \cdot x^{1-d} \Big|_1^{\infty}$$

Subcazul $d > 1$: $x^{1-d} \Big|_1^{\infty} = 0 - 1 = -1$. Deci în acest caz $c \int_1^{\infty} x^{-d}dx = \frac{c}{d-1}$.

Subcazul $d < 1$: $x^{1-d} \Big|_1^{\infty} = +\infty$. Deci în acest caz $c \int_1^{\infty} x^{-d}dx = +\infty$.

Prin urmare, $\int_{-\infty}^{\infty} p(x)dx = 1 \Rightarrow d > 1$ și $\frac{c}{d-1} = 1$.

În concluzie, p poate reprezenta o funcție densitate de probabilitate dacă

$$c > 0, d > 1 \text{ și } c = d - 1.$$

Referitor la aceste trei restricții, se observă că ultimele două o implică pe cea dintâi, deci o fac superfluă.

19. (Coeficientul de corelație pentru două variabile aleatoare:
două proprietăți)

Liviu Ciortuz, 2019, după

■ · *Sheldon Ross, A First Course in Probability, 5th ed.,
Prentice Hall, 1997, pag. 332*

Pentru două variabile aleatoare oarecare X și Y având $\text{Var}(X) \neq 0$ și $\text{Var}(Y) \neq 0$, coeficientul de corelație se definește astfel:

$$\rho(X, Y) \stackrel{\text{def.}}{=} \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}}.$$

a. Să se demonstreze că $-1 \leq \rho(X, Y) \leq 1$.

Consecință: $\text{Cov}(X, Y) \in [-\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}, +\sqrt{\text{Var}(X) \cdot \text{Var}(Y)}]$, dacă $\text{Var}(X) \neq 0$ și $\text{Var}(Y) \neq 0$.

b. Să se arate că dacă $\rho(X, Y) = 1$, atunci $Y = aX + b$, cu $a = \text{Var}(Y)/\text{Var}(X) > 0$. Similar, dacă $\rho(X, Y) = -1$, atunci $Y = aX + b$, cu $a = -\text{Var}(Y)/\text{Var}(X) < 0$.¹⁸

Indicații:

1. La punctul a, notând $\text{Var}(X) = \sigma_X$ și $\text{Var}(Y) = \sigma_Y$, pentru a demonstra inegalitatea $\rho(X, Y) \geq -1$ vă sugerăm să dezvoltăm expresia $\text{Var}\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right)$ folosind următoarele două proprietăți, valabile pentru orice variabile aleatoare X și Y : $\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$,¹⁹ și $\text{Cov}(aX, bY) = ab\text{Cov}(X, Y)$, pentru orice $a, b \in \mathbb{R}$.²⁰ Apoi veți proceda similar, dezvoltând expresia $\text{Var}\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right)$, ca să demonstrați inegalitatea $\rho(X, Y) \leq 1$.
2. La punctul b, veți ține cont de faptul că pentru o variabilă aleatoare oarecare X , avem $\text{Var}(X) = 0$ dacă și numai dacă variabila X este constantă.²¹

Răspuns:

a. Pentru a demonstra inegalitatea $\rho(X, Y) \geq -1$, procedăm conform *Indicației 1*:

$$\begin{aligned} \text{Var}\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right) &= \text{Var}\left(\frac{X}{\sigma_X}\right) + \text{Var}\left(\frac{Y}{\sigma_Y}\right) + 2\text{Cov}\left(\frac{X}{\sigma_X}, \frac{Y}{\sigma_Y}\right) \\ &= \frac{1}{\sigma_X^2} \text{Var}(X) + \frac{1}{\sigma_Y^2} \text{Var}(Y) + 2\frac{1}{\sigma_X} \frac{1}{\sigma_Y} \text{Cov}(X, Y) \\ &= 1 + 1 + 2\rho(X, Y) = 2[1 + \rho(X, Y)]. \end{aligned}$$

Întrucât $\text{Var}(X) \geq 0$ pentru orice variabilă aleatoare X , rezultă că $\text{Var}\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right) \geq 0$, deci $1 + \rho(X, Y) \geq 0$, adică $\rho(X, Y) \geq -1$.

În mod similar, putem să arătăm că

$$\text{Var}\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = 2[1 - \rho(X, Y)] \geq 0,$$

deci $\rho(X, Y) \leq 1$.

b. Dacă $\rho(X, Y) = -1$, atunci din primul calcul de la punctul a va rezulta că $\text{Var}\left(\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}\right) = 0$. Știm că aceasta se întâmplă dacă $\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y}$ este o variabilă aleatoare constantă, adică, mai precis, există $a' \in \mathbb{R}$ astfel încât $\frac{X}{\sigma_X} + \frac{Y}{\sigma_Y} = a'$ cu probabilitate 1. Prin urmare, putem scrie $Y = a'\sigma_Y - \frac{\sigma_Y}{\sigma_X}X$. Rezultă că $Y = aX + b$, unde $a = -\frac{\sigma_Y}{\sigma_X} < 0$ și $b = a'\sigma_Y$.

¹⁸Așadar, coeficientul de corelație reprezintă o „măsură“ a gradului de „dependență liniară“ dintre X și Y .
LC: Coeficientul a din relația $Y = aX + b$ nu poate lua valori în afara intervalului $[-1, 1]$ — așa cum ne-am aștepta dacă facem legătura cu ecuația unei drepte oarecare din planul euclidian — din cauza simetriei $\text{Cov}(X, Y) = \text{Cov}(Y, X)$.

¹⁹Vedeți rezolvarea problemei 21.c.

²⁰Vedeți *Consecința* de la problema 9.c.

²¹Mai precis, există $c \in \mathbb{R}$ astfel încât $P(X = c) = 1$, unde P este distribuția de probabilitate considerată la definirea variabilelor din enunțul problemei.

În mod similar, dacă $\rho(X, Y) = 1$, atunci din al doilea calcul de la punctul a va rezulta că $\text{Var}\left(\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y}\right) = 0$, deci există $a'' \in \mathbb{R}$ astfel încât $\frac{X}{\sigma_X} - \frac{Y}{\sigma_Y} = a''$ cu probabilitate 1. Așadar, $Y = -a''\sigma_Y + \frac{\sigma_Y}{\sigma_X}X$. Renotând, obținem $Y = aX + b$, cu $a = \frac{\sigma_Y}{\sigma_X} > 0$ și $b = -a''\sigma_Y$.

20.

(O proprietate: matricea de covarianță a oricărui vector de variabile aleatoare este simetrică și pozitiv semidefinită)

■ □ prelucrare de Liviu Ciortuz, după

“The Multivariate Gaussian Distribution”, Chuong B. Do, 2008

Fie variabilele aleatoare X_1, \dots, X_n , cu $X_i : \Omega \rightarrow \mathbb{R}$ pentru $i = 1, \dots, n$. Matricea de covarianță a vectorului de variabile aleatoare $X = (X_1, \dots, X_n)$ este o matrice pătratică de dimensiune $n \times n$, ale cărei elemente se definesc astfel: $[Cov(X)]_{ij} \stackrel{\text{def.}}{=} Cov(X_i, X_j)$, pentru orice $i, j \in \{1, \dots, n\}$.

Arătați că $\Sigma \stackrel{\text{not.}}{=} Cov(X)$ este matrice simetrică și pozitiv semidefinită, cea de-a doua proprietate însemnând că pentru orice vector $z \in \mathbb{R}^n$ are loc inegalitatea $z^\top \Sigma z \geq 0$. (Vectorii $z \in \mathbb{R}^n$ sunt considerați vectori-coloană, iar simbolul \top reprezintă operația de transpunere a matricelor.)

Răspuns:

Faptul că matricea Σ este simetrică decurge imediat din definiția ei: dacă $X = (X_1, \dots, X_n)$, atunci $[Cov(X)]_{i,j} \stackrel{\text{def.}}{=} Cov(X_i, X_j) = E[(X_i - E[X_i])(X_j - E[X_j])] = E[(X_j - E[X_j])(X_i - E[X_i])] = Cov(X_j, X_i) = [Cov(X)]_{j,i}$, pentru orice $i, j \in \{1, \dots, n\}$.

Apoi, pentru orice vector $z \in \mathbb{R}^n$ de forma $z = (z_1, \dots, z_n)^\top$, avem:

$$\begin{aligned} z^\top \Sigma z &= \sum_{i=1}^n z_i \left(\sum_{j=1}^n \Sigma_{ij} z_j \right) = \sum_{i=1}^n \sum_{j=1}^n (z_i \Sigma_{ij} z_j) = \sum_{i=1}^n \sum_{j=1}^n (z_i Cov[X_i, X_j] z_j) \\ &= \sum_{i=1}^n \sum_{j=1}^n (z_i E[(X_i - E[X_i])(X_j - E[X_j])] z_j) \\ &= E \left[\sum_{i=1}^n \sum_{j=1}^n z_i (X_i - E[X_i])(X_j - E[X_j]) z_j \right] \end{aligned}$$

Ultima dintre egalitățile de mai sus derivă din proprietatea de liniaritate a mediilor. Mai departe,

$$\begin{aligned} z^\top \Sigma z &= E \left[\left(\sum_{i=1}^n z_i (X_i - E[X_i]) \right) \left(\sum_{j=1}^n (X_j - E[X_j]) z_j \right) \right] \\ &= E \left[\left(\sum_{i=1}^n (X_i - E[X_i]) z_i \right) \left(\sum_{j=1}^n (X_j - E[X_j]) z_j \right) \right] \end{aligned}$$

Pentru a finaliza demonstrația, să mai observăm că ultima expresie obținută mai sus se scrie sub formă vectorială astfel:

$$E[(X - E[X])^\top \cdot z]^2,$$

ceea ce evident, reprezintă o cantitate nenegativă. Așadar, $z^\top \Sigma z \geq 0$.

21.

(Variabile aleatoare: Adevărat sau Fals?)

*CMU, 2006 fall, E. Xing, T. Mitchell, final exam, pr. 1.b**CMU, 2008 fall, Eric Xing, midterm exam, pr. 1.1, 1.2, 1.3*

a. Dacă o variabilă aleatoare continuă X are funcția densitate de probabilitate p diferită de zero pe tot domeniul de definiție, atunci probabilitatea ca X să ia o valoare oarecare x (notație: $P(X = x)$) este egală cu $p(x)$.

b. $E[X + Y] = E[X] + E[Y]$ pentru orice două variabile aleatoare X și Y .

c. $Var[X + Y] = Var[X] + Var[Y]$ pentru orice două variabile aleatoare X și Y .

d. $E[XY] = E[X] \cdot E[Y]$ pentru orice două variabile aleatoare X și Y .

Răspuns:

a. Fals (în general). Dacă variabila aleatoare continuă X are funcția densitate de probabilitate p , atunci $P(a \leq X \leq b) = \int_a^b p(x)dx$. Prin urmare, $P(X = x) = 0$ pentru orice $x \in \mathbb{R}$. Enunțul afirmă pe de o parte că $p(x) = P(X = x) = 0$ pentru un anume $x \in \mathbb{R}$, iar pe de altă parte că p este diferită de zero pe tot domeniul de definiție, ceea ce este absurd.

b. Adevărat. Demonstrația este făcută în problema 9 punctul a.

c. Fals (în general). Se poate considera situația $Y = -X$, caz în care $Var[X + Y] = 0$, dar $Var[X] + Var[Y] = E[X^2] - (E[X])^2 + E[(-X)^2] - (E[-X])^2 = 2Var[X]$. Așadar, pentru orice variabilă aleatoare X cu $Var[X] \neq 0$, luând $Y = -X$, rezultă că $Var[X + Y] = 0$ iar $Var[X] + Var[Y] \neq 0$. Un exemplu de variabilă aleatoare cu varianța nenulă este distribuția gaussiană.

Mai general, se poate demonstra că $Var[X + Y] = Var[X] + Var[Y] + 2Cov[X, Y]$ astfel:

$$\begin{aligned} Var[X + Y] &= E[(X + Y)^2] - (E[X + Y])^2 \\ &= E[X^2 + 2XY + Y^2] - (E[X] + E[Y])^2 \\ &= E[X^2] + 2E[XY] + E[Y^2] - (E[X])^2 - 2E[X] \cdot E[Y] - (E[Y])^2 \\ &= (E[X^2] - (E[X])^2) + (E[Y^2] - (E[Y])^2) + (2E[XY] - 2E[X] \cdot E[Y]) \\ &= Var[X] + Var[Y] + 2Cov[X, Y] \end{aligned}$$

Prin urmare, egalitatea din enunț este adevărată pentru orice două variabile aleatoare X și Y pentru care $Cov[X, Y] = 0$ (vedeți problema 9.c), dar este falsă în rest. Egalitatea $Cov[X, Y] = 0$ este adevărată de exemplu atunci când X și Y sunt variabile independente (vedeți problema 10).

d. Fals, în general. Afirmatia din enunț este echivalentă cu $Cov[X, Y] = 0$. Așa cum am menționat la punctul c, ea este adevărată dacă, de exemplu X

și Y sunt variabile aleatoare independente. Dacă, în schimb, vom considera de pildă cazul $Y = X$, cu X variabilă aleatoare binară care ia valoarea 1 cu probabilitatea p și valoarea 0 cu probabilitatea $1 - p$, se poate deduce imediat că $E[X^2] \neq (E[X])^2$, fiindcă $p(1 - p) \neq 0$ pentru orice $p \in (0, 1)$.

0.1.3 Distribuții probabiliste uzuale

22. (Variabile aleatoare discrete – distribuția Bernoulli – și evenimente aleatoare identic distribuite: calcul de medii)

• CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 2.2

Un iepuraș se joacă „de-a săritelele”. Poziția lui inițială coincide cu originea axei reale. După aceea, iepurașul face câte un salt de-a lungul axei, fie la stânga fie la dreapta. Pentru a determina în ce direcție să sară, iepurașul dă cu banul. Dacă obține stema, va sări spre dreapta, iar dacă obține banul, va sări spre stânga. Probabilitatea de a obține stema este p . Se presupune că toate salturile iepurașului au aceeași lungime, și anume 1.

Care va fi poziția (medie) la care ne așteptăm să fie iepurașul după ce face n salturi?

Răspuns:

Fiecare săritură a iepurașului este modelată de o variabilă aleatoare. Un salt spre dreapta înseamnă o deplasare cu $+1$ pe axă, iar un salt spre stânga -1 . Să notăm cu X_i variabila aleatoare corespunzătoare săriturii i . Aceasta este:

$$X_i : \begin{pmatrix} -1 & 1 \\ 1-p & p \end{pmatrix}$$

Media acestei variabile aleatoare este $E[X_i] = -1 \cdot (1 - p) + 1 \cdot p = 2p - 1$.

Ținând cont de proprietatea de liniaritate a mediilor (vedeți pr. 9.a), poziția iepurașului după n salturi este de așteptat să fie:

$$E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n] = n(2p - 1).$$

De exemplu, pentru $p = 1/2$, se va obține poziția 0 pe axa reală (așa cum este de așteptat dacă n este număr par), în vreme ce pentru $p = 2/3$ va rezulta poziția $n/3$ (dacă $n/3 \in \mathbb{N}$), iar pentru $p = 1/3$ poziția $-n/3$ (similar).

23. (Distribuția binomială: verificarea condițiilor de definiție pentru p.m.f.; calculul mediei și al varianței)

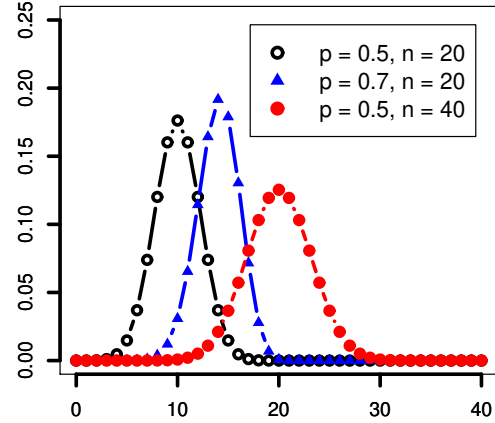
■ □ Livi Ciortuz, 2015

Distribuția binomială de parametri n și p are funcția masă de probabilitate (engl., probability mass function, p.m.f.) definită astfel:

$$b(r; n, p) = C_n^r p^r (1-p)^{n-r} \quad \forall r \in \{0, \dots, n\}.$$

Distribuția binomială: p.m.f.

Vă reamintim că $b(r; n, p)$ este probabilitatea care corespunde numărului (r) de apariții ale feței cu stema (corespondent engl., *head*) obținute la efectuarea a n aruncări independente ale unei monede, atunci când se presupune că probabilitatea de apariție a stemei la o aruncare oarecare a aceastei monede este p .



a. Verificați că funcția $b(r; n, p)$, așa cum a fost definită mai sus, reprezintă într-adevăr o funcție masă de probabilitate. Aceasta revine la a arăta că $b(r; n, p) \geq 0$ pentru orice $p \in [0, 1]$, $n \in \mathbb{N}$ și $r \in \{0, 1, \dots, n\}$, iar $\sum_{r=0}^n b(r; n, p) = 1$ pentru orice astfel de n și p , fixați.

b. Calculați media și varianța distribuției binomiale.

Răspuns:

a. Evident, $b(r; n, p) \stackrel{\text{def.}}{=} C_n^r p^r (1-p)^{n-r} \geq 0$ pentru orice $p \in [0, 1]$, $n \in \mathbb{N}$ și $r \in \{0, 1, \dots, n\}$, iar

$$\begin{aligned} \sum_{r=0}^n b(r; n, p) &= (1-p)^n + C_n^1 p (1-p)^{n-1} + \dots + C_n^{n-1} p^{n-1} (1-p) + p^n \\ &= [p + (1-p)]^n = 1 \end{aligned}$$

b. Calculul mediei se poate face pornind de la definiție:

$$\begin{aligned} E[b(r; n, p)] &\stackrel{\text{def.}}{=} \sum_{r=0}^n r \cdot b(r; n, p) = \\ &= 1 \cdot C_n^1 p (1-p)^{n-1} + 2 \cdot C_n^2 p^2 (1-p)^{n-2} + \dots + (n-1) \cdot C_n^{n-1} p^{n-1} (1-p) + n \cdot p^n \\ &= p [C_n^1 (1-p)^{n-1} + 2 \cdot C_n^2 p (1-p)^{n-2} + \dots + (n-1) \cdot C_n^{n-1} p^{n-2} (1-p) + n \cdot p^{n-1}] \\ &= np [(1-p)^{n-1} + C_{n-1}^1 p (1-p)^{n-2} + \dots + C_{n-1}^{n-2} p^{n-2} (1-p) + C_{n-1}^{n-1} p^{n-1}] \quad (8) \\ &= np [p + (1-p)]^{n-1} = np. \quad (9) \end{aligned}$$

Pentru egalitatea (8) am folosit faptul că

$$\begin{aligned} k C_n^k &= k \frac{n!}{k! (n-k)!} = \frac{n!}{(k-1)! (n-k)!} = \frac{n (n-1)!}{(k-1)! (n-1-(k-1))!} \\ &= n C_{n-1}^{k-1}, \quad \forall k = 1, \dots, n. \end{aligned}$$

Pentru calculul varianței, vom folosi formula $\text{Var}[X] = E[X^2] - E^2[X]$, care a fost demonstrată la problema 9.b. Întrucât am calculat deja $E[b(r; n, p)]$, rămâne să calculăm $E[b^2(r; n, p)]$.²² Notând $q = 1 - p$, vom avea:

$$\begin{aligned} E[b^2(r; n, p)] &\stackrel{\text{def.}}{=} \sum_{r=0}^n r^2 C_n^r p^r q^{n-r} = \sum_{r=0}^n r^2 \frac{n(n-1) \dots (n-r+1)}{r!} p^r q^{n-r} \\ &= \sum_{r=1}^n r n \frac{(n-1) \dots (n-r+1)}{(r-1)!} p^r q^{n-r} = \sum_{r=1}^n r n C_{n-1}^{r-1} p^r q^{n-r} \\ &= np \sum_{r=1}^n r C_{n-1}^{r-1} p^{r-1} q^{(n-1)-(r-1)}. \end{aligned}$$

Mai departe, notând pentru conveniență $r - 1$ cu j , urmează:

$$\begin{aligned} E[b^2(r; n, p)] &= np \sum_{j=0}^{n-1} (j+1) C_{n-1}^j p^j q^{(n-1)-j} \\ &= np \left[\sum_{j=0}^{n-1} j C_{n-1}^j p^j q^{(n-1)-j} + \sum_{j=0}^{n-1} C_{n-1}^j p^j q^{(n-1)-j} \right]. \end{aligned}$$

Prima sumă din paranteza pătrată de mai sus este chiar $E[b^2(r; n-1, p)]$, conform relației (9), iar cea de-a doua sumă este egală cu 1, conform unui calcul absolut similar cu cel de la punctul a. Prin urmare,

$$E[b^2(r; n, p)] = np[(n-1)p + 1] = n^2p^2 - np^2 + np.$$

Așadar, $\text{Var}[X] = E[b^2(r; n, p)] - (E[b(r; n, p)])^2 = n^2p^2 - np^2 + np - n^2p^2 = np(1-p)$.

Observație: O altă cale de a calcula varianța distribuției binomiale este următoarea:

- se demonstrează relativ ușor că orice variabilă aleatoare urmând distribuția binomială $b(r; n, p)$ poate fi văzută ca o sumă de n variabile independente care urmează distribuția Bernoulli de parametru p ;²³
- se știe (sau, se poate dovedi imediat) că varianța distribuției Bernoulli de parametru p este $p(1-p)$;
- ținând cont de proprietatea $\text{Var}[X_1 + X_2 + \dots + X_n] = \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n]$ atunci când X_1, X_2, \dots, X_n sunt variabile independente, conform demonstrației de la problema 21.c, rezultă că $\text{Var}[X] = np(1-p)$.

²²Rezolvarea de mai jos urmează îndeaproape linia demonstrației găsite pe site-ul www.proofwiki.org/wiki/Variance_of_Binomial_Distribution, accesat la data de 5 octombrie 2015. La rândul său, acest site menționează ca sursă “Probability: An Introduction” de Geoffrey Grimmett și Dominic Welsh, Oxford Science Publications, 1986.

²³Vedeți www.proofwiki.org/wiki/Bernoulli_Process_as_Binomial_Distribution, care citează ca sursă “Probability: An Introduction” de Geoffrey Grimmett și Dominic Welsh, Oxford Science Publications, 1986.

24.

(Distribuția Poisson:
verificarea condițiilor de definiție pentru p.m.f.;
calculul mediei și al varianței)

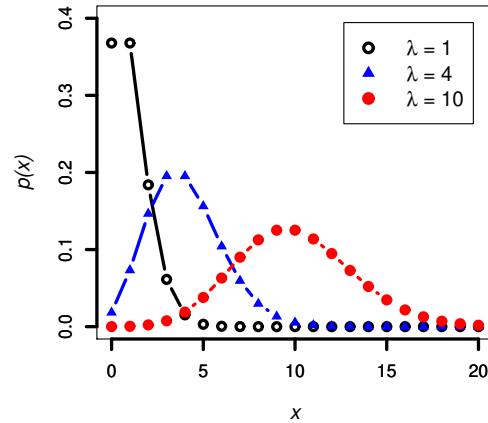
Liviu Ciortuz, 2017

Distribuția Poisson: pmf

Distribuția *Poisson* este o distribuție discretă de parametru $\lambda > 0$, a cărei funcție masă de probabilitate este dată de expresia

$$p(x | \lambda) = \frac{1}{e^\lambda} \cdot \frac{\lambda^x}{x!}, \text{ pentru orice } x \in \mathbb{N}.$$

Prin convenție, se consideră că $0! = 1$. Factorul $\frac{1}{e^\lambda}$, care nu depinde de x , este așa-numitul *factor de normalizare*.



Demonstrați mai întâi că funcția $p(\cdot)$ este într-adevăr funcție masă de probabilitate (engl., probability mass function, p.m.f.) și apoi că media acestei distribuții este λ , iar varianța ei este tot λ .

Sugestie: Țineți cont că $\sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = e^\lambda$ (limită fundamentală).

Răspuns:

Vom arăta mai întâi că $p(\cdot)$ este într-adevăr o funcție masă de probabilitate (p.m.f.). Evident, $p(x|\lambda) > 0$ pentru orice $x \in \mathbb{N}$, fiindcă $\lambda > 0$. Apoi,

$$\sum_{x \in \mathbb{N}} \frac{1}{e^\lambda} \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} e^\lambda = 1,$$

ținând cont de limita fundamentală care apare în *Sugestia* din enunț.

Pentru calcularea mediei acestei distribuții, aplicăm definiția:

$$\sum_{x=0}^{\infty} x \frac{1}{e^\lambda} \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} \sum_{x=0}^{\infty} x \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} = \frac{\lambda}{e^\lambda} \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} = \frac{\lambda}{e^\lambda} \underbrace{\sum_{x=0}^{\infty} \frac{\lambda^x}{x!}}_{e^\lambda} = \frac{\lambda}{e^\lambda} e^\lambda = \lambda.$$

În ce privește calculul varianței pentru distribuția Poisson, știm că $\text{Var}[X] = E[X^2] - E^2[X]$ (proprietate demonstrată la problema 9.b) pentru orice distribuție aleatoare X și, prin urmare, pentru că am calculat deja media acestei distribuții, trebuie să mai calculăm $\sum_{x=0}^{\infty} x^2 p(x|\lambda)$.

$$\begin{aligned} \sum_{x=0}^{\infty} x^2 p(x|\lambda) &= \sum_{x=0}^{\infty} x^2 \frac{1}{e^\lambda} \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} \sum_{x=1}^{\infty} x^2 \frac{\lambda^x}{x!} = \frac{1}{e^\lambda} \left[\sum_{x=1}^{\infty} [x(x-1) + x] \frac{\lambda^x}{x!} \right] \\ &= \frac{1}{e^\lambda} \left[\sum_{x=1}^{\infty} x(x-1) \frac{\lambda^x}{x!} + \sum_{x=1}^{\infty} x \frac{\lambda^x}{x!} \right] = \frac{1}{e^\lambda} \left[\sum_{x=2}^{\infty} \frac{\lambda^x}{(x-2)!} + \sum_{x=1}^{\infty} \frac{\lambda^x}{(x-1)!} \right] \\ &= \frac{1}{e^\lambda} \left[\lambda^2 \sum_{x=2}^{\infty} \frac{\lambda^{x-2}}{(x-2)!} + \lambda \sum_{x=1}^{\infty} \frac{\lambda^{x-1}}{(x-1)!} \right] = \frac{1}{e^\lambda} \left[\lambda^2 \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} + \lambda \sum_{x=0}^{\infty} \frac{\lambda^x}{x!} \right] \end{aligned}$$

$$= \frac{1}{e^\lambda} [\lambda^2 e^\lambda + \lambda e^\lambda] = \frac{1}{e^\lambda} e^\lambda (\lambda^2 + \lambda) = \lambda^2 + \lambda.$$

Prin urmare, varianța distribuției Poisson este $\lambda^2 + \lambda - \lambda^2 = \lambda$.

25.

(Distribuția geometrică:
numărul „așteptat“ / mediu de „observații“ necesare
pentru ca un anumit eveniment să se producă)

□ • ◦ CMU, 2012 spring, Ziv Bar-Joseph, HW1, pr. 1.4

În cazul unui zar perfect cu șase fețe, probabilitățile de apariție pentru fiecare dintre fețele zarului sunt egale.

Mickey se duce la un cazinou și dorește ca, folosindu-și cunoștințele din domeniul probabilităților, să-și evalueze șansa pe care o are de a obține la aruncarea unui astfel de zar fața 6.

Mai precis, Mickey se întreabă care este numărul mediu (sau, numărul „așteptat“; engl., expected number) de aruncări ale zarului pe care ar trebui să le efectueze până să obțină fața 6.

Justificați răspunsul în detaliu.

Răspuns:

Experimentul lui Mickey poate fi modelat cu ajutorul *distribuției geometrice*.²⁴ Știm că distribuția geometrică poate fi reprezentată de o tabelă de „repartiție“ de forma

$$\begin{pmatrix} 1 & 2 & 3 & \dots & n & \dots \\ q & pq & p^2q & \dots & p^{n-1}q & \dots \end{pmatrix},$$

unde $p, q \geq 0$ și $p + q = 1$. Se verifică imediat că $\sum_{i=1}^{\infty} p^{i-1}q = 1$ atunci când $p \in [0, 1)$. În particular, pentru $p = \frac{5}{6}$ și $q = \frac{1}{6}$, avem

$$\sum_{i=1}^{\infty} \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{i-1} = 1. \quad (10)$$

Vom nota cu E numărul mediu de aruncări ale zarului pe care ar trebui să le efectueze Mickey până să obțină fața 6.

Conform definiției mediei, E se poate exprima astfel:

$$E = 1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} \cdot \frac{5}{6} + 3 \cdot \frac{1}{6} \cdot \left(\frac{5}{6}\right)^2 \dots + n \cdot \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{n-1} \dots + \dots$$

Punând termenul generic $n \cdot \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{n-1}$ sub forma $(1 + (n-1)) \cdot \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{n-1}$ și aplicând apoi distributivitatea înmulțirii față de adunare, putem rescrie E în mod echivalent astfel:

$$E = \frac{1}{6} + \frac{5}{6} \cdot \left[\frac{1}{6} + \frac{1}{6} \cdot \frac{5}{6} + \dots + \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{n-1} + \dots \right]$$

²⁴ Simpu spus, distribuția gemetrică poate fi gândită ca modelând următorul experiment aleatoriu: Fie o monedă a cărei probabilitate de apariție a feței-stemă este p . Aruncăm moneda o dată sau de mai multe ori, până când apare stema. Notăm numărul de aruncări care au precedat apariția stemei cu k . Acest număr $k \in \{0, 1, \dots\}$ va fi [asociat cu] valoarea unei variabile aleatoare X , despre care spunem că urmează distribuția geometrică. Evident, $P(X = k) = (1 - p)^k p$.

$$+\frac{5}{6} \cdot \left[1 \cdot \frac{1}{6} + 2 \cdot \frac{1}{6} \cdot \frac{5}{6} + \dots + (n-1) \cdot \frac{1}{6} \cdot \left(\frac{5}{6}\right)^{n-1} + \dots \right].$$

Ținând cont de relația (10), egalitatea precedentă se rescrie astfel:

$$E = \frac{1}{6} + \frac{5}{6} \cdot 1 + \frac{5}{6} \cdot E.$$

Rezultă imediat că $\frac{1}{6}E = \frac{1}{6} + \frac{5}{6}$, deci $E = 6$.

26.

(O mixtură de distribuții categoriale:
calculul mediei și al varianței)

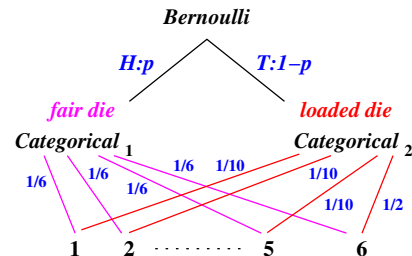
■ □ ● ○ CMU, 2010 fall, Aarti Singh, HW1, pr. 2.2.1-2

Presupunem că avem două zaruri cu șase fețe. Un zar este perfect — deci vom considera funcția sa masă de probabilitate (p.m.f.) definită prin $P_1(x) = 1/6$ pentru $x = 1, \dots, 6$ —, iar celălalt zar este măsluit și are următoarea funcție masă de probabilitate:

$$P_2(x) = \begin{cases} \frac{1}{2} & \text{pentru } x = 6; \\ \frac{1}{10} & \text{pentru } x \in \{1, 2, 3, 4, 5\}. \end{cases}$$

Pentru a decide ce zar să aruncăm, vom folosi o monedă; considerăm că probabilitatea să obținem stema la aruncarea monedei este $p \in (0, 1)$. Dacă obținem stema (engl., head), vom arunca apoi zarul perfect; în caz contrar vom arunca zarul măsluit.

Putem reprezenta grafic mixtura formată din cele două distribuții categoriale ca în figura alăturată.



a. Calculați în funcție de p media variabilei aleatoare (notată cu X) care reprezintă „mixtura” de distribuții [categoriale] descrisă mai sus.

b. Calculați în funcție de p varianța variabilei aleatoare X .

Răspuns:

a. Ținând cont de modul în care a fost definită în enunț mixtura celor două distribuții categoriale, putem calcula media variabilei X , care reprezintă această mixtură, în felul următor:

$$\begin{aligned} E[X] &= \sum_{i=1}^6 i \cdot P(i) = \sum_{i=1}^6 i \cdot [P(i|fair) \cdot p + P(i|loaded) \cdot (1-p)] \\ &= \sum_{i=1}^6 i \cdot [P_1(i) \cdot p + P_2(i) \cdot (1-p)] = \left[\sum_{i=1}^6 i \cdot P_1(i) \right] p + \left[\sum_{i=1}^6 i \cdot P_2(i) \right] (1-p) \end{aligned}$$

b. Considerăm familia mixturilor de două distribuții gaussiene unidimensionale:

$$f(x|\theta, \mu_1, \sigma_1, \mu_2, \sigma_2) = \theta \mathcal{N}(x|\mu_1, \sigma_1) + (1 - \theta) \mathcal{N}(x|\mu_2, \sigma_2),$$

unde $0 < \theta < 1$, $0 < \sigma_1$, $0 < \sigma_2$, iar μ_1, μ_2 sunt numere reale oarecare, distincte. Este imediat că pentru orice estimări bine-definite $\hat{\theta}, \hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1$ și $\hat{\sigma}_2$ obținute pe un set oarecare de instanțe x_1, \dots, x_n , următorul set de estimări

$$\hat{\theta}' = 1 - \hat{\theta}, \hat{\mu}'_1 = \hat{\mu}_2, \hat{\mu}'_2 = \hat{\mu}_1, \hat{\sigma}'_1 = \hat{\sigma}_2, \hat{\sigma}'_2 = \hat{\sigma}_1$$

vor produce întotdeauna aceeași verosimilitate ca și estimările inițiale. Și totuși $(\hat{\theta}', \hat{\mu}'_1, \hat{\mu}'_2, \hat{\sigma}'_1, \hat{\sigma}'_2) \neq (\hat{\theta}, \hat{\mu}_1, \hat{\mu}_2, \hat{\sigma}_1, \hat{\sigma}_2)$. Așadar, estimările în sensul MLE nu sunt în mod necesar unice.

c. În exemplul de la punctul a, domeniul de existență pentru parametrul distribuției pe care am ales-o este, ca și în multe alte cazuri, un interval deschis. Însă valoarea acestui parametru pentru care se atinge maximul funcției de log-verosimilitate este situată la „marginea” [inferioară a] intervalului, nu în interiorul intervalului respectiv. Această problemă apare adeseori atunci când mărimea eșantionului (adică numărul instanțelor x_i) este mică în comparație cu numărul parametrilor care trebuie estimați.

În exemplul de la punctul b, funcția de log-verosimilitate nu este concavă — pentru exemplificare, vedeți problema 17 de la capitolul de *Clusterizare*, în special graficele care reprezintă curbele de izocontur ale funcției de log-verosimilitate ale mixturilor de gaussiene, la pag. 777 și pag. 778 —, ceea ce implică faptul că pot exista mai multe puncte de maxim local.

0.1.5 Elemente de teoria informației⁹²

50.

(Exercițiu cu caracter teoretic:
proprietăți dezirabile ale entropiei)

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

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

Scopul acestui exercițiu este de a arăta că orice funcție $\psi_n(p_1, \dots, p_n)$ care satisface trei proprietăți dezirabile pentru entropie este în mod necesar de forma $-K \sum_i p_i \log p_i$ unde K este o constantă reală pozitivă. Iată care sunt aceste *proprietăți*:⁹³

⁹²Observație importantă: În toate problemele care urmează, referitor la entropie / teoria informației se va considera în mod implicit că notația ‘log’ desemnează logaritmul în baza 2. De asemenea, prin convenție, se va considera $p \cdot \log p = 0$ pentru $p = 0$.

⁹³LC: 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 ca măsură a *dezordinii*; de asemenea, $\psi_1(1) = 0$ pentru că în acest caz particular nu există niciun fel de dezordine.

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.

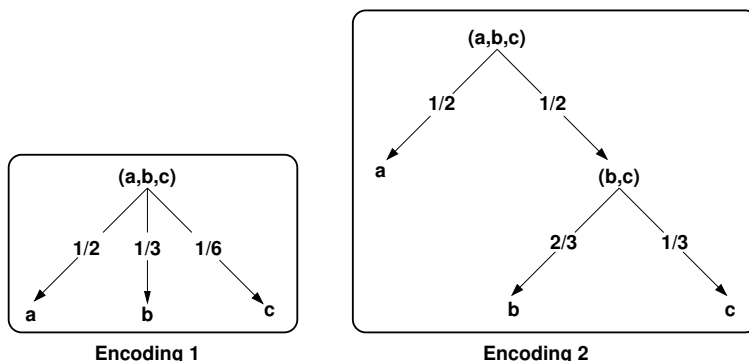
A3. Dacă faptul de a alege între mai multe evenimente posibile poate fi realizat prin mai multe alegeri succesive, atunci $\psi_n(p_1, \dots, p_n)$ trebuie să se poată scrie ca o sumă ponderată a entropiilor calculate la fiecare stadiu / alegere.

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 $\psi_2(1/2, 1/2) + 1/2 \cdot \psi_2(2/3, 1/3)$.



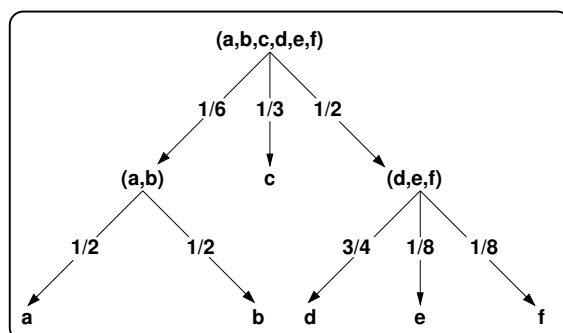
Așadar, în acest exercițiu vi se cere să arătați că dacă o funcție de n variabile $\psi_n(p_1, \dots, p_n)$ satisface proprietățile A1, A2 și A3 de mai sus, atunci există $K \in \mathbb{R}^+$ astfel încât $\psi_n(p_1, \dots, p_n) = -K \sum_i p_i \log p_i$ unde, vom vedea, K depinde de $\psi_s\left(\frac{1}{s}, \dots, \frac{1}{s}\right)$ pentru un anumit $s \in \mathbb{S}^*$.

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

$$A(s^m) = m A(s) \text{ pentru orice } s, m \in \mathbb{N}^*. \quad (49)$$

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

$$s^m \leq t^n \leq s^{m+1}. \quad (50)$$

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

$$\left| \frac{m}{n} - \frac{\log t}{\log s} \right| \leq \frac{1}{n}. \quad (51)$$

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

f. Combinând ultima inegalitate de mai sus cu inegalitatea (49), 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

$$\left| \frac{m}{n} - \frac{A(t)}{A(s)} \right| \leq \frac{1}{n} \text{ pentru } s \neq 1. \quad (52)$$

g. Combinând inegalitățile (51) și (52), arătați că

$$\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n} \text{ pentru } s \neq 1 \quad (53)$$

și, în consecință

$$A(t) = K \log t \text{ cu } K > 0 \text{ (din cauza proprietății A2)}. \quad (54)$$

Observație:

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

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

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

Conform egalității (54), 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,

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

Răspuns:

a. Facem calculele:

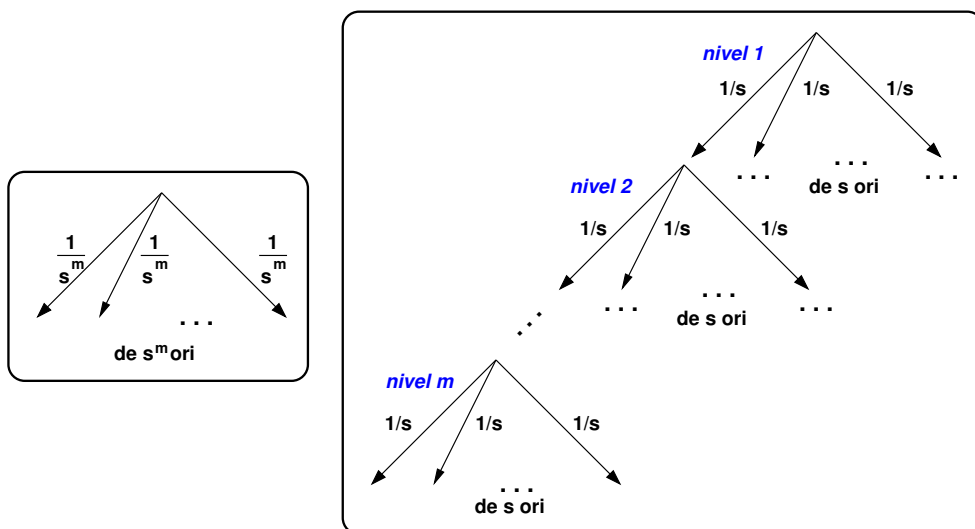
$$\begin{aligned} H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) &= \frac{1}{2} \log 2 + \frac{1}{3} \log 3 + \frac{1}{6} \log 6 = \left(\frac{1}{2} + \frac{1}{6}\right) \log 2 + \left(\frac{1}{3} + \frac{1}{6}\right) \log 3 \\ &= \frac{2}{3} + \frac{1}{2} \log 3 \\ H\left(\frac{1}{2}, \frac{1}{2}\right)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

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

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

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

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



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

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

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

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

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

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

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

f. Datorită proprietății (54), dubla inegalitate $A(s^m) \leq A(t^n) \leq A(s^{m+1})$ devine $m A(s) \leq n A(t) \leq (m+1) A(s)$. Împărțind această inegalitate prin $n A(s)$, despre care se poate spune că este nenul pentru orice $s \neq 1$,⁹⁴ rezultă:

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

g. Inegalitățile duble de mai jos rescriu convenabil proprietățile (51) și (52):

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

⁹⁴Putem considera în mod natural $A(1) = 0$. Conform proprietății A2, urmează că $A(s) > A(1) = 0$ pentru orice $s > 1$.

Însumându-le membru cu membru, rezultă

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

Din inegalitatea dublă $s^m \leq t^n \leq s^{m+1}$ rezultă că odată cu m crește și n (acesta din urmă depinzând de valorile lui s, t și m).⁹⁵ Așadar, atunci când m tinde la infinit vom avea și $n \rightarrow \infty$. Dacă trecem la limită inegalitatea $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \leq \frac{2}{n}$ pentru $n \rightarrow \infty$, rezultă $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| \rightarrow 0$, de unde avem $\left| \frac{A(t)}{A(s)} - \frac{\log t}{\log s} \right| = 0$ și deci $\frac{A(t)}{A(s)} = \frac{\log t}{\log s}$. Rezultă că $A(t) = \frac{A(s)}{\log s} \log t = K \log t$. Evident, constanta $K = \frac{A(s)}{\log s} = \frac{1}{\log s} \psi\left(\frac{1}{s}, \dots, \frac{1}{s}\right)$ nu depinde de t . Variind valorile lui t , rezultă că $A(t) = K \log t$ pentru orice $t \in \mathbb{N}^*$ astfel încât relația (50) are loc.

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

51.

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

■ • Liviu Ciortuz, pornind de la

CMU, 2005 fall, T. Mitchell, A. Moore, HW1, pr. 2

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

- 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.)

⁹⁵ Mai precis, este util să observăm o succesiune de forma $s^m \leq t^{n_1} \leq s^{m+1} \leq t^{n_2} \leq s^{m+2} \leq \dots$.

Arătați că:

- a. $H(X) \geq 0$. În particular, $H(X) = 0$ dacă și numai dacă variabila X este constantă.
- b. $H(Y | X) = - \sum_i \sum_j P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i)$.
- c. $H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$.
- Mai general: $H(X_1, \dots, X_n) = H(X_1) + H(X_2 | X_1) + \dots + H(X_n | X_1, \dots, X_{n-1})$ (regula de înlănțuire).

Răspuns:

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

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

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

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

„ \Rightarrow “ Presupunem că $H(X) = 0$, adică $\sum_i P(X = x_i) \log \frac{1}{P(X = x_i)} = 0$. Datorită faptului că fiecare termen din această sumă este mai mare sau egal cu 0, rezultă că $H(X) = 0$ doar dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $\log \frac{1}{P(X = x_i)} = 0$, adică dacă pentru $\forall i$, $P(X = x_i) = 0$ sau $P(X = x_i) = 1$. Cum însă $\sum_i P(X = x_i) = 1$ rezultă că există o singură valoare x_1 pentru X astfel încât $P(X = x_1) = 1$, iar $P(X = x) = 0$ pentru orice $x \neq x_1$. Altfel spus, variabila aleatoare discretă X este constantă.⁹⁶

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

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

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

⁹⁶Mai corect spus, X este constantă pe tot domeniul de definiție, eventual cu excepția unei mulțimi de probabilitate 0.

c. În primul rând, trebuie să demonstrăm că

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

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

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

Egalitatea $\sum_j p(y_j | x_i) = 1$ se justifică ușor ținând cont de *proprietatea de aditivitate numărabilă* din definiția funcției / distribuției de probabilitate.

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

Pentru cazul general

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

vom folosi regula de înlănțuire de la variabile aleatoare

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

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

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

La penultima egalitate am ținut cont de definiția distribuției marginale pornind de la distribuția comună, iar la ultima egalitate am folosit rezultatul de la punctul b.

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

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

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

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

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

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

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

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

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

c. Calculează entropia variabilei S .

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

Răspuns:

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

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

b. Conform definiției date, vom avea:

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

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

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

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

c. Conform definiției pentru entropie (vedeți problema 51), $H(S)$ este media ponderată (cu ajutorul probabilităților) a „surprizelor“ / cantităților de informație produse la „observarea“ tuturor valorilor variabilei S . Făcând calculele, vom obține:

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

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

S	2	3	4	5	6	7	8	9	10	11	12
$P(S \dots)$	0	0	0	1/6	1/6	1/6	1/6	1/6	1/6	0	0

În consecință, folosind definiția entropiei condiționale specifice (vedeți de asemenea problema 51), vom avea:

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

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

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

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

53.

(Probabilități marginale, entropii, entropii condiționale medii)

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

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

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

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

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

$$P(C = no) = P(C = no, T = sunny) + P(C = no, T = rainy) + P(C = no, T = snowy).$$

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

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

Răspuns:

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

b. Aplicând definiția pentru entropie (vedeți problema 51), rezultă:

$$\begin{aligned} H(C) &= 0.6 \log_2 \frac{5}{3} + 0.4 \log_2 \frac{5}{2} = \log_2 5 - 0.6 \log_2 3 - 0.4 = 0.971 \text{ biți} \\ H(T) &= 2 \cdot 0.35 \log_2 \frac{20}{7} + 0.3 \log_2 \frac{10}{3} \\ &= 0.7(2 + \log_2 5 - \log_2 7) + 0.3(1 + \log_2 5 - \log_2 3) \\ &= 1.7 + \log_2 5 - 0.7 \log_2 7 - 0.3 \log_2 3 = 1.581 \text{ biți.} \end{aligned}$$

c. Aplicând definiția pentru entropie condițională medie (vedeți de asemenea problema 51), vom avea:

$$\begin{aligned} H(C|T) &\stackrel{\text{def.}}{=} \sum_{t \in \text{Val}(T)} P(T = t) \cdot H(C|T = t) \\ &= P(T = sunny) \cdot H(C|T = sunny) + P(T = rainy) \cdot H(C|T = rainy) + \\ &\quad P(T = snowy) \cdot H(C|T = snowy) \\ &= 0.35 \cdot H\left(\frac{0.30}{0.30 + 0.05}, \frac{0.05}{0.30 + 0.05}\right) + 0.35 \cdot H\left(\frac{0.20}{0.20 + 0.15}, \frac{0.15}{0.20 + 0.15}\right) + \\ &\quad 0.30 \cdot H\left(\frac{0.10}{0.10 + 0.20}, \frac{0.20}{0.20 + 0.10}\right) \\ &= \frac{7}{20} \cdot H\left(\frac{6}{7}, \frac{1}{7}\right) + \frac{7}{20} \cdot H\left(\frac{4}{7}, \frac{3}{7}\right) + \frac{3}{10} \cdot H\left(\frac{1}{3}, \frac{2}{3}\right) \\ &= \frac{7}{20} \cdot \left(\frac{6}{7} \log_2 \frac{7}{6} + \frac{1}{7} \log_2 7\right) + \frac{7}{20} \cdot \left(\frac{4}{7} \log_2 \frac{7}{4} + \frac{3}{7} \log_2 \frac{7}{3}\right) + \frac{3}{10} \cdot \left(\frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 \frac{3}{2}\right) \\ &= \frac{7}{20} \cdot \left(\log_2 7 - \frac{6}{7} - \frac{6}{7} \log_2 3\right) + \frac{7}{20} \cdot \left(\log_2 7 - \frac{8}{7} - \frac{3}{7} \log_2 3\right) + \frac{3}{10} \cdot \left(\log_2 3 - \frac{2}{3}\right) \\ &= \frac{7}{10} \log_2 7 - \left(\frac{3}{10} + \frac{4}{10} + \frac{2}{10}\right) - \left(\frac{6}{20} + \frac{3}{20} - \frac{3}{10}\right) \cdot \log_2 3 \\ &= \frac{7}{10} \log_2 7 - \frac{3}{20} \log_2 3 - \frac{9}{10} = 0.82715 \text{ biți.} \end{aligned}$$

Similar,

$$\begin{aligned} H(T|C) &\stackrel{\text{def.}}{=} \sum_{c \in \text{Val}(C)} P(C = c) \cdot H(T|C = c) \\ &= P(C = no) \cdot H(T|C = no) + P(C = yes) \cdot H(T|C = yes) \\ &= 0.60 \cdot H\left(\frac{0.30}{0.30 + 0.20 + 0.10}, \frac{0.20}{0.30 + 0.20 + 0.10}, \frac{0.10}{0.30 + 0.20 + 0.10}\right) \end{aligned}$$

$$\begin{aligned}
& +0.40 \cdot H\left(\frac{0.05}{0.05+0.15+0.20}, \frac{0.15}{0.05+0.15+0.20}, \frac{0.20}{0.05+0.15+0.20}\right) \\
&= \frac{3}{5} \cdot H\left(\frac{1}{2}, \frac{1}{3}, \frac{1}{6}\right) + \frac{2}{5} \cdot H\left(\frac{1}{8}, \frac{3}{8}, \frac{1}{2}\right) \\
&= \frac{3}{5} \left(\frac{1}{2} + \frac{1}{3} \log_2 3 + \frac{1}{6} (1 + \log_2 3)\right) + \frac{2}{5} \left(\frac{1}{8} \cdot 3 + \frac{3}{8} (3 - \log_2 3) + \frac{1}{2}\right) \\
&= \frac{3}{5} \left(\frac{2}{3} + \frac{1}{2} \log_2 3\right) + \frac{2}{5} \left(2 - \frac{3}{8} \log_2 3\right) = \frac{6}{5} + \frac{3}{20} \log_2 3 = 1.43774 \text{ biți.}
\end{aligned}$$

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

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

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

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

Calculați entropia *distribuției* continue *exponențiale* de parametru $\lambda > 0$. Vă reamintim că definiția p.d.f.-ului acestei distribuții este următoarea:⁹⁷

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

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

Răspuns:

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

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

Prima dintre aceste două ultime integrale este 0, conform *indicației* din enunț. Pentru a putea calcula mai ușor cea de-a doua integrală (în expresia căreia apare numărul e), vom schimba baza logaritmului, și anume vom trece din baza 2 în baza e (baza logaritmului natural, \ln).⁹⁸

⁹⁷La ex. 28.a puteți vedea graficul acestei funcții de densitate pentru câteva valori ale parametrului (λ).

⁹⁸Pentru aceasta, vom folosi formula $\log_a b = \frac{\log_c b}{\log_c a}$, valabilă pentru orice $a > 0$, $b > 0$ și $c > 0$, cu $a \neq 1$ și $c \neq 1$. În calculele care urmează vom folosi și alte formule de la logaritmi.

Prin urmare,

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

Prima integrală are valoarea 1, întrucât p este p.d.f.-ul distribuției exponențiale (vedeți ex.28.a) Cea de-a doua integrală are valoarea $1/\lambda$ (vedeți rezolvarea aceluiași ex.28.a). Prin urmare,

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

55.

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

■ □ ● prelucrare de Liviu Ciortuz, după

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

Entropia relativă sau divergența Kullback-Leibler (KL) a unei distribuții p în raport cu o altă distribuție q — ambele distribuții fiind discrete — se definește astfel:⁹⁹

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

Din perspectiva teoriei informației, divergența KL specifică numărul de *biți adiționali* care sunt necesari în medie pentru a transmite valorile variabilei X atunci când presupunem că aceste valori sunt distribuite conform distribuției („model”) q , dar în realitate ele urmează o altă distribuție, p .¹⁰⁰

⁹⁹În învățarea automată, divergența KL este folosită de exemplu la fundamentarea schemei algoritmice EM. Vedeți problemele 1.b și 2 de la capitolul *Algoritmul EM*. Mai general, maximizarea verosimilității datelor — în vederea estimării parametrilor distribuțiilor probabiliste — poate fi văzută ca fiind echivalentă cu minimizarea divergenței KL (vedeți problema 135).

¹⁰⁰Atenție: Divergența KL nu este o măsură de *distanță* între două distribuții probabiliste, fiindcă în general ea nu este simetrică ($KL(p||q) \neq KL(q||p)$) și nici nu satisface inegalitatea triunghiului. Pentru „simetrizare”, se consideră $M(p, q) = \frac{1}{2}(p + q)$, apoi se definește funcția $JSD(p||q) = \frac{1}{2}KL(p||M) + \frac{1}{2}KL(q||M)$, care se numește *divergența Jensen-Shannon*. În sfârșit, se poate arăta că $\sqrt{JSD(p||q)}$ definește o măsură de distanță (metrică), adică este nenegativă, simetrică, implică egalitatea indiscernabililor și satisface inegalitatea triunghiului; ea este numită *distanța Jensen-Shannon*.

Variația informației, 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 distanță.

a. Demonstrați inegalitatea $KL(p||q) \geq 0$ și apoi arătați că egalitatea are loc dacă și numai dacă $p = q$.¹⁰¹

Indicație:

Pentru a demonstra punctul acesta puteți folosi *inegalitatea lui Jensen*.¹⁰²

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 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 pe de altă parte:

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

Arătați că această nouă definiție a câștigului de informație este echivalentă cu definiția dată anterior (vedeți problema 51). 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 . Arătați că $IG(X, Y) = 0$ dacă și numai dacă X și Y sunt independente.

Răspuns:

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

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

¹⁰¹Mai general, $KL(p||q)$ este cu atât mai mică cu cât „asemănarea“ dintre distribuțiile p și q este mai mare.

¹⁰²Vedeți problema 71.

Așadar, $KL(p \parallel q) \geq 0$, oricare ar fi distribuțiile (discrete) p și q .

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

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

Pentru a demonstra implicația inversă, se ține cont că în inegalitatea lui Jensen, în cazul funcțiilor strict convexe (cum este $-\log_2$) are loc egalitatea doar în cazul în care $x_i = x_j$ pentru orice i și j . În cazul de față, această condiție se traduce prin faptul că raportul $\frac{q(x)}{p(x)}$ este același (α) pentru orice valoare a lui x . Ținând cont că $\sum_x p(x) = 1$ și $\sum_x q(x) = \sum_x p(x) \frac{q(x)}{p(x)} = 1$, rezultă că $\alpha = \frac{q(x)}{p(x)} = 1$, deci $p(x) = q(x)$ pentru orice x , ceea ce înseamnă că distribuțiile p și q sunt identice.

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

$$\begin{aligned}
 KL(p_{X,Y} \parallel (p_X p_Y)) &\stackrel{\text{def.}}{=}_{KL} - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)p(y)}{p(x, y)} \right) \\
 &= - \sum_x \sum_y p(x, y) \log \left(\frac{p(x)}{p(x \mid y)} \right) = - \sum_x \sum_y p(x, y) [\log p(x) - \log p(x \mid y)] \\
 &= - \sum_x \sum_y p(x, y) \log p(x) - \left(- \sum_x \sum_y p(x, y) \log p(x \mid y) \right) \\
 &\stackrel{\text{pr. 51.b}}{=} - \sum_x \log p(x) \underbrace{\sum_y p(x, y)}_{=p(x)} - H[X \mid Y] = \sum_x p(x) \log p(x) - H[X \mid Y] \\
 &= H[X] - H[X \mid Y] = IG(X, Y)
 \end{aligned}$$

c. Conform punctului b, egalitatea $IG(X, Y) = 0$ este echivalentă cu egalitatea $KL(p_{X,Y} \parallel p_X p_Y) = 0$. Conform punctului a, această a doua relație este adevărată dacă și numai dacă distribuțiile $p_{X,Y}$ și $p_X p_Y$ sunt identice, or aceasta este exact definiția independenței variabilelor X și Y .

56.

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

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

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

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

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

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

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

Răspuns:

a. Pentru calculul informației mutuale putem folosi formula din problema 55.b:

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

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

	P_{X_1}	P_{X_2}	P_{X_3}	P_{X_4}	P_{X_5}	P_Y
0	1/2	3/8	1/2	1/2	0.5	1/4
1	1/2	5/8	1/2	1/2	0.5	3/4

iar probabilitățile comune sunt:

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

Se poate observa că X_1 și Y sunt independente, deci $MI(X_1, Y) = 0$, conform proprietății care a fost demonstrată la problema 55.c. Similar, $MI(X_3, Y) = 0$. În rest, efectuând calculele obținem $MI(X_2, Y) = 0.01571$, $MI(X_4, Y) = 0.3113$ și $MI(X_5, Y) = 0.3113$.

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

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

Observație: Deși $MI(X_1, Y) = 0$ — sau, echivalent spus, X_1 este independent de Y —, nu rezultă că variabila X_1 combinată cu una sau mai multe variabile X_j , cu $j \in \{2, 4, 5\}$, și formând astfel o nouă variabilă aleatoare, rămâne independentă de Y . Nouă variabilă poate avea câștig de informație nenul (în unele cazuri chiar maxim!) în raport cu Y .

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

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

Conform problemei 51.c, *formula de înlănțuire* a entropiilor pentru cazul general (adică, indiferent dacă X și Y sunt sau nu independente) este:

$$H(X, Y) = H(X) + H(Y|X) = H(Y) + H(X|Y). \quad (55)$$

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

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

Răspuns:

Conform definiției câștigului de informație (vedeți problema 51),

$$IG(X, Y) = H(X) - H(X|Y) = H(Y) - H(Y|X) \quad (56)$$

De asemenea, conform problemei 55.c,

$$IG(X, Y) = 0 \Leftrightarrow X \text{ și } Y \text{ sunt independente.} \quad (57)$$

Din relațiile (56) și (57) rezultă că

$$H(Y) = H(Y|X) \Leftrightarrow X \text{ și } Y \text{ sunt independente.} \quad (58)$$

Așadar, dacă X și Y sunt independente, coroborând relațiile (58) și (55) vom avea $H(X, Y) = H(Y) + H(X)$.

Invers, dacă $H(X, Y) = H(X) + H(Y)$, din relația (55) rezultă că $H(Y) = H(Y|X)$, ceea ce implică faptul că X și Y sunt independente, conform relației (58).

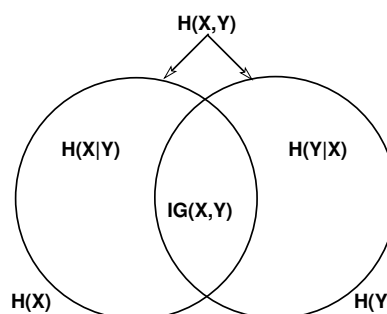
Observație: Din egalitățile (55) și (56), rezultă

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

Conform proprietății de nenegativitate a câștigului de informație ($IG(X, Y) \geq 0$; vedeți problema 55.c), rezultă

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

Această ultimă relație, precum și relațiile (55) și (56) sunt ilustrate în figura alăturată.



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

■ □ ● ○ CMU, 2011 spring, Roni Rosenfeld, HW2, pr. 3.c

Cross-entropia a două distribuții p și q , desemnată prin $CH(p, q)$, reprezintă numărul mediu de biți necesari pentru a codifica un eveniment dintr-o mulțime oarecare de posibilități, atunci când schema de codificare folosită se bazează pe o distribuție de probabilitate dată q , în loc să se bazeze pe distribuția „adevărată” p . În cazul în care distribuțiile p și q sunt discrete, această noțiune se definește formal astfel:¹⁰³

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

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

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

Observație: Ținând cont de definiția entropiei relative (cunoscută și sub numele de divergența Kullback-Leibler), vedeți pr. 55, putem scrie:

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

Cross-entropia — ca și entropia relativă; vedeți problema 55 —, 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.

¹⁰³Pentru exemple de folosire a cross-entropiei în învățarea automată, vedeți problema 15 de la capitolul *Rețele neuronale artificiale*, precum și problema 2 de la capitolul *Algoritmul EM*.

Răspuns:

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

Observație importantă: Spre deosebire de entropie (vedeți problema 129), cross-entropia nu este mărginită superior. Ea poate crește la infinit; vedeți cazul când pentru o anumită valoare x sunt adevărate simultan relațiile $p(x) \neq 0$ și $q(x) = 0$.¹⁰⁴ Prin urmare, nici entropia relativă (adică divergența Kullback-Leibler) nu este mărginită superior (vedeți *Observația* din enunț).

b. Facem calculele, folosind formula cross-entropiei:

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

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

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

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

¹⁰⁴Mai precis, $\lim_{q(x) \rightarrow 0} (-p(x) \cdot \log_2 q(x)) = -p(x)(-\infty) = +\infty$.

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

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

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

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

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

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

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.

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

Puteți face legătura dintre expresia din partea dreaptă a acestei ultime inegalități și definiția *entropiei relative* (numită de asemenea *divergența Kullback-Leibler*, vedeți problema 55) și apoi să folosiți proprietățile entropiei relative.

Răspuns:

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

- a. La problema 55.a am demonstrat inegalitatea $KL(P||Q) \geq 0$, care are loc pentru orice distribuții probabiliste discrete P și Q . Aceasta este exact proprietatea de care avem nevoie pentru a justifica inegalitatea dată în enunț la acest punct ($H(P) \leq CH(P, Q)$).¹⁰⁵

- b. Tot la problema 55.a s-a demonstrat că $KL(P||Q)$ are valoarea 0 dacă și numai dacă distribuțiile P și Q sunt identice. În contextul nostru, această proprietate se transpune imediat sub forma $H(P) = CH(P, Q) \Leftrightarrow p_i = q_i$ pentru $i = 1, \dots, n$.

¹⁰⁵Pentru o demonstrație mai directă a inegalității lui Gibbs, de această dată folosind inegalitatea lui Jensen,¹⁰⁶ vedeți nota de subsol 749 (pagina 850) de la problema 2 de la capitolul *Algoritmii EM*.

4 Arbori de decizie

Sumar

Noțiuni preliminare

- partiție a unei mulțimi: ex. 82 de la cap. *Fundamente*;
- proprietăți elementare ale funcției logaritm; formule uzuale pentru calcule cu logaritmi;
- Elemente de *teoria informației* (vedeți secțiunea corespunzătoare din cap. *Fundamente*):
 - entropie, definiție: T. Mitchell, *Machine Learning*, 1997 (desemnată în continuare simplu prin *cartea ML*), pag. 57; ex. 2.a, ex. 39.a, ex. 35.a;
 - entropie condițională specifică: ex. 14.a;
 - entropie condițională medie: ex. 2.cd, ex. 35.c;
 - câștig de informație (definiție: *cartea ML*, pag. 58): ex. 2.cd, ex. 5.a, ex. 33, ex. 39.b, ex. 35.e;
- *arbori de decizie*, văzuți ca *structură de date*: ex. 1, ex. 31
și, respectiv, ca program în logica propozițiilor: ex. 2.e, ex. 38.bc;
 - (P0) *expresivitatea arborilor de decizie* cu privire la *funcții boolene*: ex. 32;
- *spațiu de versiuni* pentru un concept (de învățat): ex. 1, ex. 31, ex. 37.

Algoritmul ID3

- pseudo-cod: *cartea ML*, pag. 56;
- *bias-ul inductiv*: *ibidem*, pag. 63-64;
- exemple simple de aplicare: ex. 2, ex. 3, ex. 5, ex. 36, ex. 37, ex. 39, ex. 40;
- ID3 ca algoritm *per se*:
 - este un *algoritm de căutare*;
spațiul de căutare — mulțimea tuturor arborilor de decizie care se pot construi cu atributele de intrare în *nodurile de test* și cu valorile atributului de ieșire în *nodurile de decizie* — este de dimensiune exponențială în raport cu numărul de atribute: ex. 1, ex. 3, ex. 31, ex. 37;
ID3 are ca *obiectiv* căutarea unui arbore / *model* care *i.* să explice cât mai bine datele (în particular, atunci când datele sunt *consistente*, modelul trebuie să fie *consistent* cu acestea), *ii.* să fie cât mai *compact*, din motive de *eficiență* la generalizare / testare și *iii.* în final să aibă o [cât mai] bună putere de *generalizare*;³⁷¹
 - ID3 ar putea fi văzut și ca algoritm de *optimizare*;³⁷²

³⁷¹LC: Alternativ, putem spune că algoritmul ID3 produce o *structură* de tip *ierarhie* (arbore) între diferite *partiționări* ale setului de instanțe de antrenament, această ierarhie fiind generată pe baza *corespondenței* dintre atributul de *ieșire* și atributele de *intrare*, care sunt adăugate la model câte unul pe rând.

³⁷²LC: Am putea să-l interpretăm pe ID3 ca fiind un algoritm care caută între diferitele *distribuții probabiliste* discrete care pot fi definite pe setul de date de antrenament una care să satisfacă cerința de *ierarhizare*, și pentru care entropia să fie *minimală* (vedeți proprietatea de structuralitate de la ex. 50 de la cap. *Fundamente*. Cerința ca arborele ID3 să fie *minimal* (ca număr de niveluri / noduri) este însă mai importantă, mai practică și mai ușor de înțeles.

- *greedy* \Rightarrow nu garantează obținerea soluției optime d.p.v. al numărului de niveluri / noduri:
ex. 4, ex. 21.a, ex. 38 (vs. ex. 37.b, ex. 3.b), ex. 46;
 - de tip *divide-et-impera* (\Rightarrow “*Iterative Dichotomizer*”), recursiv;
 - *1-step look-ahead*;
 - complexitate de timp, cf. *Weka book*:³⁷³
la antrenare, în anumite condiții: $\mathcal{O}(dm \log m)$; la testare $\mathcal{O}(d)$,
unde d este numărul de atribute, iar m este numărul de exemple;
- ID3 ca algoritm de învățare automată:
- *bias*-ul inductiv al algoritmului ID3:
[dorim ca modelul să aibă structură ierarhică, să fie compatibil / consistent cu datele dacă acestea sunt consistente (adică, necontradictorii), iar] *arborele* produs de ID3 trebuie să aibă un număr cât mai mic de niveluri / noduri;
 - algoritm de învățare de tip “*eager*”;
 - *analiza erorilor*:
la antrenare: ex. 7.a, ex. 10.a, ex. 42;³⁷⁴ [acuratețe la antrenare: ex. 6;]
la validare: CMU, 2003 fall, T. Mitchell, A. Moore, midterm, pr. 1;
la n -fold cross-validare
la cross-validare leave-one-out (CVLOO): ex. 10.b, ex. 45.bc;³⁷⁵
 - *robustețe la „zgomote” și overfitting*: ex. 10, ex. 21.bc, ex. 45, ex. 67.b;³⁷⁶
 - *zone de decizie și granițe de separare / decizie* pentru arbori de decizie cu variabile continue: ex. 10, ex. 44, ex. 45, ex. 46.
Observație: Zonele de decizie produse de algoritmul ID3 nu sunt în mod neapărat unice, fiindcă arborele de decizie creat de ID3 nu este determinat în mod unic (vedeți proprietatea (P2), mai jos).

Extensii / variante ale algoritmului ID3

- *attribute cu valori continue*: ex. 10-12, ex. 14.c, ex. 43-47; cap. *Învățare bazată pe memorare*, ex. 11.b;
- *attribute discrete cu multe valori*: ex. 13;
- *attribute cu valori nespificate / absente pentru unele instanțe*;
- *attribute cu diferite costuri asociate*: ex. 14;
- *reducerea caracterului “eager” al învățării*: ex. 16;
- *reducerea caracterului “greedy” al învățării*:
IG cu “2-step look-ahead”: ex. 17, ex. 18;
variante de tip “look-ahead” specifice atributelor continue: ex. 48;
- *folosirea altor măsuri de „impuritate” în locul câștigului de informație*:
Gini Impurity, Misclassification Impurity: ex. 15;

³⁷³“Data Mining: Practical Machine Learning Tools and Techniques with Java Implementations”, Ian Witten, Eibe Frank (3rd ed.), Morgan Kaufmann Publishers, 2011.

³⁷⁴De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

³⁷⁵De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

³⁷⁶De asemenea, ex. 4.ab de la capitolul *Clasificare bayesiană*.

- reducerea *overfitting*-ului:
 - reduced-error pruning (folosind un set de date de validare):
cartea ML, pag. 69-71; A. Cornuéjols, L. Miclet, 2nd ed., pag. 418-421;
 - rule post-pruning: cartea ML, pag. 71-72; ex. 50
 - top-down vs. bottom-up pruning, folosind un criteriu bazat pe câștigul de informație: ex. 19, ex. 49;
 - pruning folosind testul statistic χ^2 : ex. 20, ex. 51.

Proprietăți ale arborilor ID3

- (P1) arborele produs de algoritmul ID3 este *consistent* (adică, în concordanță) cu datele de antrenament, dacă acestea sunt *consistente* (adică, necontradictorii). Altfel spus, *eroarea la antrenare* produsă de algoritmul ID3 pe orice set de *date consistente* este 0: ex. 2-4, ex. 37-38;
 - (P2) arborele produs de algoritmul ID3 *nu* este în mod neapărat *unic*: ex. 3, ex. 37;
 - (P3) arborele ID3 *nu* este neapărat *optimal* (ca nr. de noduri / niveluri): ex. 4, ex. 21, ex. 38;
 - (P4) influența *atributelor identice* și, respectiv, a *instanțelor multiple* asupra arborelui ID3: ex. 8;
 - (P5) o *margine superioară* pentru *eroarea la antrenare* a algoritmului ID3, în funcție de numărul de valori ale variabilei de ieșire): ex. 7.b;
 - (P6) o aproximare simplă a numărului de *instanțe greșit clasificate* din totalul de M instanțe care au fost asignate la un nod frunză al unui arbore ID3, cu ajutorul entropiei (H) nodului respectiv: ex. 41;
 - (P7) *granițele de separare / decizie* pentru arborii ID3 cu *attribute de intrare continue* sunt întotdeauna paralele cu axele de coordonate: ex. 10, ex. 12, ex. 44, ex. 45, ex. 46 și cap. *Învățare bazată pe memorare*, ex. 11.b;
- Observație:* Următoarele trei proprietăți se referă la arbori de decizie în general, nu doar la arbori ID3.
- (P8) *adâncimea maximă* a unui arbore de decizie, când attributele de intrare sunt categoriale: numărul de attribute: ex. 52.c;
 - (P9) o *margine superioară* pentru *adâncimea* unui arbore de decizie când attributele de intrare sunt continue, iar datele de antrenament sunt (ne)separabile liniar: ex. 11;
 - (P10) o *margine superioară* pentru numărul de *noduri-frunză* dintr-un arbore de decizie, în funcție de numărul de exemple și de numărul de attribute de intrare, atunci când acestea (attributele de intrare) sunt binare: ex. 9.

Învățare automată de tip ansamblist folosind arbori de decizie: Algoritmul AdaBoost

- Noțiuni preliminare:
 - distribuție de probabilitate discretă, factor de normalizare pentru o distribuție de probabilitate, ipoteze „slabe” (engl., weak hypothesis), compas de decizie (engl., decision stump), prag de separare (engl., threshold split) pentru un

- compas de decizie, prag exterior de separare (engl., outside threshold split), eroare ponderată la antrenare (engl., weighted training error), vot majoritar ponderat (engl., weighted majority vote), overfitting, ansambluri de clasificatori (vedeți ex. 64), funcții de cost / pierdere (engl., loss function) (vedeți ex. 29 și ex. 23);
- pseudo-codul algoritmului AdaBoost + proprietăți de bază + convergența erorii la antrenare: ex. 22 și 23;
 - exemple de aplicare: ex. 24, 56, 54, 55, 57, 58.
 - **AdaBoost ca algoritm *per se*:**
algoritm iterativ, algoritm de căutare (spațiul de căutare este mulțimea combinațiilor liniare care se pot construi peste clasa de ipoteze „slabe” considerate), algoritm de optimizare secvențială (minimizează o margine superioară pentru eroarea la antrenare), algoritm greedy (dacă la fiecare iterație se alege cea mai bună ipoteză „slabă”).
 - *învățabilitate empirică γ -slabă*:
definiție: ex. 23.e
exemplificarea unor cazuri când nu există *garanție* pentru învățabilitate γ -slabă: ex. 25, 60;
 - AdaBoost ca algoritm de *optimizare secvențială* în raport cu funcția de cost / „pierdere” negativ-exponențială: ex. 26;
 - marginea de votare: ex. 27, 28 și 65;
 - *selectarea trăsăturilor* folosind AdaBoost; aplicare la clasificarea de documente: ex. 62;
 - o variantă generalizată a algoritmului AdaBoost: ex. 29 și ex. 65;
 - recapitulare (întrebări cu răspuns *adevărat* / *fals*): ex. 30 și 66.
 - **Proprietăți ale algoritmului AdaBoost:**
 - (P0) AdaBoost poate produce rezultate diferite atunci când are posibilitatea să aleagă între două sau mai multe [cele mai bune] ipoteze „slabe”: ex. 24, 54;
 - (P1) $err_{D_{t+1}}(h_t) = \frac{1}{2}$ (ex. 22.vii);
ca o *consecință*, rezultă că ipoteza h_t nu poate fi reselectată și la iterația $t + 1$; ea poate fi reselectată la o iterație ulterioară;
 - (P2) Din relația de definiție pentru distribuția D_{t+1} rezultă

$$Z_t = e^{-\alpha_t} \cdot (1 - \varepsilon_t) + e^{\alpha_t} \cdot \varepsilon_t = 2\sqrt{\varepsilon_t(1 - \varepsilon_t)}$$
și

$$\varepsilon_t \in (0, 1/2) \Rightarrow Z_t \in (0, 1)$$
(ex. 22).
 - (P3) $D_{t+1}(i) = \frac{1}{m \prod_{t'=1}^t Z_{t'}} e^{-y_i f_t(x_i)}$, unde $f_t(x_i) \stackrel{\text{def.}}{=} -\sum_{t'=1}^t \alpha_{t'} h_{t'}(x_i)$ (ex. 23.a).
Produsul $y_i f_t(x_i)$ se numește *margine algebrică*;
 - (P4) $err_S(H_t) \leq \prod_{t'=1}^t Z_{t'}$, adică eroarea la antrenare comisă de ipoteza combinată produsă de AdaBoost este majorată de produsul factorilor de normalizare (ex. 23.b);
 - (P5) AdaBoost nu optimizează în mod direct $err_S(H_t)$, ci marginea sa superioară, $\prod_{t'=1}^t Z_{t'}$; optimizarea se face în mod secvențial (greedy): la iterația t se minimizează valoarea lui Z_t ca funcție de α_t , ceea ce conduce la $\alpha_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$

(ex. 23.c);

(P5') $\varepsilon_i > \varepsilon_j \Leftrightarrow \alpha_i < \alpha_j$ (consecință imediată din relația $\alpha_t = \ln \sqrt{\frac{1 - \varepsilon_t}{\varepsilon_t}}$):

ex. 22.v;

(P6) O consecință din relația (175) și (P5): $D_{t+1}(i) = \begin{cases} \frac{1}{2\varepsilon_t} D_t(i), & i \in M \\ \frac{1}{2(1 - \varepsilon_t)} D_t(i), & i \in C \end{cases}$

(ex. 22.iv);

(P7) $err_S(H_t)$ nu neapărat descrește de la o iterație la alta; în schimb, descresc marginile sale superioare: $\prod_{t'=1}^t Z_{t'}$ și $\exp(-\sum_{t'=1}^t \gamma_{t'}^2)$ (ex. 23.d);

(P8) O condiție suficientă pentru învățabilitate γ -slabă, bazată pe marginea de votare: marginea de votare a oricărei instanțe de antrenament să fie de cel puțin 2γ , la orice iterație a algoritmului AdaBoost (ex. 28);

(P9) Orice mulțime formată din m instanțe din \mathbb{R} care sunt etichetate în mod consistent poate fi corect clasificată de către o combinație liniară formată din cel mult $2m$ compași de decizie (ex. 64.a);

(P10) Orice mulțime de instanțe distincte [și etichetate] din \mathbb{R} este γ -slab învățabilă cu ajutorul compașilor de decizie (ex. 63).

- Alte metode de învățare ansamblistă bazate pe arbori de decizie: Bagging, Random Forests;
- Alte metode de învățare automată bazate pe arbori: arbori de regresie (CART).

4.1 Arbori de decizie — Probleme rezolvate

4.1.1 Algoritmul ID3

1. (Arbori de decizie; optimalitate, relativ la numărul de noduri)

Reprezentați arborele / arborii de decizie care are / au numărul minim de noduri posibile și corespunde / corespund funcției booleene $(\neg A \vee B) \wedge \neg(C \wedge A)$ definită peste atributele booleene A, B și C .

Răspuns:

Vom determina arborele de decizie optimal (ca număr de noduri) parcurgând în mod *exhaustiv spațiul de versiuni*, adică mulțimea tuturor arborilor de decizie (construiți cu variabilele A, B și C) care sunt *consistenți* cu funcția dată. Așadar, vom examina ce se întâmplă când în nodul rădăcină se pun pe rând atributele A, B și respectiv C .

Notăm cu X funcția $(\neg A \vee B) \wedge \neg(C \wedge A)$, ale cărei valori sunt date în tabelul alăturat.

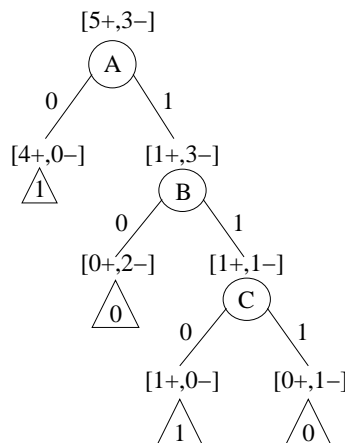
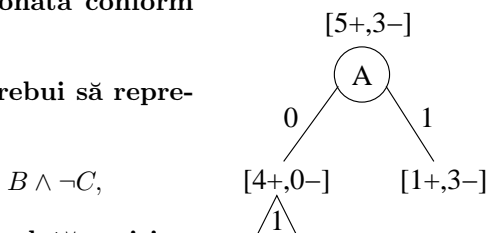
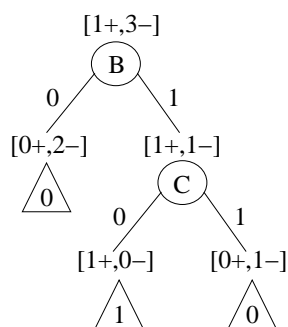
A	B	C	X
0	0	0	1
0	0	1	1
0	1	0	1
0	1	1	1
1	0	0	0
1	0	1	0
1	1	0	1
1	1	1	0

- *Cazul 1:* Dacă în nodul rădăcină se plasează atributul A , mulțimea de exemple va fi re-partiționată conform reprezentării alăturate.

Subarboarele drept al acestui arbore va trebui să reprezinte arborele de decizie pentru funcția

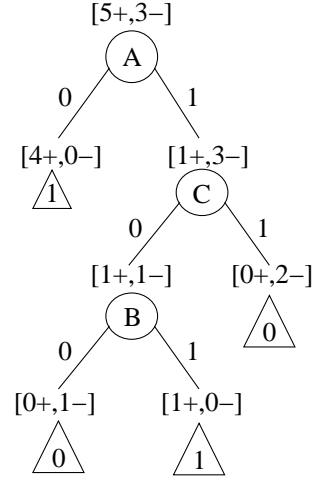
$$X_1 = X[A/1] = (\neg 1 \vee B) \wedge \neg(C \wedge 1) = B \wedge \neg C,$$

pentru care o reprezentare optimă este redată mai jos, în partea stângă:

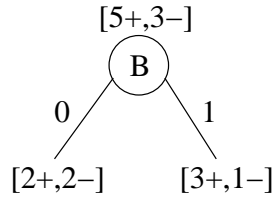


Prin urmare, un arbore optim (ca număr de noduri) care are variabila A în nodul rădăcină este cel reprezentat mai sus, în partea dreaptă.

Observație: Evident, există încă un arbore optim care are variabila A în nodul rădăcină (el corespunde unei alte reprezentări optimale a conjuncției $B \wedge \neg C$ față de cea de mai sus). Vedeți desenul alăturat.



- **Cazul 2:** Dacă în nodul rădăcină se alege atributul B , se obține partiția



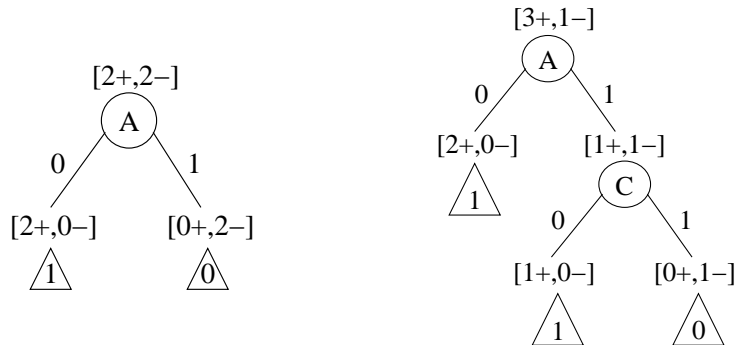
Subarboarele stâng și subarboarele drept trebuie să reprezinte arborele de decizie pentru funcțiile

$$X_2 = X[B/0] = (\neg A \vee 0) \wedge \neg(C \wedge A) = \neg A \wedge (\neg C \vee \neg A) = (\neg A \wedge \neg C) \vee \neg A = \neg A,$$

și respectiv

$$X_3 = X[B/1] = (\neg A \vee 1) \wedge \neg(C \wedge A) = 1 \wedge (\neg C \vee \neg A) = \neg C \vee \neg A,$$

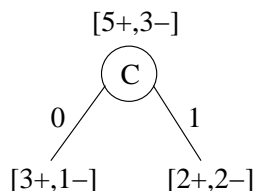
care au ca reprezentări optime arborii de mai jos:



Pentru arborele din dreapta există un arbore de decizie echivalent, obținut prin interschimbarea lui A cu C .

Prin urmare, orice arbore optim având variabila B în rădăcină are 3 niveluri și 4 noduri, așadar cu un nod (de test) mai mult decât cel determinat în primul caz.

- *Cazul 3:* În sfârșit, când în nodul rădăcină se alege atributul C , obținem partiția



Subarborele stâng și subarborele drept trebuie să reprezinte arborele de decizie pentru funcțiile

$$X_4 = X[C/0] = (\neg A \vee B) \wedge \neg(0 \wedge A) = (\neg A \vee B) \wedge \neg 0 = \neg A \vee B$$

și respectiv

$$X_5 = X[C/1] = (\neg A \vee B) \wedge \neg(1 \wedge A) = (\neg A \vee B) \wedge \neg A = \neg A.$$

Urmând un raționament similar cu cel de la cazul anterior, putem spune că orice arbore optim cu atributul C în rădăcină are 3 niveluri și 4 noduri, cu un nod (de test) mai mult decât cel determinat în primul caz.

Așadar, putem concludiona că arborii de decizie optimi corespunzători funcției date sunt cei determinați în primul caz.

2.

(Algoritmul ID3: aplicare)

■ • *CMU, 2002 spring, A. Moore, midterm example questions, pr. 2*

Ai naufragiat pe o insulă pustie, unde nu găsești niciun alt fel de hrană decât ciuperci. Despre unele dintre aceste ciuperci se știe că sunt otrăvitoare, despre altele se știe că sunt comestibile, iar despre restul nu se știe ce fel sunt. Ai rămas singur pe insulă — foștii tăi camarazi, fiind epuizați de foame, au folosit metoda ‘trial and error’... — și ai la dispoziție următoarele date:

Exemplu	<i>Ușoară</i>	<i>Mirositoare</i>	<i>ArePete</i>	<i>Netedă</i>	<i>Comestibilă</i>
<i>A</i>	1	0	0	0	1
<i>B</i>	1	0	1	0	1
<i>C</i>	0	1	0	1	1
<i>D</i>	0	0	0	1	0
<i>E</i>	1	1	1	0	0
<i>F</i>	1	0	1	1	0
<i>G</i>	1	0	0	1	0
<i>H</i>	0	1	0	0	0
<i>U</i>	0	1	1	1	?
<i>V</i>	1	1	0	1	?
<i>W</i>	1	1	0	0	?

Atunci când nu vei mai avea la dispoziție pentru a supraviețui decât ciuperci U , V , sau W , ai putea estima care dintre ele sunt comestibile, folosind arbori de decizie.

În primele trei întrebări care urmează, ne vom referi la ciupercile $A - H$:

- Care este entropia atributului *Comestibilă*?
- Doar privind datele — adică fără a face explicit calculul câștigului de informație (engl., information gain) pentru cele patru atribute — poți determina ce atribut vei alege ca rădăcină a arborelui de decizie?
- Calculează câștigul de informație pentru atributul pe care l-ai ales la întrebarea precedentă.
- Elaborează întregul arbore de decizie ID3 bazat pe datele din tabel și apoi clasifică ciupercile U, V, W.
- Exprimă cu ajutorul calculului propozițional (logica predicatelor de ordinul 0) clasificarea produsă de arborele de decizie obținut. (*Comestibilă* $\leftrightarrow \dots$)
- Există vreun risc dacă vei consuma ciuperci care au fost clasificate de arborele de decizie ca fiind comestibile? De ce da? sau, de ce nu?

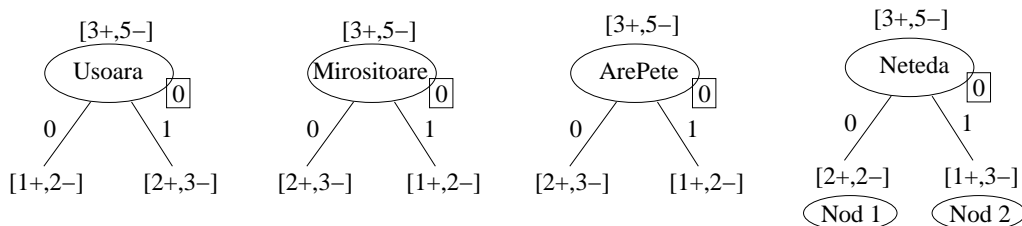
Răspuns:

a. Entropia atributului *Comestibilă* este:

$$\begin{aligned} H_{\text{Comestibilă}} &\stackrel{\text{not.}}{=} H[3+, 5-] \stackrel{\text{def.}}{=} -\frac{3}{8} \log_2 \frac{3}{8} - \frac{5}{8} \log_2 \frac{5}{8} = \frac{3}{8} \log_2 \frac{8}{3} + \frac{5}{8} \log_2 \frac{8}{5} = \\ &= \frac{3}{8} 3 - \frac{3}{8} \log_2 3 + \frac{5}{8} 3 - \frac{5}{8} \log_2 5 = 3 - \frac{3}{8} \log_2 3 - \frac{5}{8} \log_2 5 \approx \\ &\approx 0.9544 \end{aligned}$$

Observație (1): Notăția $[3+, 5-]$ simbolizează o mulțime partiționată în 3 exemple pozitive și 5 exemple negative. Vom folosi acest gen de notație peste tot în continuare, cu mici variații determinate de valorile pe care le poate lua atributul de ieșire. De exemplu, dacă vorbim despre o mulțime cu 5 obiecte roșii, 3 albastre și 4 verzi, am putea nota: $[5R, 3A, 4V]$.

b. În rădăcina arborelui de decizie se alege atributul care aduce cel mai mare câștig de informație. Adică, atributul care, intuitiv vorbind, partiționează cel mai bine datele de antrenament în raport cu atributul de ieșire. În cazul nostru, variantele pe care le avem la dispoziție pentru rădăcina arborelui (nodul 0) sunt:



Este ușor de observat că atributele *Ușoară*, *Mirositoare* și *ArePete* împart mulțimea exemplilor în mod similar: o submulțime cu 3 elemente, dintre care unul este pozitiv iar două sunt negative, și o submulțime cu 5 elemente, dintre care două sunt pozitive, iar trei sunt negative.

Dacă am considera un arbore de decizie cu un singur nod de test în care plasăm atributul *Netedă*, atunci numărul minim de erori la antrenare pe care îl putem obține este 3, utilizând următoarea clasificare:

- $Neted\grave{a} = 0$: $Comestibil\grave{a} = 1 \Rightarrow$ ciupercile E și H sunt clasificate greșit
- $Neted\grave{a} = 1$: $Comestibil\grave{a} = 0 \Rightarrow$ ciuperca C este clasificată greșit

Dacă, în schimb, vom pune în rădăcina arborelui de decizie unul dintre celelalte trei atribute, spre exemplu atributul $Ușoar\grave{a}$, și dacă vom lua votul majoritar în fiecare nod descendent din nodul rădăcină, eroarea rezultată la antrenare va fi aceeași ca mai sus ($3/8$), însă toate instanțele vor fi clasificate la fel (și anume, negativ). Dacă nu lucrăm cu vot majoritar pentru ambii descendenți, ci doar pentru cel cu entropie mai mică (în vreme ce pentru celălalt nod descendent luăm decizia contrară), se observă că pentru atributul $Ușoar\grave{a}$ vom obține 4 erori pe setul de antrenament, iar pentru atributul $Neted\grave{a}$ vom obține 3 erori.

Sumarizând, suntem înclinați să credem că ar fi o alegere sensibil mai bună să punem în rădăcină atributul $Neted\grave{a}$. Pentru o justificare numerică riguroasă a acestei alegeri folosind criteriul maximizării câștigului de informație, vedeți punctul d .

c. Pentru a obține câștigul de informație pentru atributul $Neted\grave{a}$, se fac calculele:

$$\begin{aligned} H_{0/Neted\grave{a}} &\stackrel{def.}{=} \frac{4}{8}H[2+, 2-] + \frac{4}{8}H[1+, 3-] = \frac{1}{2} \cdot 1 + \frac{1}{2} \left(\frac{1}{4} \log_2 \frac{4}{1} + \frac{3}{4} \log_2 \frac{4}{3} \right) \\ &= \frac{1}{2} + \frac{1}{2} \left(\frac{1}{4} \cdot 2 + \frac{3}{4} \cdot 2 - \frac{3}{4} \log_2 3 \right) = \frac{1}{2} + \frac{1}{2} \left(2 - \frac{3}{4} \log_2 3 \right) \\ &= \frac{1}{2} + 1 - \frac{3}{8} \log_2 3 = \frac{3}{2} - \frac{3}{8} \log_2 3 \approx 0.9056 \end{aligned}$$

$$IG_{0/Neted\grave{a}} \stackrel{def.}{=} H_{Comestibil\grave{a}} - H_{0/Neted\grave{a}} = 0.9544 - 0.9056 = 0.0488$$

Observație (2): În cele de mai sus am notat cu $H_{0/Neted\grave{a}}$ entropia partiției [multimii de exemple de antrenament] determinate de alegerea atributului $Neted\grave{a}$ în nodul 0,³⁷⁷ iar cu $IG_{0/Neted\grave{a}}$ câștigul de informație corespunzător acestei alegeri. În general, prin notația $H_{n/A}$ vom înțelege entropia partiției determinate de alegerea atributului A în nodul n .

d. Arborele de decizie ID3 se construiește pornind din rădăcină și alegând atributul pentru fiecare nod de test în modul următor:

Nodul 0 (rădăcina):

Să verificăm dacă alegerea făcută la punctul b este cea corectă:

$$\begin{aligned} H_{0/Ușoar\grave{a}} &\stackrel{def.}{=} \frac{3}{8}H[1+, 2-] + \frac{5}{8}H[2+, 3-] \\ &= \frac{3}{8} \left(\frac{1}{3} \log_2 \frac{3}{1} + \frac{2}{3} \log_2 \frac{3}{2} \right) + \frac{5}{8} \left(\frac{2}{5} \log_2 \frac{5}{2} + \frac{3}{5} \log_2 \frac{5}{3} \right) \\ &= \frac{3}{8} \left(\frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 3 - \frac{2}{3} \cdot 1 \right) + \frac{5}{8} \left(\frac{2}{5} \log_2 5 - \frac{2}{5} \cdot 1 + \frac{3}{5} \log_2 5 - \frac{3}{5} \log_2 3 \right) \\ &= \frac{3}{8} \left(\log_2 3 - \frac{2}{3} \right) + \frac{5}{8} \left(\log_2 5 - \frac{3}{5} \log_2 3 - \frac{2}{5} \right) \\ &= \frac{3}{8} \log_2 3 - \frac{2}{8} + \frac{5}{8} \log_2 5 - \frac{3}{8} \log_2 3 - \frac{2}{8} = \frac{5}{8} \log_2 5 - \frac{4}{8} \approx 0.9512 \end{aligned}$$

³⁷⁷Mai riguros, folosind terminologia din *Teoria informației*, vom spune că notația $H_{0/Neted\grave{a}}$ se referă la entropia condițională medie a atributului de ieșire $Comestibil\grave{a}$ în raport cu atributul de intrare $Neted\grave{a}$.

Urmează că

$$IG_{0/U\text{șoară}} \stackrel{\text{def.}}{=} H_{Comestibilă} - H_{0/U\text{șoară}} = 0.9544 - 0.9512 = 0.0032,$$

deci

$$IG_{0/U\text{șoară}} = IG_{0/Mirositoare} = IG_{0/ArePete} = 0.0032 < IG_{0/Netedă} = 0.0488$$

Am avut deci dreptate să alegem atributul *Netedă* la punctul b.

Observație importantă:

În loc să fi calculat efectiv aceste câștiguri de informație, pentru a determina atributul cel mai „bun”, ar fi fost *suficient* să elaborăm un *raționament* de tip *relațional*, bazat pe comparația dintre valorile entropiilor condiționale medii $H_{0/Netedă}$ și $H_{0/U\text{șoară}}$:

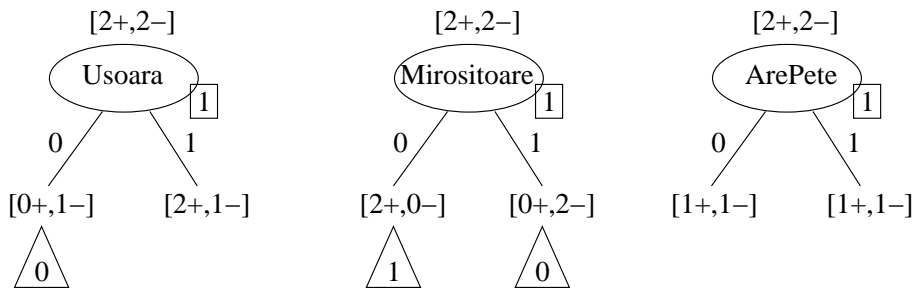
$$\begin{aligned} IG_{0/Netedă} > IG_{0/U\text{șoară}} &\Leftrightarrow H_{0/Netedă} < H_{0/U\text{șoară}} \\ &\Leftrightarrow \frac{3}{2} - \frac{3}{8} \log_2 3 < \frac{5}{8} \log_2 5 - \frac{1}{2} \Leftrightarrow 12 - 3 \log_2 3 < 5 \log_2 5 - 4 \\ &\Leftrightarrow 16 < 5 \log_2 5 + 3 \log_2 3 \Leftrightarrow 16 < 11.6096 + 4.7548 \text{ (adev.)} \end{aligned}$$

În mod *alternativ*, ținând cont de relația (197) de la problema 35, putem proceda chiar *mai simplu* relativ la calcule (nu doar aici, ci ori de câte ori *nu* avem de-a face cu un *număr mare de instanțe*):

$$\begin{aligned} H_{0/Netedă} < H_{0/U\text{șoară}} &\Leftrightarrow \frac{4^4}{2^2 \cdot 2^2} \cdot \frac{4^4}{3^3} < \frac{5^5}{2^2 \cdot 3^3} \cdot \frac{3^3}{2^2} \Leftrightarrow \frac{4^8}{3^3} < 5^5 \Leftrightarrow 4^8 < 3^3 \cdot 5^5 \\ &\Leftrightarrow 2^{16} < 3^3 \cdot 5^5 \Leftrightarrow 64 \cdot 2^{10} < 27 \cdot 25 \cdot 125 \text{ (adev.)} \end{aligned}$$

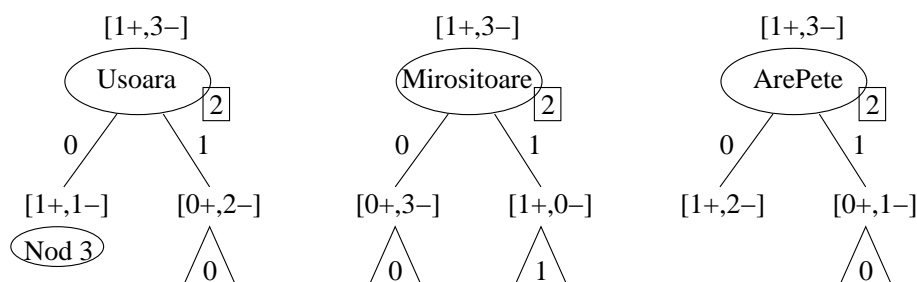
Vă sfătuim să *procedați așa* la rezolvarea problemelor propuse din acest capitol, acolo unde este cazul.

Nodul 1: Trebuie să clasificăm acele exemple care au $Netedă = 0$; avem de ales între 3 atribute - *Ușoară*, *Mirositoare* și *ArePete*.



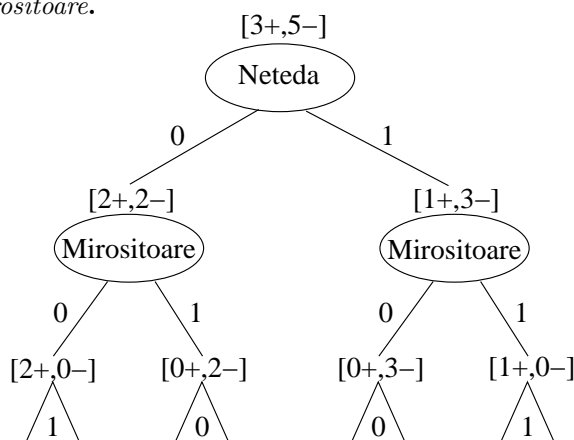
Avem $H_{1/Mirositoare} = \frac{2}{4}H[2+, 0-] + \frac{2}{4}H[0+, 2-] = 0$. Oricare ar fi valorile pentru $H_{1/ArePete}$ și $H_{1/U\text{șoară}}$, întrucât știm că entropia are întotdeauna valori nenegative, rezultă că atributul *Mirositoare* maximizează în nodul 1 câștigul de informație. În imaginea de mai sus valorile din triunghi reprezintă decizia luată de subarboarele construit în nodul-frunză respectiv.

Nodul 2: Avem de clasificat exemplele pentru care $Netedă = 1$. Atributele disponibile sunt: *Ușoară*, *Mirositoare* și *ArePete*.



Evident, $H_{2/Mirositoare} = \frac{3}{4}H[0+, 3-] + \frac{1}{4}H[1+, 0-] = \frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 0 = 0$ Așadar, pentru nodul 2 putem alege atributul *Mirositoare*.

Arborele complet arată astfel:



Parcurgând arborele construit, ciupercile U , V și W vor fi clasificate astfel:

U	$Netedă = 1, Mirositoare = 1 \Rightarrow Comestibilă = 1$
V	$Netedă = 1, Mirositoare = 1 \Rightarrow Comestibilă = 1$
W	$Netedă = 0, Mirositoare = 1 \Rightarrow Comestibilă = 0$

e. $Comestibilă \leftrightarrow (\neg Netedă \wedge \neg Mirositoare) \vee (Netedă \wedge Mirositoare)$

Același lucru poate fi exprimat și sub forma unui pseudo-cod *if ... then Comestibilă else \neg Comestibilă*:

```

IF      (Netedă = 0 AND Mirositoare = 0) OR
        (Netedă = 1 AND Mirositoare = 1)
THEN    Comestibilă;
ELSE     $\neg$ Comestibilă;

```

f. Arborele de decizie produs de către algoritmul ID3 elaborat mai sus este consistent cu datele de antrenament pe care le-am avut la dispoziție (fiindcă aceste date sunt necontradictorii). Întrucât în realitate clasificarea poate depinde și de alte trăsături / informații decât cele de care dispunem noi, nu avem garanția că arborele ID3 face identificarea corectă a etichetei / clasei pentru toate instanțele din setul de test. Așadar, nu putem fi siguri că nu ne vom îmbolnăvi dacă vom consuma ciupercile U și V , sau că ne vom îmbolnăvi dacă vom consuma ciuperca W . În multe aplicații practice, calitatea unui model de învățare automată (în cazul de față, un arbore de decizie) se verifică pe un set de *date de validare*.

3. (Algoritmul ID3, aplicat pe expresii booleene; exploatarea simetriilor operațiilor \vee, \wedge în alegerea nodurilor; analiza „optimalității” arborelui ID3, ca număr de noduri de test)

*prelucrare de Liviu Ciortuz, după
Tom Mitchell, “Machine Learning”, 1997, ex. 3.1.b*

Considerăm următoarea funcție booleană: $A \vee (B \wedge C)$. Presupunem că această funcție este deja definită — adică valoarea ei este cea cunoscută din logica propozițiilor —, însă dorim să o reprezentăm ca arbore de decizie.

a. Aplicați algoritmul ID3 [tabelei de adevăr corespunzătoare] acestei funcții. *Observație:* Dacă exploatați simetriile, veți avea nevoie doar de puține calcule, altfel vă veți complica inutil.

b. Arborele ID3 obținut la punctul precedent este optimal?

Alfel spus, puteți găsi un alt arbore de decizie, de adâncime mai mică sau cu număr mai mic de noduri (comparativ cu arborele obținut la punctul a), care să reprezinte această funcție? (Țineți cont că în fiecare nod al unui arbore de decizie se poate testa un singur atribut.)

Răspuns:

a. Observăm că funcția dată este simetrică în B și C , datorită comutativității operatorului logic \wedge . O consecință a acestui fapt este că dacă, pe parcursul algoritmului ID3, avem de ales (și) între cele două atribute este nevoie să-l studiem doar pe unul dintre ele, celălalt comportându-se identic.

A	B	C	$Y = A \vee (B \wedge C)$
0	0	0	0
0	0	1	0
0	1	0	0
0	1	1	1
1	0	0	1
1	0	1	1
1	1	0	1
1	1	1	1

Nodul 0 (rădăcină):

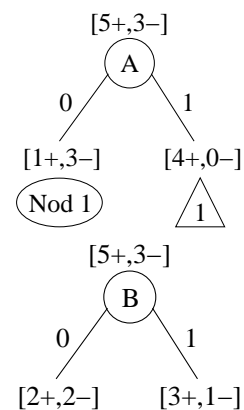
$$\begin{aligned} H_{0/A} &= \frac{4}{8}H[1+, 3-] + \frac{4}{8}H[4+, 0-] = \\ &= \frac{1}{2}H[1+, 3-] + \frac{1}{2} \cdot 0 = \frac{1}{2}H[1+, 3-] \end{aligned}$$

$$\begin{aligned} H_{0/B} &= \frac{4}{8}H[2+, 2-] + \frac{4}{8}H[3+, 1-] = \\ &= \frac{1}{2} \cdot 1 + \frac{1}{2}H[3+, 1-] = \frac{1}{2} + \frac{1}{2}H[1+, 3-] \end{aligned}$$

Este evident că $H_{0/A} < H_{0/B}$, deci vom alege atributul A în rădăcină.

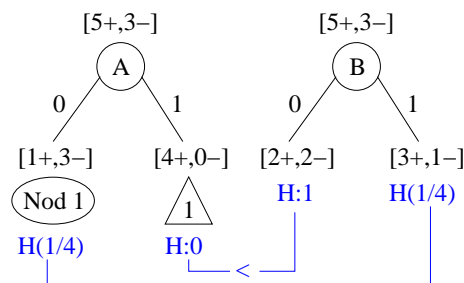
Observații importante:

1. La aceeași concluzie se putea ajunge *imediat* pe baza unui *raționament* de tip *calitativ*, și anume, comparând atent cei doi arbori („compași de decizie”) de mai sus. Mai precis, vom compara (două câte două) entropiile



condiționale specifice din nodurile descendente, precum și ponderile cu care se combină aceste entropii în scrierea entropiilor condiționale medii corespunzătoare atributelor A și B .

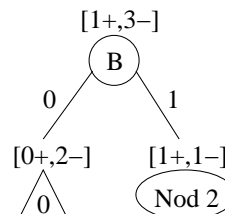
Putem pune în evidență acest fapt în figura alăturată, în care simbolul H , scris uneori însoțit de un argument (așadar, ca $H(p)$), se referă la entropia unei variabile Bernoulli de parametru p .



Mai facem *precizarea* că semnele $<$ și $=$ din figura de mai sus se referă de fapt nu [doar] la entropiile condiționale specifice, ci [și] la produsul acestora cu ponderile asociate în mod corespunzător: $\frac{4}{8}H[1+,3-] = \frac{4}{8}H[3+,1-]$ și respectiv $\frac{4}{8}H[4+,0-] < \frac{4}{8}H[2+,2-]$.

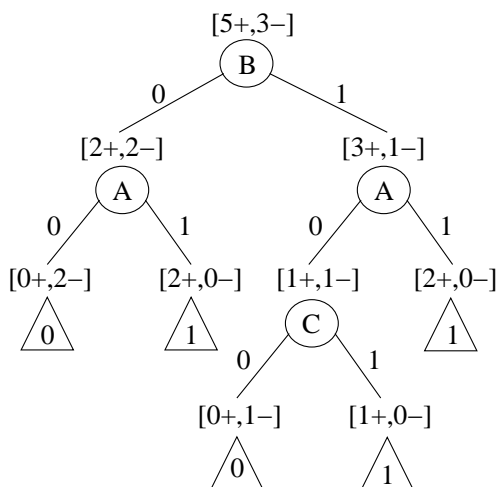
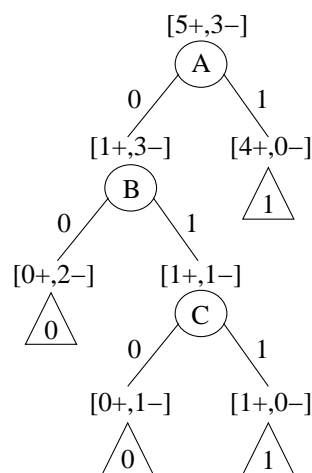
2. Pentru câteva formule de calcul convenabile pentru entropii și câștiguri de informație relative la compași de decizie, atunci când se folosește calculatorul de buzunar, dar numărul de instanțe de antrenament asociate nu este prea mare, vedeți problema 35.

Nodul 1: Avem de clasificat instanțele care au $A = 0$ și putem alege între atributele B și C . Datorită simetriei, îl putem alege pe oricare dintre ele. Pentru fixare, îl alegem pe B .



Nodul 2: La acest punct a mai rămas disponibil doar atributul C .

Arborele construit de ID3 este cel reprezentat mai jos, în partea stângă:³⁷⁸



³⁷⁸Un alt arbore ID3 este cel obținut din acesta interschimbând atributele B și C . (Vedeți *Observația* din enunț.)

b. Pentru a vedea dacă arborele construit de algoritmul ID3 este cel optimal, trebuie să reconsiderăm toate deciziile pe care le-am luat în construirea acestuia:

– La nodul 1 al arborelui avem de clasificat exemplele pentru care $A = 0$, deci funcția care trebuie reprezentată de subarborele în cauză este $f' = f[A/0] = 0 \vee (B \wedge C) = B \wedge C$, funcție care este reprezentată în mod optimal de subarborele construit de ID3. (Notăția $A/0$ semnifică faptul că variabila logică A este instanțiată la valoarea 0.) Prin urmare, nu există un arbore mai bun care să reprezinte funcția dată și să aibă în rădăcină atributul A .

– În rădăcină am ales atributul A în detrimentul celorlalte două attribute deoarece am demonstrat că aduce cel mai mare câștig de informație. Să vedem ce se întâmplă dacă alegem unul dintre attributele B sau C . După cum am discutat mai sus, datorită simetriei, pe oricare dintre cele două l-am alege, arborele rezultat ar avea aceeași formă. Pentru fixare, să-l alegem pe B .

Subarborele stâng și drept vor trebui să reprezinte funcțiile:

$$f'' = f[B/0] = A \vee (0 \wedge C) = A \vee 0 = A$$

și respectiv

$$f''' = f[B/1] = A \vee (1 \wedge C) = A \vee C.$$

Arborele minimal care poate fi construit în aceste circumstanțe este cel reprezentat mai sus în partea dreaptă. Acest arbore are 3 niveluri și 4 noduri, cu un nod în plus față de cel construit de algoritmul ID3.

Putem deci conchide că arborele construit respectând specificațiile algoritmului ID3 este cel optimal.

Observație:

Această problemă pune în evidență două modalități de parcurgere a spațiului de versiuni pentru un concept, în particular unul din logica propozițiilor, care este reprezentat cu ajutorul arborilor de decizie. Pe de o parte avem explorarea (incompletă) făcută de algoritmul ID3 care este de tip “greedy”, iar pe de altă parte avem explorarea exhaustivă. Prima strategie de explorare procedează la o căutare „orientată” a soluției (și din această cauză este mai eficientă, dar se va vedea, ca revers, că nu asigură întotdeauna găsirea optimului), iar cea de-a doua strategie de explorare, deși asigură găsirea optimului, nu este utilizabilă în cazurile (frecvente!) în care spațiul de versiuni este foarte mare.

4. (ID3, ca algoritm “greedy”;
un exemplu când arborele ID3 nu este optimal
ca număr de noduri și de niveluri)

prelucrare de Liviu Ciortuz, după

■ CMU, 2003 fall, T. Mitchell, A. Moore, midterm exam, pr. 9.a

Fie attributele binare de intrare A, B, C , atributul de ieșire Y și următoarele exemple de antrenament:

A	B	C	Y
1	1	0	0
1	0	1	1
0	1	1	1
0	0	1	0

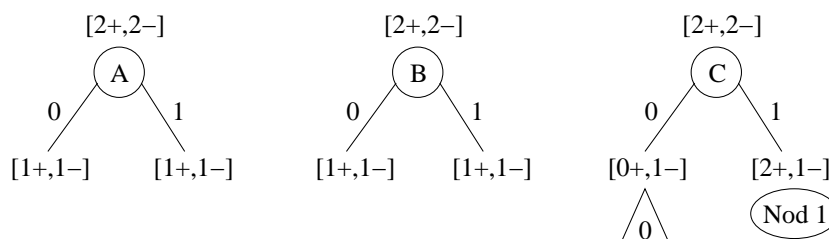
a. Determinați arborele de decizie calculat de algoritmul ID3. Este acest arbore de decizie *consistent* cu datele de antrenament?

b. Există un arbore de decizie de adâncime mai mică (decât cea a arborelui ID3) consistent cu datele de mai sus? Dacă da, ce concept (logic) reprezintă acest arbore?

Răspuns:

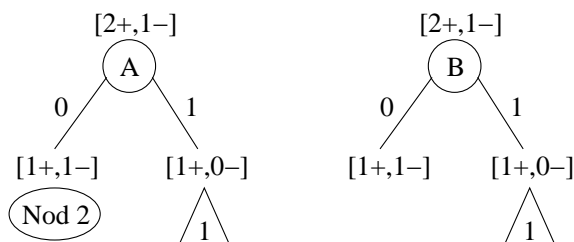
a. Se construiește arborele de decizie cu algoritmul ID3 astfel:

Nodul 0 (rădăcina):



Observăm că $H_{0/A} = H_{0/B} = \frac{2}{4}H[1+,1-] + \frac{2}{4}H[1+,1-] = H[1+,1-] = 1$, care este valoarea maximă a entropiei [condiționale medii a] unei variabile booleene. Prin urmare, $H_{0/C}$ nu poate fi decât mai mică sau egală cu $H_{0/A}$ și $H_{0/B}$. Deci vom alege în nodul rădăcină atributul C .

Nodul 1: Avem de clasificat instanțele cu $C = 1$, deci alegerea se face între atributele A și B .

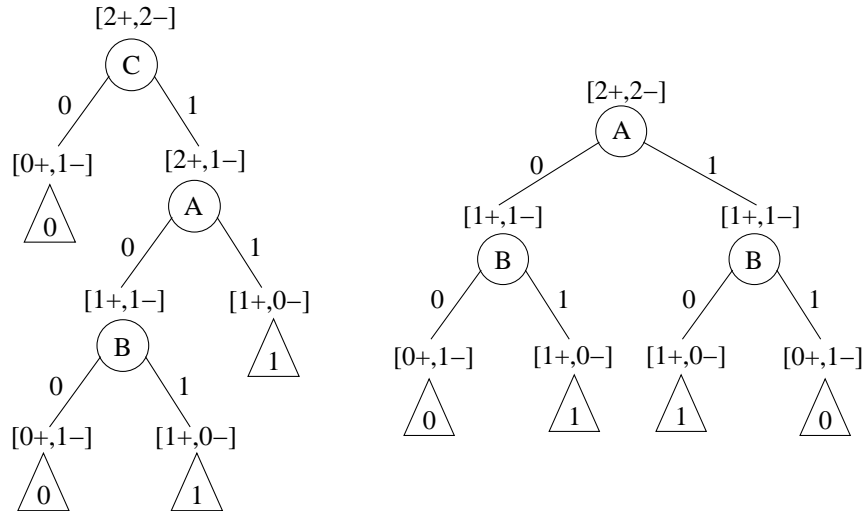


Cele două entropii condiționale medii sunt egale:

$$H_{1/A} = H_{1/B} = \frac{2}{3}H[1+,1-] + \frac{1}{3}H[1+,0-]$$

Așadar, putem alege oricare dintre cele două atribute. Pentru fixare, îl alegem pe A .

Nodul 2: La acest nod nu mai avem decât atributul B , deci îl vom pune pe acesta. Arborele complet este reprezentat în partea stângă:



Prin construcție, arborele ID3 este consistent cu datele de antrenament dacă acestea sunt consistente (i.e., necontradictorii). În cazul nostru, se verifică imediat că datele de antrenament sunt consistente.

b. Se observă că atributul de ieșire Y reprezintă de fapt funcția logică $A \text{ XOR } B$. Reprezentând această funcție ca arbore de decizie, vom obține arborele desenat mai sus în partea dreaptă. Acest arbore are cu un nivel mai puțin decât arborele construit cu algoritmul ID3. Prin urmare, arborele obținut de algoritmul ID3 pe datele din enunț *nu* este *optim* din punctul de vedere al numărului de niveluri. Aceasta este o *consecință* a caracterului “greedy” al algoritmului ID3, datorat faptului că la fiecare iterație alegem „cel mai bun” atribut în raport cu criteriul câștigului de informație. Se știe că algoritmi de tip “greedy” nu garantează obținerea optimului global.

5. (Clasificare ternară: “decision stump” produs de ID3, pe date care conțin duplicări și „zgomete“)

Presupunem că se dau șase date de antrenament (precizate în tabel) pentru o problemă de clasificare cu două atribute binare și trei clase $Y \in \{1, 2, 3\}$. Se va crea un arbore ID3, bazat pe câștigul de informație.

- a. Calculați câștigul de informație atât pentru X_1 cât și pentru X_2 . Se va folosi aproximarea $\log_2 3 = 19/12$ și se va scrie câștigul de informație sub formă de fracții.

X_1	X_2	Y
1	1	1
1	1	1
1	1	2
1	0	3
0	0	2
0	0	3

- b. Pe baza rezultatelor anterioare, ce atribut va fi folosit pentru primul nod al arborelui ID3? Desenați arborele de decizie care rezultă folosind doar acest singur nod. Etichetați corespunzător nodul, ramurile și eticheta prevăzută în fiecare frunză.

- c. Cum va clasifica acest arbore instanța determinată de $X_1 = 0$ și $X_2 = 1$?

Răspuns:

a. Redăm formula pentru calculul câștigului de informație în varianta folosită de Tom Mitchell:

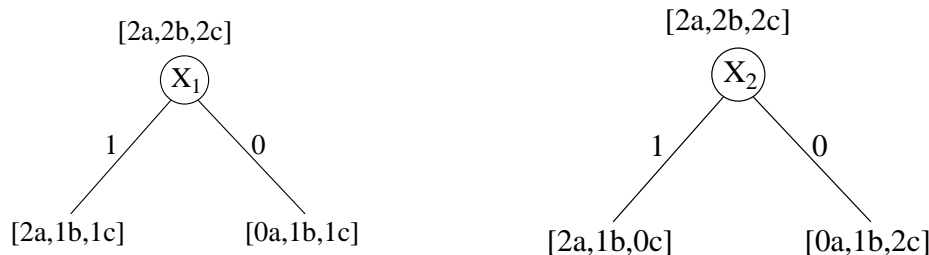
$$Gain(X_i) = Entropy(S) - \sum_{v \in Values(X_i)} \frac{|S_v|}{|S|} Entropy(S_v)$$

unde S este mulțimea celor 6 date de antrenament, iar

$$Entropy(S) = - \sum_{y \in Y} p_y \log_2 p_y$$

Notăm cu a clasa instanțelor având eticheta $Y = 1$, cu b clasa instanțelor cu $Y = 2$ și cu c clasa instanțelor cu $Y = 3$. În mulțimea S există câte 2 elemente din fiecare clasă și atunci putem scrie că $S = [2a, 2b, 2c]$.

Pentru nodul rădăcină, putem alege fie atributul X_1 , fie atributul X_2 , ceea ce determină următoarele împărțiri ale mulțimii S :



Putem calcula câștigurile de informație pentru cele două atribute:

$$Gain(X_1) = H[2a, 2b, 2c] - \left(\frac{4}{6} H[2a, 1b, 1c] + \frac{2}{6} H[0a, 1b, 1c] \right)$$

unde $H[2a, 2b, 2c]$ este o altă notație pentru entropia mulțimii compuse din două exemple de clasă a , două de clasă b și două de clasă c .

Calculăm entropiile care intervin în formulă:

$$\begin{aligned} H[2a, 2b, 2c] &= -\frac{2}{6} \log_2 \frac{2}{6} - \frac{2}{6} \log_2 \frac{2}{6} - \frac{2}{6} \log_2 \frac{2}{6} \\ &= -3 \cdot \frac{2}{6} \log_2 \frac{2}{6} = -\log_2 \frac{1}{3} = \log_2 3 = \frac{19}{12} \\ H[2a, 1b, 1c] &= -\frac{2}{4} \log_2 \frac{2}{4} - \frac{1}{4} \log_2 \frac{1}{4} - \frac{1}{4} \log_2 \frac{1}{4} = +\frac{1}{2} \cdot \log_2 2 + \frac{1}{4} \log_2 4 + \frac{1}{4} \log_2 4 \\ &= \frac{1}{2} + \frac{1}{4} \cdot 2 + \frac{1}{4} \cdot 2 = \frac{1}{2} + \frac{1}{2} + \frac{1}{2} = \frac{3}{2} \\ H[0a, 1b, 1c] &= 0 - \frac{1}{2} \cdot \log_2 2 - \frac{1}{2} \cdot \log_2 2 = \log_2 2 = 1 \end{aligned}$$

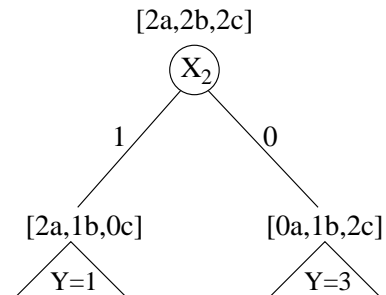
Înlocuind aceste valori numerice în formula câștigului de informație, obținem:

$$Gain(X_1) = \frac{19}{12} - \left(\frac{2}{3} \cdot \frac{3}{2} + \frac{1}{3} \cdot 1 \right) = \frac{19}{12} - \left(1 + \frac{1}{3} \right) = \frac{19}{12} - \frac{4}{3} = \frac{3}{12} = \frac{1}{4}$$

Se aplică aceleași formule și pentru atributul X_2 :

$$\begin{aligned}
 Gain(X_2) &= H[2a, 2b, 2c] - \left(\frac{3}{6} H[2a, 1b, 0c] + \frac{3}{6} H[0a, 1b, 2c] \right) \\
 H[2a, 1b, 0c] &= H[0a, 1b, 2c] \\
 &= -\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} = -\frac{2}{3} (\log_2 2 - \log_2 3) - \frac{1}{3} (-1) \log_2 3 \\
 &= -\frac{2}{3} \cdot 1 + \frac{2}{3} \cdot \frac{19}{12} + \frac{1}{3} \cdot \frac{19}{12} = \frac{19}{12} - \frac{2}{3} = \frac{11}{12} \\
 \Rightarrow Gain(X_2) &= \frac{19}{12} - \left(\frac{1}{2} \cdot \frac{11}{12} + \frac{1}{2} \cdot \frac{11}{12} \right) = \frac{19}{12} - \frac{11}{12} = \frac{8}{12} = \frac{2}{3}.
 \end{aligned}$$

b. Deoarece $Gain(X_1) = \frac{3}{12}$, $Gain(X_2) = \frac{8}{12}$, și deci $Gain(X_1) < Gain(X_2)$, se va alege atributul X_2 ca rădăcină a arborelui. Arborele de decizie construit din acest singur nod este cel din figura alăturată.



c. O instanță care are $X_1 = 0$ și $X_2 = 1$ va fi clasificată de acest arbore cu $Y = 1$ (cu probabilitate $2/3$).

6. (Algoritmul ID3: aplicare pe date inconsistente, “decision stumps”, calculul acurateții)

• ◦ CMU, 2012 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 2.ab

Tabelul de mai jos sumarizează situația celor 2201 de pasageri și membri ai echipajului de la bordul vasului Titanic, în urma naufragiului din data de 15 Aprilie 1912. Pentru fiecare combinație de valori ale celor 3 variabile (Clasă, Sex, Vârstă) am indicat în tabel câți oameni au supraviețuit și câți nu au supraviețuit. (*Observație:* Datele originale au patru valori pentru atributul Clasă; am comasat valorile II, III, și Echipaj într-o singură valoare, denumită „Inferioară“.)

Clasa	Sexul	Vârsta	Supraviețuitori		
			Nu	Da	Total
I	Masculin	Copil	0	5	5
I	Masculin	Adult	118	57	175
I	Feminin	Copil	0	1	1
I	Feminin	Adult	4	140	144
Inferioară	Masculin	Copil	35	24	59
Inferioară	Masculin	Adult	1211	281	1492
Inferioară	Feminin	Copil	17	27	44
Inferioară	Feminin	Adult	105	176	281
Total			1490	711	2201

Pentru a vă ușura calculele pe care va trebui să le faceți, am făcut noi totalurile pentru fiecare variabilă:

Clasa	Supraviețuitori		
	Nu	Da	Total
I	122	203	325
Inferioară	1368	508	1876

Sexul	Supraviețuitori		
	Nu	Da	Total
Masculin	1364	367	1731
Feminin	126	344	470

Vârsta	Supraviețuitori		
	Nu	Da	Total
Copil	52	57	109
Adult	1438	654	2092

a. Folosind un arbore de decizie, dorim să prezicem variabila de ieșire Y (Supraviețuitor), pornind de la atributele de intrare C (Clasa), S (Sexul), V (Vârsta). Utilizați criteriul câștigului de informație pentru a alege care dintre aceste trei atribute C , S sau V trebuie să fie folosit în nodul-rădăcină al arborelui de decizie.

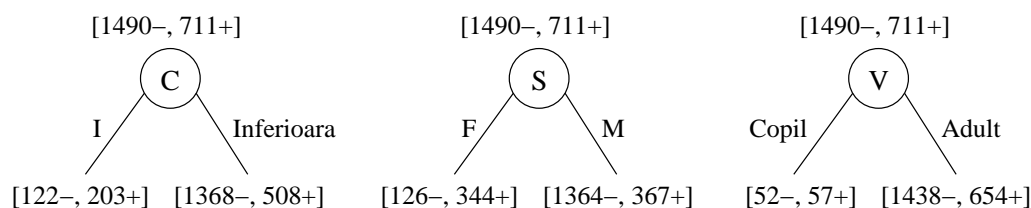
De fapt, ce vi se cere este să învățați un arbore de decizie de adâncime 1 care folosește doar atributul din rădăcină pentru a clasifica datele. (Astfel de arbori de decizie de adâncime 1 sunt adesea numiți în terminologia de limbă engleză “decision stumps”.) Parcurgeți toate etapele rezolvării, redând inclusiv calculele pentru câștigul de informație al fiecărui atribut.

b. Care este acuratețea [medie] obținută pe datele de antrenament de către arborele de decizie cu adâncime 1 de la punctul precedent?

c. Dacă ați crea un arbore de decizie care folosește toate cele trei variabile, care ar fi acuratețea lui [medie] pe datele de antrenament? (*Observație:* Nu trebuie neapărat să creați arborele de decizie pentru a afla răspunsul!)

Răspuns:

a. Totalurile care au fost furnizate în enunț pentru fiecare dintre variabilele C , S și V ne servesc foarte bine pentru a crea rapid cei trei “decision stumps”:



Analizând datele conform figurii de mai sus, se poate „intui” că atributul S va avea un câștig de informație (în raport cu atributul de ieșire Y – *Supraviețuitor*) mai bun decât al celorlalte două atribute de intrare (C și V). Intuiția se verifică făcând calculele:

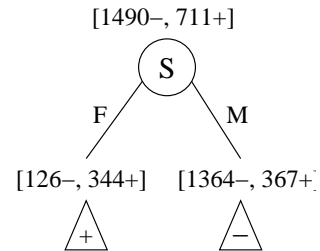
$$\begin{aligned}
 IG(Y, C) &= H[1490-, 711+] - \left(\frac{325}{2201} H[122-, 203+] + \frac{1876}{2201} H[1368-, 508+] \right) \\
 &= 0.048501 \\
 IG(Y, S) &= H[1490-, 711+] - \left(\frac{470}{2201} H[126-, 344+] + \frac{1731}{2201} H[1364-, 367+] \right) \\
 &= 0.142391
 \end{aligned}$$

$$\begin{aligned}
 IG(Y, V) &= H[1490-, 711+] - \left(\frac{109}{2201} H[52-, 57+] + \frac{2092}{2201} H[1438-, 654+] \right) \\
 &= 0.006411.
 \end{aligned}$$

Deci, într-adevăr, câștigul maxim de informație se obține pentru atributul S .

b. Arborele de decizie de adâncime 1 care are în nodul rădăcină atributul S este cel din figura alăturată. Acuratețea [medie a] acestui arbore de decizie este:

$$\frac{470}{2201} \cdot \frac{344}{470} + \frac{1731}{2201} \cdot \frac{1364}{1731} = \frac{344 + 1364}{2201} = \frac{1708}{2201} = 0.776.$$



c. Se poate constata imediat că arborele ID3 produs pe datele din această problemă va avea 8 noduri-frunză, iar în fiecare dintre aceste noduri-frunză se va asigura câte una dintre mulțimile descrise (pe linie) în coloanele 4 și 5 ale tabelului principal din enunț: $[0-, 5+]$, $[118-, 57+]$, \dots , $[17-, 27+]$, $[105-, 176+]$. Decizia care va fi luată în fiecare nod-frunză este dictată de votul majoritar, și anume: $+$, $-$, \dots , $+$ și respectiv $+$.

Putem calcula acuratețea [medie] astfel:

$$\frac{5 + 118 + 1 + 140 + 35 + 1211 + 27 + 176}{2201} = \frac{1713}{2201} = 0.778.$$

Se observă că se produce (din păcate) o creștere foarte mică în raport cu acuratețea celui mai bun "decision stump": doar 0.002.

7.

(Algoritmul ID3: cazul când există repetiții și inconsistențe în datele de antrenament; o margine superioară pentru eroarea la antrenare în funcție de numărul de valori ale variabilei de ieșire)

CMU, 2002 fall, Andrew Moore, midterm exam, pr. 1.fg

Presupunem că învățăm un arbore de decizie care să prezică atributul de ieșire Z pornind de la atributele de intrare A , B , C . Se folosesc datele de antrenament din tabelul alăturat.

a. Care va fi eroarea la antrenare pe acest set de date? Exprimați răspunsul sub forma fracției de înregistrări care vor fi clasificate eronat ($n/12$).

b. Considerăm un arbore de decizie construit pe un set arbitrar de date. Dacă atributul de ieșire este cu valori discrete și poate lua k valori distincte, care este eroarea de antrenare maximă (exprimată ca fracție)?

A	B	C	Z
0	0	0	0
0	0	1	0
0	0	1	0
0	1	0	0
0	1	1	0
0	1	1	1
1	0	0	0
1	0	1	1
1	1	0	1
1	1	0	1
1	1	1	0
1	1	1	1

Răspuns:

a. Este ușor de observat că datele de antrenament conțin „inconsistențe” (contradicții, relativ la etichetare), și anume la exemplele $(0, 1, 1)$ și $(1, 1, 1)$. Fiecare dintre aceste exemple sunt etichetate o dată cu 0 și altă dată cu 1. Prin urmare, jumătate din aceste exemple vor fi clasificate eronat de arborele învățat de către algoritmul ID3. Eroarea la antrenare va fi deci $\frac{2}{12}$.

b. Vom analiza pe rând mai multe cazuri, care sunt din ce în ce mai generale.

Cazul *i*: Mai întâi vom calcula eroarea la antrenare pentru cazul în care setul de date de antrenament este compus din k instanțe care sunt identice ca tupluri de valori pentru atributele ce le caracterizează, dar sunt clasificate pe rând cu fiecare din cele k valori posibile ale atributului de ieșire. Este evident că arborele de decizie învățat va clasifica eronat $k - 1$ instanțe. Așadar, în acest caz, eroarea la antrenare va fi

$$E = \frac{k - 1}{k}$$

Cazul *ii*: Aceeași eroare [la antrenare] ca mai sus se va înregistra dacă în locul fiecărei instanțe dintre cele considerate la cazul precedent vom avea l instanțe identice, inclusiv în ce privește eticheta. (În total sunt kl instanțe de antrenament.)

$$E = \frac{(k - 1) \cdot l}{k \cdot l} = \frac{k - 1}{k}$$

Cazul *iii*: Dacă relaxăm condiția de mai sus considerând l_1, l_2, \dots, l_k instanțe identice, iar $l = \max_{i=1}^k l_i$, este imediat că eroarea maximă se va atinge în cazul $l_1 = l_2 = \dots = l_k = l$, și va avea aceeași valoare ca mai sus. Așadar, în continuare vom putea renunța la a considera factorul de multiplicare l , fără ca prin aceasta să restrângem generalitatea raționamentului.

Cazul *iv*: Fie n exemple de antrenament (instanțe etichetate) dintre care d sunt distincte (ca tupluri de valori ale atributelor de intrare). Fie k_1, k_2, \dots, k_d numărul de instanțe etichetate pentru fiecare caz distinct în parte din cele d . Atunci vom avea:

$$k_1 \leq k, k_2 \leq k, \dots, k_d \leq k \Rightarrow n = k_1 + k_2 + \dots + k_d \leq k \cdot d \text{ deci } n \leq k \cdot d$$

Eroarea maximă la antrenare va fi dată de formula

$$E = \frac{(k_1 - 1) + (k_2 - 1) + \dots + (k_d - 1)}{n} = \frac{n - d}{n}$$

Avem:

$$E \leq \frac{k - 1}{k} \Leftrightarrow \frac{n - d}{n} \leq \frac{k - 1}{k} \Leftrightarrow k \cdot n - k \cdot d \leq k \cdot n - n \Leftrightarrow k \cdot d \geq n \text{ (adev.)}$$

Prin urmare, eroarea maximă la antrenare care poate fi atinsă atunci când atributul de ieșire poate lua k valori distincte este $\frac{k - 1}{k}$.

8. (ID3, aspecte computaționale: influența atributelor duplicate, respectiv a instanțelor de antrenament duplicate asupra arborelui ID3 rezultat)

CMU, 2009 spring, Ziv Bar-Joseph, final exam, pr. 3.1-2

Se dorește construirea unui arbore de decizie pentru n vectori, cu m atribute.

a. Să presupunem că există i și j astfel încât pentru TOȚI vectorii X din datele de antrenament, aceste atribute au valori egale (adică, $x_i = x_j$ pentru toți vectorii, unde x_i este valoarea atributului i în vectorul X). Să presupunem de asemenea că în cazul în care ambele atribute duc la același câștig de informație vom folosi atributul i . Ștergerea atributului j din datele de antrenament poate schimba arborele de decizie obținut? Explicați pe scurt.

b. Să presupunem că există în mulțimea de antrenament doi vectori egali X și Z (adică, toate atributele lui X și Z sunt exact la fel, inclusiv etichetele). Ștergerea vectorului Z din datele de antrenament poate schimba arborele de decizie obținut? Explicați pe scurt.

Răspuns:

a. Nu, îndepărtarea atributului j nu schimbă arborele de decizie, deoarece atributele i și j conduc la valori egale pentru câștigul de informație în fiecare nod al arborelui.

b. Da, în acest caz arborele de decizie se poate schimba, fiindcă entropia condițională — care se calculează în fiecare nod pentru a determina atributul cu câștigul de informație cel mai mare — depinde de numărul de instanțe de antrenament luate în considerare.

9. (Arbori de decizie: o margine superioară pentru numărul de noduri frunză, în funcție de numărul atributelor și numărul de exemple)

CMU, 2005 fall, T. Mitchell, A. Moore, midterm exam, pr. 2.d

Presupunem că învățăm un arbore de decizie folosind un set de R instanțe de antrenament descrise de M atribute de intrare având valori binare.

Care este numărul maxim posibil de noduri frunză din arborele de decizie, presupunând că fiecărui nod frunză îi este asociat măcar un exemplu de antrenament? Încercuiți unul din răspunsurile de mai jos; justificați alegerea făcută.

$R, \log_2(R), R^2, 2^R, M, \log_2(M), M^2, 2^M,$
 $\min(R, M), \min(R, \log_2(M)), \min(R, M^2), \min(R, 2^M),$
 $\min(\log_2(R), M), \min(\log_2(R), \log_2(M)), \min(\log_2(R), M^2), \min(\log_2(R), 2^M),$
 $\min(R^2, M), \min(R^2, \log_2(M)), \min(R^2, M^2), \min(R^2, 2^M),$
 $\min(2^R, M), \min(2^R, \log_2(M)), \min(2^R, M^2), \min(2^R, 2^M),$
 $\max(R, M), \max(R, \log_2(M)), \max(R, M^2), \max(R, 2^M),$
 $\max(\log_2(R), M), \max(\log_2(R), \log_2(M)), \max(\log_2(R), M^2), \max(\log_2(R), 2^M),$
 $\max(R^2, M), \max(R^2, \log_2(M)), \max(R^2, M^2), \max(R^2, 2^M),$
 $\max(2^R, M), \max(2^R, \log_2(M)), \max(2^R, M^2), \max(2^R, 2^M).$

Răspuns:

Notăm cu \max_{frunze} valoarea căutată. Trebuie luate în considerare două aspecte:

- (a) fiecare nod frunză trebuie să clasifice măcar un exemplu de antrenament
 \Rightarrow nu putem să avem mai multe frunze decât exemple de antrenament
 $\Rightarrow \max_{frunze} \leq R$
- (b) fiecare atribut poate fi testat o singură dată pe un drum de la rădăcină la o frunză oarecare \Rightarrow arborele obținut va avea adâncimea cel mult M
 $\Rightarrow \max_{frunze} \leq 2^M$

$$\left. \begin{array}{l} \max_{frunze} \leq R \\ \max_{frunze} \leq 2^M \end{array} \right\} \Rightarrow \max_{frunze} \leq \min(R, 2^M)$$

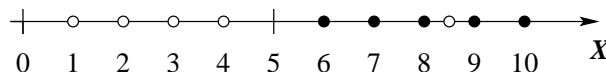
Această valoare poate fi atinsă: putem să luăm, de exemplu, cazul trivial când avem o singură instanță de antrenament. Prin urmare, avem $\max_{frunze} = \min(R, 2^M)$.

10. (Extensii ale algoritmului ID3: variabile de intrare continue; “decision stumps”; eroarea la antrenare, eroarea la CVLOO; overfitting)

■ • CMU, 2002 fall, Andrew Moore, midterm exam, pr. 3

Fie setul de date de mai jos. X este atribut de intrare și ia valori reale, iar Y este variabilă de ieșire cu valori booleene. (*Observație:* Am reprezentat acest set de date sub formă grafică, marcând valoarea / eticheta $Y = 1$ prin bulină neagră, iar valoarea / eticheta $Y = 0$ prin cerculeț alb.) Pe acest set de date se folosește algoritmul ID3 pentru învățare de arbori de decizie.

X	Y
1	0
2	0
3	0
4	0
6	1
7	1
8	1
8.5	0
9	1
10	1



Algoritmul ID3 (extins) va trebui să decidă cum divide (engl., split) intervale de valori asociate variabilei reale X . Separarea în intervale diferite va fi stabilită cu ajutorul unor *praguri* (engl., split thresholds), determinate în felul următor:

- Mai întâi se ordonează crescător acele valori ale variabilei X care apar în datele de antrenament.

- Se stabilesc apoi perechi de valori consecutive pentru care există instanțe de antrenament care sunt etichetate în mod diferit pentru o valoare, comparativ cu cealaltă valoare. Pentru fiecare pereche de valori de acest fel, va fi plasat un prag de separare la jumătatea distanței dintre cele două valori.

Inițial, se alege pragul de separare care conduce la un câștig de informație maxim. Apoi, la fiecare nouă execuție a buclei principale a algoritmului ID3 — vă readucem aminte că acest algoritm este recursiv — se va selecta câte un alt prag dintre cele rămase disponibile, aplicând același criteriu: maximizarea câștigului de informație.

De *exemplu*, pentru $X = 4$ avem o instanță de antrenament negativă, iar pentru $X = 6$ avem o instanță de antrenament pozitivă. Se poate arăta că algoritmul ID3 va decide să spliteze mai întâi la valoarea $X = 5$ (care reprezintă jumătatea distanței dintre $X = 4$ și $X = 6$) și apoi la valoarea $X = 8.25$ (care reprezintă jumătatea distanței dintre $X = 8$ și $X = 8.5$).

Fie DT^* arborele de decizie complet, obținut de algoritmul ID3 fără a face pruning, iar $DT2$ arborele de decizie produs de ID3 urmat de pruning, care are doar două noduri frunză (deci $DT2$ face o singură divizare de interval).

- Care este *eroarea la antrenare* produsă de $DT2$ respectiv DT^* (exprimată ca număr de exemple clasificate eronat din totalul de 10 exemple)?
- Care este *eroarea* produsă de $DT2$ respectiv DT^* la *cross-validare* folosind metoda *Leave-One-Out* (CVLOO)?

Răspuns:

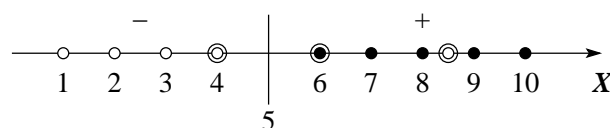
- Deoarece $DT2$ reține doar un singur split, și anume la $X = 5$, regula de decizie pe care o reprezintă este:

```
IF      X ≤ 5
THEN   Y = 0
ELSE   Y = 1
```

Este evident că această regulă produce o clasificare eronată doar pentru una dintre instanțele de antrenament, și anume $X = 8.5$. Avem deci $E_{antren.}(DT2) = 1/10$.

Întrucât exemplele nu conțin inconsistente, arborele de decizie ID3 clasifică corect toate datele de antrenament. Avem deci $E_{antren.}(DT^*) = 0/10 = 0$.

- Figura de mai jos reprezintă împărțirea axei reale în *intervale* / *zone de decizie* conform arborelui (și *pragului de decizie*) învățat de către algoritmul $DT2$, folosind întregul set de exemple date. (Am încercuit acele puncte care, după cum se va vedea mai jos, vor constitui cazuri aparte la calcularea erorii de tip CVLOO cu algoritmul $DT2$.)



În ce privește cross-validarea prin metoda “Leave-One-Out”, se poate demonstra — calculele nu sunt arătate aici³⁷⁹ — următorul fapt: pentru fiecare din cele 10 exemple ($X = 1, X = 2, \dots, X = 10$) considerate pe rând, pragul de decizie identificat de algoritmul DT2 va fi $X = 5$, cu excepția următoarelor două cazuri:

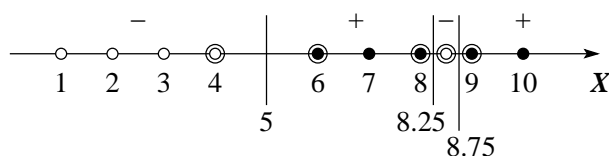
- $X = 4$: în acest caz, splitarea se va face la mijlocul intervalului $[3, 6]$, deci la 4.5. Cum $4 \leq 4.5$, rezultă că exemplul $X = 4$ va fi clasificat corect;
- $X = 6$: splitarea se va face la mijlocul intervalului $[4, 7]$, deci la 5.5. Cum $6 > 5.5$, rezultă că exemplul $X = 6$ va fi clasificat corect.

Atunci când punctul $X = 8.5$ este lăsat deoparte, pragul de decizie selectat fiind $X = 5$, va rezulta că punctul $X = 8.5$ este clasificat pozitiv (deci eronat), întrucât $8.5 > 5$.

Este imediat că pentru restul de 7 cazuri ($X = 1, 2, 3, 7, 8, 9, 10$), arborele determinat de algoritmul DT2 va clasifica corect punctul X .

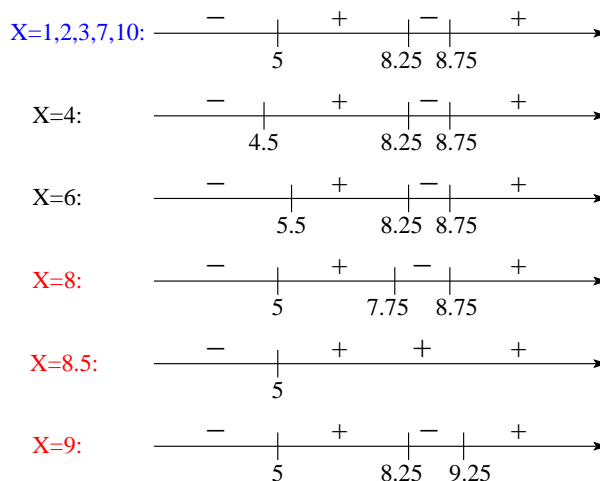
Așadar, pentru DT2 rezultă $E_{CVLOO} = 1/10$.

Figura de mai jos reprezintă împărțirea axei reale în *intervale* / *zone de decizie* conform *pragurilor de decizie* (și arborelui DT*) determinate de către algoritmul ID3 pe întregul set de date de antrenament.³⁸⁰



În această figură am încercuit exemplele care ar putea conduce la eroare la cross-validarea de tip “Leave-One-Out”:

- $X = 4$: corect clasificat, explicația este aceeași ca în cazul DT2;
- $X = 6$: idem;
- $X = 8$: split-ul trebuie făcut la mijlocul intervalului $[7, 8.5]$, adică 7.75. Cum $8 > 7.75$, rezultă că punctul $X = 8$ va fi clasificat negativ, ceea ce este eronat;
- $X = 8.5$: nu este nevoie decât de un singur split, arborele DT* învățat în acest caz fiind identic cu cel construit de algoritmul DT2. Cum $8 > 5$, rezultă că punctul $X = 8.5$ va fi clasificat pozitiv, deci eronat;



³⁷⁹ Aceste calcule nu sunt dificile, mai ales dacă se utilizează raționamente de tip „calitativ”, așa cum am arătat la problema 3.

³⁸⁰ *Observație*: Atunci când ne raportăm doar la zonele de decizie în ansamblu, ordinea în care s-au stabilit testele în arborele ID3 este irelevantă! Acest fapt este valabil și mai jos, unde discutăm despre eroarea la CVLOO cu algoritmul DT*.

• $X = 9$: intervalul de split devine $[8.5, 10]$, split-ul făcându-se la 9.25. Cum $9 \leq 9.25$, rezultă că punctul $X = 9$ va fi clasificat negativ, ceea ce este eronat.

Așadar, pentru DT^* avem $E_{CVLOO} = 3/10$.

În *concluzie*, se observă că $E_{antren.}(DT2) = 1/10 > 0 = E_{antren.}(DT^*)$, în vreme ce $E_{CVLOO}(DT2) = 1/10 < 3/10 = E_{CVLOO}(DT^*)$. Acesta este un caz tipic de manifestare a fenomenului de *overfitting* (*supra-specializare*).

11. (Arbori de decizie cu variabile de intrare continue;
o margine superioară pentru adâncimea arborilor
în cazul (ne)separabilității liniare în \mathbb{R}^2)

CMU, 2009 spring, Ziv Bar-Joseph, midterm exam, pr. 5.cd

Se consideră n vectori bidimensionali ($x = \{x_1, x_2\}$) care pot fi clasificați folosind o funcție [de regresie] liniară, adică există $w \in \mathbb{R}^2$ și $b \in \mathbb{R}$ astfel încât:

$$y = \begin{cases} +1 & \text{if } w \cdot x + b > 0 \\ -1 & \text{if } w \cdot x + b \leq 0 \end{cases}$$

a. Poate un arbore binar de decizie — atenție!, nu neapărat arborele ID3 — să clasifice corect acești vectori? Dacă nu, justificați. Dacă da, determinați adâncimea maximă (adică numărul maxim de niveluri de test) ale unui astfel de arbore de decizie, care este optim (ca număr de niveluri de test).

b. Acum să presupunem că aceste n date nu sunt separabile liniar (adică nu există $w \in \mathbb{R}^2$ și $b \in \mathbb{R}$ cu proprietatea de mai sus). Poate un arbore de decizie (binar) să clasifice corect acești vectori? Dacă nu, justificați. Dacă da, care este adâncimea maximă a unui arbore de decizie corespunzător, optim ca număr de niveluri?

Răspuns:

a. Da. O strategie posibilă pentru a construi un arbore de decizie binar care să clasifice corect aceste puncte este următoarea:

Mai întâi vom construi un arbore de decizie considerând doar atributul x_1 . În particular, făcând „înjumătățiri” succesive ale mulțimii de valori ale acestui atribut, arborele rezultat va avea o adâncime maximă de $\lceil \log_2 n \rceil$ (adică partea întreagă superioară din $\log_2 n$) niveluri. În fiecare nod frunză al acestui arbore se va găsi o mulțime de puncte, însă nu neapărat cu aceeași clasificare. Totuși, deoarece datele sunt liniar separabile, pentru fiecare dintre aceste mulțimi de puncte se poate găsi o valoare a atributului x_2 care să o împartă în două submulțimi clasificate corect.

Prin urmare, arborele construit inițial considerând doar valoarea x_1 are nevoie să i se adauge doar cel mult câte un nod în fiecare frunză, nod care să clasifice corect datele luând în considerare valoarea x_2 . Adâncimea totală a arborelui astfel construit este $1 + \lceil \log_2 n \rceil$, adică de ordinul $O(\log n)$.

b. Da. Similar punctului anterior, se construiește un arbore de decizie considerând doar atributul x_1 , ceea ce înseamnă o adâncime de $\lceil \log_2 n \rceil$. În fiecare

nod frunză al acestui arbore se va găsi o mulțime de puncte, posibil cu clasificări diferite. Datele de antrenament nu mai sunt liniar separabile, deci nu pot fi clasificate corect printr-un singur nod.

Pentru fiecare astfel de nod frunză în care punctele nu au aceeași clasificare, se aplică algoritmul de determinare a arborelui de decizie, de această dată luând în considerare doar valoarea lui x_2 . Aceasta presupune din nou o adâncime maximă a subarborelui de $\lceil \log_2 n \rceil$. În total, arborele obținut are adâncimea maximă $\lceil \log_2 n \rceil + \lceil \log_2 n \rceil$, deci tot de ordinul $O(\log n)$.

12. (Algoritmul ID3 cu attribute discrete, respectiv attribute discrete și un atribut continuu; aplicare; predicție)

■ • CMU, 2012 fall, E. Xing, A. Singh, HW1, pr. 1.1

Până în luna septembrie a anului 2012, 800 de planete extrasolare (numite în continuare exoplanete) au fost identificate în galaxia noastră. Niște nave spațiale super-secrete au fost trimise pentru a survola toate aceste exoplanete, cu scopul de a stabili dacă ele sunt locuibile sau nu de către oameni. Evident, a trimite câte o navetă spațială la fiecare dintre aceste exoplanete este extrem de costisitor. De aceea, în această problemă vă propunem să elaborați un arbore de decizie pentru a prezice dacă o exoplanetă este locuibilă sau nu, folosind doar trăsături / caracteristici (engl., features) observabile cu ajutorul telescoapelor terestre.

a. În tabelul de mai jos vi se dau anumite date în legătură cu toate cele 800 de planete survolate până acum. Trăsăturile observate cu ajutorul telescoapelor sunt *Size* ("Big" sau "Small") și *Orbit* ("Near" sau "Far").

Fiecare linie din tabel indică valori ale acestor două trăsături, caracterul habitabil ("Yes" sau "No"), precum și de câte ori a fost identificată fiecare combinație de valori [pentru cele trei trăsături]. De exemplu, au fost identificate 20 de planete mari ("Big"), care sunt situate pe orbite apropiate ("Near") de soarele / steaua lor și sunt locuibile.

Size	Orbit	Habitable	Count
Big	Near	Yes	20
Big	Far	Yes	170
Small	Near	Yes	139
Small	Far	Yes	45
Big	Near	No	130
Big	Far	No	30
Small	Near	No	11
Small	Far	No	255

Elaborați și desenați arborele de decizie învățat de către algoritmul ID3 pe aceste date. (Folosiți criteriul câștigului de informație; nu aplicați pruning-ul.) La fiecare nod din arbore, scrieți numărul de planete locuibile și respectiv nelocuibile din datele de antrenament care sunt asociate cu nodul respectiv.

b. Pentru doar 9 dintre aceste exoplanete, a fost măsurată o a treia trăsătură, *Temperature* (exprimată în grade Kelvin), după cum se arată în tabelul de mai jos.

Refaceți toți pașii de la punctul *a*, de data aceasta folosind toate cele trei trăsături de intrare. Pentru trăsătura *Temperature* (văzută ca atribut numeric cu valori continue), la fiecare iterație va trebui să faceți maximizarea în raport cu toate pragurile adecvate pentru separare binară (engl., binary thresholding splits). Iată un *exemplu* de test pentru o astfel de separare binară: $T \leq 250$ vs. $T > 250$.

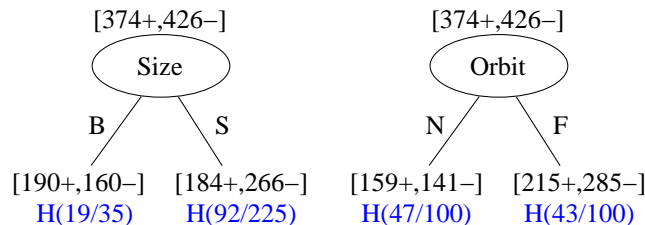
Size	Orbit	Temperature	Habitable
Big	Far	205	No
Big	Near	205	No
Big	Near	260	Yes
Big	Near	380	Yes
Small	Far	205	No
Small	Far	260	Yes
Small	Near	260	Yes
Small	Near	380	No
Small	Near	380	No

c. Conform arborelui de decizie pe care l-ați obținut la punctul *b*, cum va fi clasificată o planetă având trăsăturile (Big, Near, 280), locuibilă sau nelocuibilă?

Indicație: Este posibil să aveți nevoie de următoarele valori pentru entropia ($H(p)$) unei variabile aleatoare Bernoulli de parametru p : $H(1/3) = 0.9182$, $H(2/5) = 0.9709$, $H(92/225) = 0.9759$, $H(43/100) = 0.9858$, $H(16/35) = 0.9946$, $H(47/100) = 0.9974$.

Răspuns:

a. „Compașii de decizie“ corespunzători nodului rădăcină sunt înfățișați în desenul următor:



Sub fiecare nod descendent am notat entropia nodului respectiv, făcând referire la distribuția Bernoulli și dând (de fiecare dată) parametrului acestei distribuții valoarea corespunzătoare. Pentru nodul descendent corespunzător lui Orbit = Near am ținut cont și de faptul că entropia distribuției Bernoulli, ca funcție de parametru p , este simetrică față de valoarea $1/2$.³⁸¹

Entropiile condiționale medii corespunzătoare acestor doi „compași de decizie“ sunt:

$$\begin{aligned}
 H(\text{Habitable}|\text{Size}) &= \frac{35}{80} \cdot H\left(\frac{19}{35}\right) + \frac{45}{80} \cdot H\left(\frac{92}{225}\right) \\
 &= \frac{35}{80} \cdot 0.9946 + \frac{45}{80} \cdot 0.9759 = 0.9841
 \end{aligned}$$

³⁸¹ Este util să revedeți explicațiile date în *observația importantă* de la pagina 420.

$$\begin{aligned}
 H(Habitable|Orbit) &= \frac{3}{8} \cdot H\left(\frac{47}{100}\right) + \frac{5}{8} \cdot H\left(\frac{43}{100}\right) \\
 &= \frac{3}{8} \cdot 0.9974 + \frac{5}{8} \cdot 0.9858 = 0.9901
 \end{aligned}$$

Suntem acum în măsură să desemnăm atributul care va fi plasat în nodul rădăcină al arborelui care va fi construit de algoritmul ID3 pe datele din enunț: este atributul *Size*, întrucât $H(Habitable|Size) < H(Habitable|Orbit)$.

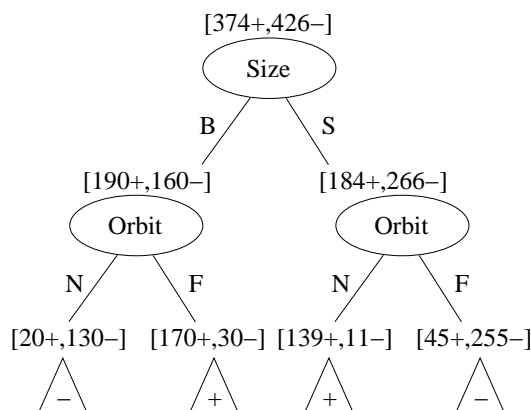
Doar cu titlu de informare, precizăm și câștigurile de informație realizate de cele două atribute de intrare în raport cu atributul de ieșire:

$$IG(Habitable; Size) = H(Habitable) - H(Habitable|Size) = 0.9969 - 0.9841 = 0.0128$$

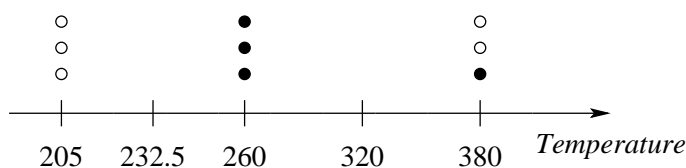
și, similar,

$$IG(Habitable; Orbit) = 0.0067.$$

Întrucât nu avem decât două atribute de intrare, arborele de decizie va fi:

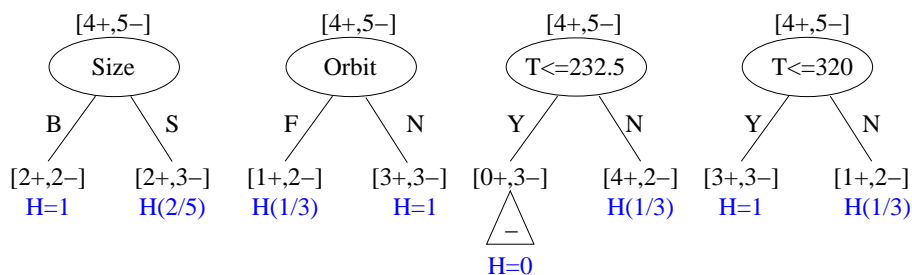


b. Pragurile de separare corespunzătoare atributului *Temperature* sunt determinate conform imaginii următoare:



Nivelul 1 (rădăcina):

„Compașii de decizie“ corespunzători acestui nivel sunt:



Observând cu atenție partițiile formate,³⁸² vom constata că entropia condițională medie pentru testul $Temperature \leq 232.5$ (în raport cu atributul de ieșire $Habitable$) are o valoare mai mică decât fiecare dintre entropiile condiționale medii $H(Habitable|Orbit)$ și $H(Habitable|Temperature \leq 320)$ (care, de fapt, sunt egale între ele).³⁸³

Relația dintre $H(Habitable|Temperature \leq 232.5)$ și $H(Habitable|Size)$ se determină cu ajutorul calculelor:

$$H(Habitable|Size) = \frac{4}{9} + \frac{5}{9} \cdot H\left(\frac{2}{5}\right) = \frac{4}{9} + \frac{5}{9} \cdot 0.9709 = 0.9838$$

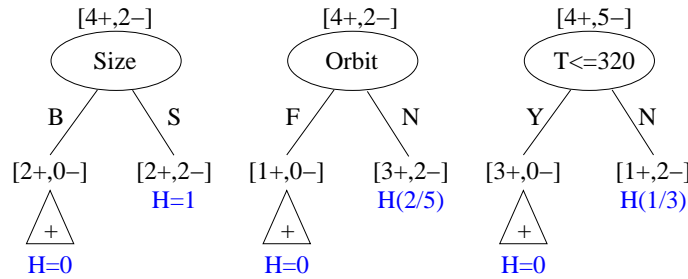
$$H(Habitable|Temp \leq 232.5) = \frac{2}{3} \cdot H\left(\frac{1}{3}\right) = \frac{2}{3} \cdot 0.9182 = 0.6121.$$

Deși nu mai este necesar, indicăm și valorile câștigurilor de informație:

$$\begin{aligned} IG(Habitable; Size) &= 0.0072 \\ IG(Habitable; Orbit) &= 0.0183 \\ IG(Habitable; Temp \leq 232.5) &= 0.3788 \\ IG(Habitable; Temp \leq 320) &= 0.0183. \end{aligned}$$

Prin urmare, vom reține pentru nivelul 0 testul $Temperature \leq 232.5$.

Nivelul 2 (mai exact, aici ne limităm la datele cu $Temperature > 232.5$: „Compașii de decizie” corespunzători [completării] acestui nivel sunt:



Este imediat că, pe aceste date, $H(Habitable|Temp \leq 320) < H(Habitable|Size)$ și, de asemenea, $H(Habitable|Temp \leq 320) < H(Habitable|Orbit)$.³⁸⁴ Prin urmare, aici va fi ales testul $Temp \leq 320$.

³⁸²Se compară două câte două entropiile condiționale specifice, precum și ponderile datelor respective în [raport cu] numărul total de instanțe asociate nodului părinte.

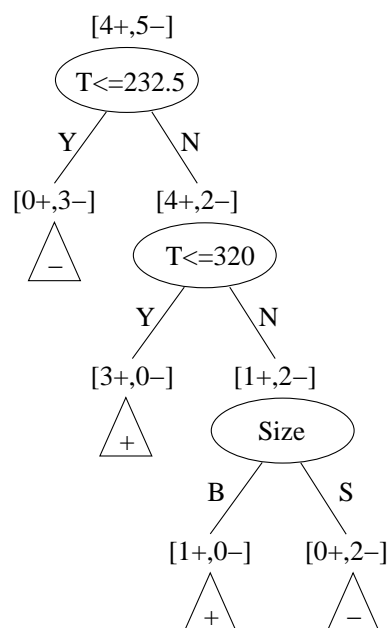
³⁸³Observați de exemplu că $H(Habitable|Orbit = F) > H(Habitable|Temperature \leq 232.5)$ și $H(Habitable|Orbit = N) > H(Habitable|Temperature > 232.5)$, iar ponderile corespunzătoare acestor entropii condiționale specifice în calculul entropiilor condiționale medii $H(Habitable|Orbit = N)$ și $Temperature \leq 232.5$ sunt egale două câte două (și anume, cu 3/9 și respectiv 6/9).

³⁸⁴Mai exact, ar fi trebuit să scriem: $H(Habitable|Temp > 232.5, Temp \leq 320) < H(Habitable|Temp > 232.5, Size)$ și respectiv $H(Habitable|Temp > 232.5, Temp \leq 320) < H(Habitable|Temp > 232.5, Orbit)$.

Nivelul 3 (mai exact, aici ne limităm la datele cu $Temperature > 320$):

Se poate observa că, pentru aceste date, atributul *Size* are putere discriminativă maximă. Așadar, arborele de decizie final va fi cel din figura alăturată.

c. Conform arborelui de decizie obținut la punctul anterior, o exoplanetă având trăsăturile (Big, Near, 280) va fi clasificată ca fiind locuibilă.



13.

(Extensii ale algoritmului ID3:
cazul atributelor discrete cu număr mare de valori)

CMU, 2008(?) spring, HW2, pr. 1

Se dorește antrenarea unui arbore de decizie care să clasifice exemple cu două attribute de intrare X_1, X_2 și un atribut de ieșire Y care are valorile 1 și 2. Primul atribut, X_1 , este binar, pe când al doilea atribut, X_2 , are 6 valori posibile A, B, C, D, E, F . Se dau următoarele 12 exemple de antrenament, câte 6 din fiecare clasă:

$Y = 1$	(1, A)	(0, E)	(1, B)	(1, B)	(1, F)	(0, D)
$Y = 2$	(0, A)	(0, C)	(1, E)	(0, F)	(0, B)	(1, D)

Vă reamintim formula câștigului de informație la partajarea mulțimii de exemple S în funcție de valorile atributului A :

$$Gain(S, A) = H(S) - H(S | A), \text{ cu } H(S | A) = \sum_{v \in \text{valori}(A)} \frac{|S_v|}{|S|} \cdot H(S_v),$$

unde $H(S)$ este entropia mulțimii de exemple S , iar $H(S | A)$ este entropia condițională medie a mulțimii S în raport cu atributul A , calculată așa cum se vede mai sus, ca sumă ponderată a entropiilor submulțimilor lui S determinate de valorile atributului A .

a. Determinați atributul ales în rădăcină folosind câștigul de informație. Folosiți aproximarea $\log_2 3 = 1.585$.

b. Determinați atributul ales în nodul rădăcină folosind o măsură numită *gain ratio impurity*, adică alegeți acel atribut care maximizează raportul

$$\frac{Gain(S, A)}{-\sum_v P(A = v) \cdot \log_2 P(A = v)} = \frac{H(S) - H(S | A)}{H(A)}.$$

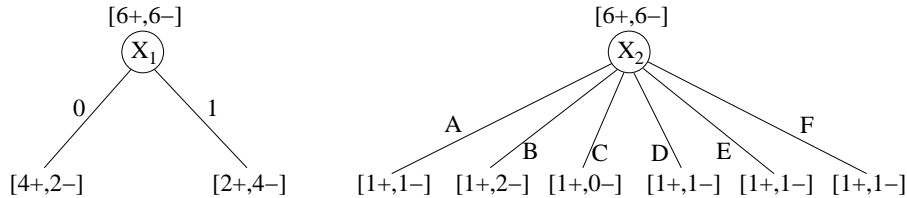
c. Având în vedere rezultatele de la punctele precedente, analizați utilitatea folosirii măsurii *gain ratio impurity* în cazurile în care atributele au numere diferite de valori posibile.

Răspuns:

a. Calculăm în primul rând entropia nodului rădăcină (sau a atributului de ieșire), care are 6 exemple pozitive și 6 negative:

$$H(Y) = H[6+, 6-] = -\frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12} = 2 \cdot \frac{1}{2} \log_2 2 = 1$$

În rădăcină poate fi ales fie atributul X_1 , fie atributul X_2 , obținând următoarele împărțiri:



Dacă vom pune atributul X_1 în nodul rădăcină, câștigul de informație va fi:

$$\begin{aligned} \text{Gain}(S, X_1) &= H(Y) - \frac{6}{12}H[4+, 2-] - \frac{6}{12}H[2+, 4-] = \\ &= 1 - \frac{6}{12} \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) - \frac{6}{12} \left(-\frac{2}{6} \log_2 \frac{2}{6} - \frac{4}{6} \log_2 \frac{4}{6} \right) \\ &= 1 - 2 \cdot \frac{1}{2} \left(\frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 \frac{3}{2} \right) = 1 - \left(\log_2 3 - \frac{2}{3} \right) \\ &= \frac{5}{3} - \log_2 3 \approx 0.0817 \end{aligned}$$

Altminteri, punând atributul X_2 în nodul rădăcină, câștigul de informație este:

$$\begin{aligned} \text{Gain}(S, X_2) &= H(Y) - 4 \cdot \frac{2}{12}H[1+, 1-] - \frac{3}{12}H[1+, 2-] - \frac{1}{12}H[1+, 0-] \\ &= 1 - \frac{2}{3} \cdot 1 - \frac{3}{12} \left(-\frac{1}{3} \log_2 \frac{1}{3} - \frac{2}{3} \log_2 \frac{2}{3} \right) - \frac{1}{12} \cdot 0 \\ &= 1 - \frac{2}{3} - \frac{1}{4} \left(\frac{1}{3} \log_2 3 + \frac{2}{3} \log_2 \frac{3}{2} \right) = 1 - \frac{2}{3} - \frac{1}{4} \left(\log_2 3 - \frac{2}{3} \right) \\ &= \frac{1}{2} - \frac{1}{4} \log_2 3 \approx 0.1037 \end{aligned}$$

Folosind drept criteriu de optimizat în fiecare nod câștigul de informație, în rădăcină vom pune atributul X_2 , întrucât $\text{Gain}(S, X_1) < \text{Gain}(S, X_2)$.

b. Dacă vom pune atributul X_1 în nodul rădăcină, *gain ratio impurity* va fi:

$$\begin{aligned} \frac{\text{Gain}(S, X_1)}{-\sum_{i \in \{0,1\}} P(X_1 = i) \cdot \log_2 P(X_1 = i)} &= \frac{\text{Gain}(S, X_1)}{-\frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12}} = \frac{\text{Gain}(S, X_1)}{1} \\ &\approx 0.0817 \end{aligned}$$

În schimb, plasând atributul X_2 în nodul rădăcină, *gain ratio impurity* va fi:

$$\begin{aligned} \frac{Gain(S, X_2)}{-\sum_{j \in \{A, \dots, F\}} P(X_2 = j) \cdot \log_2 P(X_2 = j)} &= \frac{Gain(S, X_2)}{-4 \cdot \frac{2}{12} \log_2 \frac{2}{12} - \frac{3}{12} \log_2 \frac{3}{12} - \frac{1}{12} \log_2 \frac{1}{12}} \\ &= \frac{Gain(S, X_2)}{\frac{2}{3} \log_2 6 + \frac{1}{4} \log_2 4 + \frac{1}{12} \log_2 12} = \frac{Gain(S, X_2)}{\frac{4}{3} + \frac{3}{4} \log_2 3} \approx 0.0411 \end{aligned}$$

Prin urmare, în nodul rădăcină vom pune atributul X_1 , pentru care măsura *gain ratio impurity* are valoarea cea mai mare.

c. Câștigul de informație favorizează alegerea atributelor care au un număr mare de valori, indiferent dacă ele determină sau nu partiționarea în mod semnificativ a datelor de antrenament. În schimb, măsura *gain ratio impurity* ia în considerare, prin cantitatea de la numitor (vedeți definiția), numărul de valori ale atributului respectiv, mai exact mărimea mulțimilor de instanțe asiguate nodurilor-fii, care au fost generate ca urmare a alegerii respectivului atribut. Valoarea de la numitor va crește odată cu numărul de noduri-fii, și totodată cu numărul de noduri-fii care au asiguate puține exemple. Prin urmare, această măsură penalizează attributele cu multe valori, evitând favorizarea de care se face vinovat câștigul de informație într-o atare situație.

14.

(Extensii ale algoritmului ID3: cazul când se ia în considerare costul calculării atributelor; attribute continue; “decision stumps”)

CMU, 2008 spring, T. Mitchell, W. Cohen, HW1, pr. 1

Se dă următorul set de date. Acesta reprezintă fișele a 12 pacienți ipotetici, ținând cont de sex, vârstă (peste 60 ani sau nu), dacă suferă sau nu de diabet, dacă au pulsul mărit (sau nu) și EKG-ul anormal (sau nu). Pacienții sunt clasificați în final după cum prezintă (sau nu) aritmie.

Pacient	Sex	Peste60	Diabetic	Puls	EKG	AreAritmie
1	M	1	1	0	0	0
2	M	0	0	1	1	1
3	M	0	1	1	0	0
4	M	1	0	0	1	1
5	M	1	1	1	0	1
6	M	0	1	1	0	1
7	F	0	0	1	0	0
8	F	1	1	1	1	1
9	F	0	1	0	1	1
10	F	1	0	0	0	0
11	F	1	1	0	0	0
12	F	1	0	1	1	1

a. Calculați entropia condițională specifică $H(\text{AreAritmie} \mid \text{Sex} = F)$.

b. Dacă pentru selectarea atributelor se folosește măsura $\frac{Gain^2(S, A)}{Cost(A)}$ în schimbul câștigului informational, care va fi atributul pus în nodul rădăcină? Se

consideră că $Cost(\text{Sex}) = Cost(\text{Peste60}) = 1$, $Cost(\text{Diabetic}) = 3$, $Cost(\text{Puls}) = 2$ și $Cost(\text{EKG}) = 5$.

Observație: În calcule se vor utiliza următoarele aproximări: $\log_2 3 = 1.585$, $\log_2 5 = 2.322$ și $\log_2 7 = 2.807$.

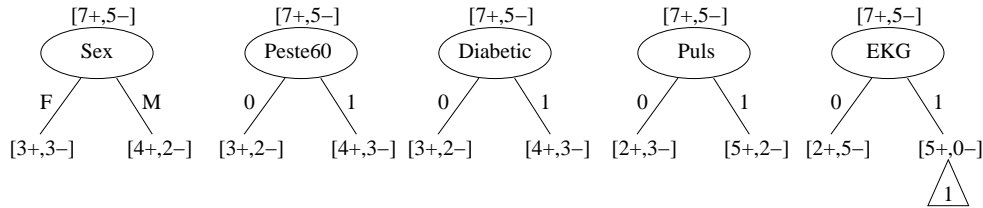
c. Să presupunem că, pentru un alt set de pacienți, se cunoaște vârsta lor exactă. Pentru exemplele pozitive, vârstele sunt: $\{40, 60, 62, 64, 70, 74, 75, 82\}$, iar pentru exemplele negative sunt: $\{33, 35, 42, 45, 49, 52, 58, 59, 80\}$. Să presupunem că toate celelalte atribute sunt predictorii „slabi”, prin urmare dorim ca arborele să aibă un singur nod, rădăcina, care să împartă exemplele cu valori continue ale atributului vârstă în două: $vârsta < k$ și $vârsta \geq k$. Care va fi valoarea aleasă pentru k , bazat pe câștigul de informație?

Răspuns:

a. Entropia condițională cerută este:

$$H(\text{AreAritmie} \mid \text{Sex} = F) = H[3+, 3-] = 1$$

b. În rădăcina arborelui de decizie se alege atributul A pentru care raportul $\frac{Gain^2(S, A)}{Cost(A)}$ este maxim. În cazul nostru, variantele pe care le avem sunt:



Se observă direct că $Gain(S, \text{EKG})$ este mai mare decât $Gain(S, A)$ pentru orice atribut $A \neq \text{EKG}$. Însă și $Cost(\text{EKG})$ este mai mare decât $Cost(A)$ pentru orice $A \neq \text{EKG}$. Așadar, trebuie să facem calculele în detaliu.

Entropia atributului de ieșire, AreAritmie, este:

$$\begin{aligned} H(\text{AreAritmie}) &= H[7+, 5-] = \frac{7}{12} \cdot \log_2 \frac{12}{7} + \frac{5}{12} \cdot \log_2 \frac{12}{5} \\ &= \frac{7}{12} \cdot \log_2 12 - \frac{7}{12} \cdot \log_2 7 + \frac{5}{12} \cdot \log_2 12 - \frac{5}{12} \cdot \log_2 5 \\ &= \log_2(3 \cdot 4) - \frac{7}{12} \cdot \log_2 7 - \frac{5}{12} \cdot \log_2 5 \\ &= \log_2 3 + \underbrace{\log_2 4}_{=2} - \frac{7}{12} \cdot \log_2 7 - \frac{5}{12} \cdot \log_2 5 \approx 0.98 \end{aligned}$$

Vom calcula câștigul de informație pentru fiecare din cele 5 atribute — se observă că pentru atributele Peste60 și Diabetic, câștigurile de informație sunt egale —, și apoi vom face raportul $\frac{Gain^2}{Cost}$:

Pentru atributul Sex:

$$\begin{aligned} Gain(S, \text{Sex}) &= H(\text{AreAritmie}) - \frac{6}{12} H[3+, 3-] - \frac{6}{12} H[4+, 2-] \\ &= H(\text{AreAritmie}) - \frac{1}{2} \cdot 1 - \frac{1}{2} \left(-\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} \right) \end{aligned}$$

$$= 0.98 - \frac{1}{2} - \frac{1}{2} \left(\log_2 3 - \frac{2}{3} \right) = 0.98 - \frac{1}{2} - \frac{1}{2} \log_2 3 + \frac{1}{3} \approx 0.02$$

Prin urmare, $\frac{Gain^2(S, Sex)}{Cost(Sex)} \approx \frac{0.02^2}{1} = 0.0004$.

Pentru attributele Peste60 și Diabetic:

$$Gain(S, Peste60) = Gain(S, Diabetic)$$

$$\begin{aligned} &= H(\text{AreAritmie}) - \frac{5}{12}H[3+, 2-] - \frac{7}{12}H[4+, 3-] \\ &= 0.98 - \frac{5}{12} \left(\frac{3}{5} \log_2 \frac{5}{3} + \frac{2}{5} \log_2 \frac{5}{2} \right) - \frac{7}{12} \left(\frac{4}{7} \log_2 \frac{7}{4} + \frac{3}{7} \log_2 \frac{7}{3} \right) \\ &= 0.98 - \frac{5}{12} \left(\log_2 5 - \frac{3}{5} \log_2 3 - \frac{2}{5} \cdot 1 \right) - \frac{7}{12} \left(\log_2 7 - \frac{4}{7} \cdot 2 - \frac{3}{7} \log_2 3 \right) \\ &= 0.98 - \frac{5}{12} \log_2 5 - \frac{7}{12} \log_2 7 + \frac{1}{2} \log_2 3 + \frac{5}{6} \approx 0.0009 \end{aligned}$$

Deci $\frac{Gain^2(S, Peste60)}{Cost(Peste60)} \approx \frac{0.0009^2}{1} = 81 \cdot 10^{-8}$ și $\frac{Gain^2(S, Diabetic)}{Cost(Diabetic)} \approx \frac{0.0009^2}{3} = 27 \cdot 10^{-8}$, ambele fiind niște valori foarte mici.

Pentru atributul Puls:

$$\begin{aligned} Gain(S, Puls) &= H(\text{AreAritmie}) - \frac{5}{12}H[2+, 3-] - \frac{7}{12}H[5+, 2-] \\ &= 0.98 - \frac{5}{12} \left(\frac{2}{5} \log_2 \frac{5}{2} + \frac{3}{5} \log_2 \frac{5}{3} \right) - \frac{7}{12} \left(\frac{5}{7} \log_2 \frac{7}{5} + \frac{2}{7} \log_2 \frac{7}{2} \right) \\ &= 0.98 - \frac{7}{12} \log_2 7 + \frac{1}{4} \log_2 3 + \frac{1}{3} \approx 0.072 \end{aligned}$$

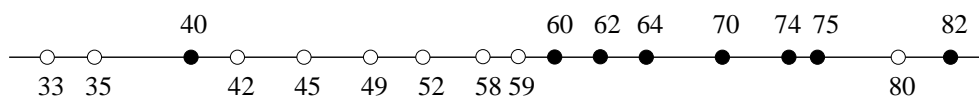
Prin urmare, $\frac{Gain^2(S, Puls)}{Cost(Puls)} \approx \frac{0.072^2}{2} = 0.002592$.

Pentru atributul EKG:

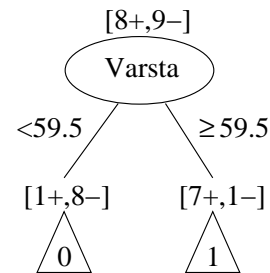
$$\begin{aligned} Gain(S, EKG) &= H(\text{AreAritmie}) - \frac{7}{12}H[2+, 5-] - \frac{5}{12}H[5+, 0-] \\ &= H(\text{AreAritmie}) - \frac{7}{12} \left(\frac{2}{7} \log_2 \frac{7}{2} + \frac{5}{7} \log_2 \frac{7}{5} \right) - \frac{5}{12} \cdot 0 \\ &= 0.98 + \frac{5}{12} \log_2 5 - \frac{7}{12} \log_2 7 + \frac{1}{6} \approx 0.476 \end{aligned}$$

Deci $\frac{Gain^2(S, EKG)}{Cost(EKG)} \approx \frac{0.476^2}{5} = 0.0453152$. Este evident că aceasta este cea mai mare valoare, de aceea în nodul rădăcină va fi ales atributul *EKG*, deși costul acestuia este cel mai mare.

c. Putem reprezenta exemplele astfel:



Se observă că alegerea cea mai bună este $k = 59.5$, arborele de decizie fiind cel din figura alăturată.



15. (Alte criterii posibile pentru selecția atributelor în ID3: Gini impurity și Misclassification impurity)

prelucrare de Liviu Ciortuz, după

■ □ ● CMU, 2003 fall, T. Mitchell, A. Moore, HW1, pr. 4

Entropia este o mărime care cuantifică gradul de neomogenitate (engl., impurity) al unui set de instanțe în raport cu etichetele asignate. Algoritmul ID3 folosește entropia drept criteriu de partiționare (engl., splitting criterion), calculând *câștigul de informație* pentru a decide care este atributul care trebuie testat în nodul curent. Există, însă, și alte măsuri de neomogenitate care pot fi folosite, de asemenea, drept criterii de partiționare. În această problemă, vom investiga două astfel de măsuri.

Presupunem că nodul curent (n) din arborele de decizie aflat în curs de elaborare are asignate instanțe din k clase: c_1, c_2, \dots, c_k . Definim

$$\text{Gini Impurity: } i(n) = 1 - \sum_{i=1}^k P^2(c_i)$$

și

$$\text{Misclassification Impurity: } i(n) = 1 - \max_{i=1}^k P(c_i),$$

unde am notat cu $P(c_i)$ probabilitatea [sau: frecvența de apariție a instanțelor aparținând] clasei c_i în ansamblul instanțelor asignate la nodul curent.

a. Presupunem $k = 2$. Așadar, în acest caz nodul n are două clase: c_1 și c_2 . Desenați un grafic în care cele trei măsuri de neomogenitate — *Entropia*, *Gini Impurity* și *Misclassification Impurity* — sunt reprezentate în funcție de $P(c_1)$.

b. Acum putem da definiția unui nou criteriu de partiționare, bazat pe măsurile de neomogenitate *Gini* și *Misclassification*. În literatura de specialitate, acest nou criteriu este denumit uneori *Drop-of-Impurity* (pentru care propunem ca traducere în limba română termenul de *diminuarea neomogenității*). El reprezintă diferența dintre neomogenitatea nodului curent pe de o parte și suma ponderată a neomogenităților fiilor pe de altă parte. În cazul partiționării atributelor binare, definim *Drop-of-Impurity* ca fiind:

$$\Delta i(n) = i(n) - P(n_l) i(n_l) - P(n_r) i(n_r),$$

unde n_l și n_r reprezintă fiul-stânga și, respectiv, fiul-dreapta, care au fost derivați din nodul n după partiționare.

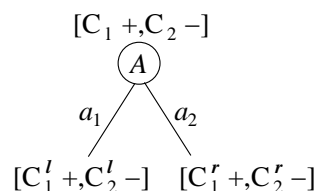
Folosind mai întâi *Gini Impurity* și apoi *Misclassification Impurity*, calculați *Drop-of-Impurity* pentru următorul set de instanțe asignate nodului pentru care se

testează atributul A luând valorile a_1 și a_2 . Am notat cu C variabila de ieșire (desemnând clasa) și cu c_1 și c_2 cele două valori ale ei.

A	a_1	a_1	a_1	a_2	a_2	a_2
C	c_1	c_1	c_2	c_2	c_2	c_2

c. Se poate crea un set de date de antrenament (sau: se poate modifica setul de date de mai sus) astfel încât, pe noul set, *Drop-of-Impurity* bazat pe *Misclassification* să fie 0 dar bazat pe *Entropy* și, respectiv, *Gini* să fie diferit de 0?

Sugestie: Puteți folosi următoarea *proprietate*, care este ușor de demonstrat: Dacă într-un set de date avem C_1 instanțe din clasa (sau: cu eticheta) c_1 și C_2 instanțe din clasa c_2 , cu $C_1 < C_2$, iar după partiționarea în funcție de valorile atributului A această relație se păstrează, adică $C_1^l < C_2^l$ și $C_1^r < C_2^r$ (evident, cu $C_1 = C_1^l + C_1^r$ și $C_2 = C_2^l + C_2^r$), unde l și r desemnează nodul-fiu stâng și respectiv nodul-fiu drept, atunci *Drop-of-Impurity* pentru *Misclassification* va fi 0. Însă, în aceleași condiții, *Drop-of-Impurity* bazat pe *Gini* sau pe *Entropy* va avea, în general, valori nenule.



(Evident, *proprietatea* de mai sus se menține dacă în locul relației $<$ vom considera peste tot relația $>$.)

Răspuns:

a. Vom scrie mai întâi expresiile celor trei funcții, iar apoi vom trasa graficele lor:

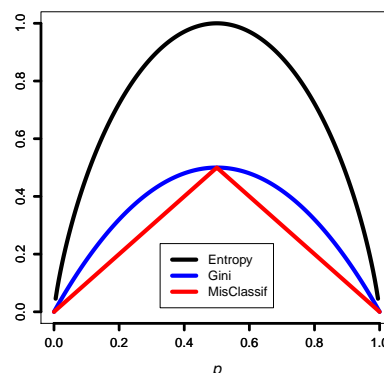
$$Entropy(p) = -p \log_2 p - (1-p) \log_2 (1-p)$$

$$Gini(p) = 1 - p^2 - (1-p)^2 = 2p(1-p)$$

$$MisClassif(p) =$$

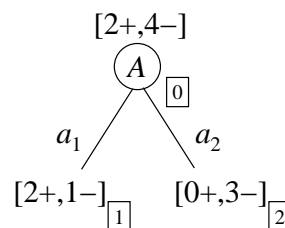
$$= \begin{cases} 1 - (1-p), & \text{pentru } p \in [0; 1/2) \\ 1 - p, & \text{pentru } p \in [1/2; 1] \end{cases}$$

$$= \begin{cases} p, & \text{pentru } p \in [0; 1/2) \\ 1 - p, & \text{pentru } p \in [1/2; 1] \end{cases}$$



Se observă că toate cele trei funcții iau valori în intervalul $[0; 1]$, sunt simetrice în raport cu punctul $p = 1/2$, sunt strict crescătoare pe intervalul $[0; 1/2]$ și strict descrescătoare pe intervalul $[1/2; 1]$, maximum fiecăreia dintre ele fiind obținut pentru $p = 1/2$.

b. Partiționarea datelor de antrenament în funcție de valorile atributului A se face așa cum se arată în figura alăturată. (În această figură și, de asemenea, în calculele de mai jos, pentru conveniență / simplitate, am asociat semnul $+$ etichetei c_1 și semnul $-$ etichetei c_2 .)



Aplicând formula $\Delta i(n) = i(n) - P(n_l) i(n_l) - P(n_r) i(n_r)$, vom obține pentru *Drop-of-Impurity* următoarele valori:

Gini: $p = 2/6 = 1/3 \Rightarrow$

$$\left. \begin{aligned} i(0) &= 2 \cdot \frac{1}{3} \left(1 - \frac{1}{3}\right) = \frac{2}{3} \cdot \frac{2}{3} = \frac{4}{9} \\ i(1) &= 2 \cdot \frac{2}{3} \left(1 - \frac{2}{3}\right) = \frac{4}{3} \cdot \frac{1}{3} = \frac{4}{9} \\ i(2) &= 0 \end{aligned} \right\} \Rightarrow \Delta i(0) = \frac{4}{9} - \frac{3}{6} \cdot \frac{4}{9} = \frac{4}{9} - \frac{2}{9} = \frac{2}{9}.$$

Misclassification: $p = 1/3 < 1/2 \Rightarrow$

$$\left. \begin{aligned} i(0) &= p = \frac{1}{3} \\ i(1) &= 1 - \frac{2}{3} = \frac{1}{3} \\ i(2) &= 0 \end{aligned} \right\} \Rightarrow \Delta i(0) = \frac{1}{3} - \frac{1}{2} \cdot \frac{1}{3} = \frac{1}{6}.$$

c. Într-adevăr, putem justifica ușor *proprietatea* din enunț:

$$\begin{aligned} \Delta i(n) &= \frac{C_1}{C_1 + C_2} - \left(\frac{C_1^l + C_2^l}{C_1 + C_2} \cdot \frac{C_1^l}{C_1^l + C_2^l} + \frac{C_1^r + C_2^r}{C_1 + C_2} \cdot \frac{C_1^r}{C_1^r + C_2^r} \right) \\ &= \frac{C_1}{C_1 + C_2} - \frac{C_1^l + C_1^r}{C_1 + C_2} = \frac{C_1}{C_1 + C_2} - \frac{C_1}{C_1 + C_2} = 0. \end{aligned}$$

Dacă în setul de date din enunț una dintre instanțele (a_1, c_1) se modifică în (a_1, c_2) și se adaugă o instanță (a_2, c_1) , atunci *Drop-of-Impurity* va avea următoarele valori:

$$\text{Entropy: } \Delta i(0) = H[5+, 2-] - \left(\frac{3}{7} H[2+, 1-] + \frac{4}{7} H[3+, 1-] \right) = 0.006 \neq 0;$$

$$\begin{aligned} \text{Gini: } 2 \left\{ \frac{2}{7} \left(1 - \frac{2}{7}\right) - \left[\frac{3}{7} \cdot \frac{1}{3} \left(1 - \frac{1}{3}\right) + \frac{4}{7} \cdot \frac{1}{4} \left(1 - \frac{1}{4}\right) \right] \right\} &= 2 \left\{ \frac{10}{49} - \left[\frac{2}{21} + \frac{3}{28} \right] \right\} \\ &= 2 \left(\frac{10}{49} - \frac{17}{84} \right) \neq 0; \end{aligned}$$

$$\text{Misclassification: } \Delta i(0) = \frac{2}{7} - \left(\frac{3}{7} \cdot \frac{1}{3} + \frac{4}{7} \cdot \frac{1}{4} \right) = 0.$$

16. (Algoritmul ID3 ca metodă de învățare de tip “eager”: posibilitatea suplimentării datelor de antrenament)

CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW1, pr. 3.4

Presupunem că, pornind de la un set de date de antrenament D , obținem un arbore de decizie ID3, notat cu T . Ulterior, cineva ne mai dă un set suplimentar de date de antrenament, D' . Putem proceda într-unul din următoarele două moduri:

– Putem rula din nou ID3, de această dată pe datele de antrenament $D \cup D'$, obținând arborele T_1 . (Dezavantaj: dacă $|D|$ este foarte mare, această procedură poate fi costisitoare ca timp.)

– Putem extinde T , arborele ID3 obținut pe mulțimea de antrenament D , ținând cont de datele D' . Arborele nou, T_2 , ar putea să nu fie la fel de bun ca T_1 (vedeți cazul de mai sus), dar el este consistent cu datele din $D \cup D'$ în cazul în care aceste mulțimi de antrenament nu conțin zgomote. În special dacă $|D'|$ este mic, această metodă este acceptabilă din punct de vedere practic.

Propuneți o procedură pentru obținerea efectivă a arborelui T_2 .

Răspuns:

Instanțele din D' se atașează nodurilor frunză ale arborelui T , conform procedurii de clasificare de la arbori de decizie. Pentru fiecare dintre aceste noduri frunză, dacă instanțele atașate nodului respectiv nu sunt toate etichetate identic, se aplică algoritmul ID3 (folosind mulțimea de attribute neutilizate pe drumul care unește nodul rădăcină al arborelui cu acest nod frunză). În acest mod se obține un alt arbore de decizie T_2 care este consistent cu datele din $D \cup D'$.

17. (Reducerea caracterului “greedy” al algoritmului ID3 prin calcularea câștigului de informație în maniera “2-step look-ahead”)

CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW4, pr. 1.2-6

Învățarea automată a arborilor de decizie depinde mult de utilizarea unui mecanism “greedy” de selecție a atributelor.

a. Dacă un set de date are A attribute booleene, calculați (ca expresie în funcție de A) numărul total de apeluri la funcția care calculează câștigul de informație pentru elaborarea întregului arbore de decizie.

Pentru *simplitate*, se va presupune că toate attributele sunt necesare pentru clasificarea unei instanțe, iar setul de date de antrenament conține toate instanțele posibile (adică toate combinațiile posibile de perechi atribut-valoare, inclusiv pentru atributul de ieșire).

b. Este posibil să îmbunătățim algoritmul ID3, făcându-l să se comporte mai puțin “greedy”, prin explorarea / prospectarea în avans (engl., look-ahead) a spațiului de căutare. La o explorare cu 2 pași înainte (engl., 2-step look-ahead), calculul câștigului de informație pentru un atribut a_i pe mulțimea de instanțe D va fi făcut cu ajutorul formulei

$$IG_{2\text{-step}}(D, a_i) = \max_{a_l, a_r} \left\{ \frac{n_l}{n_l + n_r} IG(D_l, a_l) + \frac{n_r}{n_l + n_r} IG(D_r, a_r) \right\},$$

unde

- a_l și a_r sunt attributele din nodurile descendente din nodul marcat cu atributul a_i ,
- D_l și D_r sunt seturile de instanțe asignate nodurilor descendente din a_i , iar n_l și n_r reprezintă cardinalul mulțimii D_l și respectiv D_r ;
- $IG(D_l, a_l)$ este câștigul de informație calculat (în sens clasic) pentru atributul a_l pe setul D_l ; similar, $IG(D_r, a_r)$ este câștigul de informație pentru atributul a_r pe setul D_r .

- b1. Explicați pe scurt de ce sunt necesari factorii $\frac{n_l}{n_l+n_r}$ și $\frac{n_r}{n_l+n_r}$ în formula de mai sus.
- b2. Dacă se folosesc A atribute booleene, câte apeluri la funcția $IG(\dots, \dots)$ sunt necesare pentru a stabili atributul din nodul rădăcină?
- b3. Vom evalua acum cât de costisitoare este această explorare în avans a spațiului de căutare.
Dacă $A = 10$ și se face aceeași presupuziție ca la punctul a, calculați câte niveluri complete din arborele ID3 standard se pot calcula cu același efort de calcul — exprimat ca număr de apeluri la funcția IG — ca la stabilirea atributului rădăcină în varianta *2-step look-ahead*.
- b4. Metoda de învățare a arborilor de decizie folosind *2-step look-ahead* crează o clasă de modele / ipoteze mai largă decât algoritmul ID3 simplu? Altfel spus, putem să calculăm în acest mod funcții de clasificare pe care nu le putem reprezenta cu arborii de decizie standard?

Răspuns:

a. Datorită presupunerii făcute pentru *simplitate*, arborele de decizie final va fi un arbore binar complet cu $A+1$ niveluri (incluzând și nodurile de decizie), notate în mod convențional de la 0 la A . Pe fiecare nivel $i = \overline{0, A-1}$ al arborelui binar se găsesc 2^i noduri. În fiecare dintre aceste noduri poate fi ales un atribut din cele $A-i$ rămase disponibile. Deci pentru determinarea nivelului i se fac $2^i \cdot (A-i)$ apeluri ale funcției IG .

Desigur, pe penultimul nivel, $A-1$, nu mai există decât un singur atribut rămas disponibil, prin urmare nu este necesar să se calculeze câștigul de informație. Așadar, numărul total de apeluri ale funcției IG pentru elaborarea întregului arbore de decizie este:

$$N_{IG} = A + 2(A-1) + 4(A-2) + \dots + 2^{A-2}(A - (A-2)) = \sum_{i=0}^{A-2} 2^i(A-i)$$

b1. n_l și n_r reprezintă numărul de instanțe asignate nodurilor descendente din a_i , deci factorii $\frac{n_l}{n_l+n_r}$ și $\frac{n_r}{n_l+n_r}$ reprezintă ponderile acestor sub-multimi în raport cu reuniunea lor. Cei doi factori vor pondera corespunzător câștigurile de informație de pe cele două ramuri din arborele de decizie. Dacă nu facem astfel de ponderări, este posibil să avem un câștig de informație mare pe o mulțime mică sau invers. În consecință, simpla însumare a celor două câștiguri de informație nu ar emula în mod veridic câștigul de informație pe întreg ansamblul lui D .

b2. Pentru stabilirea atributului din nodul rădăcină făcând o explorare cu 2 pași înainte, se calculează pentru fiecare dintre cele A atribute câștigurile de informație corespunzătoare celor 2 descendenți, adică:

$$N_{IG_{2\text{-step}}}(\text{rădăcină}) = A \cdot 2(A-1) = 2A^2 - 2A$$

b3. Dacă $A = 10$, atunci pentru stabilirea atributului rădăcină în varianta *2-step look-ahead* se apelează funcția IG de $N_{IG_{2\text{-step}}}(\text{rădăcină}) = 200 - 20 = 180$ de ori. Trebuie determinat numărul x de niveluri complete din arborele ID3

standard care se pot calcula folosind maxim 180 de apeluri ale funcției IG , adică argmax_x astfel încât $N_{IG}(x) = \sum_{i=0}^{x-1} 2^i(10-i) \leq 180$.

$$\begin{aligned} x=1 &\Rightarrow N_{IG}(1) = 2^0(10-0) = 10 \\ x=2 &\Rightarrow N_{IG}(2) = 10 + 2^1(10-1) = 28 \\ x=3 &\Rightarrow N_{IG}(3) = 28 + 2^2(10-2) = 60 \\ x=4 &\Rightarrow N_{IG}(4) = 60 + 2^3(10-3) = 116 \\ x=5 &\Rightarrow N_{IG}(5) = 116 + 2^4(10-4) = 212 > 180 \end{aligned}$$

Așadar, se pot calcula 4 niveluri complete din arborele ID3 standard cu același efort de calcul ca la stabilirea atributului din rădăcină în varianta (mai puțin “greedy”) a algoritmului ID3 cu *2-step look-ahead*.

b4. Nu. Clasa de modele / ipoteze pe care lucrează algoritmul ID3 cu *2-step look-ahead* este aceeași ca la algoritmul ID3 standard. Însă este foarte posibil ca arborele de decizie construit să fie mai bun dacă se face cătarea în maniera *2-step look-ahead*.

18.

(Îmbunătățirea algoritmului ID3, folosind IG cu “2-step look-ahead”)

Liviu Ciortuz, folosind date de la CMU, 2008 fall, Eric Xing, HW2, pr. 3

Se consideră variabilele booleene X_1 , X_2 și X_3 , precum și clasificarea $Y = \{0, 1\}$. Fie setul de date de antrenament din tabelul de mai jos.

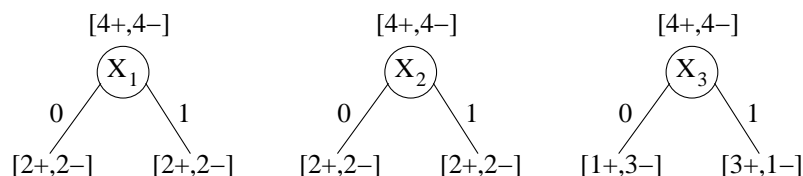
	X_1	X_2	X_3	Y
1	0	0	0	0
2	0	0	1	0
3	0	1	0	1
4	0	1	1	1
5	1	0	1	1
6	1	0	1	1
7	1	1	0	0
8	1	1	0	0

a. Elaborați arborele de decizie cu algoritmul ID3 standard.

b. Elaborați arborele de decizie cu algoritmul ID3 folosind câștigul de informație cu *2-step look-ahead*, așa cum este acesta definit în problema 17.

Răspuns:

a. Pentru nodul rădăcină vom alege unul dintre atributele X_1 , X_2 și X_3 . Corespunzător, vom obține următoarele partajări ale setului de date de antrenament:



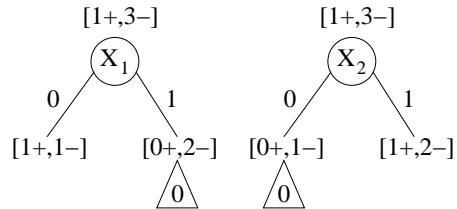
În cele ce urmează vom scrie câștigul de informație (IG) omițând primul argument, întrucât se consideră cunoscut din context. Vom calcula câștigul de informație corespunzător fiecărui atribut:

$$IG(X_1) = IG(X_2) = H[4+, 4-] - \left(\frac{4}{8}H[2+, 2-] + \frac{4}{8}H[2+, 2-] \right) = 1 - \left(\frac{1}{2} + \frac{1}{2} \right) = 0$$

$$IG(X_3) = H[4+, 4-] - \left(\frac{4}{8}H[1+, 3-] + \frac{4}{8}H[3+, 1-] \right)$$

Însă $H[1+, 3-] = H[3+, 1-]$, iar valoarea corespunzătoare este sub-unitară, deci $IG(X_3) > 0$. Așadar, în nodul rădăcină se alege atributul X_3 .

Pentru nodul corespunzător ramurii $X_3 = 0$ se poate alege dintre atributele care au mai rămas, adică X_1 sau X_2 . Se calculează câștigurile de informație corespunzătoare:



$$\begin{aligned} IG(X_1 | X_3 = 0) &= H[1+, 3-] - \left(\frac{2}{4}H[1+, 1-] + \frac{2}{4}H[0+, 2-] \right) \\ &= 2 - \frac{3}{4}\log_2 3 - \left(\frac{1}{2} + 0 \right) = \frac{3}{2} - \frac{3}{4}\log_2 3 = 0.311 \end{aligned}$$

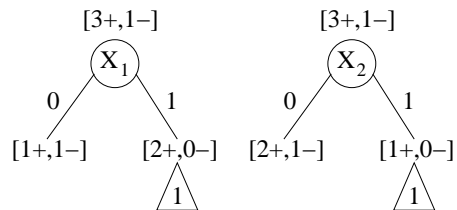
$$\begin{aligned} IG(X_2 | X_3 = 0) &= H[1+, 3-] - \left(\frac{1}{4}H[0+, 1-] + \frac{3}{4}H[1+, 2-] \right) \\ &= 2 - \frac{3}{4}\log_2 3 - \left(0 + \frac{3}{4}H[1+, 2-] \right) \\ &= 2 - \frac{3}{4}\log_2 3 - \frac{3}{4} \left(\frac{1}{3}\log_2 3 + \frac{2}{3}\log_2 \frac{3}{2} \right) \\ &= 2 - \frac{3}{4}\log_2 3 - \frac{3}{4} \left(\log_2 3 - \frac{2}{3} \right) = \frac{5}{2} - \frac{3}{2}\log_2 3 = 0.122 \end{aligned}$$

Cum $IG(X_1 | X_3 = 0) > IG(X_2 | X_3 = 0)$, se alege atributul X_1 .

În cele de mai sus, s-au folosit notațiile în maniera condițională $IG(X_1 | X_3 = 0)$ și $IG(X_2 | X_3 = 0)$ doar pentru a identifica în mod neambiguu care este nodul din arbore pentru care se calculează câștigul de informație respectiv.

Pentru nodul corespunzător ramurii $X_3 = 1$ ce pornește din nodul rădăcină se poate alege din nou unul dintre atributele X_1 și X_2 .

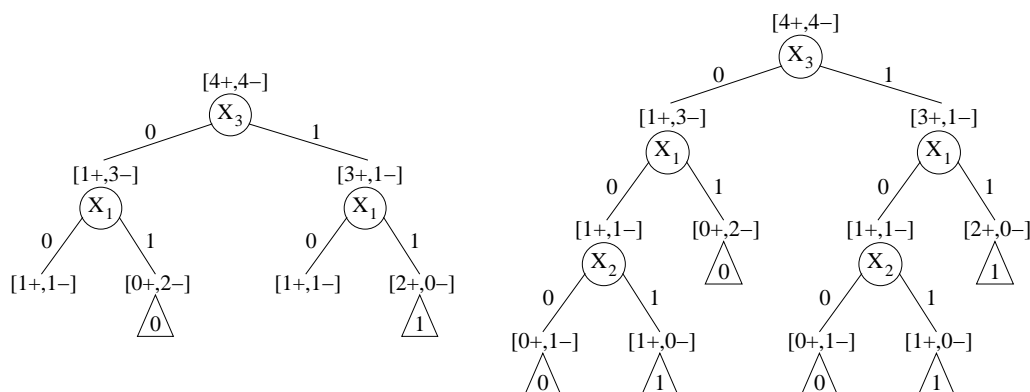
Pentru aceasta, se compară câștigurile de informație corespunzătoare. Se observă că acestea sunt egale cu cele din cazul precedent. Așadar,



$$IG(X_1 | X_3 = 1) = IG(X_1 | X_3 = 0) = \frac{3}{2} - \frac{3}{4}\log_2 3 = 0.311$$

$$IG(X_2 | X_3 = 1) = IG(X_2 | X_3 = 0) = \frac{5}{2} - \frac{3}{2}\log_2 3 = 0.122$$

Prin urmare, se alege tot atributul X_1 . Arborele de decizie construit până în acest moment este cel reprezentat mai jos în partea stângă:



În cele două noduri care au rămas, nu poate fi ales decât atributul X_2 . Deci arborele de decizie complet construit de algoritmul ID3 este cel reprezentat mai sus în partea dreaptă.

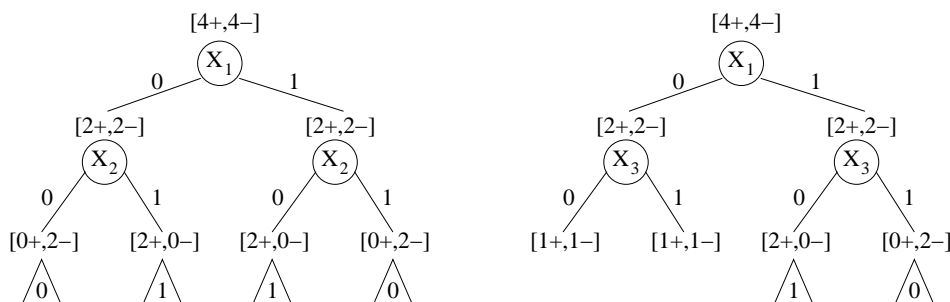
b. Dacă se folosește câștigul de informație cu *2-step look-ahead* pentru algoritmul ID3, atunci la alegerea fiecărui nod se iau în considerare [și] toate combinațiile posibile de noduri de pe nivelul următor, utilizându-se formula:

$$IG_{2\text{-step}}(D, a_i) = \max_{a_l, a_r} \left\{ \frac{n_l}{n_l + n_r} IG(D_l, a_l) + \frac{n_r}{n_l + n_r} IG(D_r, a_r) \right\}$$

Pentru nodul rădăcină se poate alege, la fel ca la punctul precedent, unul dintre atributele X_1 , X_2 și X_3 . Așadar, câștigul de informație corespunzător atributului X_1 va fi calculat astfel:

$$IG_{2\text{-step}}(X_1) = \max \left\{ \begin{array}{l} \frac{4}{8} IG(X_2 | X_1 = 0) + \frac{4}{8} IG(X_2 | X_1 = 1) \\ \frac{4}{8} IG(X_2 | X_1 = 0) + \frac{4}{8} IG(X_3 | X_1 = 1) \\ \frac{4}{8} IG(X_3 | X_1 = 0) + \frac{4}{8} IG(X_2 | X_1 = 1) \\ \frac{4}{8} IG(X_3 | X_1 = 0) + \frac{4}{8} IG(X_3 | X_1 = 1) \end{array} \right.$$

Pentru a determina această valoare, va trebui să calculăm cele 4 câștiguri de informație (în sens clasic) implicate în formulă, și este util să reprezentăm două din cele 4 situații posibile:



$$\begin{aligned}
IG(X_2 \mid X_1 = 0) &= IG(X_2 \mid X_1 = 1) = IG(X_3 \mid X_1 = 1) \\
&= H[2+, 2-] - \left(\frac{1}{2}H[0+, 2-] + \frac{1}{2}H[2+, 0-] \right) = 1 - 0 = 1 \\
IG(X_3 \mid X_1 = 0) &= H[2+, 2-] - \left(\frac{1}{2}H[1+, 1-] + \frac{1}{2}H[1+, 1-] \right) = 1 - \left(\frac{1}{2} + \frac{1}{2} \right) = 0
\end{aligned}$$

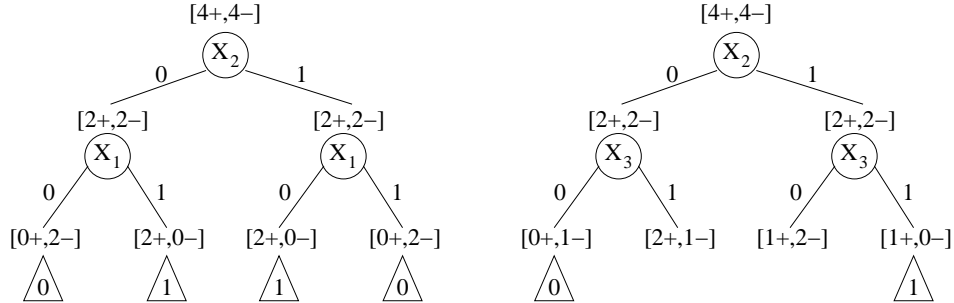
Prin urmare,

$$IG_{2\text{-step}}(X_1) = \max \left\{ \begin{array}{l} \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1 \\ \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1 \\ \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 1 \\ \frac{1}{2} \cdot 0 + \frac{1}{2} \cdot 1 \end{array} \right\} = \max \left\{ 1, 1, \frac{1}{2}, \frac{1}{2} \right\} = 1.$$

Câștigul de informație corespunzător atributului X_2 plasat în nodul rădăcină este:

$$IG_{2\text{-step}}(X_2) = \max \left\{ \begin{array}{l} \frac{4}{8}IG(X_1 \mid X_2 = 0) + \frac{4}{8}IG(X_1 \mid X_2 = 1) \\ \frac{4}{8}IG(X_1 \mid X_2 = 0) + \frac{4}{8}IG(X_3 \mid X_2 = 1) \\ \frac{4}{8}IG(X_3 \mid X_2 = 0) + \frac{4}{8}IG(X_1 \mid X_2 = 1) \\ \frac{4}{8}IG(X_3 \mid X_2 = 0) + \frac{4}{8}IG(X_3 \mid X_2 = 1) \end{array} \right\}$$

Vom reprezenta din nou două dintre situațiile posibile:



$$\begin{aligned}
IG(X_1 \mid X_2 = 0) &= IG(X_1 \mid X_2 = 1) = H[2+, 2-] - 0 = 1 \\
IG(X_3 \mid X_2 = 0) &= IG(X_3 \mid X_2 = 1) = H[2+, 2-] - \left(\frac{1}{4}H[0+, 1-] + \frac{3}{4}H[2+, 1-] \right) = \\
&= 1 - \frac{3}{4} \left(\log_2 3 - \frac{2}{3} \right) = \frac{3}{2} - \frac{3}{4} \log_2 3 = 0.311
\end{aligned}$$

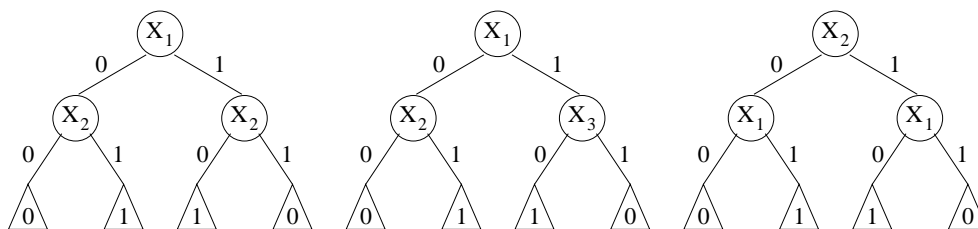
Prin urmare,

$$IG_{2\text{-step}}(X_2) = \max \left\{ \begin{array}{l} \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 1 \\ \frac{1}{2} \cdot 1 + \frac{1}{2} \cdot 0.311 \\ \frac{1}{2} \cdot 0.311 + \frac{1}{2} \cdot 1 \\ \frac{1}{2} \cdot 0.311 + \frac{1}{2} \cdot 0.311 \end{array} \right\} = \max \{ 1, 0.655, 0.655, 0.311 \} = 1$$

Pentru calculul câștigului de informație corespunzător atributului X_3 plasat în nodul rădăcină, au fost calculate la punctul a cele 4 câștiguri de informație în sens clasic implicate în formulă, deci rezultă:

$$\begin{aligned}
 IG_{2\text{-step}}(X_3) &= \max \begin{cases} \frac{4}{8}IG(X_1 | X_3 = 0) + \frac{4}{8}IG(X_1 | X_3 = 1) \\ \frac{4}{8}IG(X_1 | X_3 = 0) + \frac{4}{8}IG(X_2 | X_3 = 1) \\ \frac{4}{8}IG(X_2 | X_3 = 0) + \frac{4}{8}IG(X_1 | X_3 = 1) \\ \frac{4}{8}IG(X_2 | X_3 = 0) + \frac{4}{8}IG(X_2 | X_3 = 1) \end{cases} \\
 &= \max \begin{cases} \frac{1}{2} \left(\frac{3}{2} - \frac{3}{4} \log_2 3 \right) + \frac{1}{2} \left(\frac{3}{2} - \frac{3}{4} \log_2 3 \right) \\ \frac{1}{2} \left(\frac{3}{2} - \frac{3}{4} \log_2 3 \right) + \frac{1}{2} \left(\frac{5}{2} - \frac{3}{2} \log_2 3 \right) \\ \frac{1}{2} \left(\frac{5}{2} - \frac{3}{2} \log_2 3 \right) + \frac{1}{2} \left(\frac{3}{2} - \frac{3}{4} \log_2 3 \right) \\ \frac{1}{2} \left(\frac{5}{2} - \frac{3}{2} \log_2 3 \right) + \frac{1}{2} \left(\frac{5}{2} - \frac{3}{2} \log_2 3 \right) \end{cases} = \max \begin{cases} \frac{3}{2} - \frac{3}{4} \log_2 3 \\ 2 - \frac{9}{8} \log_2 3 \\ 2 - \frac{9}{8} \log_2 3 \\ \frac{5}{2} - \frac{3}{2} \log_2 3 \end{cases} \\
 &= \max\{0.311, 0.216, 0.216, 0.122\} = 0.311
 \end{aligned}$$

Comparând $IG_{2\text{-step}}(X_1)$, $IG_{2\text{-step}}(X_2)$ și $IG_{2\text{-step}}(X_3)$, obținem valoarea maximă 1 fie pentru X_1 , fie pentru X_2 în nodul rădăcină. Având în vedere calculele realizate pentru a obține aceste valori, nu mai sunt necesare operații suplimentare pentru determinarea arborelui de decizie construit de ID3 cu metoda *2-step look-ahead*. Există de fapt 3 arbori de decizie optimi (ca număr de niveluri și / sau noduri) pentru aceste date de antrenament, și anume:



Este important de remarcat faptul că algoritmul ID3 cu *2-step look-ahead* identifică toate aceste trei soluții optime, în vreme ce algoritmul ID3 standard nu identifică niciuna dintre ele.

19.

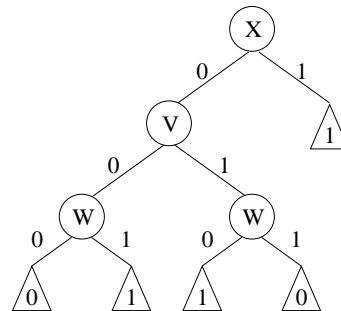
(O strategie de pruning pentru arborele ID3:
eliminarea nodurilor cu $IG < \epsilon$;
explorare top-down vs. bottom-up)

prelucrare de Liviu Ciortuz, după

■ □ ● ○ CMU, 2006 spring, Carlos Guestrin, midterm exam, pr. 4

Puteți constata ușor că, aplicând algoritmul ID3 pe datele din tabelul de mai jos, se va obține arborele de decizie alăturat.

V	W	X	Y
0	0	0	0
0	1	0	1
1	0	0	1
1	1	0	0
1	1	1	1



a. Pentru un astfel de arbore de decizie, o strategie simplă de trunchiere (engl., pruning) în vederea contracarării fenomenului de “overfitting” constă în a parcurge arborele de sus în jos, începând deci cu nodul-rădăcină și identificând fiecare nod de test pentru care câștigul de informație (sau un alt criteriu fixat în avans) are o valoare mai mică decât o valoare pozitivă, mică, fixată de la început, ε . Orice astfel de nod de test este imediat înlocuit — împreună cu subarboarele corespunzător lui — cu un nod de decizie, conform etichetei majoritare a instanțelor asignate nodului de test. Această strategie se numește “top-down pruning”.

Care este arborele de decizie obținut aplicând această strategie pe arborele de mai sus, dacă se consideră $\varepsilon = 0.0001$? Care este eroarea la antrenare pentru noul arbore?

b. O altă posibilitate de a face pruning este să parcurgem arborele de decizie începând cu părinții nodurilor-frunză și să eliminăm în mod recursiv acele noduri de test pentru care câștigul de informație (sau un alt criteriu ales) este mai mic decât ε . Aceasta este strategia de “bottom-up pruning”.

Observație: Spre deosebire de strategia top-down, în varianta de pruning de tip bottom-up nu vor fi eliminate noduri (cu $IG < \varepsilon$) pentru care există descendenți al căror câștig de informație este mai mare sau egal cu ε .

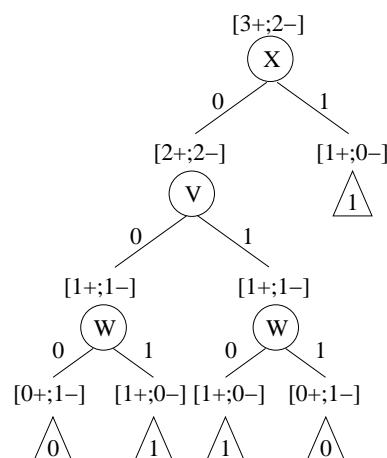
Ce arbore se obține făcând “bottom-up pruning” pe arborele dat mai sus, dacă se consideră $\varepsilon = 0.0001$? Care este eroarea la antrenare pentru arborele rezultat?

c. Stabiliți în ce situații ar fi indicat să alegem strategia “bottom-up pruning” în loc de “top-down pruning” și viceversa. Comparați acuratețea la antrenare și complexitatea computațională a celor două strategii de pruning.

d. Cât este înălțimea — adică, numărul de niveluri de test — pentru arborele returnat de ID3 urmat de “bottom-up pruning”? Puteți găsi un arbore de decizie având o înălțime mai mică, dar care clasifică perfect setul de antrenament? Ce concluzie putem trage despre calitatea [outputului] algoritmului ID3?

Răspuns:

Înainte de a rezolva efectiv punctele a și b , vom augmenta arborele de decizie dat cu informațiile referitoare la numărul de instanțe (pozitive și, respectiv, negative) asignate fiecărui nod de test. Obținem astfel figura alăturată.



a. Câștigul de informație al atributului X plasat în nodul-rădăcină este:

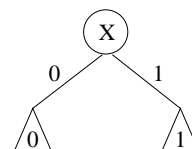
$$H[3+; 2-] - 1/5 \cdot 0 - 4/5 \cdot 1 = 0.971 - 0.8 = 0.171 > \varepsilon.$$

Prin urmare, acest nod nu va fi eliminat din arbore.

Câștigul de informație al atributului V este:

$$H[2+; 2-] - 1/2 \cdot 1 - 1/2 \cdot 1 = 1 - 1 = 0 < \varepsilon.$$

Așadar, nodul reprezentat de atributul V va fi eliminat și vom obține arborele de decizie [trunchiat], reprezentat în figura alăturată. Menționăm că am fi putut la fel de bine (din punctul de vedere al numărului de erori la antrenare) să alegem decizia $Y = 1$ în nodul-fiu stâng, însă în acel caz arborele s-ar fi redus de fapt la un singur nod de decizie (cu outputul $Y = 1$).



Eroarea la antrenare produsă de acest arbore este $2/5$.

b. Câștigul de informație al celor două noduri marcate cu atributul W în arborele dat în enunț este același, și anume 1 (se poate verifica imediat). Prin urmare, aplicând strategia de “bottom-up pruning”, arborele de decizie rămâne identic cu cel inițial. Evident, eroarea la antrenare pentru acest arbore este 0, întrucât datele de antrenament nu conțin inconsistențe.

c. Din cauza faptului că la top-down pruning, odată cu un nod de test pentru care câștigul de informație (IG) este mai mic decât valoarea ε se elimină întregul subarbore care are ca rădăcină acel nod de test, această strategie este mai rapidă decât (sau, în cel mai rău caz, la fel de rapidă / lentă ca și) pruning-ul de tip bottom-up.

Așa cum s-a menționat în *observația* din enunț și s-a exemplificat apoi la punctele a și b , dezavantajul pruning-ului de tip top-down este că odată ce este eliminat un nod cu IG mai mic decât ε , este posibil ca între descendenții săi (eliminați) să fie și noduri care au câștigul de informație mai mare sau egal cu ε . Pruning-ul bottom-up nu are acest dezavantaj; el este deci mai „conservativ”.

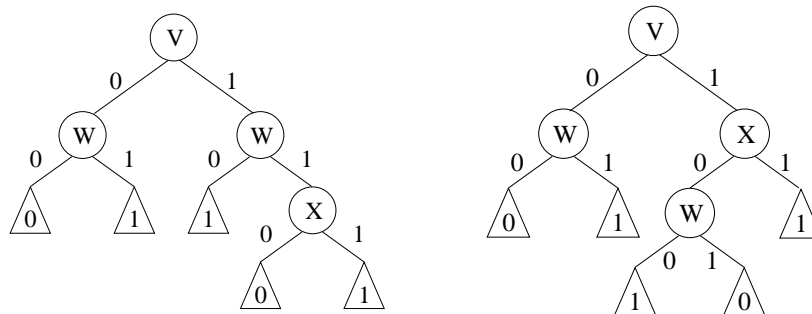
O problemă care poate însă să apară și în cazul pruning-ului de tip bottom-up este faptul că în nodurile apropiate de nodurile de decizie (acestea din urmă

fiind nodurile-frunză), câștigul de informație se calculează uneori pe mulțimi mici de exemple, deci testul $IG \geq \varepsilon$ nu este neapărat semnificativ din punct de vedere statistic. Așadar, este *recomandabil* ca în astfel de situații să se folosească un test statistic, de exemplu *testul* χ^2 . (A se vedea problemele 20 și 51.)

În ce privește comparația dintre acuratețile arborilor obținuți prin aplicarea celor două variante de pruning: se observă că arborele mai simplu (cel de la punctul a) are o eroare la antrenare mai mare decât arborele mai complex (cel de la punctul b), însă este mai probabil ca acesta din urmă să producă overfitting.

d. Din rezolvarea dată la punctul b, rezultă imediat că înălțimea arborelui obținut prin aplicarea algoritmului ID3 urmat de pruning de tip bottom-up este 3. (Nu se iau în considerare nodurile frunză.) Vom demonstra — exact ca la problema 1 — că nu există un arbore de decizie consistent cu datele de antrenament, care să aibă adâncimea strict mai mică decât 3:

Se observă ușor din tabelul dat în enunț că $Y = (V \text{ XOR } W) \vee X$, așadar variabilele V și W au rol simetric în definirea funcției reprezentate de variabila Y . Punând atributul X în nodul rădăcină, arborele de decizie minimal care se poate obține este cel dat în enunț (sau, echivalent, arborele care obține din acesta interschimbând V și W). Dacă, în schimb, punem atributul V în nodul rădăcină, se pot obține doi arbori de decizie minimali, așa cum se arată grafic mai jos. (Și similar, dacă în nodul rădăcină punem atributul W .)



20. (ID3 cu post-pruning: folosirea testului statistic χ^2 pentru limitarea overfitting-ului)
 prelucrare de Livi Ciortuz, după
 ■ • CMU, 2010 fall, Ziv Bar-Joseph, HW2, pr. 2.1

În acest exercițiu veți face pruning asupra unui arbore ID3 (după ce s-a făcut antrenarea pe toate datele disponibile), folosind o metodă statistică de testare / verificare a ipotezelor.

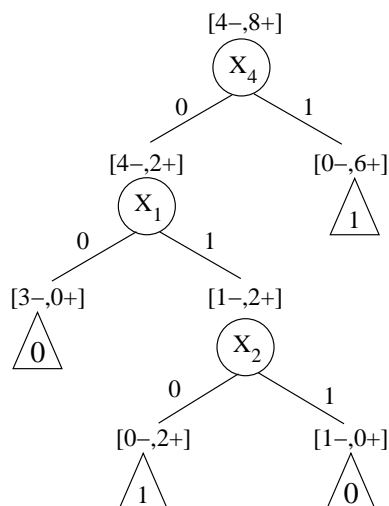
După ce a fost învățat arborele de decizie, vizităm fiecare nod intern (inclusiv nodul rădăcină) și testăm dacă atributul care a fost pus în nodul respectiv nu este cumva necorelat cu eticheta / clasa specificată de către atributul de ieșire.

Pentru aceasta, mai întâi presupunem că atributul din nodul respectiv este independent de atributul de ieșire (aceasta este așa-numita „ipoteză nulă”), iar apoi folosim testul χ^2 al lui Pearson pentru a genera o „statistică”, care

poate constitui temeiul pentru respingerea „ipotezei nule“. Dacă ipoteza nulă nu poate fi respinsă, eliminăm sub-arborele din nodul respectiv (de fapt, îl înlocuim cu un nod de decizie).

Pentru a ilustra aceste chestiuni, considerăm arborele de decizie de mai jos (partea dreaptă); el a fost construit pornind de la datele din tabelul alăturat lui, folosind algoritmul ID3.

X_1	X_2	X_3	X_4	Class
1	1	0	0	0
1	0	1	0	1
0	1	0	0	0
1	0	1	1	1
0	1	1	1	1
0	0	1	0	0
1	0	0	0	1
0	1	0	1	1
1	0	0	1	1
1	1	0	1	1
1	1	1	1	1
0	0	0	0	0



a. Pentru fiecare nod intern din arborele de decizie vom crea o *tabelă de contingență* (engl., contingency table) pentru exemplele de antrenare care sunt asignate nodului respectiv. Tabela aceasta va avea coloanele etichetate cu cele c clase / valori ale variabilei de ieșire. Similar, valorile atributului testat în nodul respectiv (având în total r valori) vor fi asignate liniilor tabelului. Dacă acceptăm o ușoară simplificare în forma de exprimare, vom putea spune că un element oarecare $O_{i,j}$ din tabela de contingență reprezintă numărul de „observații“ (adică, instanțe de antrenament asignate nodului respectiv) pentru care valoarea atributului testat este i , iar eticheta / clasa este j .

Calculați tabelele de contingență pentru cele trei noduri de test ale arborelui de decizie dat mai sus. Apoi, pornind de la datele conținute în fiecare dintre aceste matrice de contingență, estimați (în sensul verosimilității maxime) probabilitățile pentru valorile variabilei de ieșire (*Class*), precum și pentru valorile atributului din nodul corespunzător.

b. Pentru a aplica testul statistic χ^2 , avem nevoie să calculăm pentru fiecare nod intern al arborelui de decizie încă o tabelă (pe care o vom nota cu E), în care să consemnăm *numărul așteptat* de apariții (engl., expected counts) ale instanțelor de antrenament la nodul respectiv, pentru fiecare pereche de indici i, j având semnificația de mai sus. Acest număr așteptat este numărul (mediu) de instanțe de antrenament pe care le-am „observa“ în nodul respectiv dacă atributul selectat și clasa (variabila de ieșire) ar fi independente.

Derivați o formulă pentru calculul fiecărui element (notat $E_{i,j}$) din această a doua tabelă.

Întrebări ajutătoare: Care este probabilitatea ca exemplele de antrenament asignate nodului respectiv să aibă o anumită etichetă (j)? Ținând cont de această probabilitate, precum și de presupuziția de independență formulată

prin ipoteza „nulă”, care este numărul de exemple cu o anumită valoare (i) pentru atributul selectat în nodul respectiv, care ar trebui (i.e., „ne așteptăm”) să aibă acea etichetă / clasă (j)?

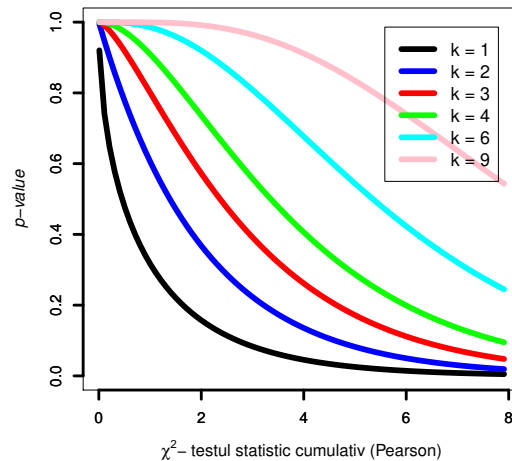
Folosind formula pe care tocmai ați derivat-o, calculați matricea E pentru fiecare dintre cele trei noduri de test ale arborelui de decizie dat mai sus.

c. Date fiind cele două tabele pentru nodul considerat, puteți calcula acum testul statistic χ^2 -pătrat:

$$\chi^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}}$$

Puteți introduce valoarea calculată χ^2 precum și numărul de grade de libertate $(r-1)(c-1)$ într-un program³⁸⁵ sau într-un calculator on-line³⁸⁶ pentru a calcula o așa-numită p -valoare (engl., p -value).³⁸⁷

În general, dacă $p < 0.05$, se va considera că nu avem suficientă *evidență* în favoarea ipotezei „nule”, care afirmă că atributul selectat și clasa (variabila de ieșire) sunt independente, și o vom respinge. Într-o astfel de situație, spunem că testul [din nodul respectiv] este *semnificativ* din punct de vedere statistic.



Pentru fiecare dintre cele trei noduri interne din arborele de decizie dat mai sus, găsiți p -valoarea corespunzătoare și precizați dacă testul din nodul respectiv este sau nu semnificativ d.p.v. statistic. Cât de multe noduri interne vor fi eliminate din arbore dacă la prunung impunem condiția $p \geq 0.05$ [pentru a elimina un nod de test din arbore și a-l înlocui cu un nod de decizie]?

Răspuns:

a. Pentru fiecare nod din arborele ID3 vom alcătui matricea de contingență asociată, pornind de la partiționările mulțimilor de exemple care au fost asigurate (de către algoritmul ID3) descendenților nodului respectiv. Apoi, din fiecare matrice de contingență vom estima (în sensul verosimilității maxime)

³⁸⁵Folosiți `1-chi2cdf(x,df)` în MATLAB sau `CHIDIST(x,df)` în Excel.

³⁸⁶<http://faculty.vassar.edu/lowry/tabs.html#csq> este un astfel de calculator.

³⁸⁷Statistica χ^2 este aproximată de distribuția χ^2 cu numărul corespunzător (k) de grade de libertate. (Distribuția χ^2 reprezintă suma pătratelor a k variabile gaussiene standard independente.)

p -valoarea despre care este vorba mai sus reprezintă probabilitatea ca distribuția χ^2 să ia valori mai mari sau egale cu valoarea considerată (i.e., valoarea calculată pentru statistica χ^2). Așadar, p -valoarea pentru testul χ^2 se calculează făcând diferența dintre 1 și valoarea funcției de distribuție cumulative (c.d.f.) pentru distribuția χ^2 cu k de grade de libertate. Vedeți site-ul

https://en.m.wikipedia.org/wiki/Chi-square_distribution#/Table_of_.CF.872_value_vs_p-value (accesat la 5.09.2015).

probabilitățile pentru valorile variabilei de ieșire (*Class*), precum și probabilitățile pentru valorile atributului din acel nod. (Atenție la condiționarea probabilităților!)

$$\begin{array}{c|cc} O_{X_4} & Class = 0 & Class = 1 \\ \hline X_4 = 0 & 4 & 2 \\ X_4 = 1 & 0 & 6 \end{array} \xrightarrow{N=12} \begin{cases} P(X_4 = 0) = \frac{6}{12} = \frac{1}{2}, P(X_4 = 1) = \frac{1}{2} \\ P(Class = 0) = \frac{4}{12} = \frac{1}{3}, P(Class = 1) = \frac{2}{3} \end{cases}$$

$$\begin{array}{c|cc} O_{X_1|X_4=0} & Class = 0 & Class = 1 \\ \hline X_1 = 0 & 3 & 0 \\ X_1 = 1 & 1 & 2 \end{array} \xrightarrow{N=6} \begin{cases} P(X_1 = 0 | X_4 = 0) = \frac{3}{6} = \frac{1}{2} \\ P(X_1 = 1 | X_4 = 0) = \frac{1}{2} \\ P(Class = 0 | X_4 = 0) = \frac{4}{6} = \frac{2}{3} \\ P(Class = 1 | X_4 = 0) = \frac{1}{3} \end{cases}$$

$$\begin{array}{c|cc} O_{X_2|X_4=0, X_1=1} & Class = 0 & Class = 1 \\ \hline X_2 = 0 & 0 & 2 \\ X_2 = 1 & 1 & 0 \end{array} \xrightarrow{N=3} \begin{cases} P(X_2 = 0 | X_4 = 0, X_1 = 1) = \frac{2}{3} \\ P(X_2 = 1 | X_4 = 0, X_1 = 1) = \frac{1}{3} \\ P(Class = 0 | X_4 = 0, X_1 = 1) = \frac{1}{3} \\ P(Class = 1 | X_4 = 0, X_1 = 1) = \frac{2}{3} \end{cases}$$

b. Considerăm i o valoare arbitrar aleasă (dar fixată) pentru atributul de intrare A care este testat în nodul curent, iar j o valoare arbitrar aleasă (de asemenea, fixată) pentru atributul de ieșire $Class$ (renotat cu C). Ținând cont de presupuziția de independență stipulată de ipoteza „nulă”, putem scrie:

$$P(A = i, C = j) = P(A = i) \cdot P(C = j)$$

Probabilitățile $P(A = i)$ și $P(C = j)$ pot fi estimate — în sensul verosimilității maxime (MLE) —, cu ajutorul celor N instanțe de antrenament asignate nodului respectiv. Instanțele pentru care atributul A are valoarea i sunt tocmai cele din linia i a matricei. În mod similar, instanțele care au eticheta / clasa j sunt cele din coloana j a matricei de count-uri observate. Așadar,

$$P(A = i) = \frac{\sum_{k=1}^c O_{i,k}}{N} \text{ și } P(C = j) = \frac{\sum_{k=1}^r O_{k,j}}{N}$$

În consecință,

$$P(A = i, C = j) = \frac{(\sum_{k=1}^c O_{i,k})(\sum_{k=1}^r O_{k,j})}{N^2},$$

iar valoarea așteptată — repetăm, în condițiile presupuziției de independență — pentru numărul de instanțe având atributul $A = i$ și clasa $C = j$ va fi dată de formula

$$E_{i,j} = N \cdot P(A = i, C = j) = \frac{(\sum_{k=1}^c O_{i,k})(\sum_{k=1}^r O_{k,j})}{N}$$

c. Folosind probabilitățile calculate la punctul precedent și ținând cont de presupuziția de independență, calculăm numărul de observații așteptate în fiecare nod, pentru a completa matricele E :

E_{X_4}	$Class = 0$	$Class = 1$	$E_{X_1 X_4}$	$Class = 0$	$Class = 1$
$X_4 = 0$	2	4	$X_1 = 0$	2	1
$X_4 = 1$	2	4	$X_1 = 1$	2	1

$E_{X_2 X_4, X_1=1}$	$Class = 0$	$Class = 1$
$X_2 = 0$	$\frac{2}{3}$	$\frac{4}{3}$
$X_2 = 1$	$\frac{1}{3}$	$\frac{2}{3}$

Ca să *exemplificăm* cum am procedat, detaliem mai jos calculul pentru primul element din matricea E_{X_4} :

$$N = 12, P(X_4 = 0) = \frac{1}{2} \text{ și } P(Class = 0) = \frac{1}{3} \Rightarrow$$

$$N \cdot P(X_4 = 0, Class = 0) = N \cdot P(X_4 = 0) \cdot P(Class = 0) = 12 \cdot \frac{1}{2} \cdot \frac{1}{3} = 2$$

Acum putem aplica pentru fiecare nod din arborele de decizie formula de calcul a valorilor / statisticilor χ^2 care ne-a fost dată în enunț:

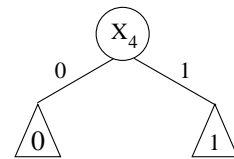
$$\chi^2_{X_4} = \frac{(4-2)^2}{2} + \frac{(2-4)^2}{4} + \frac{(0-2)^2}{2} + \frac{(6-4)^2}{4} = 2 + 2 + 1 + 1 = 6$$

$$\chi^2_{X_1|X_4=0} = \frac{(3-2)^2}{2} + \frac{(0-1)^2}{1} + \frac{(1-2)^2}{2} + \frac{(2-1)^2}{1} = 3$$

$$\chi^2_{X_2|X_4=0, X_1=1} = \frac{\left(0 - \frac{2}{3}\right)^2}{\frac{2}{3}} + \frac{\left(2 - \frac{4}{3}\right)^2}{\frac{4}{3}} + \frac{\left(1 - \frac{1}{3}\right)^2}{\frac{1}{3}} + \frac{\left(0 - \frac{2}{3}\right)^2}{\frac{2}{3}} = \frac{4}{9} \cdot \frac{27}{4} = 3$$

Accesând pagina web indicată în enunț, am obținut p -valorile următoare: 0.0143, 0.0833 și 0.0833.

În consecință, cu un grad de încredere de cel puțin 95%, nodurile situate pe nivelurile 1 și 2 din arborele ID3 pot fi eliminate. Pentru nodul rădăcină, ipoteza „nulă“ nu se verifică, adică variabilele X_4 și $Class$ nu sunt independente. Arborele obținut în urma prunării este cel din figura alăturată.



Observație: Este de remarcat faptul că arborele ID3 furnizat în enunț are (întâmplător) atât pentru atributul X_4 (din nodul rădăcină) cât și pentru nodul X_1 (de pe primul nivel) același câștig de informație, 0.4591. În urma testului χ^2 , se va elimina însă doar nodul care-l conține pe X_1 . Valoarea statisticii χ^2 asociate nodului X_1 este mult mai mică (3, față de 6, cât este pentru nodul care-l conține pe X_4), deci vom avea suficientă „evidență“ pentru a trage concluzia, cu un grad de încredere de 95%, că variabila de ieșire, $Class$, și $X_1|X_4 = 0$ sunt independente.³⁸⁸ Dacă am fi aplicat o metodă de prunare

³⁸⁸Observați și faptul că nodul care conține atributul X_1 are mai puține instanțe asociate (și anume, jumătate) față de cele asociate nodului rădăcină.

bazată pe câștigul de informație, aceasta n-ar fi putut să trateze în mod diferit cele două noduri, adică să-l păstreze pe unul și să-l elimine pe celălalt.

21. (Adevărat sau Fals?)

a. *Liviu Ciortuz*

Algoritmul ID3 garantează obținerea arborelui de decizie optimal (ca număr de niveluri sau de noduri).

CMU, 2002 spring, A. Moore, midterm example questions, pr. 1.c

b. Întrucât arborii de decizie pot învăța să clasifice instanțe într-un număr discret de clase (deci nu învață funcții cu valori reale), este imposibil ca ei să manifeste fenomenul de overfitting.

CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, midterm, pr. 1.4

c. Fie A și B doi algoritmi de clasificare automată. Algoritmul A este mai bun decât algoritmul B dacă eroarea la antrenare a algoritmului A este mai mică decât eroarea la antrenare a algoritmului B . Justificați.

CMU, 2005 spring, C. Guestrin, T. Mitchell, midterm, pr. 2.c

d. Presupunem că avem m instanțe și că vom folosi jumătate dintre ele pentru antrenarea unui clasificator oarecare (nu neapărat ID3) și jumătate pentru testare. Diferența dintre eroarea la antrenare și eroarea la testare descrește pe măsură ce numărul m crește.

Răspuns:

a. Fals. ID3 nu garantează obținerea arborelui optim (relativ la numărul de niveluri și / sau noduri), ci încearcă să găsească o soluție convenabilă, însă fără să caute în mod exhaustiv în tot spațiul soluțiilor. Mai exact, căutarea soluției se face în manieră “greedy”, maximizând un anumit criteriu (e.g., câștigul de informație) la fiecare iterație. O astfel de căutare nu garantează obținerea optimului.

b. Fals. Arborii de decizie manifestă fenomenul de overfitting. Prima parte a afirmației din enunț este adevărată — arborii de decizie pot învăța să clasifice instanțe într-un număr discret de clase —, însă a doua parte este falsă, deci avem o implicație de forma: $T \rightarrow F \equiv \neg T \vee F \equiv F$.

c. Fals, fiindcă la testare algoritmul B poate să aibă o eroare mai mică decât algoritmul A . Într-un astfel de caz, se spune că algoritmul A este “overfit” (rom., supra-antrenat).

d. Adevărat. Pe măsură ce dispunem de tot mai multe date de antrenament, dacă datele de antrenament sunt inconsistente, eroarea la antrenare va crește, fiindcă va fi din ce în ce mai greu ca modelul învățat să se adapteze la „zgometele” din date. Similar, eroarea la testare va descrește, fiindcă producem un clasificator care este din ce în ce mai puțin afectat de “overfitting” pe datele de antrenament. Cele două erori vor converge la așa-numita *eroare adevărată* (engl., true error) fiindcă diferențele statistice dintre datele de antrenament și datele de testare vor dispărea.

2 Clasificare bayesiană

Sumar

Noțiuni preliminare

- probabilități și probabilități condiționate;
- formula lui Bayes: ex. 5.b;
cap. *Fundamente*, ex. 6, ex. 7, ex. 83, ex. 84;
- independența [condițională a] evenimentelor aleatoare:
cap. *Fundamente*, ex. 4, ex. 80, ex. 81;
- independența [condițională a] variabilelor aleatoare: ex. 9, ex. 10, ex. 12,
ex. 31-38; vedeți și cap. *Fundamente*, ex. 15, ex. 27, ex. 88.b, ex. 97, ex. 95;
- distribuții probabiliste comune, marginale și condiționale: ex. 8, ex. 10, ex. 12,
ex. 31; vedeți și cap. *Fundamente*, ex. 13, ex. 14;
- distribuția gaussiană: de la cap. *Fundamente*, ex. 29, ex. 30 (pentru cazul
unidimensional), ex. 32 (pentru cazul bidimensional), ex. 20, ex. 31, ex. 33,
ex. 34 (pentru cazul multidimensional);
- estimarea parametrilor pentru distribuții de tip Bernoulli, categorial și gaus-
sian (ultimul doar pentru cazul clasificării bayesiene de tip gaussian);³⁰⁵
- ipoteze MAP vs. ipoteze ML:
formulare [ca soluții la] probleme de optimizare:³⁰⁶ ex. 25;
exemplificare: ex. 1, ex. 2, ex. 3, ex. 24, ex. 37;
exemplificare în cazul arborilor de decizie: ex. 4;
- regresia logistică, chestiuni introductive:³⁰⁷ de la cap. *Metode de regresie*,
ex. 12.

Algoritmi de clasificare bayesiană

- Algoritmul Bayes Naiv și algoritmul Bayes Optimal:³⁰⁸
formulare ca probleme de optimizare / estimare în sens MAP: cartea ML,
pag. 167;
pseudo-cod: cartea ML, pag. 177; vedeți și slide-urile lui Tom Mitchell;
exemple de aplicare: ex. 5, ex. 7, ex. 8, ex. 9, ex. 26, ex. 27, ex. 28;

³⁰⁵De la cap. *Fundamente*, pentru estimarea parametrului unei distribuții Bernoulli vedeți ex. 40 și ex. 113.a, pentru estimarea parametrilor unei distribuții categoriale vedeți ex. 115, iar pentru estimarea parametrilor unei distribuții gaussiene vedeți ex. 45, ex. 46, ex. 122 (pentru cazul unidimensional) și ex. 48 (pentru cazul multidimensional).

³⁰⁶Vedeți cartea ML, pag. 156-157.

³⁰⁷Vedeți draftul capitolului suplimentar pentru cartea ML a lui T. Mitchell, *Generative and discriminative classifiers: Naive Bayes and logistic regression* (în special secțiunea 3).

³⁰⁸La secțiunea aceasta, precum și la următoarea secțiune, considerăm (implicit) că toate variabilele de intrare sunt de tip Bernoulli sau, mai general, de tip categorial. După aceea vom considera și variabile de intrare de tip continuu, în genere de tip gaussian. Variabila de ieșire se consideră întotdeauna de tip Bernoulli / categorial.

- aplicarea / adaptarea algoritmului Bayes Naiv pentru clasificare de texte:³⁰⁹ ex. 6, ex. 29;
folosirea regulii “add-one” [a lui Laplace] pentru „netezirea” parametrilor: ex. 6, ex. 30;
- calculul ratei medii a erorilor pentru algoritmi Bayes Naiv și Bayes Optimal: ex. 10, ex. 11, ex. 31, ex. 32, ex. 33, ex. 34, ex. 38;
- evidențierea grafică a neconcordanței predicțiilor făcute de clasificatorii Bayes Naiv și Bayes Optimal: ex. 12.

Proprietăți ale algoritmilor Bayes Naiv și Bayes Optimal

- (P0) dacă proprietatea de independență condițională a atributelor de intrare în raport cu variabila de ieșire se verifică, atunci rezultatele produse de către cei doi algoritmi (Bayes Naiv și Bayes Optimal) în faza de testare coincid;
- (P1) numărul de parametri necesari de estimat din date: liniar pentru Bayes Naiv ($2d + 1$) și exponențial pentru Bayes Optimal ($2^{d+1} - 1$):³¹⁰ ex. 7.e, ex. 28.ab, ex. 33.ac;
- (P2) complexitatea algoritmului Bayes Naiv:
 - complexitatea de spațiu: $\mathcal{O}(dn)$
 - complexitatea de timp:
 - la antrenare: $\mathcal{O}(dn)$
 - la testare: $\mathcal{O}(d')$,
 - unde n este numărul de exemple, iar d este numărul de atribute de intrare [LC: d' este numărul de atribute de intrare din instanța de test];
- (P3) algoritmul Bayes Optimal poate produce eroare [la clasificare] din cauza faptului că ia decizia în sensul unui vot majoritar. Algoritmul Bayes Naiv are și el această „sursă” de eroare; în plus el poate produce eroare și din cauza faptului că lucrează cu presupuziția de independență condițională (care nu este satisfăcută în mod neapărat);
- (P4) acuratețea [la clasificare a] algoritmului Bayes Naiv scade atunci când unul sau mai multe atribute de intrare sunt duplicate: ex. 10.d, ex. 31.def;
- (P5) în cazul „învățării” unei funcții booleene (oarecare), rata medie a erorii produse la antrenare de către algoritmul Bayes Optimal (spre deosebire de Bayes Naiv!) este 0: ex. 33.d;
- (P6) complexitatea de eșantionare: de ordin logaritmic pentru Bayes Naiv și de ordin exponențial pentru Bayes Optimal: ex. 13;
- (P7) corespondența dintre regula de decizie a algoritmului Bayes Naiv (când toate variabilele de intrare sunt de tip Bernoulli) și regula de decizie a *regresiei logistice* și, în consecință, liniaritatea granițelor de decizie: ex. 14.
- *comparații* între algoritmul Bayes Naiv și alți algoritmi de clasificare automată: ex. 36, ex. 38.

³⁰⁹Atenție: Noi am folosit aici versiunea de bază a algoritmului Bayes Naiv; varianta “bag of words” (vedeți cartea Machine Learning a lui Tom Mitchell, pag. 183) diferă ușor de aceasta.

³¹⁰Numărul de parametri indicați în paranteze se referă la cazul când atât atributele de intrare cât și atributul de ieșire sunt de tip Bernoulli.

Algoritmii Bayes Naiv și Bayes Optimal cu variabile de intrare *de tip gaussian*

- Aplicare: G[N]B: ex. 15, ex. 40 și ex. 48; GJB: ex. 45, ex. 46 și ex. 47; GNB vs. GJB: ex. 21.
- Numărul de parametri necesari de estimat din date: ex. 42.
- Proprietăți:
 - (P0') presupunem că variabila de ieșire este booleană, i.e. ia valorile 0 sau 1; dacă pentru orice atribut de intrare, variabilele condiționale $X_i|Y = 0$ și $X_i|Y = 1$ au distribuții gaussiene de varianțe egale ($\sigma_{i0} = \sigma_{i1}$), atunci regula de decizie GNB (Gaussian Naive Bayes) este echivalentă (ca formă) cu cea a regresiei logistice, deci separarea realizată de către algoritmul GNB este de formă liniară: demonstrație: ex. 17; exemplificare în \mathbb{R} : ex. 40.a; exemplificare în \mathbb{R}^2 : ex. 41.c;
 - (P1') similar, presupunem că variabila de ieșire este booleană; dacă variabilele de intrare (notație: $X = (X_1, \dots, X_d)$) au distribuțiile [comune] condiționale $X|Y = 0$ și $X|Y = 1$ de tip gaussian [multidimensional], cu matricele de covarianță egale ($\Sigma_0 = \Sigma_1$), atunci regula de decizie a algoritmului "full" / Joint Gaussian Bayes este și ea echivalentă (ca formă) cu cea a regresiei logistice, deci separarea realizată este tot de formă liniară: ex. 18, ex. 20.a.i – ii;
 - (P2') când variabilele de intrare satisfac condiții mixte de tip (P0') sau (P7), atunci concluzia – separare liniară – se menține: ex. 44.b;
 - (P3') dacă în condițiile de la propozițiile (P0')-(P2') presupuziția de independență condițională este satisfăcută, iar numărul de instanțe de antrenament tinde la infinit, atunci rezultatul de clasificare obținut de către algoritmul Bayes Naiv gaussian este identic cu cel al regresiei logistice: ex. 22.a.
Atunci când presupuziția de independență condițională nu este satisfăcută, iar numărul de instanțe de antrenament tinde la infinit, regresia logistică se comportă mai bine decât algoritmul Bayes Naiv [gaussian]: ex. 22.b;
 - (P4') nu există o corespondență 1-la-1 între parametrii calculați de regresia logistică și între parametrii calculați de algoritmul Bayes Naiv [gaussian]: ex. 23.a;
 - (P5') atunci când varianțele distribuțiilor gaussiene care corespund probabilităților condiționale $P(X_i|Y = k)$ depind și de eticheta k , separatorul decizional determinat de algoritmul Bayes Naiv gaussian nu mai are (în mod necesar) forma regresiei logistice: ex. 40.bc, ex. 43;
similar, pentru algoritmul Bayes Optimal gaussian, atunci când $\Sigma_0 \neq \Sigma_1$, ecuația separatorului decizional este de ordin pătratic: ex. 19, ex. 20.a.iii – vi, ex. 42.e (separatorul decizional este un cerc), ex. 42.f (o hiperbolă), ex. 47 (o reuniune de două drepte);
 - (P6') parametrii algoritmilor Bayes Naiv gaussian și Bayes Optimal gaussian se pot estima în timp liniar în raport cu numărul de instanțe din setul de date de antrenament: ex. 41.bd, ex. 23.b.

2.1 Clasificare bayesiană — Probleme rezolvate

2.1.1 Ipoteze de probabilitate maximă a posteriori (MAP)

1. (Formula lui Bayes; medii ale unor variabile aleatoare discrete; [ipoteze MAP;] măsuri statistice folosite în clasificare)

■ • CMU, 2009 fall, Geoff Gordon, HW1, pr. 2

O anumită boală afectează una din 500 de persoane în medie. Identificarea persoanelor care au această boală se poate face cu ajutorul unei analize a sângelui, care costă 100 de dolari de persoană. Această analiză indică în cazul unui rezultat *pozitiv* faptul că se *poate* ca persoana respectivă să sufere de acea boală.

Testul / analiza are o *sensibilitate* (engl., *sensitivity* sau *recall*) perfectă — adică raportul dintre numărul instanțelor pozitive identificate ca atare de acel test și numărul total de instanțe pozitive este 1 —, ceea ce înseamnă că pentru orice persoană care are boala respectivă, rezultatul testului este pozitiv cu probabilitate de 100%. Pe de altă parte, testul are o *specificitate* — raportul dintre numărul instanțelor negative identificate ca atare de acel test și numărul total de instanțe negative — de 99%, adică o persoană care nu suferă de acea boală va avea cu probabilitate de 1% rezultatul testului pozitiv.

a. Se testează o persoană selectată în mod aleatoriu, iar rezultatul este pozitiv. Care este probabilitatea ca persoana respectivă să sufere de acea boală?

b. Există și un al doilea test, care costă 10.000 de dolari și are atât sensibilitatea cât și specificitatea de 100%. Dacă am cere ca toate persoanele detectate pozitiv la testul precedent să fie supuse acestui test mult mai scump, care ar fi costul mediu pentru testarea / analiza unui individ?

c. O companie farmaceutică încearcă să reducă prețul celui de-al doilea test (care este perfect), adică are atât *sensibilitatea* cât și *specificitatea* de 100%. Cât ar trebui să fie prețul acesta pentru ca primul test să nu mai fie necesar? (Adică, la ce preț va rezulta că este mai ieftin să se utilizeze doar testul al doilea, decât să se facă ambele teste, ca la punctul b?)

Răspuns:

Definim următoarele variabile aleatoare:

B : ia valoarea 1 / adevărat pentru persoanele care suferă de această boală și 0 / fals în caz contrar

T_1 : rezultatul primului test, care poate fi + (în caz de boală) sau –

T_2 : rezultatul celui de-al doilea test, care poate fi tot + sau –.

Folosind aceste variabile aleatoare, datele problemei se pot rescrie astfel:

$$\begin{aligned}
P(B) &= \frac{1}{500} \\
P(T_1 = + \mid B) &= 1 \\
P(T_1 = + \mid \bar{B}) &= \frac{1}{100} \\
P(T_2 = + \mid B) &= 1 \\
P(T_2 = + \mid \bar{B}) &= 0
\end{aligned}$$

a. Probabilitatea ca o persoană oarecare să sufere de boala respectivă, știind că rezultatul primului test este pozitiv, este $P(B \mid T_1 = +)$ și se calculează cu ajutorul formulei lui Bayes:

$$\begin{aligned}
P(B \mid T_1 = +) &= \frac{P(T_1 = + \mid B) \cdot P(B)}{P(T_1 = + \mid B) \cdot P(B) + P(T_1 = + \mid \bar{B}) \cdot P(\bar{B})} \\
&= \frac{1 \cdot \frac{1}{500}}{1 \cdot \frac{1}{500} + \frac{1}{100} \cdot \frac{499}{500}} = \frac{100}{599} \approx 0.1669
\end{aligned}$$

Observație: Remarcați faptul că $P(\bar{B} \mid T_1 = +) = 0.8331$. Este imediat care dintre cele două probabilități, $P(B \mid T_1 = +)$ și $P(\bar{B})$, este mai mare. În mod echivalent, pentru a stabili acest fapt era suficient să comparăm $P(B \mid T_1 = +)$ cu $1/2$, sau să stabilim care dintre produsele $P(T_1 = + \mid B) \cdot P(B)$ și $P(T_1 = + \mid \bar{B}) \cdot P(\bar{B})$ este mai mare. Aceste observații sunt utile pentru că ele fac legătura cu noțiunea de *ipoteză de probabilitate maximă a posteriori* (vedeți pr. 3).

b. Vom calcula costul mediu al testării unui individ folosind o nouă variabilă aleatoare, notată cu C , care reprezintă costul total de testare al unei persoane. Notând cu c_1 și c_2 costurile celor două teste, putem scrie:

$$C = \begin{cases} c_1 & \text{dacă persoana este testată doar cu primul test} \\ c_1 + c_2 & \text{dacă persoana este testată cu ambele teste} \end{cases}$$

O persoană este testată cu al doilea test doar dacă are rezultatul pozitiv la primul test, deci probabilitățile pentru variabila aleatoare C sunt:

$$P(C = c_1) = P(T_1 = -) \text{ și } P(C = c_1 + c_2) = P(T_1 = +)$$

Costul mediu cerut de problemă este media variabilei aleatoare C , deci se poate calcula astfel:

$$\begin{aligned}
E[C] &= c_1 \cdot P(C = c_1) + (c_1 + c_2) \cdot P(C = c_1 + c_2) \\
&= c_1 \cdot P(T_1 = -) + (c_1 + c_2) \cdot P(T_1 = +)
\end{aligned}$$

Știm că $P(T_1 = -) = 1 - P(T_1 = +)$, iar din formula probabilității totale avem:

$$\begin{aligned}
P(T_1 = +) &= P(T_1 = + \mid B) \cdot P(B) + P(T_1 = + \mid \bar{B}) \cdot P(\bar{B}) \\
&= 1 \cdot \frac{1}{500} + \frac{1}{100} \cdot \frac{499}{500} = \frac{599}{50000} = 0.01198
\end{aligned}$$

Așadar, vom obține:

$$\begin{aligned}
 E[C] &= c_1 \cdot (1 - P(T_1 = +)) + (c_1 + c_2) \cdot P(T_1 = +) \\
 &= c_1 - c_1 \cdot P(T_1 = +) + c_1 \cdot P(T_1 = +) + c_2 \cdot P(T_1 = +) \\
 &= c_1 + c_2 \cdot P(T_1 = +) \\
 &= 100 + 10000 \cdot \frac{599}{50000} \\
 &= 219.8 \approx 220\$
 \end{aligned}$$

c. Notăm cu c_n noul preț pentru al doilea test (T'_2). Acest preț trebuie să fie mai mic sau egal cu costul mediu de aplicare al ambelor teste (T_1 și T'_2), deci:

$$\begin{aligned}
 c_n \leq E[C'] &= c_1 \cdot P(C = c_1) + (c_1 + c_n) \cdot P(C = c_1 + c_n) \\
 &= c_1 + c_n \cdot P(T_1 = +) = 100 + c_n \cdot \frac{599}{50000}
 \end{aligned}$$

Rezolvând ecuația $c_n = 100 + c_n \cdot 0.01198$, obținem $c_n \approx 101.2125$.

Așadar, dacă al doilea test ar costa cel mult 101.21 dolari, atunci primul test n-ar mai fi necesar.

2.

(Formula lui Bayes;
[ipoteze MAP;] inferențe statistice)

prelucrare de Liviu Ciortuz, după

■ • CMU, 2009 fall, Geoff Gordon, HW1, pr. 1
("Monty's haunted house" problem)

Fără să știi cum s-a întâmplat, ai nimerit într-o casă plină de fantome. Acum ești blocat în fața unui perete care are 3 uși (pentru conveniență, le vom nota cu numerele 1, 2, 3). Apare o fantomă care îți spune: „Scăparea ta este să ieși din casă printr-una din aceste uși. Însă doar una dintre ele dă în afară; celelalte două sunt păzite de câte un monstru care te va ucide imediat dacă încerci să ieși pe acolo. Trebuie să alegi o ușă!”

Decizi să alegi la întâmplare una din cele trei uși, să zicem ușa 1.

Observație: Probabilitatea *a priori* ca ușa aceasta să dea în afară este $1/3$. (La fel este și în cazul celorlalte două uși.) Prin adăgarea altor informații — vedeți continuarea problemei — probabilitatea *a posteriori* a aceluiași eveniment se poate modifica, deci într-un caz fericit ea poate crește.

Într-adevăr, tocmai când ai pus mâna pe clanță ca să o deschizi, fantoma îți spune: „Așteaptă puțin! Îți voi mai da o informație.” Zicând aceasta, fantoma întredeschide o altă ușă (să zicem ușa 2) și îți arată că în spatele ei se află un monstru groaznic. Apoi fantoma te întreabă: „Vrei acum să alegi ultima ușă (adică ușa 3) sau consideri că este mai bine să rămâi la alegerea pe care ai făcut-o inițial?”

Fantoma te mai ajută spunându-ți că în alegerea ei, ea a urmat o *strategie* bazată pe două *principii*:

P1. După ce tu ai făcut alegerea inițială, fantoma a ales una din celelalte două uși, mai precis o ușă în spatele căruia se află un monstru. (Este evident că întotdeauna există o astfel de ușă, indiferent de alegerea ta.)

P2. În cazul în care care ambele uși între care are de ales fantoma au în spate câte un monstru, ea procedează după *una* din următoarele trei *variante* (pe care o alege a priori și ți-o aduce la cunoștință):

- a. Fantoma alege una din cele două uși cu probabilitate egală ($1/2$).
 - b. Dacă, așa cum a fost cazul mai sus, tu ai ales ușa 1, fantoma alege ușa 2 (cu probabilitate 1).
 - c. Dacă, tot așa, tu ai ales ușa 1, fantoma alege ușa 3 (cu probabilitate 1).
- (*Notă:* Alte variante nu interesează pentru rezolvarea care se cere mai jos.)

Se cere ca, pentru fiecare din aceste 3 variante în parte, să determini probabilitățile ca ieșirea să se afle în spatele ușii 1, respectiv în spatele ușii 3, dacă fantoma a deschis ușa 2.

Indicație: Pentru rezolvare, vom folosi două *variabile aleatoare*:

- O (de la engl. outside), cu valori în mulțimea $\{1, 2, 3\}$, indicând unde este ieșirea cea bună;
- G (de la engl. ghost), pentru a desemna ușa aleasă de fantomă.

Pentru fiecare din variantele de strategie ale fantomei (a, b, c), folosind formula lui Bayes se calculează $P(O = 1 \mid G = 2)$ și $P(O = 3 \mid G = 2)$.

Răspuns:

Pentru variabila aleatoare O avem probabilități *a priori* egale pentru toate ieșirile:

$$P(O = 1) = P(O = 2) = P(O = 3) = \frac{1}{3}. \quad (161)$$

Folosind formula lui Bayes combinată cu formula probabilității totale, probabilitățile (*a posteriori*) cerute vor putea fi calculate astfel:

$$P(O = 1 \mid G = 2) = \frac{P(G = 2 \mid O = 1) \cdot P(O = 1)}{P(G = 2 \mid O = 1) \cdot P(O = 1) + P(G = 2 \mid O = 3) \cdot P(O = 3)}$$

$$P(O = 3 \mid G = 2) = \frac{P(G = 2 \mid O = 3) \cdot P(O = 3)}{P(G = 2 \mid O = 3) \cdot P(O = 3) + P(G = 2 \mid O = 1) \cdot P(O = 1)}.$$

Remarcăm că la fiecare din numitorii celor două fracții de mai sus ar fi trebuit să mai scriem $P(G = 2 \mid O = 2) \cdot P(O = 2)$, însă acesta este 0 fiindcă $P(G = 2 \mid O = 2) = 0$, conform principiului P1. Deci $P(O = 3 \mid G = 2) = 1 - P(O = 1 \mid G = 2)$.

Pentru a exprima succint probabilitățile condiționate necesare pentru calculul probabilităților $P(O = 1 \mid G = 2)$ și $P(O = 3 \mid G = 2)$ folosind formulele de mai sus, completăm următorul tabel:³¹¹

G	O	$P(G \mid O)$		
		varianta a	varianta b	varianta c
2	1	$1/2$	1	0
3	1	$1/2$	0	1
2	2	0	0	0
3	2	1	1	1
2	3	1	1	1
3	3	0	0	0

³¹¹Se observă că pe fiecare dintre cele trei coloane din tabel, ultimele patru linii sunt identice, ceea ce este natural, conform principiilor P1 și P2.

În aceste condiții putem calcula ușor probabilitățile cerute în fiecare din cele trei variante:

Varianta a:

$$P(O = 1 | G = 2) = \frac{\frac{1}{2} \cdot \frac{1}{3}}{\frac{1}{2} \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} = \frac{1}{3} \text{ și } P(O = 3 | G = 2) = 1 - \frac{1}{3} = \frac{2}{3},$$

deci vei alege ușa a treia, fiindcă ea are probabilitatea mai mare de a te duce afară.

Varianta b:

$$P(O = 1 | G = 2) = \frac{1 \cdot \frac{1}{3}}{1 \cdot \frac{1}{3} + 1 \cdot \frac{1}{3}} = \frac{1}{2} \text{ și } P(O = 3 | G = 2) = 1 - \frac{1}{2} = \frac{1}{2},$$

deci vei alege la întâmplare oricare din cele două uși rămase (ușa 1 și ușa 3), ele având aceeași probabilitate de salvare.

Varianta c:

$$P(O = 1 | G = 2) = 0 \text{ și } P(O = 3 | G = 2) = 1 - 0 = 1,$$

deci vei alege ușa a treia, care este cu siguranță ieșirea cea bună.

Observații:

1. Echivalent, pentru a determina maximul dintre $P(O = 1 | G = 2)$ și $P(O = 3 | G = 2)$ ar fi fost suficient, conform formulei lui Bayes, să comparăm $P(G = 2 | O = 1) \cdot P(O = 1)$ și $P(G = 2 | O = 3) \cdot P(O = 3)$. Mai mult, ținând cont de relația (161), aceasta revine la a compara $P(G = 2 | O = 1)$ și $P(G = 2 | O = 3)$. Răspunsul poate fi citit imediat din tabelul de mai sus (vedeți prima linie și penultima linie): în cazul variantei *a*, $O = 3$ este varianta pentru care se obține probabilitatea [a posteriori] maximă. Altfel spus, pentru variantei *a*, $O = 3$ este *ipoteza de probabilitate maximă a posteriori* (engl., maximum a posteriori probability (MAP) hypothesis). Absolut similar se poate proceda și pentru variantele *b* și *c*.³¹²

2. Alternativ, pentru variantele *b* și *c* putem răspunde la întrebare și făcând un raționament care nu folosește formula lui Bayes. Ieșirea cea bună se poate găsi în spatele uneia dintre cele trei uși (vedeți figura alăturată). Cum fantoma a deschis deja ușa cu numărul 2, una dintre aceste situații (și anume, a doua din figură) este eliminată, fiindcă în spatele ei este un monstru. În continuare, putem raționa în felul următor:

1	2	3
	M	M
M		M
M	M	

Varianta b: Întrucât fantoma alege ușa 2 cu probabilitate 1, vom putea afirma că ambele variante – 1 și 3 – au probabilități egale, și anume $\frac{1}{2}$. Într-adevăr,

³¹² *Observație:* În cazuri precum cel de mai sus ($P(O = 1) = P(O = 2) = P(O = 3) = 1/3$), ipoteza MAP coincide cu *ipoteza de verosimilitate maximă* (engl., maximum likelihood (ML) hypothesis).

- fie ușa 1, cea aleasă de mine, dă înspre afară, iar atunci fantoma trebuie, conform principiului P2, să aleagă ușa 2;
- fie ușa 3 dă înspre afară, iar atunci, din nou conform principiului P2, fantoma trebuie să aleagă ușa 2;
- conform principiului P1, nu există o a treia posibilitate;
- nu dispun de alte informații pentru a decide între cele două situații de mai sus.

Varianta c: Știind că fantoma nu a deschis ușa 3 (care ar fi opțiunea corespunzătoare principiului P2), ci a ales ușa 2, înseamnă că nu a putut face altfel, deci ușa 3 reprezintă ieșirea.

3. (Formula lui Bayes; inferențe statistice; exemplificarea noțiunii de ipoteză / ipoteze MAP (“Maximum A posteriori Probability”))

■ ● ○ CMU, 2012 spring, Ziv Bar-Joseph, HW1, pr. 1.5

Mickey dă cu zarul de mai multe ori, sperând să obțină un 6. Secvența celor 10 rezultate obținute de el în urma acestor aruncări este următoarea: 1, 3, 4, 2, 3, 3, 2, 5, 1, 6. Mickey se întreabă dacă nu cumva zarul este măsluit (având tendința să producă de mai multe ori fața 3 decât ar fi normal dacă zarul ar fi perfect).

Concepeți o analiză simplă bazată pe teorema lui Bayes care să-i furnizeze lui Mickey informația care-l interesează: [în ce măsură putem spune că] zarul este măsluit [sau nu]?

Explicați raționamentul dumneavoastră.

Veți presupune că în general fiecare set de 100 de zaruri conține 5 zaruri măsluite (engl., unfair) în așa fel încât este favorizată apariția feței 3, rezultând următoarea distribuție de probabilitate a celor șase fețe, (1, 2, 3, 4, 5, 6): $P = [0.1, 0.1, 0.5, 0.1, 0.1, 0.1]$.

Răspuns:

Acesta este un exercițiu [tipic] de punere în evidență a noțiunii de ipoteză de *probabilitate maximă a posteriori* (engl., Maximum A posteriori Probability, MAP).

Vă reamintim definiția:

$$h_{MAP} = \operatorname{argmax}_{h \in H} P(h|D),$$

unde D este setul de date cu care se lucrează, iar H este mulțimea de ipoteze considerate.

În cazul nostru, $D = \{1, 3, 4, 2, 3, 3, 2, 5, 1, 6\}$, iar $H = \{FD, LD\}$, unde am notat cu FD zarul corect / cinstit (engl., fair dice) și cu LD zarul măsluit (engl., loaded dice).

Folosind formula lui Bayes, definiția de mai sus, poate fi „rafinată” astfel:

$$h_{MAP} \stackrel{\text{def.}}{=} \operatorname{argmax}_{h \in H} P(h|D) \stackrel{F.B.}{=} \operatorname{argmax}_{h \in H} \frac{P(D|h) \cdot P(h)}{P(D)} = \operatorname{argmax}_{h \in H} P(D|h) \cdot P(h),$$

ultima egalitate având loc datorită faptului că $P(D)$ este o cantitate pozitivă care nu depinde de h .³¹³

Așadar, în cazul de față, a determina ipoteza de probabilitate maximă a posteriori (h_{MAP}) revine la a determina maximul dintre două produse: $P(D|FD) \cdot P(FD)$ și $P(D|LD) \cdot P(LD)$.

Facem observația că pentru a calcula $P(D|FD)$ și $P(D|LD)$, vom ține cont de faptul că aruncările zarului au fost independente unele de altele. Prin urmare, notând $D = \{x_1, x_2, \dots, x_{10}\}$, vom putea scrie:

$$\begin{aligned} P(D|FD) \cdot P(FD) &= P(x_1, x_2, \dots, x_{10}|FD) \cdot P(FD) \stackrel{i.i.d.}{=} \left(\prod_{i=1}^{10} P(x_i|FD) \right) \cdot P(FD) \\ &= \left(\frac{1}{6} \right)^{10} \cdot \frac{95}{100} = \frac{1}{2^{10} \cdot 3^{10}} \cdot \frac{19}{20}. \end{aligned}$$

Similar,

$$\begin{aligned} P(D|LD) \cdot P(LD) &= P(x_1, x_2, \dots, x_{10}|LD) \cdot P(LD) \stackrel{i.i.d.}{=} \left(\prod_{i=1}^{10} P(x_i|LD) \right) \cdot P(LD) \\ &= \left(\frac{1}{10} \cdot \frac{1}{2} \cdot \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{10} \right) \cdot \frac{5}{100} \\ &= \frac{1}{10^7 \cdot 2^3} \cdot \frac{1}{20} = \frac{1}{2^{10} \cdot 5^7} \cdot \frac{1}{20}. \end{aligned}$$

Așadar, a vedea care dintre produsele $P(D|FD) \cdot P(FD)$ și $P(D|LD) \cdot P(LD)$ este mai mare revine la a compara fracțiile $\frac{19}{3^{10}}$ și $\frac{1}{5^7}$. În loc să facem ridicările la putere (3^{10} și 5^7), este mai convenabil să logaritmăm, folosind ca bază un număr supra-unitar:

$$\ln \frac{19}{3^{10}} = \ln 19 - \ln 3^{10} = \ln 19 - 10 \ln 3 = 2.9444 - 10.9861 = -8.0417$$

și

$$\ln \frac{1}{5^7} = -\ln 5^7 = -7 \ln 5 = -11.2661.$$

Conchidem că ipoteza de probabilitate maximă a posteriori este FD , deci că zarul lui Mickey *nu* este măsluit.

Observație: Se obișnuiește ca, în loc să se lucreze cu cele două produse ($P(D|FD) \cdot P(FD)$ și $P(D|LD) \cdot P(LD)$) în mod separat, așa cum am procedat noi mai sus, să se facă raportul lor,

$$\frac{P(D|LD) \cdot P(LD)}{P(D|FD) \cdot P(FD)} = \frac{P(LD|D)}{P(FD|D)}.$$

Acest raport se numește *raportul de șanse* (engl., odds ratio). (Dacă acest raport este supra-unitar, înseamnă că ipoteza LD este mai plauzibilă decât ipoteza FD .) Mai departe, aplicând logaritmul — fiindcă la calcule probabilitatea $P(D|\dots)$ se exprimă ca produs de n factori —, se obține ceea ce în limba

³¹³Folosind formula probabilității totale, atunci când H este o mulțime discretă, putem exprima $P(D)$ ca $\sum_{h' \in H} P(D|h') \cdot P(h')$. Ar fi util să calculați această probabilitate *a priori* în cazul datelor din acest exercițiu.

engleză se numește *log-odds ratio*. (Evident, dacă *log-odds ratio* are valoare pozitivă, ipoteza *LD* este mai plauzibilă.) Pe datele noastre,

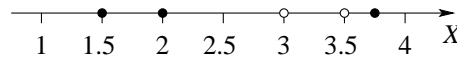
$$\ln \frac{P(LD|D)}{P(FD|D)} = -11.2661 - (-8.0417) = -3.2244 < 0.$$

În concluzie, folosind un astfel de raport (de „potrivire“), vom putea spune nu doar care ipoteză este mai plauzibilă, ci și în ce măsură acea ipoteză (în cazul nostru, *FD*) este mai plauzibilă decât cealaltă ipoteză (*LD*).

4. (Arbori ID3 cu decizii probabiliste, ca ipoteze ML și respectiv MAP)

■ CMU, 2009 spring, T. Mitchell, midterm, pr. 2

Se consideră următorul set de date de antrenament din spațiul real unidimensional:



Este vorba de 5 date caracterizate de atributul real X , împărțite în două clase: clasa 0 constituită din mulțimea $\{3, 3.5\}$, și clasa 1 – mulțimea $\{1.5, 2, 3.75\}$.

Pe acest set de date se va aplica algoritmul ID3 pentru construirea unor arbori de decizie. Deoarece atributul are valori reale, testele vor fi de forma $X > t$, unde t reprezintă o valoare-prag.

Se notează cu DT^* algoritmul care construiește arborele de decizie cu număr minim de noduri necesare pentru clasificarea perfectă a datelor de antrenament, și cu $DT1$ algoritmul care construiește un arbore de decizie cu un singur nod de test.

Indicație: Veți presupune că atunci când algoritmul de învățare de arbori de decizie găsește două praguri astfel încât testele $x < t_1$ și $x < t_2$ produc același câștig de informație, el (adică, algoritmul) va alege pragul cel mai din stânga. (Așadar, dacă $t_1 < t_2$, atunci el va alege testul $(x < t_1)$).

a. Care este eroarea la antrenare a lui DT^* pe datele specificate? Dar eroarea la cross-validare folosind metoda “Leave-One-Out” (CVLOO)?

b. Care este eroarea la antrenare a lui $DT1$ pe datele specificate? Dar eroarea la cross-validare folosind metoda “Leave-One-Out”?

În continuare se consideră o nouă clasă de arbori de decizie, care au *etichete probabiliste*. Fiecare nod frunză specifică probabilitatea fiecărei etichete posibile, probabilitate scrisă sub forma raportului dintre datele cu acea etichetă din nodul respectiv și toate datele din acel nod.

De *exemplu*, un arbore de decizie neavând niciun nod de test, construit pe datele specificate mai sus, clasifică astfel: $P(Y = 1) = 3/5$ și $P(Y = 0) = 2/5$. Un arbore de decizie cu un singur nod de test (în raport cu valoarea / pragul 2.5) conține probabilitățile: $P(Y = 1) = 1$ dacă $X \leq 2.5$, și $P(Y = 1) = 1/3$ dacă $X > 2.5$.

c. Pentru setul de date de mai sus, determinați arborele de decizie de tip ML (engl., maximum likelihood), adică acel arbore cu decizii probabiliste care maximizează *verosimilitatea* datelor de antrenament:

$$T_{ML} = \operatorname{argmax}_T P_T(D), \text{ unde}$$

$$P_T(D) \stackrel{\text{def.}}{=} P(D|T) \stackrel{\text{indep.}}{=} \prod_{i=1}^5 P(Y = y_i | X = x_i, T),$$

cu y_i eticheta / clasa instanței $x_i \in \{1.5, 2, 3, 3.5, 3.75\}$.

d. Se consideră o distribuție a priori $P(T)$ care penalizează numărul de teste / split-uri din arborele de decizie T , și anume:

$$P(T) \propto \left(\frac{1}{4}\right)^{\text{plits}(T)^2}$$

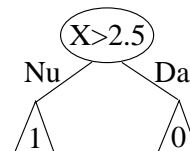
unde $\text{plits}(T)$ reprezintă numărul nodurilor de test din arborele T , iar simbolul \propto înseamnă „este proporțional cu”.

Pentru același set de date, folosind această distribuție a priori $P(T)$, găsiți arborele de decizie de tip MAP (engl., Maximum A posteriori Probability):

$$T_{MAP} = \operatorname{argmax}_T P_T(T|D)$$

Răspuns:

a. Dacă se aplică algoritmul DT1, întrucât câștigurile de informație calculate pentru testele $X > 2.5$ și $X > 3.625$ sunt 0.419 și respectiv 0.171, rezultatul este:



Eroarea la antrenare este $1/5$, deoarece punctul 3.75 este clasificat greșit.

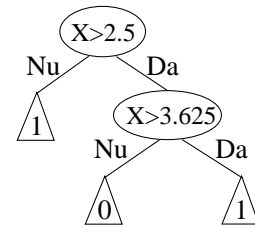
Eroarea la cross-validare cu metoda “Leave-One-Out” se determină astfel:

- $x_1 = 1.5$. Testul se face relativ la pragul 2.5 — ținând cont de *Indicația* din enunț —, deci punctul 1.5 este clasificat corect.
- $x_2 = 2$. Testul se face relativ la pragul 2.25 — ținând din nou cont de *Indicația* din enunț —, deci punctul $x_2 = 2$ este clasificat corect.
- $x_3 = 3$. Testul se face relativ la pragul 2.75, fiindcă se verifică imediat că IG-ul acestui prag este $1/2$, deci mai mare decât IG-ul pragului 3.625, și anume 0,311. Prin urmare, punctul $x_3 = 3$ va fi clasificat fie corect (în cazul în care se consideră că decizia arborelui DT1 este 0 pentru $X > 2.75$), fie eronat (în cazul în care se consideră că decizia arborelui DT1 este 1 pentru $X > 2.75$). Însă în al doilea caz arborele DT1 s-ar reduce la un singur nod (care este nod frunză), ceea ce este contrar definiției sale, așa că vom reține doar primul caz.

- $x_4 = 3.5$. Testul se face din nou relativ la pragul 2.5 — justificarea este similară cu cea de la punctul precedent —, iar punctul $x_4 = 3.5$ este clasificat fie corect (în cazul în care se consideră că decizia arborelui $DT1$ este 0 pentru $X > 2.5$), fie eronat (în cazul în care se consideră că decizia arborelui $DT1$ este 1 pentru $X > 2.5$). Însă (din nou!) în al doilea caz arborele $DT1$ s-ar reduce la un singur nod (frunză), așa că vom reține doar primul caz.
- $x_5 = 3.75$. Testul se face tot relativ la pragul 2.5, iar punctul $x_5 = 3.75$ este clasificat greșit.

Deci eroarea la cross-validare cu metoda “Leave-One-Out” este de $1/5$, punctul 3.75 fiind clasificat eronat.

b. Dacă se aplică algoritmul DT^* , rezultatul este:



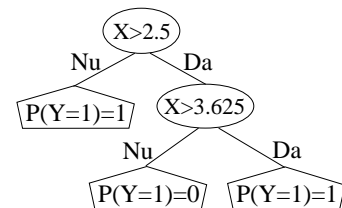
Eroarea la antrenare este bineînțeles 0, deoarece datele de antrenament nu conțin inconsistențe.

La cross-validare cu metoda “Leave-One-Out”, punctele care ar putea genera probleme sunt:

- $x_2 = 2$. Cele două teste se fac la pragurile 2.25 și la 3.625 (ordinea nu contează), deci punctul $x_2 = 2$ este corect clasificat.
- $x_3 = 3$. Primul test se face la pragul 2.75, deci punctul $x_3 = 3$ este corect clasificat.
- $x_4 = 3.5$. Al doilea test se face la pragul 3.375, deci punctul $x_4 = 3.5$ este clasificat greșit.
- $x_5 = 3.75$. Se face un singur test (la pragul 2.5), iar punctul $x_5 = 3.75$ este clasificat greșit.

Deci eroarea la cross-validare pentru arborele DT^* folosind metoda “Leave-One-Out” este $2/5$.

c. Un arbore de decizie care maximizează probabilitățile datelor de antrenament va fi unul care clasifică în mod perfect aceste date. Deci un arbore de decizie ML poate fi obținut din arborele construit de DT^* , extins cu etichete probabiliste în frunze, conform figurii alăturate.



d. Pentru a determina arborele de decizie MAP folosind distribuția a priori $P(T)$, va trebui să comparăm probabilitățile a posteriori ale celor 3 arbori de decizie posibili pe setul de date de antrenament. Vom scrie aceste probabilități a posteriori folosind formula lui Bayes:

$$P(T_j | D) = \frac{P(D | T_j) \cdot P(T_j)}{P(D)} = \frac{\prod_{i=1}^5 P(Y = y_i | T_j, X = x_i) \cdot P(T_j)}{P(D)}$$

Evident, la compararea propriu-zisă a probabilităților $P(T_j | D)$ cu $j = 0, 1, 2$, nu vom avea nevoie de numitorul $P(D)$. Așadar,

- pentru 0 noduri de test:

După cum s-a precizat și în enunț, acest arbore va avea $P(Y = 1) = 3/5$ și $P(Y = 0) = 2/5$. Deci

$$P(T_0 | D) \propto \left(\frac{3}{5}\right)^3 \cdot \left(\frac{2}{5}\right)^2 \cdot \left(\frac{1}{4}\right)^0 = \frac{3^3 \cdot 2^2}{5^5} = \frac{108}{3125} = 0.0336.$$

- pentru un nod de test:

Acest arbore va avea probabilitățile: $P(Y = 1) = 1$ dacă $X < 2.5$, și $P(Y = 1) = 1/3$ dacă $X \geq 2.5$. Prin urmare,

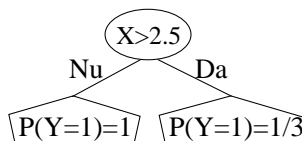
$$P(T_1 | D) \propto 1^2 \cdot \left(\frac{2}{3}\right)^2 \cdot \frac{1}{3} \cdot \left(\frac{1}{4}\right)^1 = \frac{1}{27} = 0.037$$

- pentru două noduri de test:

Este vorba de arborele construit la punctul c. Acesta clasifică perfect datele, dar

$$P(T_2) \propto \left(\frac{1}{4}\right)^4 \Rightarrow P(T_2 | D) \propto 1 \cdot \left(\frac{1}{4}\right)^4 = \frac{1}{256} = 0.0039$$

Probabilitatea a posteriori maximă o are arborele T_1 , deci acesta va fi arborele MAP. Reprezentarea grafică a acestui arbore este:



2.1.2 Algoritmii Bayes Naiv și Bayes Optimal

5.

(Algoritm Bayes Naiv: aplicare; comparație cu estimarea MLE)

prelucrare de Liviu Ciortuz, după

CMU, 2004 fall, T. Mitchell Z. Bar-Joseph, final exam, pr. 2

Fie setul de date alăturat, cu trei variabile booleene de intrare a , b , c și o variabilă booleană de ieșire K .

a. Estimați probabilitățile $P(K = 1 | a = 1, b = 1)$ și $P(K = 1 | a = 1, b = 1, c = 0)$ în sensul verosimilității maxime (engl., Maximum Likelihood Estimation, MLE).

b. Cum clasifică algoritmul Bayes Naiv instanța ($a = 1, b = 1$)?

Dar instanța ($a = 1, b = 1, c = 0$)?

a	b	c	K
1	0	1	1
1	1	1	1
0	1	1	0
1	1	0	0
1	0	1	0
0	0	0	1
0	0	0	1
0	0	1	0

Răspuns:

a. $P(K = 1 \mid a = 1, b = 1) = \frac{1}{2}$, fiindcă avem două instanțe în care ambele atribute a, b sunt adevărate, dintre care una este clasificată $K = 1$, iar cealaltă $K = 0$.

$P(K = 1 \mid a = 1, b = 1, c = 0) = 0$, fiindcă există o singură instanță cu $a = 1, b = 1, c = 0$ în setul de antrenament și ea este clasificată $K = 0$.

b. Cazul ($a = 1, b = 1$):

$$\begin{aligned}\hat{k}_{MAP} &= \operatorname{argmax}_{k \in \{0,1\}} P(K = k \mid a = 1, b = 1) = \\ &= \operatorname{argmax}_{k \in \{0,1\}} \frac{P(a = 1, b = 1 \mid K = k) \cdot P(K = k)}{P(a = 1, b = 1)} = \\ &= \operatorname{argmax}_{k \in \{0,1\}} P(a = 1, b = 1 \mid K = k) \cdot P(K = k) = \\ &= \operatorname{argmax}_{k \in \{0,1\}} P(a = 1 \mid K = k) \cdot P(b = 1 \mid K = k) \cdot P(K = k)\end{aligned}$$

Avem:

$$p_0 = P(a = 1 \mid K = 0) \cdot P(b = 1 \mid K = 0) \cdot P(K = 0) = \frac{2}{4} \cdot \frac{2}{4} \cdot \frac{4}{8} = \frac{1}{8}$$

$$p_1 = P(a = 1 \mid K = 1) \cdot P(b = 1 \mid K = 1) \cdot P(K = 1) = \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{4}{8} = \frac{1}{16}$$

Prin urmare, $p_0 > p_1$. Așadar, clasificatorul Bayes Naiv va prezice $K = 0$ pentru instanța ($a = 1, b = 1$) cu probabilitatea

$$\begin{aligned}P(K = 0 \mid a = 1, b = 1) &= \frac{P(a = 1, b = 1 \mid K = 0) \cdot P(K = 0)}{P(a = 1, b = 1 \mid K = 0) \cdot P(K = 0) + P(a = 1, b = 1 \mid K = 1) \cdot P(K = 1)} \\ &= \frac{p_0}{p_0 + p_1} = \frac{\frac{1}{8}}{\frac{1}{8} + \frac{1}{16}} = \frac{2}{3}\end{aligned}$$

Observație: Se constată că probabilitatea $P(K = 1 \mid a = 1, b = 1)$ estimată în sensul MLE (adică, $1/2$; vedeți punctul a) diferă de probabilitatea (*a posteriori*) cu care algoritmul Bayes Naiv clasifică instanța ($a = 1, b = 1$) ca aparținând clasei $K = 1$ (adică, $2/3$). Explicația rezidă în faptul că presupuziția de independență condițională asumată de către algoritmul Bayes Naiv nu este validă. Într-adevăr, $P(a = 0 \mid K = 1) = 2/4$ și $P(a = 0 \mid b = 1, K = 1) = 0$, deci $P(a = 0 \mid K = 1) \neq P(a = 0 \mid b = 1, K = 1)$.

Cazul ($a = 1, b = 1, c = 0$):

$$\begin{aligned}\hat{k}_{MAP} &= \operatorname{argmax}_{k \in \{0,1\}} P(K = k \mid a = 1, b = 1, c = 0) = \\ &= \operatorname{argmax}_{k \in \{0,1\}} \frac{P(a = 1, b = 1, c = 0 \mid K = k) \cdot P(K = k)}{P(a = 1, b = 1, c = 0)} = \\ &= \operatorname{argmax}_{k \in \{0,1\}} P(a = 1, b = 1, c = 0 \mid K = k) \cdot P(K = k) = \\ &= \operatorname{argmax}_{k \in \{0,1\}} P(a = 1 \mid K = k) \cdot P(b = 1 \mid K = k) \cdot P(c = 0 \mid K = k) \cdot P(K = k)\end{aligned}$$

Avem:

$$\begin{aligned}
 p_0 &= P(a=1|K=0) \cdot P(b=1|K=0) \cdot P(c=0|K=0) \cdot P(K=0) = \\
 &= \frac{2}{4} \cdot \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{4}{8} = \frac{1}{32} \\
 p_1 &= P(a=1|K=1) \cdot P(b=1|K=1) \cdot P(c=0|K=1) \cdot P(K=1) = \\
 &= \frac{2}{4} \cdot \frac{1}{4} \cdot \frac{2}{4} \cdot \frac{4}{8} = \frac{1}{32}
 \end{aligned}$$

Întrucât $p_0 = p_1$, clasificatorul Bayes Naiv va prezice $K = 0$ sau $K = 1$ cu aceeași probabilitate ($1/2$).

Notă: Observația de mai sus este valabilă și pentru cazul ($a = 1, b = 1, c = 0$).

6.

(Algoritmul Bayes Naiv: aplicație la filtrarea emailurilor spam)

■ • CMU, 2009 spring, Ziv Bar-Joseph, midterm, pr. 2

Circa $2/3$ dintre emailurile tale sunt spam, așadar te-ai decis să descarci de pe internet un filtru spam open-source care utilizează un clasificator Bayes Naiv.

Presupunem că ai strâns următoarele emailuri spam și non-spam (engl., regular), și de asemenea că doar trei cuvinte sunt discriminative pentru această clasificare, deci fiecare email este reprezentat ca un vector de 3 componente binare, fiecare dintre ele indicând dacă respectivul cuvânt este conținut (sau nu) în email.

'study'	'free'	'money'	Category	count
1	0	0	Regular	1
0	0	1	Regular	1
1	0	0	Regular	1
1	1	0	Regular	1
0	1	0	Spam	4
0	1	1	Spam	4

a. Descoperi că filtrul spam open-source folosește o probabilitate a priori $P(\text{spam}) = 0.1$. Explică în mod succint de ce crezi că această alegere este rezonabilă.

b. Calculează următorii parametri ai modelului prin metoda estimării de verosimilitate maximă (MLE), folosind netezire (engl., smoothing) de tip “add-one” (regula lui Laplace).

$$P(\text{study}|\text{spam}) =$$

$$P(\text{study}|\text{regular}) =$$

$$P(\text{free}|\text{spam}) =$$

$$P(\text{free}|\text{regular}) =$$

$$P(\text{money}|\text{spam}) =$$

$$P(\text{money}|\text{regular}) =$$

c. Folosind probabilitatea a priori și probabilitățile condiționate de mai sus, calculează probabilitatea ca mesajul $s = \text{"money for psychology study"}$ să fie spam, adică $P(\text{spam} | s)$.

d. Care ar trebui să fie valoarea probabilității a priori $P(\text{spam})$ în cazul în care dorim ca mesajul de mai sus să aibă aceeași probabilitate de fi spam respectiv non-spam (i.e., el va fi clasificat ca spam cu probabilitatea 0.5)?

Răspuns:

a. Diferența dintre $P_{MLE}(\text{Category} | \text{Spam}) = 2/3$ și probabilitatea a priori indicată în enunț (0.1) se explică prin faptul că se preferă trecerea prin filtru a unor emailuri spam, decât să fie marcate ca spam unele emailuri non-spam și astfel să nu ajungă în Inbox.

b. Dacă nu am folosi regula de tip "add-one" a lui Laplace (pentru „netezirea” probabilităților), parametrii modelului ar avea valorile (obținute prin metoda estimării de verosimilitate maximă – MLE) care apar mai jos în partea stângă. Folosind regula lui Laplace, parametrii primesc valorile indicate în partea dreaptă. Observați că aparițiile lui 2 de la numitorul fracțiilor corespund numărului de valori pentru fiecare dintre atributele / variabilele de intrare.

$$\begin{array}{ll}
 P(\text{study} | \text{spam}) = \frac{0}{8} = 0 & P(\text{study} | \text{spam}) \stackrel{\text{Laplace}}{=} \frac{0+1}{8+2} = \frac{1}{10} \\
 P(\text{study} | \text{regular}) = \frac{3}{4} & P(\text{study} | \text{regular}) \stackrel{\text{Laplace}}{=} \frac{3+1}{4+2} = \frac{2}{3} \\
 P(\text{free} | \text{spam}) = \frac{8}{8} = 1 & P(\text{free} | \text{spam}) \stackrel{\text{Laplace}}{=} \frac{8+1}{8+2} = \frac{9}{10} \\
 P(\text{free} | \text{regular}) = \frac{1}{4} & P(\text{free} | \text{regular}) \stackrel{\text{Laplace}}{=} \frac{1+1}{4+2} = \frac{1}{3} \\
 P(\text{money} | \text{spam}) = \frac{4}{8} = \frac{1}{2} & P(\text{money} | \text{spam}) \stackrel{\text{Laplace}}{=} \frac{4+1}{8+2} = \frac{1}{2} \\
 P(\text{money} | \text{regular}) = \frac{1}{4} & P(\text{money} | \text{regular}) \stackrel{\text{Laplace}}{=} \frac{1+1}{4+2} = \frac{1}{3}
 \end{array}$$

c. Avem mesajul $s = \text{"money for psychology study"}$, deci trebuie să calculăm $P(\text{spam} | s) = P(\text{spam} | \text{study}, \neg \text{free}, \text{money})$.

$$\begin{aligned}
 P(\text{spam} | s) &\stackrel{F. Bayes}{=} \\
 &= \frac{P(\text{study}, \neg \text{free}, \text{money} | \text{spam}) \cdot P(\text{spam})}{P(\text{study}, \neg \text{free}, \text{money} | \text{spam})P(\text{spam}) + P(\text{study}, \neg \text{free}, \text{money} | \text{reg})P(\text{reg})}
 \end{aligned}$$

Calculăm probabilitățile folosind ipoteza de independență condițională:

$$\begin{aligned}
 P(\text{study}, \neg \text{free}, \text{money} | \text{spam}) &= P(\text{study} | \text{spam}) \cdot P(\neg \text{free} | \text{spam}) \cdot P(\text{money} | \text{spam}) \\
 &= \frac{1}{10} \cdot \frac{1}{10} \cdot \frac{1}{2} = \frac{1}{200} \\
 P(\text{study}, \neg \text{free}, \text{money} | \text{reg}) &= P(\text{study} | \text{reg}) \cdot P(\neg \text{free} | \text{reg}) \cdot P(\text{money} | \text{reg}) \\
 &= \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{1}{3} = \frac{4}{27}
 \end{aligned}$$

$$\text{Înlocuind valorile în formulă, obținem: } P(\text{spam} | s) = \frac{\frac{1}{200} \cdot \frac{1}{10}}{\frac{1}{200} \cdot \frac{1}{10} + \frac{4}{27} \cdot \frac{9}{10}} \approx 0.0037$$

Aceasta este o probabilitate mică. Se observă însă că dacă nu am fi folosit regula lui Laplace, probabilitatea ca emailul s să fie spam ar fi fost 0. Aceasta se datorează faptului că niciunul dintre emailurile spam din datele de antrenament nu conține cuvântul *study*, care apare însă în emailul s .

d. Dacă notăm cu p probabilitatea a priori cerută, $P(\text{spam})$, știind că $P(\text{spam} | s) = 0.5$, putem scrie:

$$0.5 = \frac{\frac{1}{200} \cdot p}{\frac{1}{200} \cdot p + \frac{4}{27} \cdot (1-p)} \Leftrightarrow \frac{1}{2} = \frac{\frac{p}{200}}{\frac{p}{200} + \frac{4}{27} - \frac{4p}{27}} \Leftrightarrow \frac{2p}{200} = \frac{p}{200} + \frac{4}{27} - \frac{4p}{27}$$

$$\Leftrightarrow 54p = 27p + 800 - 800p \Leftrightarrow p = \frac{800}{827} \approx 0.9673$$

7. (Algoritmul Bayes Naiv și algoritmul Bayes Optimal; comparație relativ la numărului de parametri)

■ • CMU, 2008 fall, Eric Xing, HW1, pr. 2

Fie următoarea problemă de clasificare: X_1 și X_2 sunt variabile aleatoare observabile, Y este eticheta clasei asignate fiecărei instanțe observate, conform tabelului alăturat.

În acest exercițiu veți compara rezultatele care se obțin în urma antrenării pe acest set de date de către doi algoritmi de clasificare: Bayes Naiv și Bayes Comun (engl., Joint Bayes), care este numit adeseori și *Bayes Optimal* (engl., Optimal Bayes).

X_1	X_2	Y	Nr. apariții
0	0	0	2
0	0	1	18
1	0	0	4
1	0	1	1
0	1	0	4
0	1	1	1
1	1	0	2
1	1	1	18

Comentariu: Pentru cel de-al doilea algoritm, denumirea *Bayes Comun* se datorează faptului că acest clasificator lucrează cu distribuția comună a variabilelor / atributelor de intrare, în vreme ce denumirea *Bayes Optimal* este justificată de faptul că nicio altă distribuție probabilistă asupra datelor de intrare nu poate conduce la o eroare medie la clasificare mai mică decât cea obținută de către acest clasificator. (În cele ce urmează, vom opta pentru cea de-a doua denumire.)

a. Clasificați instanța $X_1 = 0, X_2 = 0$ folosind clasificatorul Bayes Naiv.

b. Clasificați instanța $X_1 = 0, X_2 = 0$ folosind clasificatorul Bayes Optimal.

c. Notăm cu P_{NB} și respectiv P_{JB} valoarea probabilității $P(Y = 1 | X_1 = 0, X_2 = 0)$ calculate pentru clasificatorul Bayes Naiv, respectiv pentru clasificatorul Bayes Optimal. De ce diferă cele două valori? *Sugestie:* Calculați $P(X_1, X_2 | Y)$.

X_1	X_2	Y	Nr. apariții
0	0	0	3
0	0	1	9
1	0	0	3
1	0	1	9
0	1	0	3
0	1	1	9
1	1	0	3
1	1	1	9

d. Care ar fi situația pentru P_{NB} și P_{JB} de la întrebarea precedentă în situația în care datele observate ar proveni din tabelul alăturat?

e. De câți parametri independenți (i.e., probabilități estimate) este nevoie în total pentru a construi clasificatorul Bayes Naiv? Dar în cazul clasificatorului Bayes Optimal?

Răspundeți la aceste întrebări și în cazul general, când se folosesc n variabile binare observate. Comentați rezultatul.

Răspuns:

a. Clasificatorul Bayes Naiv face predicția pentru valoarea lui Y după formula de mai jos:

$$\hat{y}_{NB} = \underset{y \in \{0,1\}}{\operatorname{argmax}} P(X_1 = 0 | Y = y) \cdot P(X_2 = 0 | Y = y) \cdot P(Y = y)$$

Avem:

$$\begin{aligned} p_0 &\stackrel{\text{not.}}{=} P(X_1 = 0 | Y = 0) \cdot P(X_2 = 0 | Y = 0) \cdot P(Y = 0) \\ &\stackrel{MLE}{=} \frac{6}{12} \cdot \frac{6}{12} \cdot \frac{12}{50} = \frac{3}{50} = \frac{6}{100} \\ p_1 &\stackrel{\text{not.}}{=} P(X_1 = 0 | Y = 1) \cdot P(X_2 = 0 | Y = 1) \cdot P(Y = 1) \\ &\stackrel{MLE}{=} \frac{19}{38} \cdot \frac{19}{38} \cdot \frac{38}{50} = \frac{19}{100} \end{aligned}$$

Întrucât $p_0 < p_1$, clasificatorul Bayes Naiv va prezice $Y = 1$ pentru instanța $(X_1 = 0, X_2 = 0)$.

b. Deoarece clasificatorul Bayes Optimal nu lucrează cu presupunerea de independență condițională a atributelor de intrare în raport cu atributul de ieșire, el va face predicția folosind formula de mai jos:

$$\hat{y}_{JB} = \underset{y \in \{0,1\}}{\operatorname{argmax}} P(X_1 = 0, X_2 = 0 | Y = y) \cdot P(Y = y)$$

Probabilitățile din formulă sunt, ca și în cazul clasificatorului Bayes Naiv, cele estimate din distribuția datelor de antrenament.

Așadar, avem:

$$\begin{aligned} p'_0 &\stackrel{\text{not.}}{=} P(X_1 = 0, X_2 = 0 | Y = 0) \cdot P(Y = 0) \stackrel{MLE}{=} \frac{2}{12} \cdot \frac{12}{50} = \frac{2}{50} \\ p'_1 &\stackrel{\text{not.}}{=} P(X_1 = 0, X_2 = 0 | Y = 1) \cdot P(Y = 1) \stackrel{MLE}{=} \frac{18}{38} \cdot \frac{38}{50} = \frac{18}{50} \end{aligned}$$

Fiindcă $p'_0 < p'_1$, clasificatorul Bayes Optimal va prezice tot $Y = 1$.

c. Vom calcula cele două valori pentru $P(Y = 1 | X_1 = 0, X_2 = 0)$:

$$\begin{aligned} P_{NB} &\stackrel{\text{not.}}{=} P(Y = 1 | X_1 = 0, X_2 = 0) \\ &\stackrel{F. Bayes}{=} \frac{P(X_1 = 0, X_2 = 0 | Y = 1) \cdot P(Y = 1)}{P(X_1 = 0, X_2 = 0 | Y = 1)P(Y = 1) + P(X_1 = 0, X_2 = 0 | Y = 0)P(Y = 0)} \\ &\stackrel{indep. cdt.}{=} \frac{P(X_1 = 0 | Y = 1) \cdot P(X_2 = 0 | Y = 1) \cdot P(Y = 1)}{P(X_1 = 0 | Y = 0) \cdot P(X_2 = 0 | Y = 0) \cdot P(Y = 0) + P(X_1 = 0 | Y = 1) \cdot P(X_2 = 0 | Y = 1) \cdot P(Y = 1)} \\ &= \frac{p_1}{p_0 + p_1} = \frac{\frac{19}{100}}{\frac{6}{100} + \frac{19}{100}} = \frac{19}{25} \end{aligned}$$

$$\begin{aligned}
P_{JB} &\stackrel{not.}{=} P(Y = 1 \mid X_1 = 0, X_2 = 0) \\
&\stackrel{F. Bayes}{=} \frac{P(X_1 = 0, X_2 = 0 \mid Y = 1) \cdot P(Y = 1)}{P(X_1 = 0, X_2 = 0 \mid Y = 1)P(Y = 1) + P(X_1 = 0, X_2 = 0 \mid Y = 0)P(Y = 0)} \\
&= \frac{p'_1}{p'_0 + p'_1} = \frac{\frac{18}{50}}{\frac{2}{50} + \frac{18}{50}} = \frac{18}{20}
\end{aligned}$$

Cele două valori diferă deoarece presupunerea de independență condițională a variabilelor X_1 și X_2 în raport cu Y făcută de clasificatorul Bayes Naiv este falsă. Acest lucru se poate vedea ușor din valorile estimate pentru $P(X_1 = 0, X_2 = 0 \mid Y = 0)$, $P(X_1 = 0 \mid Y = 0)$ și $P(X_2 = 0 \mid Y = 0)$:

$$\left. \begin{aligned}
P(X_1 = 0, X_2 = 0 \mid Y = 0) &\stackrel{MLE}{=} \frac{2}{12} \\
P(X_1 = 0 \mid Y = 0) \cdot P(X_2 = 0 \mid Y = 0) &\stackrel{MLE}{=} \frac{6}{12} \cdot \frac{6}{12} = \frac{1}{4}
\end{aligned} \right\} \Rightarrow$$

$$\begin{aligned}
&\Rightarrow P(X_1 = 0, X_2 = 0 \mid Y = 0) \neq P(X_1 = 0 \mid Y = 0) \cdot P(X_2 = 0 \mid Y = 0) \Rightarrow \\
&\Rightarrow P(X_1, X_2 \mid Y) \neq P(X_1 \mid Y) \cdot P(X_2 \mid Y)
\end{aligned}$$

Așadar, variabilele X_1 și X_2 nu sunt independente condițional în raport cu variabila de ieșire Y .

d. Vom calcula cele două valori pentru $P(Y = 1 \mid X_1 = 0, X_2 = 0)$ în cazul noilor date:

$$\begin{aligned}
P_{NB} &= P(Y = 1 \mid X_1 = 0, X_2 = 0) \\
&= \frac{P(X_1 = 0, X_2 = 0 \mid Y = 1) \cdot P(Y = 1)}{P(X_1 = 0, X_2 = 0 \mid Y = 1)P(Y = 1) + P(X_1 = 0, X_2 = 0 \mid Y = 0)P(Y = 0)} \\
&= \frac{P(X_1 = 0 \mid Y = 1) \cdot P(X_2 = 0 \mid Y = 1) \cdot P(Y = 1)}{P(X_1 = 0 \mid Y = 0) \cdot P(X_2 = 0 \mid Y = 0) \cdot P(Y = 0) + P(X_1 = 0 \mid Y = 1) \cdot P(X_2 = 0 \mid Y = 1) \cdot P(Y = 1)} \\
&= \frac{\frac{18}{36} \cdot \frac{18}{36} \cdot \frac{36}{48}}{\frac{6}{12} \cdot \frac{6}{12} \cdot \frac{12}{48} + \frac{18}{36} \cdot \frac{18}{36} \cdot \frac{36}{48}} = \frac{\frac{9}{48}}{\frac{3}{48} + \frac{9}{48}} = \frac{9}{12} = \frac{3}{4} \\
P_{JB} &= P(Y = 1 \mid X_1 = 0, X_2 = 0) \\
&= \frac{P(X_1 = 0, X_2 = 0 \mid Y = 1) \cdot P(Y = 1)}{P(X_1 = 0, X_2 = 0 \mid Y = 1)P(Y = 1) + P(X_1 = 0, X_2 = 0 \mid Y = 0)P(Y = 0)} \\
&= \frac{\frac{9}{36} \cdot \frac{36}{48}}{\frac{3}{12} \cdot \frac{12}{48} + \frac{9}{36} \cdot \frac{36}{48}} = \frac{\frac{9}{48}}{\frac{3}{48} + \frac{9}{48}} = \frac{9}{12} = \frac{3}{4}
\end{aligned}$$

Așadar, în acest caz avem $P_{NB} = P_{JB}$.

De fapt, se poate constata ușor că în cazul distribuției probabiliste date la acest punct al problemei se verifică independența condițională a variabilelor X_1 și X_2 în raport cu Y . Prin urmare, predicțiile făcute de cei doi clasificatori, Bayes Naiv și Bayes Optimal, vor coincide întotdeauna.

e. În contextul problemei noastre, clasificatorul Bayes Naiv are nevoie de estimările următoarelor probabilități:

$$\begin{aligned}
P(Y = 0) &\Rightarrow P(Y = 1) = 1 - P(Y = 0) \\
P(X_1 = 0 \mid Y = 0) &\Rightarrow P(X_1 = 1 \mid Y = 0) = 1 - P(X_1 = 0 \mid Y = 0) \\
P(X_1 = 0 \mid Y = 1) &\Rightarrow P(X_1 = 1 \mid Y = 1) = 1 - P(X_1 = 0 \mid Y = 1) \\
P(X_2 = 0 \mid Y = 0) &\Rightarrow P(X_2 = 1 \mid Y = 0) = 1 - P(X_2 = 0 \mid Y = 0) \\
P(X_2 = 0 \mid Y = 1) &\Rightarrow P(X_2 = 1 \mid Y = 1) = 1 - P(X_2 = 0 \mid Y = 1)
\end{aligned}$$

Avem nevoie, prin urmare, doar de 5 valori pentru a construi complet clasificatorul Bayes Naiv.

În cazul general, dacă avem n variabile de intrare, avem nevoie de estimări pentru probabilitățile $P(Y)$, $P(X_i \mid Y)$ și $P(X_i \mid \neg Y)$ pentru $i = \overline{1, n}$, deci $2n + 1$ valori.

Pentru clasificatorul Bayes Optimal avem nevoie de:

$$\begin{aligned}
P(Y = 0) & \quad P(Y = 1) = 1 - P(Y = 0) \\
P(X_1 = 0, X_2 = 0 \mid Y = 0) & \quad P(X_1 = 1, X_2 = 1 \mid Y = 0) \text{ se poate determina} \\
P(X_1 = 0, X_2 = 1 \mid Y = 0) & \quad \text{din celelalte 3 valori, așa cum vom arăta mai} \\
P(X_1 = 1, X_2 = 0 \mid Y = 0) & \quad \text{jos.} \\
P(X_1 = 0, X_2 = 0 \mid Y = 1) & \quad \text{Similar, } P(X_1 = 1, X_2 = 1 \mid Y = 1) \text{ se poate} \\
P(X_1 = 0, X_2 = 1 \mid Y = 1) & \quad \text{determina din celelalte 3 valori.} \\
P(X_1 = 1, X_2 = 0 \mid Y = 1) &
\end{aligned}$$

Notăm evenimentul $X_1 = 0$ cu A , $X_2 = 0$ cu B și $Y = 0$ cu C . Știm că:

$$\begin{aligned}
\Omega &= (A \wedge B) \vee (\neg A \wedge B) \vee (A \wedge \neg B) \vee (\neg A \wedge \neg B) \Rightarrow \\
\Omega \wedge C &= ((A \wedge B) \vee (\neg A \wedge B) \vee (A \wedge \neg B) \vee (\neg A \wedge \neg B)) \wedge C \Rightarrow \\
C &= ((A \wedge B) \wedge C) \vee ((\neg A \wedge B) \wedge C) \vee ((A \wedge \neg B) \wedge C) \vee ((\neg A \wedge \neg B) \wedge C)
\end{aligned}$$

De asemenea, deoarece toate evenimentele din partea dreaptă a egalității sunt disjuncte două câte două, putem scrie egalitatea de mai sus și cu probabilități:

$$\begin{aligned}
P(C) &= P((A \wedge B) \wedge C) + P((\neg A \wedge B) \wedge C) + P((A \wedge \neg B) \wedge C) + P((\neg A \wedge \neg B) \wedge C) \\
\Rightarrow 1 &= \frac{P((A \wedge B) \wedge C)}{P(C)} + \frac{P((\neg A \wedge B) \wedge C)}{P(C)} + \frac{P((A \wedge \neg B) \wedge C)}{P(C)} + \frac{P((\neg A \wedge \neg B) \wedge C)}{P(C)} \\
\Rightarrow 1 &= P(A, B \mid C) + P(\neg A, B \mid C) + P(A, \neg B \mid C) + P(\neg A, \neg B \mid C) \\
\Rightarrow P(\neg A, \neg B \mid C) &= 1 - (P(A, B \mid C) + P(\neg A, B \mid C) + P(A, \neg B \mid C))
\end{aligned}$$

Prin urmare, știind 3 valori o putem afla și pe a patra. La fel și pentru $\neg C$.

Pentru a avea un clasificator Bayes Optimal complet avem deci nevoie de 7 valori diferite.

În cazul general, pentru n variabile de intrare, avem nevoie de probabilitățile $P(Y)$, $P(\tilde{X}_1, \dots, \tilde{X}_n \mid Y)$ și $P(\tilde{X}_1, \dots, \tilde{X}_n \mid \neg Y)$, unde

$$\tilde{X}_i \in \{X_i, \neg X_i\} \quad \forall i \in \overline{1, n} \quad \text{și} \quad (\tilde{X}_1, \dots, \tilde{X}_n) \neq (\neg X_1, \dots, \neg X_n).$$

Avem, deci, $2(2^n - 1) + 1 = 2^{n+1} - 1$ valori.

Se observă că algoritmul Bayes Naiv folosește un număr liniar de parametri (în raport cu n , numărul de atribute de intrare), în vreme ce algoritmul Bayes Optimal folosește un număr exponențial de parametri (în raport cu același n).

8. (Calculul parametrilor pentru clasificatorul Bayes Naiv pornind de la distribuția comună a variabilelor; comparație între algoritmi Bayes Naiv și Bayes Optimal)

prelucrare de L. Ciortuz, după

■ • CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, midterm, pr. 2.1

Fie P o distribuție de probabilitate comună peste variabilele aleatoare booleene x_1, x_2 și y .

- a. Exprimați $P(y = 0 \mid x_1, x_2)$ în funcție de $P(x_1, x_2, y = 0)$ și $P(x_1, x_2, y = 1)$.

x_1	x_2	y	$P(x_1, x_2, y)$
0	0	0	0.15
0	0	1	0.25
0	1	0	0.05
0	1	1	0.08
1	0	0	0.10
1	0	1	0.02
1	1	0	0.20
1	1	1	0.15

În cele ce urmează vom considera distribuția P , definită conform tabelului alăturat.

- b. Pornind de la această distribuție comună, calculați probabilitățile necesare pentru clasificare bayesiană naivă. Completați tabelele de mai jos:

y	$P(y)$	$P(x_1 \mid y)$	$x_1 = 0$	$x_1 = 1$	$P(x_2 \mid y)$	$x_2 = 0$	$x_2 = 1$
$y = 0$		$y = 0$			$y = 0$		
$y = 1$		$y = 1$			$y = 1$		

- c. Cât este probabilitatea $P(y = 1 \mid x_1 = 1, x_2 = 0)$ calculată de către clasificatorul Bayes Naiv?

- d. Cât este probabilitatea $P(y = 1 \mid x_1 = 1, x_2 = 0)$ calculată de către clasificatorul Bayes Optimal? (Vă readucem aminte că acest algoritm este similar cu algoritmul Bayes Naiv, însă nu folosește presupuziția de independență condițională a atributelor.)

- e. Răspunsurile la precedentele două întrebări ar trebui să fie diferite. Care este motivul? Justificați.

Răspuns:

- a. Aplicând definiția probabilității condiționate și apoi proprietatea de aditivitate numărabilă din definiția funcției de probabilitate, obținem:

$$P(y = 0 \mid x_1, x_2) = \frac{P(x_1, x_2, y = 0)}{P(x_1, x_2)} = \frac{P(x_1, x_2, y = 0)}{P(x_1, x_2, y = 0) + P(x_1, x_2, y = 1)}$$

- b. Probabilitățile cerute în enunț sunt $P(x_1 \mid y)$, $P(x_2 \mid y)$ și $P(y)$.

$P(y)$ este o probabilitate marginală a distribuției comune, deci se calculează astfel:

$$\begin{aligned} P(y = 0) &= P(x_1 = 0, x_2 = 0, y = 0) + P(x_1 = 0, x_2 = 1, y = 0) + \\ &\quad + P(x_1 = 1, x_2 = 0, y = 0) + P(x_1 = 1, x_2 = 1, y = 0) \\ &= 0.15 + 0.05 + 0.1 + 0.2 = 0.5 \\ P(y = 1) &= 1 - P(y = 0) = 0.5 \end{aligned}$$

$P(x_1 | y)$ se calculează folosind din nou definiția probabilității condiționate și formula probabilității totale:

$$P(x_1 = 0 | y = 0) = \frac{P(x_1 = 0, y = 0)}{P(y = 0)} = \frac{P(x_1 = 0, y = 0)}{P(x_1 = 0, y = 0) + P(x_1 = 1, y = 0)}$$

Probabilitățile implicate în formulă sunt probabilități marginale ale distribuției comune și se calculează astfel:

$$P(x_1 = 0, y = 0) = P(x_1 = 0, x_2 = 0, y = 0) + P(x_1 = 0, x_2 = 1, y = 0) = 0.15 + 0.05 = 0.2$$

$$P(x_1 = 1, y = 0) = P(x_1 = 1, x_2 = 0, y = 0) + P(x_1 = 1, x_2 = 1, y = 0) = 0.1 + 0.2 = 0.3$$

Prin urmare,

$$P(x_1 = 0 | y = 0) = \frac{0.2}{0.2 + 0.3} = 0.4,$$

iar

$$P(x_1 = 1 | y = 0) = 1 - P(x_1 = 0 | y = 0) = 0.6$$

Analog se calculează și celelalte probabilități corespunzătoare lui $P(x_1 | y)$:

$$P(x_1 = 0 | y = 1) = \frac{P(x_1 = 0, y = 1)}{P(y = 1)} = \frac{0.25 + 0.08}{0.5} = 0.66$$

$$P(x_1 = 1 | y = 1) = 1 - P(x_1 = 0 | y = 1) = 0.34$$

$P(x_2 | y)$ se calculează în același mod.

Putem completa tabelele următoare cu valorile numerice ale probabilităților calculate la acest punct:

y	$P(y)$	$P(x_1 y)$	$x_1 = 0$	$x_1 = 1$	$P(x_2 y)$	$x_2 = 0$	$x_2 = 1$
$y = 0$	0.5	$y = 0$	0.40	0.60	$y = 0$	0.50	0.50
$y = 1$	0.5	$y = 1$	0.66	0.34	$y = 1$	0.54	0.46

c. Clasificatorul Bayes Naiv face presupunerea de *independență condițională* a variabilelor, deci:

$$\begin{aligned}
 P(y = 1 | x_1 = 1, x_2 = 0) &\stackrel{\text{Bayes}}{=} \frac{P(x_1 = 1, x_2 = 0 | y = 1) \cdot P(y = 1)}{P(x_1 = 1, x_2 = 0)} \\
 &= \frac{P(x_1 = 1, x_2 = 0 | y = 1) \cdot P(y = 1)}{P(x_1 = 1, x_2 = 0 | y = 1)P(y = 1) + P(x_1 = 1, x_2 = 0 | y = 0)P(y = 0)} \stackrel{\text{indep. cdt.}}{=} \\
 &= \frac{P(x_1 = 1 | y = 1) \cdot P(x_2 = 0 | y = 1) \cdot P(y = 1)}{P(x_1 = 1 | y = 1)P(x_2 = 0 | y = 1)P(y = 1) + P(x_1 = 1 | y = 0)P(x_2 = 0 | y = 0)P(y = 0)} \\
 &= \frac{0.34 \cdot 0.54 \cdot 0.5}{0.34 \cdot 0.54 \cdot 0.5 + 0.6 \cdot 0.5 \cdot 0.5} \approx 0.3796
 \end{aligned}$$

d. Clasificatorul Bayes Optimal nu face niciun fel de presupunere, deci folosind o formulă similară cu cea obținută la punctul a, vom avea:

$$\begin{aligned}
 P(y = 1 | x_1 = 1, x_2 = 0) &= \frac{P(x_1 = 1, x_2 = 0, y = 1)}{P(x_1 = 1, x_2 = 0, y = 1) + P(x_1 = 1, x_2 = 0, y = 0)} \\
 &= \frac{0.02}{0.02 + 0.1} = 0.1(6).
 \end{aligned}$$

e. Valorile calculate de clasificatorul Bayes Naiv și de clasificatorul Bayes Optimal pentru $P(y = 1 \mid x_1 = 1, x_2 = 0)$ sunt diferite deoarece presupunerea de independență condițională făcută de clasificatorul Bayes Naiv nu este adevărată. Într-adevăr, se observă că variabilele x_1 și x_2 nu sunt independente condițional în raport cu variabila y :

$$P(x_1 = 1, x_2 = 0 \mid y = 1) = \frac{P(x_1 = 1, x_2 = 0, y = 1)}{P(y = 1)} = \frac{0.02}{0.5} = 0.04$$

$$P(x_1 = 1 \mid y = 1) \cdot P(x_2 = 0 \mid y = 1) = 0.34 \cdot 0.54 = 0.1836 \neq 0.04$$

9. (Clasificare bayesiană: un caz particular)

CMU, 2010 spring, E. Xing, A. Singh, T. Mitchell, midterm, pr. 2.2

Se consideră variabilele aleatoare X_1 , X_2 , X_3 și X_4 . Aceste variabile sunt independente condițional două câte două în raport cu variabila Y , cu excepția perechii X_3 , X_4 . (Așadar, dacă am aplica algoritmul Bayes Naiv, acesta ar produce erori de clasificare.)

Cum am putea modifica regula de decizie a algoritmului Bayes Naiv pentru a ține cont de această particularitate a datelor?

Răspuns:

Întrucât are loc egalitatea

$$\operatorname{argmax}_y P(Y = y \mid X_1, X_2, X_3, X_4) =$$

$$\operatorname{argmax}_y P(X_1 \mid Y = y) \cdot P(X_2 \mid Y = y) \cdot P(X_3, X_4 \mid Y = y) \cdot P(Y = y),$$

urmează că regula de decizie a algoritmului Bayes Naiv în cazul dat este:

$$\operatorname{argmax}_y P_{MLE}(X_1 \mid Y = y) \cdot P_{MLE}(X_2 \mid Y = y) \cdot P_{MLE}(X_3, X_4 \mid Y = y) \cdot P_{MLE}(Y = y).$$

Pentru a justifica egalitatea de mai sus, se folosește regula de înlănțuire condițională:

$$P(A_1, A_2, A_3 \mid B) = P(A_3 \mid B) \cdot P(A_2 \mid A_3, B) \cdot P(A_1 \mid A_2, A_3, B),$$

care se demonstrează ușor. Varianta necondițională a regulii de înlănțuire este:

$$P(A_1, A_2, A_3) = P(A_3) \cdot P(A_2 \mid A_3) \cdot P(A_1 \mid A_2, A_3).$$

Apoi, regula de înlănțuire condițională se particularizează pentru evenimentele $A_1 = (X_1 = x_1)$, $A_2 = (X_2 = x_2)$, $A_3 = (X_3 = x_3, X_4 = x_4)$ și $B = (Y = y)$. În fine, se va ține cont că $P(X_1 = x_1 \mid X_2 = x_2, X_3 = x_3, X_4 = x_4, Y = y) = P(X_1 = x_1 \mid Y = y)$ și $P(X_2 = x_2 \mid X_3 = x_3, X_4 = x_4, Y = y) = P(X_2 = x_2 \mid Y = y)$ datorită proprietății de independență condițională din enunț.

10.

(Algoritmul Bayes Naiv:
calculul ratei medii a erorii la antrenare)

CMU, 2006 fall, T. Mitchell, E. Xing, midterm, pr. 6

Considerăm o problemă de clasificare binară în care fiecare exemplu X are două atribute binare $X_1, X_2 \in \{0, 1\}$ și eticheta $Y \in \{0, 1\}$. Vom presupune că X_1 și X_2 sunt independente condițional în raport cu Y ,³¹⁴ și că $P(Y = 0) = P(Y = 1) = 0,5$. De asemenea, probabilitățile condiționate sunt date în tabelele următoare:

$P(X_1 Y)$	$Y = 0$	$Y = 1$
$X_1 = 0$	0,7	0,2
$X_1 = 1$	0,3	0,8

$P(X_2 Y)$	$Y = 0$	$Y = 1$
$X_2 = 0$	0,9	0,5
$X_2 = 1$	0,1	0,5

a. Calculați predicția \hat{Y} făcută de clasificatorul Bayes Naiv pentru fiecare din cele patru combinații posibile de valori ale variabilelor X_1 și X_2 . Completați următorul tabel:

X_1	X_2	$P(X_1, X_2, Y = 0)$	$P(X_1, X_2, Y = 1)$	$\hat{Y}(X_1, X_2)$
0	0	$0,7 \cdot 0,9 \cdot 0,5$	$0,2 \cdot 0,5 \cdot 0,5$	0
0	1			
1	0			
1	1			

b. Presupunând că se folosesc o infinitate de exemple, calculați *rata medie a erorii* (engl., the expected error rate) făcute de acest clasificator *la antrenare*, folosind formula:

$$P(Y \neq \hat{Y}(X_1, X_2)) = \sum_{X_1=0}^1 \sum_{X_2=0}^1 P(X_1, X_2, Y = 1 - \hat{Y}(X_1, X_2))$$

c. Care din următorii doi clasificatori are rata medie a erorii la antrenare mai mică:

- clasificatorul Bayes Naiv care prezice Y având ca input doar X_1 ;
- clasificatorul Bayes Naiv care prezice Y având ca input doar X_2 .

d. Presupunem că definim un nou atribut X_3 , care este o copie a lui X_2 . Care este rata medie a erorii la antrenare a clasificatorului Bayes Naiv care prezice Y folosind toate atributele X_1, X_2, X_3 ? (Se presupune că datele de antrenament sunt în număr infinit.)

e. Explicați de ce rata erorii de la punctul d diferă față de cea de la punctul a.

Răspuns:

a. Predicția \hat{Y} făcută de clasificatorul Bayes Naiv pentru fiecare din cele patru combinații posibile de valori ale variabilelor X_1 și X_2 este înregistrată în ultima coloană a tabelului de mai jos. Calculele necesare pentru justificare

³¹⁴ Așadar, în această situație rezultatele algoritmilor Bayes Naiv și Bayes Optimal vor coincide.

- folosind formula $P(X_1, X_2, Y) = P(X_1, X_2|Y) \cdot P(Y) = P(X_1|Y) \cdot P(X_2|Y) \cdot P(Y)$
 — sunt conținute în coloanele a treia și a patra.

X_1	X_2	$P(X_1, X_2, Y = 0)$	$P(X_1, X_2, Y = 1)$	$\hat{Y}(X_1, X_2)$
0	0	$0,7 \cdot 0,9 \cdot 0,5 = 0,315$	$0,2 \cdot 0,5 \cdot 0,5 = \mathbf{0,05}$	0
0	1	$0,7 \cdot 0,1 \cdot 0,5 = \mathbf{0,035}$	$0,2 \cdot 0,5 \cdot 0,5 = 0,05$	1
1	0	$0,3 \cdot 0,9 \cdot 0,5 = \mathbf{0,135}$	$0,8 \cdot 0,5 \cdot 0,5 = 0,2$	1
1	1	$0,3 \cdot 0,1 \cdot 0,5 = \mathbf{0,015}$	$0,8 \cdot 0,5 \cdot 0,5 = 0,2$	1

Observație (1):

În acest tabel putem vedea valorile distribuției comune $P(X_1, X_2, Y)$. Spre deosebire de problema 8, unde distribuția comună era dată iar distribuțiile marginale condiționale erau calculate pornind de la aceasta, aici se procedează invers, ținând cont [și] de presupuziția de independență condițională.

b. Rata medie a erorii este:

$$\begin{aligned}
 P(Y \neq \hat{Y}(X_1, X_2)) &= \\
 &= \sum_{X_1=0}^1 \sum_{X_2=0}^1 P(X_1, X_2, Y = 1 - \hat{Y}(X_1, X_2)) \\
 &= P(X_1 = 0, X_2 = 0, Y = 1 - 0) + P(X_1 = 0, X_2 = 1, Y = 1 - 1) \\
 &\quad + P(X_1 = 1, X_2 = 0, Y = 1 - 1) + P(X_1 = 1, X_2 = 1, Y = 1 - 1) \\
 &= 0,05 + 0,035 + 0,135 + 0,015 = 0,235
 \end{aligned}$$

Pentru justificarea penultimei egalități, vedeți tabelul de la punctul precedent.

c. Făcând prezicerea doar cu X_1 ca atribut de intrare — folosind formula de multiplicare $P(X_1, Y) = P(X_1|Y) \cdot P(Y)$ —, obținem:

X_1	$P(X_1, Y = 0)$	$P(X_1, Y = 1)$	$\hat{Y}_1(X_1, X_2)$
0	$0,7 \cdot 0,5 = 0,35$	$0,2 \cdot 0,5 = \mathbf{0,1}$	0
1	$0,3 \cdot 0,5 = \mathbf{0,15}$	$0,8 \cdot 0,5 = 0,4$	1

Rata medie a erorii în acest caz va fi:

$$\begin{aligned}
 P(Y \neq \hat{Y}_1(X_1, X_2)) &= \\
 &= \sum_{X_1=0}^1 \sum_{X_2=0}^1 P(X_1, X_2, Y = 1 - \hat{Y}_1(X_1, X_2)) \\
 &= P(X_1 = 0, X_2 = 0, Y = 1 - 0) + P(X_1 = 0, X_2 = 1, Y = 1 - 0) \\
 &\quad + P(X_1 = 1, X_2 = 0, Y = 1 - 1) + P(X_1 = 1, X_2 = 1, Y = 1 - 1) \\
 &= 0,05 + 0,05 + 0,135 + 0,015 = 0,1 + 0,15 = 0,25
 \end{aligned}$$

Similar, dacă luăm în considerare doar variabila X_2 , avem:

X_2	$P(X_2, Y = 0)$	$P(X_2, Y = 1)$	$\hat{Y}_2(X_1, X_2)$
0	$0,9 \cdot 0,5 = 0,45$	$0,5 \cdot 0,5 = \mathbf{0,25}$	0
1	$0,1 \cdot 0,5 = \mathbf{0,05}$	$0,5 \cdot 0,5 = 0,25$	1

Acum, rata medie a erorii va fi:

$$\begin{aligned}
 P(Y \neq \hat{Y}_2(X_1, X_2)) &= \\
 &= \sum_{X_1=0}^1 \sum_{X_2=0}^1 P(X_1, X_2, Y = 1 - \hat{Y}_2(X_1, X_2)) \\
 &= P(X_1 = 0, X_2 = 0, Y = 1 - 0) + P(X_1 = 0, X_2 = 1, Y = 1 - 1) \\
 &\quad + P(X_1 = 1, X_2 = 0, Y = 1 - 0) + P(X_1 = 1, X_2 = 1, Y = 1 - 1) \\
 &= 0,05 + 0,035 + 0,2 + 0,015 = 0,25 + 0,05 = 0,3
 \end{aligned}$$

Prin urmare, rata medie a erorii la antrenare este mai mică pentru clasificatorul Bayes Naiv care prezice Y având ca input doar X_1 (decât pornind doar de la X_2).³¹⁵

Observație (2):

Atât în cazul lui X_1 cât și în cazul lui X_2 , rata medie a erorii (pentru algoritmul Bayes Naiv) putea fi calculată folosind direct probabilitățile din cele două tabele de mai sus. Justificarea ține de faptul că distribuțiile calculate de Bayes Naiv în aceste două tabele sunt distribuții marginale în raport cu distribuția comună (reală!) din tabelul de la punctul *a*. La punctul *d* veți vedea că acolo trebuie procedat altfel, fiindcă în cazul respectiv distribuția calculată de către Bayes Naiv *nu* mai coincide cu distribuția reală a datelor!

d. Clasificatorul Bayes Naiv care prezice valoarea lui Y în funcție de toate cele trei variabilele $X_1, X_2, X_3 = X_2$ va lua deciziile conform tabelului următor:³¹⁶

X_1	X_2	X_3	$P(X_1, X_2, X_3, Y = 0)$	$P(X_1, X_2, X_3, Y = 1)$	$\hat{Y}_3(X_1, X_2)$
0	0	0	$0,7 \cdot 0,9 \cdot 0,9 \cdot 0,5 = 0,2835$	$0,2 \cdot 0,5 \cdot 0,5 \cdot 0,5 = 0,025$	0
0	1	1	$0,7 \cdot 0,1 \cdot 0,1 \cdot 0,5 = 0,0035$	$0,2 \cdot 0,5 \cdot 0,5 \cdot 0,5 = 0,025$	1
1	0	0	$0,3 \cdot 0,9 \cdot 0,9 \cdot 0,5 = 0,1215$	$0,8 \cdot 0,5 \cdot 0,5 \cdot 0,5 = 0,1$	0
1	1	1	$0,3 \cdot 0,1 \cdot 0,1 \cdot 0,5 = 0,0015$	$0,8 \cdot 0,5 \cdot 0,5 \cdot 0,5 = 0,1$	1

Observație (3):

Este util să observați că distribuția de probabilitate comună calculată aici de către algoritmul Bayes Naiv nu mai coincide cu distribuția „reală” (vedeți tabelul de la punctul *a*). Rata erorii se calculează în raport cu distribuția „reală” a datelor, nu cu cea calculată de către Bayes Naiv (deși pentru a identifica situațiile în care $\hat{Y} \neq Y$ se folosește ultimul tabel)!

Rata medie a erorii produse la antrenare va fi:

$$\begin{aligned}
 P(Y \neq \hat{Y}_3(X_1, X_2)) &= \\
 &= \sum_{X_1=0}^1 \sum_{X_2=0}^1 P(X_1, X_2, Y = 1 - \hat{Y}_3(X_1, X_2)) \\
 &= P(X_1 = 0, X_2 = 0, Y = 1 - 0) + P(X_1 = 0, X_2 = 1, Y = 1 - 1) \\
 &\quad + P(X_1 = 1, X_2 = 0, Y = 1 - 0) + P(X_1 = 1, X_2 = 1, Y = 1 - 1) \\
 &= 0,05 + 0,035 + 0,2 + 0,015 = 0,3
 \end{aligned}$$

e. Diferența dintre cele două rate medii ale erorilor care au fost calculate la punctele *a* și *d* se datorează faptului că presupunerea de independență condițională a variabilelor nu este adevărată. Într-adevăr, X_2 nu este independent condițional față de X_3 deoarece cele două variabile au tot timpul valori identice.

Observații importante:

1. Este imediat că atunci când datele satisfac presupuziția de independență condițională, algoritmi Bayes Naiv și Bayes Optimal produc aceleași rezultate

³¹⁵Însă ambii clasificatori au rata medie a erorii mai mare decât clasificatorul Bayes Naiv care folosește atât X_1 cât și X_2 , ceea ce era de așteptat întrucât în condițiile date Bayes Naiv are același comportament ca și Bayes Optimal.

³¹⁶Este bine de observat că în această situație presupuziția de independență condițională este încălcată. Așadar, Bayes Naiv nu va mai furniza aceleași rezultate ca Bayes Optimal.

și au aceeași rată medie a erorilor.

2. Se poate arăta ușor că algoritmul Bayes Optimal nu produce în mod neapărat o rată medie a erorilor nulă, chiar dacă lucrează cu distribuția reală a datelor. De *exemplu*, în condițiile definite inițial de problema noastră, algoritmul Bayes Optimal produce rata medie 0.235, ca și algoritmul Bayes Naiv (vedeți punctul *b*). Erorile produse de către algoritmul Bayes Optimal se datorează faptului că el aplică operatorul $\arg \max$, echivalent cu luarea unui vot majoritar (impus deci minorității).

11. (Cât de naiv / prost este algoritmul Bayes Naiv?)

■ • CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW2, pr. 1.2

În mod evident, clasificatorul Bayes Naiv lucrează cu o presupuziție foarte restrictivă (engl., strong presupposition). Însă ne putem întreba dacă acest clasificator nu este totuși destul de folositor chiar și în cazul în care respectiva presupuziție nu este satisfăcută.

În consecință, în acest exercițiu ne propunem să folosim un exemplu simplu pentru a explora limitările algoritmului Bayes Naiv.

Fie X_1 și X_2 variabile aleatoare binare de tip Bernoulli de parametru $p = 0.5$, iar Y o funcție deterministă în raport cu valorile lui X_1 și X_2 , luând valori în mulțimea $\{1, 2\}$.

a. Definiți Y astfel încât (pe setul de date respectiv) algoritmul Bayes Naiv să aibă rata medie a erorii de 50%.

Pe acest caz, observați cum se corelează valorile lui X_1 și X_2 când valoarea lui Y este fixată. (Altfel spus, observați cât de (in)dependente sunt în acest context valorile lui X_1 și X_2 , dată fiind valoarea lui Y .)

X_1	X_2	Y
0	0	
0	1	
1	0	
1	1	

b. Există în total $2^4 = 16$ moduri în care poate fi definită funcția Y . Însă, datorită simetriei (relativ la valorile lui Y), problema se reduce la doar 4 cazuri,

X_1	X_2	Y	X_1	X_2	Y	X_1	X_2	Y	X_1	X_2	Y
0	0	1	0	0	1	0	0	1	0	0	1
0	1	1	0	1	1	0	1	2	0	1	2
1	0	1	1	0	1	1	0	1	1	0	2
1	1	1	1	1	2	1	1	2	1	1	1

dintre care un caz corespunde punctului *a* de mai sus. În fiecare din acele trei cazuri rămase după rezolvarea de la punctul *a*, arătați că rata erorii înregistrate de algoritmul Bayes Naiv este 0.

Răspuns:

a. Considerăm Y definit conform tabelului de mai jos.

Observație: Dacă se consideră valoarea lui Y fixată (fie 1, fie 2), atunci putem să stabilim o regulă astfel încât dacă îl cunoaștem pe X_1 să-l determinăm pe X_2 (și invers).³¹⁷ Altfel spus, X_1 este unic determinat de X_2 (și invers), dată fiind o valoare fixată a lui Y . Deci condiția de independență condițională este încălcată. Mai mult, în acest caz avem maximul posibil de „dependență” între cele două variabile (în raport cu Y).

X_1	X_2	Y
0	0	1
0	1	2
1	0	2
1	1	1

Dorim să calculăm rata erorii înregistrate de algoritmul Bayes Naiv pe datele din tabelul de mai sus. Bayes Naiv estimează valoarea lui Y astfel:

$$\hat{y} = \operatorname{argmax}_{y \in \{1,2\}} P(X_1 | Y = y) \cdot P(X_2 | Y = y) \cdot P(Y = y)$$

Pentru $X_1 = 0, X_2 = 0$, algoritmul compară următoarele două valori:

$$\begin{aligned} p_1 &= P(X_1 = 0 | Y = 1) \cdot P(X_2 = 0 | Y = 1) \cdot P(Y = 1) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8} \\ p_2 &= P(X_1 = 0 | Y = 2) \cdot P(X_2 = 0 | Y = 2) \cdot P(Y = 2) = \frac{1}{2} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{8} \end{aligned}$$

Cum $p_1 = p_2$, algoritmul va alege una dintre ele cu o probabilitate de 0.5. Deoarece valoarea lui Y din tabel este 1, înseamnă că algoritmul va alege greșit în 50% din cazuri.

Pentru celelalte 3 cazuri, $(X_1 = 0, X_2 = 1)$, $(X_1 = 1, X_2 = 0)$ și $(X_1 = 1, X_2 = 1)$, se observă ușor că se obțin de asemenea valori egale, iar algoritmul va alege pentru Y una dintre valorile 1 sau 2 cu o probabilitate de 0.5.

Deci pentru această definiție a lui Y rata erorii este de 50%.

b. Vom calcula rata erorii pentru fiecare dintre cele 3 moduri de definire a lui Y care nu a fost studiat.

Cazul 1:	<table><tr><th>X_1</th><th>X_2</th><th>Y</th></tr><tr><td>0</td><td>0</td><td>1</td></tr><tr><td>0</td><td>1</td><td>1</td></tr><tr><td>1</td><td>0</td><td>1</td></tr><tr><td>1</td><td>1</td><td>1</td></tr></table>	X_1	X_2	Y	0	0	1	0	1	1	1	0	1	1	1	1	Este similar cu cazul:	<table><tr><th>Y</th></tr><tr><td>2</td></tr><tr><td>2</td></tr><tr><td>2</td></tr><tr><td>2</td></tr></table>	Y	2	2	2	2
X_1	X_2	Y																					
0	0	1																					
0	1	1																					
1	0	1																					
1	1	1																					
Y																							
2																							
2																							
2																							
2																							

- Pentru $X_1 = 0, X_2 = 0$, algoritmul compară:

$$\begin{aligned} p_1 &= P(X_1 = 0 | Y = 1) \cdot P(X_2 = 0 | Y = 1) \cdot P(Y = 1) = \frac{2}{4} \cdot \frac{2}{4} \cdot 1 = \frac{1}{4} \\ p_2 &= P(X_1 = 0 | Y = 2) \cdot P(X_2 = 0 | Y = 2) \cdot P(Y = 2) = 0 \cdot 0 \cdot 0 = 0 \end{aligned}$$

Cum $p_1 > p_2$ algoritmul alege pentru Y valoarea 1, ceea ce este corect.

- Pentru celelalte 3 cazuri, $(X_1 = 0, X_2 = 1)$, $(X_1 = 1, X_2 = 0)$ și $(X_1 = 1, X_2 = 1)$, se observă că se obțin aceleași valori pentru p_1 și p_2 ca mai sus, deci algoritmul alege (în mod corect) pentru Y valoarea 1.

Așadar, am obținut că rata erorii este în acest caz 0.

³¹⁷Pentru $Y = 1$, regula este: X_2 are aceeași valoare ca și X_1 . Pentru $Y = 2$, regula este: X_1 și X_2 au valori complementare.

Cazul 2:	X_1	X_2	Y	Cazuri similare:	Y	Y	Y	Y	Y	Y	Y
	0	0	1		1	1	2	2	2	2	1
	0	1	1		1	2	1	2	2	1	2
	1	0	1		2	1	1	2	1	2	2
	1	1	2		1	1	1	1	2	2	2

- Pentru $X_1 = 0, X_2 = 0$:

$$\left. \begin{aligned} p_1 &= P(X_1 = 0 \mid Y = 1) \cdot P(X_2 = 0 \mid Y = 1) \cdot P(Y = 1) = \frac{2}{3} \cdot \frac{2}{3} \cdot \frac{3}{4} = \frac{1}{3} \\ p_2 &= P(X_1 = 0 \mid Y = 2) \cdot P(X_2 = 0 \mid Y = 2) \cdot P(Y = 2) = 0 \cdot 0 \cdot \frac{1}{4} = 0 \end{aligned} \right\} \Rightarrow \hat{y} = 1$$

- Pentru $X_1 = 0, X_2 = 1$:

$$\left. \begin{aligned} p_1 &= P(X_1 = 0 \mid Y = 1) \cdot P(X_2 = 1 \mid Y = 1) \cdot P(Y = 1) = \frac{2}{3} \cdot \frac{1}{3} \cdot \frac{3}{4} = \frac{1}{6} \\ p_2 &= P(X_1 = 0 \mid Y = 2) \cdot P(X_2 = 1 \mid Y = 2) \cdot P(Y = 2) = 0 \cdot 1 \cdot \frac{1}{4} = 0 \end{aligned} \right\} \Rightarrow \hat{y} = 1$$

- Pentru $X_1 = 1, X_2 = 0$:

$$\left. \begin{aligned} p_1 &= P(X_1 = 1 \mid Y = 1) \cdot P(X_2 = 0 \mid Y = 1) \cdot P(Y = 1) = \frac{1}{3} \cdot \frac{2}{3} \cdot \frac{3}{4} = \frac{1}{6} \\ p_2 &= P(X_1 = 1 \mid Y = 2) \cdot P(X_2 = 0 \mid Y = 2) \cdot P(Y = 2) = 1 \cdot 0 \cdot \frac{1}{4} = 0 \end{aligned} \right\} \Rightarrow \hat{y} = 1$$

- Pentru $X_1 = 1, X_2 = 1$:

$$\left. \begin{aligned} p_1 &= P(X_1 = 1 \mid Y = 1) \cdot P(X_2 = 1 \mid Y = 1) \cdot P(Y = 1) = \frac{1}{3} \cdot \frac{1}{3} \cdot \frac{3}{4} = \frac{1}{12} \\ p_2 &= P(X_1 = 1 \mid Y = 2) \cdot P(X_2 = 1 \mid Y = 2) \cdot P(Y = 2) = 1 \cdot 1 \cdot \frac{1}{4} = \frac{1}{4} \end{aligned} \right\} \Rightarrow \hat{y} = 2$$

Deci rata erorii este 0 pentru această definiție a lui Y .

Cazul 3:	X_1	X_2	Y	Cazuri similare:	Y	Y	Y
	0	0	1		2	1	2
	0	1	2		1	1	2
	1	0	1		2	2	1
	1	1	2		1	2	1

- Pentru $X_1 = 0, X_2 = 0$:

$$\left. \begin{aligned} p_1 &= \frac{1}{2} \cdot 1 \cdot \frac{1}{2} = \frac{1}{4} \\ p_2 &= \frac{1}{2} \cdot 0 \cdot \frac{1}{2} = 0 \end{aligned} \right\} \Rightarrow p_1 > p_2 \Rightarrow \hat{y} = 1 \text{ (corect)}$$

- Pentru $X_1 = 0, X_2 = 1$:

$$\left. \begin{aligned} p_1 &= \frac{1}{2} \cdot 0 \cdot \frac{1}{2} = 0 \\ p_2 &= \frac{1}{2} \cdot 1 \cdot \frac{1}{2} = \frac{1}{4} \end{aligned} \right\} \Rightarrow p_1 < p_2 \Rightarrow \hat{y} = 2 \text{ (corect)}$$

- Pentru $X_1 = 1, X_2 = 0$:

$$\left. \begin{aligned} p_1 &= \frac{1}{2} \cdot 1 \cdot \frac{1}{2} = \frac{1}{4} \\ p_2 &= \frac{1}{2} \cdot 0 \cdot \frac{1}{2} = 0 \end{aligned} \right\} \Rightarrow p_1 > p_2 \Rightarrow \hat{y} = 1 \text{ (corect)}$$

- Pentru $X_1 = 1, X_2 = 1$:

$$\left. \begin{aligned} p_1 &= \frac{1}{2} \cdot 0 \cdot \frac{1}{2} = 0 \\ p_2 &= \frac{1}{2} \cdot 1 \cdot \frac{1}{2} = \frac{1}{4} \end{aligned} \right\} \Rightarrow p_1 < p_2 \Rightarrow \hat{y} = 2 \text{ (corect)}$$

Prin urmare, rata erorii este 0 și în acest caz.

<i>Cazul 4:</i> (cel de la punctul a)	X_1	X_2	Y	Este similar cu cazul:	Y
	0	0	1		2
	0	1	2		1
	1	0	2		1
	1	1	1		2

În concluzie, doar pentru 2 moduri (cazul 4) de definire a lui Y rata erorii este de 50%; pentru celelalte 14 moduri (cazurile 1, 2, 3) rata erorii este 0.

12. (O reprezentare grafică a neconcordanței deciziilor luate de algoritmi Bayes Naiv și Bayes Optimal)
- CMU, 2009 fall, Geoff Gordon, HW4, pr. 1*
CMU, 2009 fall, Carlos Guestrin, HW1, pr. 4.1.5

Pentru un task de clasificare se consideră atributele X_1 , X_2 și X_3 și eticheta Y . Toate acestea sunt variabile aleatoare binare. X_1 și X_2 sunt independente condițional în raport cu Y , iar X_3 este o copie a lui X_2 (așadar, întotdeauna $X_2 = X_3$).

Sunt date următoarele probabilități condiționate:

$$\begin{aligned} P(X_1 = T \mid Y = T) &= p, & P(X_1 = T \mid Y = F) &= 1 - p, \\ P(X_2 = F \mid Y = T) &= q, & P(X_2 = F \mid Y = F) &= 1 - q, \\ P(Y = T) &= 0.5. \end{aligned}$$

Se dă instanța de test $X_1 = T, X_2 = X_3 = F$. Vrem să clasificăm această instanță, adică să prezicem valoarea lui Y pentru ea.

- a. Arătați că dacă se folosește clasificatorul Bayes Naiv, atunci eticheta pentru instanța aceasta de test este T — ceea ce revine la $P(Y = T \mid X_1 = T, X_2 = F, X_3 = F) \geq 0.5$ — dacă $p \geq \frac{(1-q)^2}{q^2 + (1-q)^2}$.

- b. Ce devine inegalitatea de la punctul precedent dacă în schimbul clasificatorului Bayes Naiv se utilizează clasificatorul Bayes Optimal?

- c. Desenați cele două curbe de decizie obținute la punctele a și b. Pe axa Ox marcați valorile lui q , iar pe axa Oy marcați valorile lui p . Atât p cât și q variază în intervalul $[0, 1]$. Indicați pe grafic zona în care clasificatorul Bayes Naiv produce un output (Y) diferit de cel al algoritmului Bayes Optimal.

Atenție! Acest exercițiu *nu* studiază în ce condiții cei doi algoritmi clasifică corect (ori dimpotrivă, eronat) instanța $X_1 = T, X_2 = X_3 = F$, ci doar când anume (în funcție de p și q) produc ei clasificări diferite pentru această instanță. Se va vedea (grafic!) că algoritmul Bayes Naiv nu se comportă deloc rău în comparație cu algoritmul Bayes Optimal.

Răspuns:

Se lucrează cu variabile aleatoare binare, iar pentru claritatea calculelor vom nota prin X faptul că valoarea variabilei aleatoare X este T , iar prin $\neg X$ faptul că $X = F$.

Folosind aceste notații, putem transcrie datele din enunț astfel:

$$\begin{aligned} P(X_1 | Y) &= p, & P(X_1 | \neg Y) &= 1 - p, \\ P(\neg X_2 | Y) &= q, & P(\neg X_2 | \neg Y) &= 1 - q, \\ P(Y) &= 0.5. \end{aligned}$$

a. Eticheta pentru instanța aceasta de test este T dacă $P(Y | X_1, \neg X_2, \neg X_3) \geq 0.5$. Vom calcula această probabilitate utilizând regula lui Bayes precum și ipoteza de independență condițională făcută de algoritmul Bayes Naiv:

$$\begin{aligned} P(Y | X_1, \neg X_2, \neg X_3) &\stackrel{\text{form. Bayes}}{=} \\ &= \frac{P(X_1, \neg X_2, \neg X_3 | Y)P(Y)}{P(X_1, \neg X_2, \neg X_3 | Y)P(Y) + P(X_1, \neg X_2, \neg X_3 | \neg Y)P(\neg Y)} \stackrel{\text{indep. cdt.}}{=} \\ &= \frac{P(X_1 | Y)P(\neg X_2 | Y)P(\neg X_3 | Y)P(Y)}{P(X_1 | Y)P(\neg X_2 | Y)P(\neg X_3 | Y)P(Y) + P(X_1 | \neg Y)P(\neg X_2 | \neg Y)P(\neg X_3 | \neg Y)P(\neg Y)} \\ &= \frac{p \cdot q \cdot q \cdot 0.5}{p \cdot q \cdot q \cdot 0.5 + (1 - p) \cdot (1 - q) \cdot (1 - q) \cdot 0.5} = \frac{pq^2}{pq^2 + (1 - p)(1 - q)^2} \end{aligned}$$

Deci eticheta este T dacă $\frac{pq^2}{pq^2 + (1 - p)(1 - q)^2} \geq 0.5$, ceea ce înseamnă că

$$\begin{aligned} pq^2 &\geq 0.5(pq^2 + (1 - p)(1 - q)^2) \Leftrightarrow pq^2 - 0.5pq^2 \geq 0.5(1 - p)(1 - q)^2 \\ \Leftrightarrow pq^2 &\geq (1 - q)^2 - p(1 - q)^2 \Leftrightarrow p(q^2 + (1 - q)^2) \geq (1 - q)^2 \Leftrightarrow p \geq \frac{(1 - q)^2}{q^2 + (1 - q)^2} \end{aligned}$$

b. Dacă în locul clasificatorului Bayes Naiv se folosește clasificatorul Bayes Optimal, nu se mai folosește presupunerea de independență condițională. În locul acesteia se folosesc informațiile furnizate în enunț, și anume că X_1 și X_2 sunt independente condițional în raport cu Y , iar X_3 este o copie a lui X_2 .

$$\begin{aligned} P(Y | X_1, \neg X_2, \neg X_3) &\stackrel{\text{form. Bayes}}{=} \\ &= \frac{P(X_1, \neg X_2, \neg X_3 | Y)P(Y)}{P(X_1, \neg X_2, \neg X_3 | Y)P(Y) + P(X_1, \neg X_2, \neg X_3 | \neg Y)P(\neg Y)} \\ &= \frac{P(X_1 | Y)P(\neg X_2, \neg X_3 | Y)P(Y)}{P(X_1 | Y)P(\neg X_2, \neg X_3 | Y)P(Y) + P(X_1 | \neg Y)P(\neg X_2, \neg X_3 | \neg Y)P(\neg Y)} \\ &= \frac{P(X_1 | Y)P(\neg X_2 | Y)P(Y)}{P(X_1 | Y)P(\neg X_2 | Y)P(Y) + P(X_1 | \neg Y)P(\neg X_2 | \neg Y)P(\neg Y)} \\ &= \frac{p \cdot q \cdot 0.5}{p \cdot q \cdot 0.5 + (1 - p) \cdot (1 - q) \cdot 0.5} = \frac{pq}{pq + (1 - p)(1 - q)} \end{aligned}$$

S-a folosit egalitatea

$$P(X_1, \neg X_2, \neg X_3 | Y) = P(X_1 | \neg X_2, \neg X_3, Y) \cdot P(\neg X_2, \neg X_3 | Y) = P(X_1 | Y) \cdot P(\neg X_2 | Y).$$

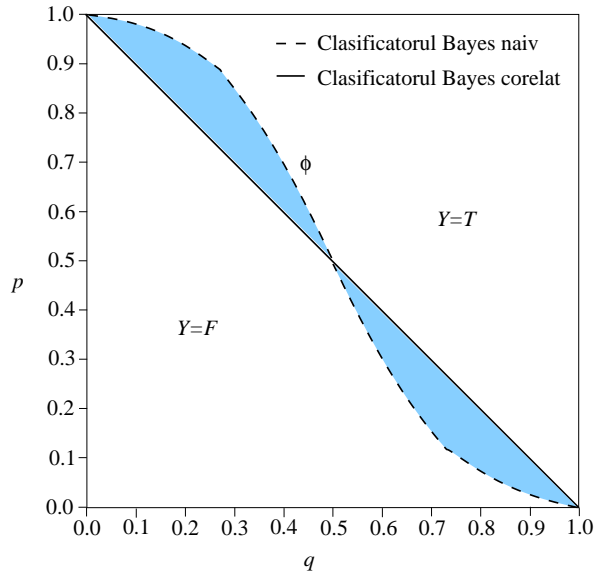
În acest caz eticheta este T dacă $\frac{pq}{pq + (1-p)(1-q)} \geq 0.5$, adică

$$\begin{aligned} pq &\geq 0.5(pq + (1-p)(1-q)) \Leftrightarrow pq - 0.5pq \geq 0.5(1-p-q+pq) \\ &\Leftrightarrow pq \geq 1-p-q+pq \Leftrightarrow p \geq 1-q. \end{aligned}$$

c. Dreapta de ecuație $p = 1 - q$ este ușor de reprezentat.

Notând $\phi(q) = \frac{(1-q)^2}{q^2 + (1-q)^2}$, se observă imediat că $\phi(q) + \phi(1-q) = 1 \Leftrightarrow \phi(1-q) = 1 - \phi(q)$. De aici, cu ajutorul unui raționament geometric simplu, se ajunge imediat la concluzia că graficul funcției ϕ este simetric față de punctul de coordonate $(1/2, 1/2)$. Așadar, va fi suficient să studiem graficul lui ϕ pe intervalul $[0, 1/2]$. Proprietățile funcției ϕ necesare elaborării graficului sunt ușor de studiat. Adicional, se poate arăta imediat că $\phi(q) \geq 1 - q$ pentru orice $q \in [0, 1/2]$ și $\phi(q) \leq 1 - q$ pentru orice $q \in [1/2, 1]$.

În figura de mai sus am reprezentat curbele $p = 1 - q$ și $p = \phi(q)$, precum și zonele de decizie pentru cei doi clasificatori obținuți la punctele a și b . Se observă ușor zona în care rezultatul produs de clasificatorul Bayes Naiv (Y) este diferit / „eronat“ în raport cu algoritmul Bayes Optimal.



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

■ • CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW2, pr. 1.1

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 variabila aleatoare comună $\bar{X} = (X_1, \dots, X_{d-1}, X_d)$, 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 $P(E_1 \cup \dots \cup E_n) \leq \sum_{i=1}^n P(E_i)$.³¹⁸
- $(1 - \frac{1}{k})^k \leq \frac{1}{e}$ 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 pentru fiecare din variabilele X_i ($i = \overline{1, d}$), atunci probabilitatea să nu fi întâlnit ambele valori pentru fiecare dintre aceste variabile 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

³¹⁸ Aceasta se numește *proprietatea de subaditivitate* a probabilităților.

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 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) .

Răspuns:

a. Dacă s-au făcut N observații și nu s-au întâlnit ambele valori ale variabilei X_i , înseamnă că ea are aceeași valoare în toate aceste observații (adică ea este fie 0 în toate observațiile, fie 1 în toate observațiile). Pentru fiecare dintre aceste două cazuri probabilitatea este $1/2^N$, deci:

$$\begin{aligned} P(\text{doar una dintre valorile variabilei } X_i \text{ a apărut în } N \text{ observații}) \\ = \left(\frac{1}{2}\right)^N + \left(\frac{1}{2}\right)^N = \frac{2}{2^N} = \frac{1}{2^{N-1}} \end{aligned}$$

b. Se cere să se calculeze probabilitatea să nu fi întâlnit ambele valori pentru fiecare dintre variabilele X_1, \dots, X_{d-1}, X_d . Pentru acesta vom folosi prima inegalitate din indicația dată în enunț:

$P(\text{nu toate valorile variabilelor } X_i, i = \overline{1, d}, \text{ au apărut în } N_{NB} \text{ observații})$

$$\begin{aligned} &\leq \sum_{i=1}^d P(\text{numai una dintre valorile variabilei } X_i \text{ a apărut în } N_{NB} \text{ observații}) \\ &= \sum_{i=1}^d \frac{1}{2^{N_{NB}-1}} = d \cdot \frac{1}{2^{N_{NB}-1}} = d \cdot \frac{1}{2^{1+\log_2 \frac{d}{\varepsilon}-1}} = d \cdot \frac{1}{2^{\log_2 \frac{d}{\varepsilon}}} = d \cdot \frac{1}{\frac{d}{\varepsilon}} = d \cdot \frac{\varepsilon}{d} = \varepsilon. \end{aligned}$$

c. Pentru algoritmul Bayes Optimal s-au făcut N observații. Trebuie să calculăm probabilitatea ca să nu se fi întâlnit niciodată instanța \bar{x} .

Cum există în total 2^d posibile observații, probabilitatea ca instanța \bar{x} să nu fie obținută la *una* dintre observații este $1 - \frac{1}{2^d}$. Observațiile fiind independente, probabilitatea ca \bar{x} să nu fie obținută *după* N observații este $\left(1 - \frac{1}{2^d}\right)^N$.

Așadar,

$$\begin{aligned} P(\text{instanța } \bar{x} \text{ n-a fost întâlnită în } N \text{ observații}) &= \left(1 - \frac{1}{2^d}\right)^N \\ &= \left[\left(1 - \frac{1}{2^d}\right)^{2^d} \right]^{N/2^d} \leq \left(\frac{1}{e}\right)^{N/2^d} = e^{-N/2^d} \end{aligned}$$

d. Vom calcula probabilitatea ca să nu se fi întâlnit toate instanțele variabilei \bar{X} utilizând din nou prima inegalitate din *indicație*:

$$\begin{aligned} &P(\text{nu toate instanțele variabilei } \bar{X} \text{ au fost întâlnite în } N_{JB} \text{ observații}) \\ &\leq \sum_{\bar{x}} P(\text{instanța } \bar{x} \text{ n-a fost întâlnită în } N_{JB} \text{ observații}) \end{aligned}$$

$$\leq \sum_{\bar{x}} e^{-N_{JB}/2^d} = 2^d \cdot e^{-N_{JB}/2^d} = 2^d \cdot e^{-\ln \frac{2^d}{\varepsilon}} = 2^d \cdot \frac{1}{e^{\ln \frac{2^d}{\varepsilon}}} = \frac{2^d}{\frac{2^d}{\varepsilon}} = \varepsilon.$$

e. Vom înlocui datele numerice în formulele $N_{NB} = 1 + \log_2 \frac{d}{\varepsilon}$, $N_{JB} = 2^d \ln \frac{2^d}{\varepsilon}$.

$$\begin{aligned} \varepsilon = 0.1, d = 2 &\Rightarrow \begin{cases} N_{NB} = 1 + \log_2 \frac{2}{0.1} = 1 + \log_2 20 \approx 5.32 \\ N_{JB} = 2^2 \cdot \ln \frac{2^2}{0.1} = 4 \cdot \ln 40 \approx 14.75 \end{cases} \\ \varepsilon = 0.1, d = 5 &\Rightarrow \begin{cases} N_{NB} = 1 + \log_2 \frac{5}{0.1} = 1 + \log_2 50 \approx 6.64 \\ N_{JB} = 2^5 \cdot \ln \frac{2^5}{0.1} = 32 \cdot \ln 320 \approx 184.58 \end{cases} \\ \varepsilon = 0.1, d = 10 &\Rightarrow \begin{cases} N_{NB} = 1 + \log_2 \frac{10}{0.1} = 1 + \log_2 100 \approx 7.64 \\ N_{JB} = 2^{10} \cdot \ln \frac{2^{10}}{0.1} = 1024 \cdot \ln 10240 \approx 9455.67 \end{cases} \end{aligned}$$

Acum se observă ușor [diferența dintre] numărul de date de antrenament necesare pentru cei doi algoritmi: de ordin logaritmic pentru Bayes Naiv, respectiv de ordin exponențial pentru Bayes Optimal.

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

■ □ ● ○ CMU, 2005 fall, T. Mitchell, A. Moore, HW2, pr. 2
 CMU, 2009 fall, Carlos Guestrin, HW1, pr. 4.1.2
 CMU, 2009 fall, Geoff Gordon, HW4, pr. 1.2-3
 CMU, 2012 fall, T. Mitchell, Z. Bar-Joseph, HW2, pr. 3.a

a. [Bayes Naiv și Regresia Logistică: relația dintre regulile de decizie]³¹⁹

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 $z = -(w_0 + w_1 X_1 + \dots + w_d X_d)$, cu parametrii $w_0, w_1, \dots, w_d \in \mathbb{R}$,³²⁰ adică

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

și, prin urmare

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

³¹⁹Pentru o introducere în chestiunea regresiei logistice, vedeți Tom Mitchell, *Generative and Discriminative Classifiers: Naive Bayes and Logistic Regression*, draft pentru un capitol suplimentar pentru o nouă ediție a cărții *Machine Learning*, 2016. Puteți vedea de asemenea problema 12 de la capitolul *Metode de regresie* din prezenta culegere.

³²⁰LC: 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 $w_0 + w_1 X_1 + \dots + w_d X_d = 0$.

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 parametrii distribuției $P(Y|X)$,³²¹ 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:

$$P(X_i|Y = 1) = \theta_{i1}^{X_i} (1 - \theta_{i1})^{(1-X_i)}, \quad (162)$$

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 (162) 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 $X_i = 0$.

b. [Relaxarea presupozității de independență condițională]³²²

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 corespondentul discriminativ al clasificatorului nostru generativ).

³²¹LC: 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. Vedeți de exemplu problemele 12 și 13 de la capitolul *Metode de regresie*.

³²²Vedeți de exemplu problema 9.

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:

$$P(X_1, X_2|Y = k) = (\beta_{11k})^{X_1 X_2} (\beta_{10k})^{X_1(1-X_2)} (\beta_{01k})^{(1-X_1)X_2} (\beta_{00k})^{(1-X_1)(1-X_2)} \quad (163)$$

pentru $k \in \{0, 1\}$, cu excepția următoarelor cazuri: **i.** $\beta_{11k} = 0$ și $X_1 X_2 = 0$, **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 $(1 - X_1)(1 - X_2) = 0$.

Răspuns:

a. Mai întâi vom rescrie probabilitatea $P(Y = 1|X = x)$ ca o fracție, folosind formula lui Bayes, compusă cu formula probabilității totale, apoi vom împărți atât numărătorul cât și numitorul fracției astfel obținute cu expresia de la numărător:³²³

$$\begin{aligned} P(Y = 1|X = x) &\stackrel{FB}{=} \frac{P(X = x|Y = 1) P(Y = 1)}{\sum_{y' \in \{0, 1\}} P(X = x|Y = y') P(Y = y')} \\ &= \frac{1}{1 + \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)}}. \end{aligned}$$

După aceea, folosind formula $e^{\ln a} = a$ (valabilă pentru orice $a > 0$), vom forța punerea fracției de mai sus într-o formă apropiată de cea a funcției sigmoideale ($\sigma(x) = 1/(1 + e^{-x})$):³²⁴

$$\begin{aligned} P(Y = 1|X = x) &= \frac{1}{1 + \exp\left(\ln \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)}\right)} \\ &= \frac{1}{1 + \exp\left(\ln \frac{P(X_1 = x_1, \dots, X_d = x_d|Y = 0)P(Y = 0)}{P(X_1 = x_1, \dots, X_d = x_d|Y = 1)P(Y = 1)}\right)}. \end{aligned}$$

Mai departe, ținând cont de presupuziția de independență condițională și de proprietățile funcției logaritm, obținem:³²⁵

$$P(Y = 1|X = x) = \frac{1}{1 + \exp\left(\ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \ln \frac{P(X_i = x_i|Y = 0)}{P(X_i = x_i|Y = 1)}\right)}.$$

Vom nota probabilitățile a priori $P(Y = 1)$ și $P(Y = 0)$ cu π și respectiv $1 - \pi$. Apoi, conform indicației 2, vom scrie $P(X_i|Y = 1)$ ca $\theta_{i1}^{X_i}(1 - \theta_{i1})^{(1-X_i)}$ și $P(X_i|Y = 0)$ ca $\theta_{i0}^{X_i}(1 - \theta_{i0})^{(1-X_i)}$. În consecință,

³²³Cu excepția cazului când $P(X = x|Y = 1) P(Y = 1) = 0$.

³²⁴Cu excepția cazului când $P(X = x|Y = 0) P(Y = 0) = 0$.

³²⁵Cu excepția cazurilor când $P(X = x_i|Y = 0) = 0$ sau $P(X = x_i|Y = 1) = 0$ pentru $i = 1, \dots, d$.

$$\begin{aligned}
P(Y = 1|X = x) &= \frac{1}{1 + \exp \left(\ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \ln \frac{\theta_{i0}^{X_i} (1-\theta_{i0})^{(1-X_i)}}{\theta_{i1}^{X_i} (1-\theta_{i1})^{(1-X_i)}} \right)} \\
&= \frac{1}{1 + \exp \left(\ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \left(X_i \ln \frac{\theta_{i0}}{\theta_{i1}} + (1-X_i) \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \right) \right)} \\
&= \frac{1}{1 + \exp \left(\ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} + \sum_{i=1}^d X_i \left(\ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \right) \right)}.
\end{aligned}$$

Pentru a pune această ultimă expresie sub forma dorită, adică $P(Y = 1|X = x) = 1/(1 + \exp(w_0 + \sum_{i=1}^d w_i X_i))$, vom alege valorile parametrilor w_i în mod natural:

$$w_0 = \ln \frac{1-\pi}{\pi} + \sum_{i=1}^d \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \quad \text{și} \quad w_i = \ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \text{ pentru } i = 1, \dots, d.$$

b. Vom începe ca și la punctul precedent, prin a pune probabilitatea $P(Y = 1|X = x)$ sub o formă apropiată de cea a funcției sigmoideale:³²⁶

$$\begin{aligned}
P(Y = 1|X) &\stackrel{FB}{=} \frac{P(X|Y = 1)P(Y = 1)}{P(X|Y = 1)P(Y = 1) + P(X|Y = 0)P(Y = 0)} \\
&= \frac{1}{1 + \frac{P(X|Y = 0)P(Y = 0)}{P(X|Y = 1)P(Y = 1)}} = \frac{1}{1 + \exp \left(\ln \frac{P(X|Y = 0)P(Y = 0)}{P(X|Y = 1)P(Y = 1)} \right)}.
\end{aligned}$$

Până aici nu avem încă nicio diferență în raport cu calculul de la punctul a. Însă acum vom ține cont că toate variabilele X_i ci $i = 1, \dots, d$ sunt independente condițional două câte două în raport cu Y , cu excepția perechii X_1, X_2 :³²⁷

$$\begin{aligned}
P(Y = 1|X) &= \frac{1}{1 + \exp \left(\ln \frac{P(X_1, X_2|Y = 0) \prod_{i=3}^d P(X_i|Y = 0)P(Y = 0)}{P(X_1, X_2|Y = 1) \prod_{i=3}^d P(X_i|Y = 1)P(Y = 1)} \right)} \\
&= \frac{1}{1 + \exp \left(\ln \frac{1-\pi}{\pi} + \sum_{i=3}^d \ln \frac{P(X_i|Y = 0)}{P(X_i|Y = 1)} + \ln \frac{P(X_1, X_2|Y = 0)}{P(X_1, X_2|Y = 1)} \right)}.
\end{aligned}$$

Ca și la punctul a, vom ține cont că

$$\ln \frac{P(X_i|Y = 0)}{P(X_i|Y = 1)} = \ln \frac{\theta_{i0}^{X_i} (1-\theta_{i0})^{(1-X_i)}}{\theta_{i1}^{X_i} (1-\theta_{i1})^{(1-X_i)}}$$

atunci când condițiile asociate cu relația (162) sunt satisfăcute. Mai departe, folosind indicația 4, vom putea înlocui $P(X_1, X_2|Y = 0)$ și $P(X_1, X_2|Y = 1)$ în

³²⁶Cu excepția cazurilor când $P(X|Y = 1)P(Y = 1) = 0$ sau $P(X|Y = 0)P(Y = 0) = 0$.

³²⁷Cu excepția cazurilor când $P(Y = 0) = 0$ sau $P(Y = 1) = 0$, respectiv $P(X_1 X_2|Y = 0) = 0$ sau $P(X_1 X_2|Y = 1) = 0$ și încă $P(X_i|Y = 0) = 0$ sau $P(X_i|Y = 1) = 0$ pentru $i = 3, \dots, d$.

funcție de β_{ijk} , obținând (atunci când condițiile asociate cu relația (163) sunt satisfăcute):

$$\ln \frac{P(X_1, X_2 | Y = 0)}{P(X_1, X_2 | Y = 1)} = \ln \frac{(\beta_{110})^{X_1 X_2} (\beta_{100})^{X_1(1-X_2)} (\beta_{010})^{(1-X_1)X_2} (\beta_{000})^{(1-X_1)(1-X_2)}}{(\beta_{111})^{X_1 X_2} (\beta_{101})^{X_1(1-X_2)} (\beta_{011})^{(1-X_1)X_2} (\beta_{001})^{(1-X_1)(1-X_2)}}.$$

Așadar, va rezulta:

$$P(Y = 1 | X) = \frac{1}{1 + \exp \left(w_0 + \sum_{i=3}^d w_i X_i + w_1 X_1 + w_2 X_2 + w_{1,2} X_1 X_2 \right)},$$

unde

$$\begin{aligned} w_0 &= \ln \frac{1 - \pi}{\pi} + \sum_{i=3}^d \ln \frac{1 - \theta_{i1}}{1 - \theta_{i0}} + \ln \frac{\beta_{000}}{\beta_{001}} \\ w_1 &= \ln \frac{\beta_{100}}{\beta_{101}} + \ln \frac{\beta_{001}}{\beta_{000}} \\ w_2 &= \ln \frac{\beta_{010}}{\beta_{011}} + \ln \frac{\beta_{001}}{\beta_{000}} \\ w_{1,2} &= \ln \frac{\beta_{110}}{\beta_{111}} + \ln \frac{\beta_{101}}{\beta_{100}} + \ln \frac{\beta_{011}}{\beta_{010}} + \ln \frac{\beta_{000}}{\beta_{001}} \\ w_i &= \ln \frac{\theta_{i0}}{\theta_{i1}} + \ln \frac{1 - \theta_{i1}}{1 - \theta_{i0}} \text{ pentru } i = 3, \dots, d. \end{aligned}$$

2.1.3 Clasificare bayesiană

[cu attribute de intrare] de tip gaussian

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

prelucrare de Liviu Ciortuz, după

■ • CMU, 2001 fall, Andrew Moore, midterm, pr. 3.a

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).³²⁸ Centralizați rezultatele, completând tabelul următor:

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

³²⁸Vedeți secțiunea corespunzătoare din capitolul de *Fundamente*, în speță (pentru acest caz) problemele 45.a și 46.a.

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.

Răspuns:

a. Pentru a estima mediile μ_A și μ_B , vom folosi formula demonstrată la problema 45.a de la capitolul de *Fundamente*:

$$\mu_{MLE} = \frac{\sum_{i=1}^n x_i}{n},$$

unde n este numărul instanțelor de antrenament. Așadar, $\mu_A = \frac{\sum_{i=1}^2 X_i}{2} = \frac{0+2}{2} = 1$, iar $\mu_B = \frac{\sum_{i=3}^7 X_i}{5} = \frac{3+4+5+6+7}{5} = 5$.

Similar, pentru calculul varianțelor σ_A și σ_B , vom folosi formula care a fost demonstrată la problema 46.a de la capitolul de *Fundamente*:

$$\sigma_{MLE}^2 = \frac{\sum_{i=1}^n (x_i - \mu_{MLE})^2}{n}.$$

Așadar, $\sigma_A^2 = \frac{1}{2}[(0-1)^2 + (2-1)^2] = 1$, iar $\sigma_B^2 = \frac{1}{5}[(3-5)^2 + (4-5)^2 + 0^2 + (6-5)^2 + (7-5)^2] = \frac{1}{5} \cdot 2 \cdot [4+1] = 2$.

Pentru calculul probabilităților $P(Y = A)$ și $P(Y = B)$ se folosește formula clasică (numărul de cazuri favorabile împărțit la numărul de cazuri posibile; vedeți problema 40 sau problema 113 de la capitolul de *Fundamente*), fiindcă Y este variabilă de tip Bernoulli. Așadar, $P(Y = A) = 2/7$ și $P(Y = B) = 5/7$.

Centralizând aceste estimări, obținem:

$\mu_A = 1$	$\sigma_A^2 = 1$	$P(Y = A) = 2/7$
$\mu_B = 5$	$\sigma_B^2 = 2$	$P(Y = B) = 5/7$

b. Folosind formula de multiplicare [a probabilităților], calculăm $p(X = 2, Y = A) = p(X = 2|Y = A) \cdot P(Y = A) = \frac{2\alpha}{7}$ și $p(X = 2, Y = B) = p(X = 2|Y = B) \cdot P(Y = B) = \frac{5\beta}{7}$.

Probabilitatea $p(X = 2)$ se poate obține aplicând formula probabilității totale:
 $p(X = 2) = p(X = 2|Y = A) \cdot P(Y = A) + p(X = 2|Y = B) \cdot P(Y = B) = \frac{1}{7}(2\alpha + 5\beta)$.

Probabilitatea condiționată $p(Y = A|X = 2)$ se calculează folosind definiția:

$$p(Y = A|X = 2) = \frac{p(Y = A, X = 2)}{p(X = 2)} = \frac{2\alpha}{2\alpha + 5\beta}.$$

c. Algoritmul Bayes [Naiv] gaussian va asocia punctului $X = 2$ eticheta $Y = A$ dacă $p(Y = A|X = 2) \geq p(Y = B|X = 2) \Leftrightarrow \frac{2}{7}\alpha \geq \frac{5}{7}\beta$.

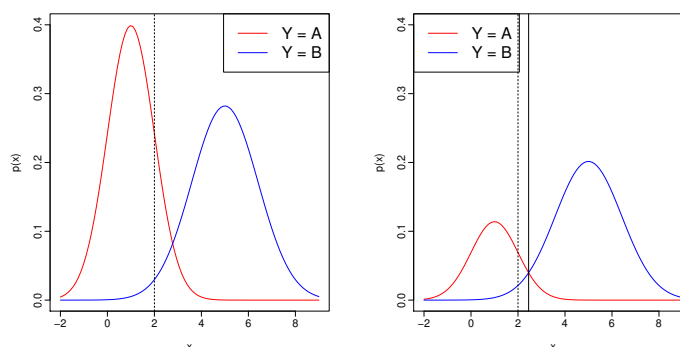
Folosind valorile estimate pentru parametrii μ_A, μ_B, σ_A și σ_B la punctul a , vom

$$\text{putea scrie: } \alpha = \frac{1}{\sqrt{2\pi}} e^{-\frac{(2-1)^2}{2}} = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}} = 0.24197 \text{ și } \beta = \frac{1}{\sqrt{2\pi} \cdot \sqrt{2}} e^{-\frac{(2-5)^2}{2 \cdot 2}} = \frac{1}{2\sqrt{\pi}} e^{-\frac{9}{4}} = 0.029732. \text{ Deci,}$$

$$\frac{2}{7}\alpha \geq \frac{5}{7}\beta \Leftrightarrow \frac{2}{7} \cdot 0.24197 \geq \frac{5}{7} \cdot 0.029732 \Leftrightarrow 0.06913 \geq 0.02123 \text{ (adevărat).}$$

Prin urmare, algoritmul Bayes [Naiv] gaussian va asocia punctului $X = 2$ eticheta $Y = A$, cu probabilitatea $\frac{0.06913}{0.06913 + 0.02123} = 0.76499$.

Observație: În cele două grafice alăturate am reprezentat $\mathcal{N}(\mu_A, \sigma_A^2)$ și $\mathcal{N}(\mu_B, \sigma_B^2)$, p.d.f.-urile gaussienele asociate claselor A și respectiv B (în partea stângă), precum și funcțiile obținute din aceste două p.d.f.-uri prin înmulțirea cu factorii de selecție $2/7$ și respectiv $5/7$ (în partea dreaptă).



Se poate constata relativ ușor că există două puncte de intersecție ($x_1 = -8.451$ și $x_2 = 2.451$) pentru graficele funcțiilor $\frac{2}{7}\mathcal{N}(\mu_A, \sigma_A^2)$ și $\frac{5}{7}\mathcal{N}(\mu_B, \sigma_B^2)$. Toate instanțele de test x situate între aceste puncte de intersecție ($x_1 < x < x_2$) vor aparține clasei A (acolo curba roșie este situată deasupra celei albastre). Instanțele situate fie la stânga lui x_1 fie la dreapta lui x_2 vor aparține clasei B (acolo curba albastră este situată deasupra celei roșii). *Separatorul decizional* este de tip pătratic, fiind constituit din punctele x_1 și x_2 .³²⁹

³²⁹LC: Mulțumesc studentului MSc Dinu Sergiu pentru această observație.

16.

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

□ Sergiu Dinu, Liviu Ciortuz, 2020

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, 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

$$\hat{Y}_{GB}(X = x) = A \Leftrightarrow p \cdot \mathcal{N}(x|\mu_A, \sigma_A^2) \geq (1 - p) \cdot \mathcal{N}(x|\mu_B, \sigma_B^2) \quad (164)$$

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 $\sigma_A = \sigma_B \stackrel{\text{not.}}{=} \sigma$ și $\mu_A \neq \mu_B$, regula de decizie (164) 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 (164) 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

$$(\mu_A - \mu_B)^2 + (\sigma_A^2 - \sigma_B^2) \ln \left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B} \right)^2 < 0. \quad (165)$$

Răspuns:

a. Putem scrie următoarea succesiune de echivalențe:

$$\begin{aligned}
 p \cdot \mathcal{N}(x|\mu_A, \sigma_A^2) &\geq (1-p) \cdot \mathcal{N}(x|\mu_B, \sigma_B^2) \Leftrightarrow \\
 p \cdot \frac{1}{\sqrt{2\pi}\sigma_A} \exp\left(-\frac{(x-\mu_A)^2}{2\sigma_A^2}\right) &\geq (1-p) \cdot \frac{1}{\sqrt{2\pi}\sigma_B} \exp\left(-\frac{(x-\mu_B)^2}{2\sigma_B^2}\right) \Leftrightarrow \\
 \exp\left(\frac{1}{2}\left[\left(\frac{x-\mu_B}{\sigma_B}\right)^2 - \left(\frac{x-\mu_A}{\sigma_A}\right)^2\right]\right) &\geq \frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B} \Leftrightarrow \\
 \left(\frac{x-\mu_B}{\sigma_B}\right)^2 - \left(\frac{x-\mu_A}{\sigma_A}\right)^2 &\geq \ln\left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B}\right) \Leftrightarrow \\
 (\sigma_A^2 - \sigma_B^2)x^2 + 2(\sigma_B^2\mu_A - \sigma_A^2\mu_B)x + \sigma_A^2\mu_B^2 - \sigma_B^2\mu_A^2 - \sigma_A^2\sigma_B^2 \ln\left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B}\right) &\geq 0.
 \end{aligned} \tag{166}$$

Membrul stâng al acestei ultime inegalități este o funcție polinomială de gradul al doilea. Prin urmare, discriminantul ei „pe jumătate“ (adică, $\Delta' \stackrel{\text{not.}}{=} \Delta/2$) este:

$$\Delta' = (\sigma_B^2\mu_A - \sigma_A^2\mu_B)^2 - (\sigma_A^2 - \sigma_B^2) \left[\sigma_A^2\mu_B^2 - \sigma_B^2\mu_A^2 - \sigma_A^2\sigma_B^2 \ln\left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B}\right) \right]$$

Însă,

$$\begin{aligned}
 &(\sigma_B^2\mu_A - \sigma_A^2\mu_B)^2 - (\sigma_A^2 - \sigma_B^2)(\sigma_A^2\mu_B^2 - \sigma_B^2\mu_A^2) \\
 &= \cancel{\sigma_B^4\mu_A^2} + \cancel{\sigma_A^4\mu_B^2} - 2\sigma_A^2\sigma_B^2\mu_A\mu_B - \cancel{\sigma_A^4\mu_B^2} + \sigma_A^2\sigma_B^2\mu_A^2 + \sigma_A^2\sigma_B^2\mu_B^2 - \cancel{\sigma_B^4\mu_A^2} \\
 &= \sigma_A^2\sigma_B^2(\mu_A - \mu_B)^2.
 \end{aligned}$$

Așadar,

$$\Delta' = \sigma_A^2\sigma_B^2 \left[(\mu_A - \mu_B)^2 + (\sigma_A^2 - \sigma_B^2) \ln\left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B}\right) \right],$$

ceea ce ne conduce imediat la concluzia din enunț.

Consecință: În acest caz (a), separatorul decizional este format din punctele x_1 și x_2 . Dacă $\sigma_A^2 > \sigma_B^2$, atunci zona de decizie corespunzătoare clasei A este $(-\infty, x_1] \cup [x_2, +\infty)$, iar zona de decizie corespunzătoare clasei B este intervalul (x_1, x_2) . Similar, dacă $\sigma_A^2 < \sigma_B^2$, atunci zona de decizie corespunzătoare clasei A este (x_1, x_2) , iar zona de decizie corespunzătoare clasei B este intervalul $(-\infty, x_1] \cup [x_2, +\infty)$.

b. Dacă $\sigma_A = \sigma_B = \sigma$, iar $\sigma > 0$, atunci este ușor de observat că inegalitatea (166) devine

$$(\mu_A - \mu_B)x + \mu_B^2 - \mu_A^2 - \sigma^2 \ln \frac{1-p}{p} \geq 0,$$

ceea ce fundamentează concluzia din enunț.

Consecință: În acest caz (b), separatorul decizional este punctul x_0 . Zonele de decizie se stabilesc imediat conform relațiilor din enunț.

c. Observăm că dacă vom considera cazul particular $\mu_A = \mu_B$, relația (165) va căpăta o formă mai simplă:

$$(\sigma_A^2 - \sigma_B^2) \ln\left(\frac{1-p}{p} \cdot \frac{\sigma_A}{\sigma_B}\right) < 0. \tag{167}$$

Această inegalitate este satisfăcută, de pildă atunci când σ_A și σ_B (care se consideră, prin convenție, strict pozitive) sunt în relația $\sigma_A > \sigma_B$, iar $\frac{1-p}{p} < \frac{\sigma_B}{\sigma_A}$. Se poate constata imediat că a doua inegalitate este echivalentă cu $p > \frac{\sigma_B}{\sigma_A + \sigma_B}$. Se observă că această ultimă fracție are valori în intervalul $(1/2, 1)$, fiindcă $\sigma_A > \sigma_B$.

Similar, dacă $\sigma_A < \sigma_B$, atunci pentru a satisface inegalitatea (167) se impune condiția $\frac{1-p}{p} > \frac{\sigma_B}{\sigma_A}$, care revine la $p < \frac{\sigma_A}{\sigma_A + \sigma_B}$. Această ultimă fracție ia valori în intervalul $(0, 1/2)$, întrucât $\sigma_A < \sigma_B$.

Consecință: În acest caz (c), nu există separator decizional. Una din cele două zone de decizie este vidă, iar cealaltă zonă coincide cu mulțimea tuturor numerelor reale, \mathbb{R} .

17. (Algoritmul Bayes Naiv gaussian:
deducerea [formeii liniare a] regulei de decizie în cazul
matricelor de covarianță diagonale și identice,
i.e., $\sigma_{i0} = \sigma_{i1}$, pentru $i = 1, \dots, d$)
■ ● ○ CMU, 2009 spring, Ziv Bar-Joseph, HW2, pr. 2

Considerăm un model de tip Bayes Naiv cu două clase ($Y \in \{0, 1\}$), definit peste spațiul real \mathbb{R}^d al atributelor de intrare X_1, \dots, X_d . Presupunem că în acest model distribuția [comună] condiționată $X|Y = 0$, unde $X = (X_1, \dots, X_d) \in \mathbb{R}^d$, poate fi definită ca un vector de distribuții gaussiene unidimensionale independente și-l vom desemna prin notația

$$\text{Gaussian}(\mu_0 = \langle \mu_{10}, \dots, \mu_{d0} \rangle, \sigma = \langle \sigma_1, \dots, \sigma_d \rangle)$$

și analog pentru $X|Y = 1$:

$$\text{Gaussian}(\mu_1 = \langle \mu_{11}, \dots, \mu_{d1} \rangle, \sigma = \langle \sigma_1, \dots, \sigma_d \rangle).$$

Observați că intrările X_1, \dots, X_d au — la condiționare în raport cu clasa — medii diferite dar varianțe (de fapt, matrice de covarianță, diagonale) identice pentru ambele clase.

În acest exercițiu vă vom arăta că, în modelul specificat mai sus, probabilitatea condiționată $P(Y = 1|X = x)$, unde $x = (x_1, \dots, x_d) \in \mathbb{R}^d$, se poate rescrie ca valoare a unei funcții sigmoideale / „logistice“, $f(x) = \frac{1}{1 + e^{-(w_0 + w \cdot x)}}$, cu parametrii $w_0 \in \mathbb{R}$ și $w = (w_1, \dots, w_d) \in \mathbb{R}^d$ aleși în mod convenabil.³³⁰

a. Folosiți regula lui Bayes (compusă cu formula probabilității totale) pentru a rescrie $P(Y = 1|X = x)$ sub forma unei fracții. Împărțiți atât numărătorul cât și numitorul fracției astfel obținute cu expresia de la numărător.

b. La punctul a ar fi trebuit să ajungeți la un rezultat de forma

$$\frac{1}{1 + f(x, X, Y)},$$

³³⁰În consecință, [se poate arăta imediat că] regula de decizie a clasificatorului Bayes Naiv pentru acest model este de tip *liniar*.

unde f este o anumită funcție având argumentele x, X și Y . Folosind formula $e^{\ln a} = a$ (valabilă pentru orice $a > 0$), putem forța punerea acestei fracții într-o formă apropiată de funcția sigmoidală:

$$\frac{1}{1 + e^{\ln f(x, X, Y)}}.$$

Vă cerem să rescrieți sub o astfel de formă rezultatul de la punctul a .

c. Explicitați presupuziția Bayes „naivă” în cadrul modelului probabilist dat. Apoi folosiți-o pentru a converti exponentul care apare în fracția rezultată la punctul b la o sumă de forma

$$\ln g(Y) + \sum_{i=1}^d \ln h(x_i, X_i, Y).$$

Precizați expresiile funcțiilor g și h .

d. Acum folosiți specificul modelului dat — și anume, de tip gaussian, având pentru componentele condiționate (și anume, variabilele aleatoare condiționate $X_i|Y = 1$, respectiv $X_i|Y = 0$, pentru $i = 1, \dots, n$) medii diferite dar varianțe egale —, pentru a aduce expresia de la punctul c la o formă mai convenabilă.

e. Rescriind $P(Y = 1|X = x)$ conform expresiilor obținute la punctele b și d , rezultatul ar trebui să semene cu un alt model pe care l-am întâlnit la curs. Care anume? Exprimați *parametrii* acelui model în raport cu $P(Y = 1)$, μ_{i0} , μ_{i1} și σ_i , cu $i = 1, \dots, d$.

Răspuns:

a. Procedând conform cerințelor, vom scrie:

$$\begin{aligned} P(Y = 1|X = x) &= \frac{P(X = x|Y = 1) P(Y = 1)}{\sum_{y \in \{0,1\}} P(X = x|Y = y) P(Y = y)} \\ &= \frac{1}{1 + \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)}}. \end{aligned}$$

Considerând $f(x, X, Y) = \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)}$, se observă că ultima expresie obținută are într-adevăr forma $\frac{1}{1 + f(x, X, Y)}$.

b. Folosind formula $a = e^{\ln a}$, obținem:

$$P(Y = 1|X = x) = \frac{1}{1 + \exp\left(\ln \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)}\right)}.$$

c. Conform presupuziției specifice algoritmului Naive Bayes, vom considera că atributele X_i sunt independente condițional două câte două în raport cu variabila de ieșire. Așadar, vom avea egalitățile următoare:

$$P(X = x|Y = 1) = \prod_{i=1}^d P(X_i = x_i|Y = 1)$$

$$P(X = x|Y = 0) = \prod_{i=1}^d P(X_i = x_i|Y = 0)$$

Prin urmare, vom scrie argumentul funcției $\exp(\cdot)$ din fracția de la punctul b astfel:

$$\ln \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)} = \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \ln \frac{P(X_i = x_i|Y = 0)}{P(X_i = x_i|Y = 1)}.$$

În consecință, funcțiile g și h care au fost cerute în enunț vor fi:

$$g(Y) = \frac{P(Y = 0)}{P(Y = 1)} \quad h(x_i, X_i, Y) = \frac{P(X_i = x_i|Y = 0)}{P(X_i = x_i|Y = 1)} \text{ pentru } i = 1, \dots, d.$$

d. Ținând cont de specificul gaussian al modelului din enunț, vom putea scrie rezultatul de la punctul b astfel:

$$\begin{aligned} & \ln \frac{P(X = x|Y = 0)P(Y = 0)}{P(X = x|Y = 1)P(Y = 1)} \\ &= \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \ln \left(\frac{\frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - \mu_{i0})^2}{2\sigma_i^2}\right)}{\frac{1}{\sqrt{2\pi}\sigma_i} \exp\left(-\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2}\right)} \right) \\ &= \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \left(\frac{(x_i - \mu_{i1})^2}{2\sigma_i^2} - \frac{(x_i - \mu_{i0})^2}{2\sigma_i^2} \right) \\ &= \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \frac{2x_i(\mu_{i0} - \mu_{i1}) + (\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2} \\ &= \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \left(\frac{x_i(\mu_{i0} - \mu_{i1})}{\sigma_i^2} + \frac{(\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2} \right) \\ &= \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} + \sum_{i=1}^d \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} x_i. \end{aligned}$$

e. Notând în expresia obținută la punctul d

$$w_0 = \ln \frac{P(Y = 0)}{P(Y = 1)} + \sum_{i=1}^d \frac{\mu_{i1}^2 - \mu_{i0}^2}{2\sigma_i^2} \quad \text{și} \quad w_i = \frac{\mu_{i0} - \mu_{i1}}{\sigma_i^2} \text{ pentru } i = 1, \dots, d$$

și apoi revenind la expresia de la punctul b, vom putea scrie:

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

Așadar, expresia probabilității a posteriori este exact cea a funcției sigmoidale, $\frac{1}{1 + e^{-x}}$. Rezultatul seamănă foarte bine cu modelul de *regresie logistică*.³³¹

³³¹Pentru o introducere la modelul regresiei logice, vedeți problema 12 de la capitolul *Metode de regresie*.

18.

(Algoritmul Bayes Optimal gaussian:
raportul față de regresia logistică
în cazul $\Sigma_0 = \Sigma_1$)

■ □ ● ○ CMU, 2011 spring, Tom Mitchell, HW2, pr. 2.2

În cele ce urmează, vom considera:

1. Y , o variabilă booleană care urmează o distribuție de tip Bernoulli, cu parametrul $\pi = P(Y = 1)$, ceea ce implică $P(Y = 0) = 1 - \pi$;
2. $X = (X_1, X_2, \dots, X_d)^\top$, un vector de variabile aleatoare care *nu* sunt independente condițional în raport cu variabila Y , probabilitatea [comună și] condițională $P(X|Y = k)$ urmând o distribuție gaussiană multidimensională, $\mathcal{N}(\mu_k, \Sigma)$, unde $k \in \{0, 1\}$.

Rețineți faptul că μ_k , media acestei distribuții multidimensionale, este un vector-coloană (altfel spus, o matrice $d \times 1$) și există două astfel de medii, câte una pentru fiecare dintre cele două valori ale variabilei Y . De asemenea, Σ , matricea de covarianță are dimensiunea $d \times d$, iar ea nu depinde de valorile variabilei Y .³³²

În rezolvarea problemei, veți folosi funcția de densitate [de probabilitate] a distribuției gaussiene multidimensionale în notație matriceală:³³³

$$\mathcal{N}(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right),$$

unde simbolul \top desemnează operația de transpunere a matricelor.

Răspundeți la următoarea întrebare:

Este oare distribuția $P(Y|X)$ corespunzătoare acestui clasificator de tip Bayes Optimal gaussian (nu naiv!) de aceeași formă cu cea a *regresiei logistice*?

Sugestie: Rafinați expresia distribuției probabiliste $P(Y|X)$.

Observație: La problema 17 am arătat că în cazul (particular!) în care intrările X_1, X_2, \dots, X_d sunt independente condițional în raport cu ieșirea Y — ceea ce, conform problemei 31 de la capitolul de *Fundamente*, este echivalent cu a spune că matricea Σ este diagonală —, răspunsul la întrebarea pusă în enunț este pozitiv.

Răspuns:

Vom demara calculele în maniera standard (adică, similar cu prima parte a rezolvării problemelor 14 și 17):

$$\begin{aligned} P(Y = 1|X) &= \frac{P(X|Y = 1) P(Y = 1)}{P(X|Y = 1) P(Y = 1) + P(X|Y = 0) P(Y = 0)} \\ &= \frac{1}{1 + \frac{P(Y = 0) P(X|Y = 0)}{P(Y = 1) P(X|Y = 1)}} = \frac{1}{1 + \exp \left(\ln \frac{P(Y = 0) P(X|Y = 0)}{P(Y = 1) P(X|Y = 1)} \right)} \\ &= \frac{1}{1 + \exp \left(\ln \frac{P(Y = 0)}{P(Y = 1)} + \ln \frac{P(X|Y = 0)}{P(X|Y = 1)} \right)}. \end{aligned}$$

³³²Altfel spus, presupunând că Σ_k , pentru $k \in \{0, 1\}$, desemnează matricea de covarianță a distribuției gaussiene multidimensionale $\mathcal{N}(\mu_k, \Sigma_k)$, atunci considerăm că $\Sigma_0 = \Sigma_1$.

³³³Vedeți problema 34 de la capitolul de *Fundamente*.

Acum ne vom concentra atenția asupra termenului $\ln \frac{P(X|Y=0)}{P(X|Y=1)}$, ținând cont de faptul că $X|Y=0 \sim \mathcal{N}(\mu_0, \Sigma)$ și $X|Y=1 \sim \mathcal{N}(\mu_1, \Sigma)$:

$$\begin{aligned}
 & \ln \frac{P(X|Y=0)}{P(X|Y=1)} \\
 &= \ln \frac{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}}}{\frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}}} + \ln \exp \left(\frac{1}{2} [(X - \mu_1)^\top \Sigma^{-1} (X - \mu_1) - (X - \mu_0)^\top \Sigma^{-1} (X - \mu_0)] \right) \\
 &= \frac{1}{2} [(X - \mu_1)^\top \Sigma^{-1} (X - \mu_1) - (X - \mu_0)^\top \Sigma^{-1} (X - \mu_0)] \\
 &= \frac{1}{2} [-X^\top \Sigma^{-1} \mu_1 - \mu_1^\top \Sigma^{-1} X + \mu_1^\top \Sigma^{-1} \mu_1 + X^\top \Sigma^{-1} \mu_0 + \mu_0^\top \Sigma^{-1} X - \mu_0^\top \Sigma^{-1} \mu_0] \\
 &= \frac{1}{2} [\mu_1^\top \Sigma^{-1} \mu_1 - \mu_0^\top \Sigma^{-1} \mu_0 + X^\top \Sigma^{-1} (\mu_0 - \mu_1) + (\mu_0^\top - \mu_1^\top) \Sigma^{-1} X] \\
 &= \frac{1}{2} \mu_1^\top \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_0^\top \Sigma^{-1} \mu_0 + (\mu_0 - \mu_1)^\top \Sigma^{-1} X.
 \end{aligned}$$

Remarcați faptul că $((\mu_0^\top - \mu_1^\top) \Sigma^{-1} X)^\top = ((\mu_0 - \mu_1)^\top \Sigma^{-1} X)^\top = X^\top (\Sigma^{-1})^\top (\mu_0 - \mu_1) = X^\top (\Sigma^\top)^{-1} (\mu_0 - \mu_1) = X^\top \Sigma^{-1} (\mu_0 - \mu_1)$, **întrucât matricea Σ^{-1} este simetrică** (ca urmare a faptului că Σ însăși, ca matrice de covarianță, este simetrică; vedeți problema 20 de la capitolul de *Fundamente*).

Prin urmare,

$$\begin{aligned}
 P(Y=1|X) &= \frac{1}{1 + \exp \left(\ln \frac{1-\pi}{\pi} + \frac{1}{2} \mu_1^\top \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_0^\top \Sigma^{-1} \mu_0 + (\mu_0 - \mu_1)^\top \Sigma^{-1} X \right)} \\
 &= \frac{1}{1 + \exp(w_0 + w^\top X)},
 \end{aligned}$$

cu $w_0 = \ln \frac{1-\pi}{\pi} + \frac{1}{2} \mu_1^\top \Sigma^{-1} \mu_1 - \frac{1}{2} \mu_0^\top \Sigma^{-1} \mu_0$ și $w = \Sigma^{-1} (\mu_0 - \mu_1)$. Evident, w_0 este un număr real (constant), iar w un vector-coloană (mai precis, o matrice de dimensiune $d \times 1$).

În concluzie, distribuția probabilistă $P(Y|X)$ are (și în acest caz!) aceeași formă cu cea din modelul regresiei logistice.

19.

(Clasificatorul Bayes Optimal Gaussian:
natura separatorului decizional în cazul în care $\Sigma_0 \neq \Sigma_1$)

■ ● ○ *Stanford, 2014 fall, Andrew Ng, midterm, pr. 2.b*

Fie setul de date de antrenament $\{(x_1, y_1), \dots, (x_n, y_n)\}$, cu x_i vectori-coloană din \mathbb{R}^d și $y_i \in \{0, 1\}$ pentru orice i . Algoritmul Bayes Optimal gaussian estimează următorii parametri: $\phi \in (0, 1)$, vectorii-coloană μ_0 și μ_1 din \mathbb{R}^d , precum și matricele de covarianță Σ_0 și Σ_1 de dimensiune $d \times d$, corespunzătoare următoarelor distribuții:

$$\begin{aligned}
 p(y) &= \phi^y (1 - \phi)^{1-y}, \text{ unde } \phi = p(y=1) \\
 p(x|y=0) &= \frac{1}{(2\pi)^{d/2}|\Sigma_0|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_0)^\top \Sigma_0^{-1} (x - \mu_0) \right)
 \end{aligned}$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2}|\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu_1)^\top \Sigma_1^{-1}(x-\mu_1)\right)$$

Vă readucem aminte *regula de decizie* a acestui clasificator:

$$y = 1 \text{ dacă } p(y=1|x) \geq p(y=0|x) \text{ și } y = 0 \text{ în caz contrar.} \quad (168)$$

Arătați că atunci când $\Sigma_0 \neq \Sigma_1$, *separatorul decizional* determinat de algoritmul Bayes Optimal gaussian este de ordin pătratic, adică inegalitatea dintre probabilitățile condiționate a posteriori, $p(y=1|x) \geq p(y=0|x)$, este echivalentă cu o inegalitate numerică de forma

$$x^\top Ax + B^\top x + C \geq 0, \quad (169)$$

unde A este o anumită matrice de dimensiune $d \times d$, cu $A \neq 0$ (matricea nulă de dimensiune $d \times d$), $B \in \mathbb{R}^d$ (un anumit vector-coloană) și $C \in \mathbb{R}$ (o anumită constantă). Veți specifica în mod clar valorile pentru A , B și C .

Răspuns:

Elaborând relația dintre probabilitățile care determină decizia luată de către clasificatorul Bayes Optimal gaussian pentru o instanță oarecare de test x (vedeți relația (168)), putem obține următoarele echivalențe:

$$\begin{aligned} p(y=1|x) \geq p(y=0|x) &\Leftrightarrow \ln p(y=1|x) \geq \ln p(y=0|x) \\ &\Leftrightarrow \ln p(y=1|x) - \ln p(y=0|x) \geq 0 \Leftrightarrow \ln \frac{p(y=1|x)}{p(y=0|x)} \geq 0 \\ \stackrel{F. \text{ Bayes}}{\Leftrightarrow} \ln \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} &\geq 0 \Leftrightarrow \ln \frac{p(y=1)}{p(y=0)} + \ln \frac{p(x|y=1)}{p(x|y=0)} \geq 0 \\ &\Leftrightarrow \ln \frac{\phi}{1-\phi} - \ln \frac{|\Sigma_1|^{1/2}}{|\Sigma_0|^{1/2}} - \frac{1}{2} \left((x-\mu_1)^\top \Sigma_1^{-1}(x-\mu_1) - (x-\mu_0)^\top \Sigma_0^{-1}(x-\mu_0) \right) \geq 0 \\ &\Leftrightarrow -\frac{1}{2} \left(x^\top (\Sigma_1^{-1} - \Sigma_0^{-1})x - 2(\mu_1^\top \Sigma_1^{-1} - \mu_0^\top \Sigma_0^{-1})x + \mu_1^\top \Sigma_1^{-1} \mu_1 - \mu_0^\top \Sigma_0^{-1} \mu_0 \right) \\ &\quad + \ln \frac{\phi}{1-\phi} - \ln \frac{|\Sigma_1|^{1/2}}{|\Sigma_0|^{1/2}} \geq 0 \\ &\Leftrightarrow x^\top \left(\frac{1}{2} (\Sigma_0^{-1} - \Sigma_1^{-1}) \right) x + \left(\mu_1^\top \Sigma_1^{-1} - \mu_0^\top \Sigma_0^{-1} \right) x \\ &\quad + \ln \frac{\phi}{1-\phi} + \ln \frac{|\Sigma_0|^{1/2}}{|\Sigma_1|^{1/2}} + \frac{1}{2} \left(\mu_0^\top \Sigma_0^{-1} \mu_0 - \mu_1^\top \Sigma_1^{-1} \mu_1 \right) \geq 0. \end{aligned} \quad (170)$$

Coroborând rezultatul pe care tocmai l-am obținut cu relația (168) din enunț, observăm că putem considera

$$\begin{aligned} A &= \frac{1}{2} (\Sigma_0^{-1} - \Sigma_1^{-1}) \\ B^\top &= \mu_1^\top \Sigma_1^{-1} - \mu_0^\top \Sigma_0^{-1} \Leftrightarrow \\ B &= (\mu_1^\top \Sigma_1^{-1} - \mu_0^\top \Sigma_0^{-1})^\top = (\Sigma_1^{-1})^\top \mu_1 - (\Sigma_0^{-1})^\top \mu_0 = \Sigma_1^{-1} \mu_1 - \Sigma_0^{-1} \mu_0 \\ C &= \ln \frac{\phi}{1-\phi} + \ln \frac{|\Sigma_0|^{1/2}}{|\Sigma_1|^{1/2}} + \frac{1}{2} \left(\mu_0^\top \Sigma_0^{-1} \mu_0 - \mu_1^\top \Sigma_1^{-1} \mu_1 \right). \end{aligned} \quad (171)$$

Se poate constata imediat că ipoteza $\Sigma_0 \neq \Sigma_1$ din enunț implică faptul că $\Sigma_0^{-1} \neq \Sigma_1^{-1}$, deci $A \neq 0$. Prin urmare, granița de separare (engl., decision boundary) determinată în acest caz de către algoritmul Bayes Optimal gaussian este de ordin pătratic.

20.

(Clasificatorul Bayes Optimal gaussian:
aplicare pe date din \mathbb{R}^2 ;
raportul cu regresia logistică)

prelucrare de Liviu Ciortuz, după

• MIT, 2001 fall, Tommi Jaakkola, HW2, pr. 2.abe

a. Precizați separatorii decizionali determinați de algoritmul Bayes Optimal gaussian pentru fiecare din cazurile de mai jos. Veți menționa în mod explicit, în fiecare caz în parte, tipul — liniar, pătratic (cerc, elipsă, parabolă, hiperbolă etc.), etc. — și ecuația acestor separatori decizionali.

(i.)

$$P_0 = 0.5, P_1 = 0.5$$

$$\mu_0 = (1, 1)^\top$$

$$\mu_1 = (-1, -1)^\top$$

$$\Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

(ii.)

$$P_0 = 0.995, P_1 = 0.005$$

$$\mu_0 = (1, 1)^\top$$

$$\mu_1 = (-1, -1)^\top$$

$$\Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

(iii.)

$$P_0 = 0.5, P_1 = 0.5$$

$$\mu_0 = (1, 1)^\top$$

$$\mu_1 = (-1, -1)^\top$$

$$\Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$$

(iv.)

$$P_0 = 0.5, P_1 = 0.5$$

$$\mu_0 = (0, 0)^\top$$

$$\mu_1 = (0, 0)^\top$$

$$\Sigma_0 = \begin{pmatrix} 4 & 0 \\ 0 & 2 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$$

(v.)

$$P_0 = 0.5, P_1 = 0.5$$

$$\mu_0 = (1, 1)^\top$$

$$\mu_1 = (1, 1)^\top$$

$$\Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}$$

(vi.)

$$P_0 = 0.5, P_1 = 0.5$$

$$\mu_0 = (-2, 1)^\top$$

$$\mu_1 = (1, 1)^\top$$

$$\Sigma_0 = \begin{pmatrix} 1 & 0 \\ 0 & 1/2 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 2 & 0 \\ 0 & 4 \end{pmatrix}$$

Sugestie: În cazul general, ecuația separatorului decizional determinat de clasificatorul Bayes Optimal gaussian este $x^\top Ax + B^\top x + C = 0$ (vedeți realția (169) de la problema 19), unde valorile pentru A , B și C au fost determinate prin relațiile (171).

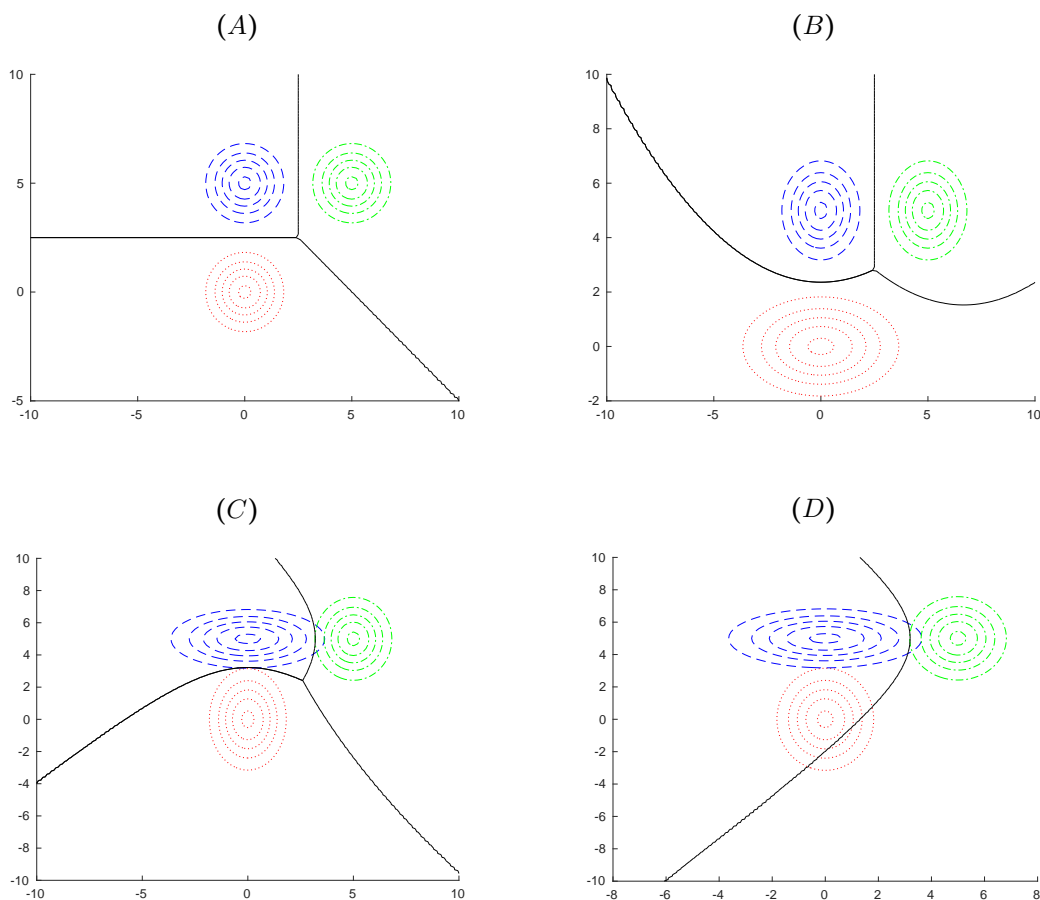
b. Arătați că este posibil ca întreg planul euclidian să fie alocat de către algoritmul Bayes Optimal gaussian³³⁴ unei singure clase, situație în care, practic, nu avem separator decizional.

c. Care dintre *separatorii decizionali* de la punctul a poate corespunde și unui model de tip *regresie logistică*? Justificați răspunsul dumneavoastră.

d. Care dintre *zonele de decizie* din figurile de mai jos poate să corespundă unui model de tip *regresie logistică multinomială (softmax)*? Justificați răspunsul dumneavoastră.³³⁵

³³⁴LC: De fapt, este suficient să ne gândim la cazul $\Sigma_0 = \Sigma_1$. (Menționăm că acest caz particular de clasificare bayesiană gaussiană este numit în literatura de specialitate *analiză gaussiană discriminativă*.)

³³⁵Pentru o introducere în regresia logistică multinomială, vedeți problema 17 de la capitolul *Metode de regresie*.

**Răspuns:**

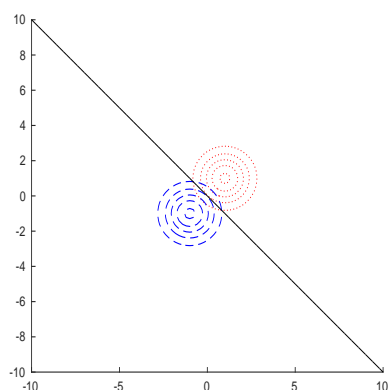
a. Folosind formulele (169), se calculează coeficienții A , B și C din ecuația $x^\top A x + B^\top x + C = 0$, pentru fiecare dintre cazurile date:³³⁶

i. $A = 0$ – matricea nulă de dimensiune 2×2 , $B^\top = (2, 2)$, $C = 0$.

Ecuția separatorului decizional este

$$2x_1 + 2x_2 = 0 \Leftrightarrow x_2 = -x_1.$$

Această ecuație reprezintă o dreaptă care trece prin originea sistemului de coordonate. Observați că această dreaptă este mediatoarea segmentului de dreaptă care unește punctele reprezentate de μ_1 și μ_2 (mediile celor două gaussiene).



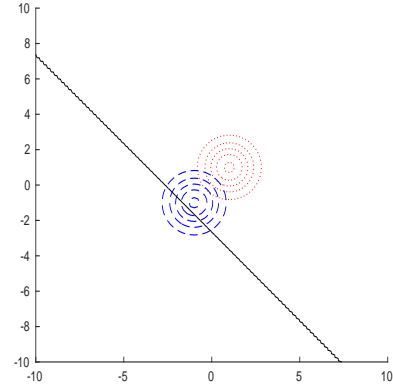
³³⁶ Am omis calculele efective, fiindcă nu sunt dificile. Invităm studentul să verifice că este capabil să obțină el însuși aceste rezultate.

ii. $A = 0$, $B^\top = (2, 2)$, $C = \ln 199$.

Ecuția separatorului decizional este

$$2x_1 + 2x_2 + \ln 199 = 0 \Leftrightarrow x_2 = -x_1 - \underbrace{\ln \sqrt{199}}_{\approx 2.6466}.$$

Această ecuație reprezintă o dreaptă paralelă cu cea de la punctul precedent (i). Ea este situată acum mult mai aproape de punctul μ_1 — datorită faptului că în acest caz probabilitatea P_0 este mult mai mare decât probabilitatea P_1 .

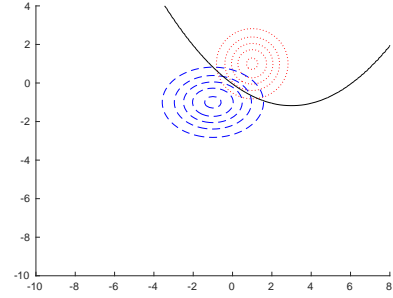


iii. $A = \begin{pmatrix} 1/4 & 0 \\ 0 & 0 \end{pmatrix}$, $B^\top = (-3/2, -2)$, $C = 1/4$.

În acest caz, spre deosebire de cazurile precedente, ecuația separatorului decizional este de ordin pătratic:

$$\frac{1}{4}x_1^2 - \frac{3}{2}x_1 - 2x_2 + \frac{1}{4} = 0 \Leftrightarrow x_2 = \frac{1}{8}x_1^2 - \frac{3}{4}x_1 + \frac{1}{8}.$$

Această ecuație reprezintă o parabolă, al cărei vârf are coordonatele $(3, -1)$.



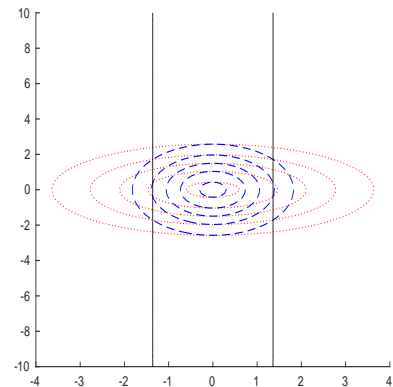
iv. $A = \begin{pmatrix} -3/8 & 0 \\ 0 & 0 \end{pmatrix}$, $B^\top = 0$ — vectorul nul din \mathbb{R}^2 , $C = \ln 2$.

Ecuția separatorului decizional este

$$-\frac{3}{8}x_1^2 + \ln 2 = 0 \Leftrightarrow x_1^2 = \frac{8}{3} \ln 2.$$

Așadar, în acest caz, separatorul decizional este tot de ordin pătratic, dar este reprezentat de două drepte, având ecuațiile $x_1 = -\sqrt{\frac{8}{3} \ln 2} \approx -1.3595$ și respectiv $x_1 = \sqrt{\frac{8}{3} \ln 2} \approx 1.3595$.

Ținând cont de relația (170) care a fost demonstrată la problema 19, deducem că zona de decizie corespunzătoare clasei 1 este situată între aceste două drepte, iar zona de decizie corespunzătoare clasei 0 este în exteriorul lor.

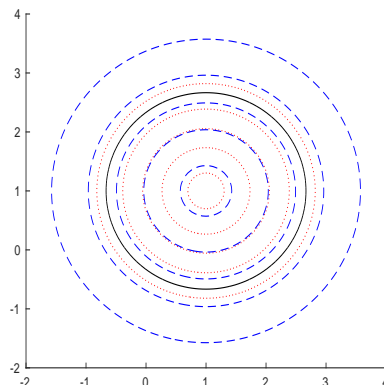


v. $A = \begin{pmatrix} 1/4 & 0 \\ 0 & 1/4 \end{pmatrix}$, $B^\top = (-1/2, -1/2)$, $C = -\ln 2 + 1/2$.

Ecuția separatorului decizional este

$$\begin{aligned}\frac{1}{4}(x_1^2 + x_2^2) - \frac{1}{2}(x_1 + x_2) - \ln 2 + \frac{1}{2} &= 0 \Leftrightarrow \\ \frac{1}{4}(x_1^2 - 2x_1 + x_2^2 - 2x_2) &= \ln 2 - \frac{1}{2} \Leftrightarrow \\ (x_1 - 1)^2 + (x_2 - 1)^2 &= 4\left(\ln 2 - \frac{1}{2}\right) + 2.\end{aligned}$$

Această ecuație — tot de ordin pătratic — reprezintă un cerc cu centrul în punctul $(1, 1)$ și raza de $2\sqrt{\ln 2} \approx 1.3862$.

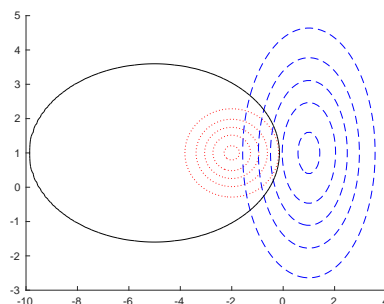


vi. $A = \begin{pmatrix} 1/4 & 0 \\ 0 & 7/8 \end{pmatrix}$, $B^\top = (5/2, -7/4)$, $C = -\ln(1/4) + 21/8$.

Ecuția separatorului decizional este

$$\begin{aligned}\frac{1}{4}x_1^2 + \frac{7}{8}x_2^2 + \frac{5}{2}x_1 - \frac{7}{4}x_2 - 2\ln 2 + \frac{21}{8} &= 0 \Leftrightarrow \\ \frac{1}{4}(x_1^2 + 10x_1) + \frac{7}{8}(x_2^2 - 2x_2) &= 2\ln 2 - \frac{21}{8} \Leftrightarrow \\ \frac{(x_1 + 5)^2}{4} + \frac{(x_2 - 1)^2}{\frac{8}{7}} &= 2\ln 2 - \frac{21}{8} + \frac{25}{4} + \frac{7}{8}.\end{aligned}$$

Această ecuație — tot de ordin pătratic — reprezintă o elipsă care are centrul de simetrie în punctul $(-5, 1)$.



b. Pornind de la regula de decizie $y_{GNB} = \operatorname{argmax}_{y \in \{0,1\}} P(Y = y|X = x)$, putem raționa astfel:

$$\begin{aligned}P(Y = 1|X = x) &\geq P(Y = 0|X = x) \stackrel{F. Bayes}{\Leftrightarrow} \\ P(X = x|Y = 1)P(Y = 1) &\geq P(X = x|Y = 0)P(Y = 0).\end{aligned}$$

Dacă alegem ca distribuțiile gaussiene corespunzătoare claselor $Y = 1$ și $Y = 0$ să fie identice — adică să aibă aceiași parametri μ și Σ —, în inegalitatea de mai sus factorii $P(X = x|Y = 1)$ și $P(X = x|Y = 0)$ vor fi egali. Prin urmare, atunci când $P(Y = 1) \neq P(Y = 0)$, va rezulta că întreg planul / spațiul va fi alocat unei singure clase, și anume acelei clase $y \in \{0, 1\}$ care are probabilitatea a priori ($P(Y = y)$) mai mare.

c. Se știe că regresia logistică determină separatori decizionali liniari.³³⁷ Prin urmare, doar separatorii din cazurile *i* și *ii* pot corespunde (și) unor modele de regresie logistică.

d. Justificarea este similară cu cea de la punctul c: doar în cazul (A) separatorii decizionali pot corespunde (și) unui model de regresie logistică multinomială / softmax.

³³⁷Vedeți problema 12 de la capitolul *Metode de regresie*, în special relațiile (123), din care rezultă:

$$P(Y = 1|X = x) \geq P(Y = 0|X = x) \Leftrightarrow \sigma(z) \geq 1/2 \Leftrightarrow z \geq 0,$$

unde $\sigma(z) \stackrel{def.}{=} \frac{1}{1+e^{-z}}$, $z \stackrel{not.}{=} w_0 + \sum_{i=1}^d w_i x_i$ și $x \stackrel{not.}{=} (x_1, \dots, x_d)$. Așadar, în cazul regresiei logistice separatorul decizional are ecuația de forma $w_0 + \sum_{i=1}^d w_i x_i = 0$, deci este de tip liniar.

21.

(Clasificatori de tip Bayes gaussian, cazul datelor cu atribute de intrare corelate: exemplificare în \mathbb{R}^2 ; comparație)

■ □ ● ○ CMU, 2014 fall, W. Cohen, Z. Bar-Joseph, HW2, pr. 5.c

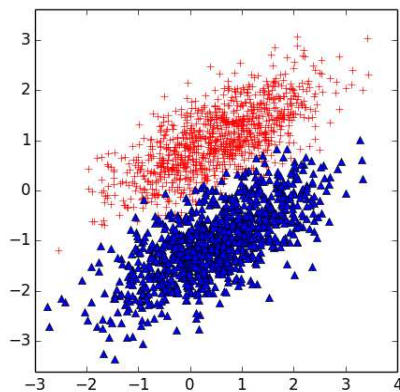
În cazul bidimensional, putem vizualiza modul în care se comportă algoritmul Bayes Naiv gaussian atunci când atributele de intrare (adică, trăsăturile; engl., features) sunt corelate (engl., correlated).³³⁸

Fie setul de date din figura (A) de mai jos, în care instanțele roșii sunt din clasa 0, iar instanțele albastre din clasa 1. Distribuțiile comune condiționale $((X_1, X_2)|Y = 0)$ și $((X_1, X_2)|Y = 1)$ sunt de tip gaussian bidimensional. Elipsele din figurile (B), (C) și (D) reprezintă curbe de izocontur pentru [diverse] distribuții condiționale asociate celor două clase. Centrele elipselor corespund mediilor, iar curbele de izocontur sunt situate la o distanță de două deviații standard față de medii.³³⁹

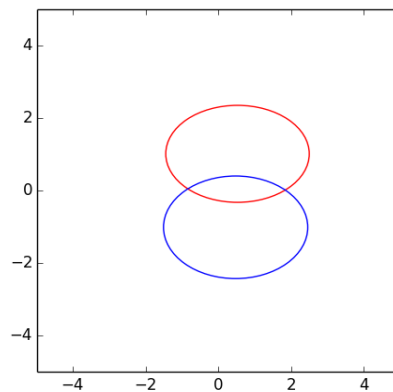
a. Care anume dintre perechile de elipse din figurile (B), (C) și (D) corespunde cel mai probabil distribuțiilor condiționale care au generat datele din figura (A)?

b. Care anume dintre aceste elipse corespunde [cel mai probabil] estimărilor de parametri făcute de către algoritmul Bayes Naiv gaussian?

c. Dacă presupunem că probabilitățile a priori pentru cele două clase sunt egale, care dintre modelele (B), (C) și (D) va obține o acuratețe mai mare pe datele de antrenament? Care va fi natura separatorului decizional în acest caz?



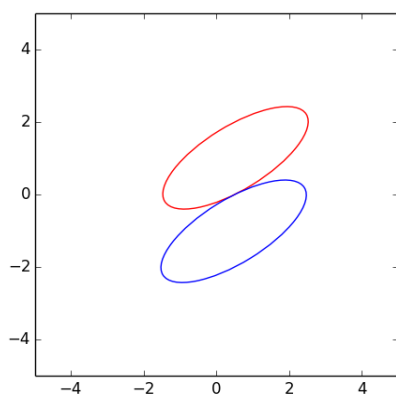
(A) Date



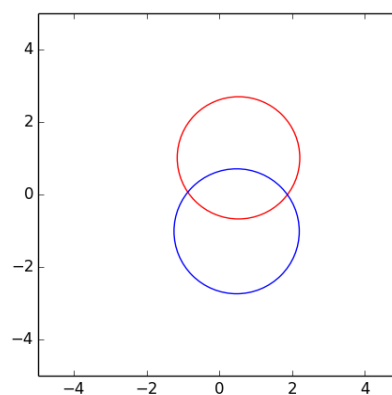
(B)

³³⁸Pentru definiția noțiunii de *corelare* pentru două variabile aleatoare, vedeți problema 19 de la capitolul de *Fundamente*.

³³⁹Mai exact, o astfel de curbă de izocontur este constituită din punctele x din planul euclidian pentru care $\Sigma^{-1/2}(x - \mu) = 2 \Leftrightarrow (x - \mu)^\top \Sigma^{-1}(x - \mu) = 4$, unde μ este media distribuției gaussiene considerate. Vedeți *Pattern Classification*, R. Duda, P. Hart, D. Stork, 2nd ed. (Wiley-Interscience, 2000), Appendix A, pag. 625.



(C)



(D)

Răspuns:

a. Se observă în figura (A) că datele au fost generate de distribuții gaussiene bidimensionale având matrice de covarianță nediagonale și identice, însă având medii diferite. În plus, dacă notăm o instanță oarecare cu $x = (x_1, x_2)$ este evident că x_1 și x_2 nu sunt independente, ci există o anumită corelare între ele (mai precis, x_1 și x_2 sunt într-o relație de dependență de tip liniar). În consecință, curbele de izocontur pentru aceste distribuții sunt elipse identice ca mărime, având axele de simetrie neparalele cu axele sistemului de coordonate. Evident, doar desenul (C) corespunde acestei situații.

b. Clasificatorul Bayes Naiv gaussian presupune independența condițională a celor două atribute (X_1 și X_2) în raport cu eticheta / variabila de ieșire (Y). Această independență corespunde unor matrice de covarianță diagonale, respectiv unor curbe de izocontur reprezentate de elipse ale căror axe de simetrie sunt paralele cu axele sistemului de coordonate. Atât elipsele din desenul (B) cât și cele din desenul (D) satisfac aceste specificații, însă în cazul (D) elipsele sunt chiar cercuri, în vreme ce datele din figura (A) sunt dispuse în elipse având deviații standard diferite pe cele două axe de simetrie. Așadar, cazul (B) corespunde estimărilor făcute de clasificatorul Bayes Naiv gaussian pe aceste date.

c. Evident, cazul (C) furnizează cea mai mică eroare la antrenare, deci cea mai mare acuratețe. În acest caz, se lucrează cu algoritmul Bayes Optimal gaussian. Întrucât matricele de covarianță sunt egale, separatorul decizional este de tip liniar. (Pentru justificare riguroasă, vedeți rezultatul teoretic demonstrat la problema 18.)

22. (Algoritmul Bayes Naiv [gaussian] vs. regresia logistică — comparații)

□ ● ○ CMU, 2012 fall, T. Mitchell, Z. Bar-Joseph, HW2, pr. 1.f, 3.b-d

Fie un set de date caracterizate de atributele X_1, \dots, X_n (pe care le vom considera ca fiind fie toate Bernoulli, fie toate gaussiene) și de eticheta Y .³⁴⁰

³⁴⁰ Așa este cazul problemei 14, în care variabilele condiționate $X_i|Y$ sunt de tip Bernoulli și sunt independente

Modelul Bayes Naiv [eventual de tip gaussian] va fi identificat în continuare cu abrevierea NB, iar modelul regresiei logistice cu LR.

- Presupunem că datele satisfac presupoziția de independență condițională de tip Bayes Naiv. Atunci când numărul de exemple de antrenament tinde la infinit, care dintre cei doi clasificatori va produce rezultate mai bune, NB sau LR? Justificați.
- Presupunem acum că datele *nu* satisfac presupoziția de independență condițională de tip Bayes Naiv. Ne punem aceeași întrebare ca mai sus: atunci când numărul de exemple de antrenament tinde la infinit, care dintre cei doi clasificatori va produce rezultate mai bune, NB sau LR? Justificați.
- Este oare posibil să calculăm distribuția $P(X)$ cu ajutorul parametrilor estimați de către algoritmul Bayes Naiv? Explicați în mod succint.
- Este oare posibil să calculăm distribuția $P(X)$ cu ajutorul parametrilor w calculați de către regresia logistică? Explicați în mod succint.

Răspuns:

a. Regresia logistică este un clasificator probabilist de tip *discriminativ*, adică aproximează / modelează $P(Y|X)$ cu ajutorul funcției logistice de argument $w \cdot X$ (unde w este vectorul de parametri, $w \in \mathbb{R}^d$, sau $w \in \mathbb{R}^{d+1}$ dacă extindem fiecare instanță X cu componenta $X_0 = 1$). În contrast cu acesta, NB este un clasificator probabilist de tip *generativ*, deci calculează distribuțiile $P(X|Y = y)$ (care în cazul GNB sunt de tip gaussian multidimensional) și de asemenea $P(Y)$, estimând parametrii acestor distribuții.

Atunci când numărul de instanțe tinde la infinit, pe de o parte aproximația calculată de regresia logistică pentru distribuția $P(Y|X)$ tinde la distribuția reală $P(Y|X)$, iar pe de altă parte estimările făcute de clasificatorul Bayes Naiv pentru distribuțiile $P(Y)$ și $P(X|Y)$ vor tinde la distribuțiile reale corespunzătoare, dat fiind (în cazul lui $P(X|Y)$) că datele de antrenament satisfac presupoziția de independență condițională. Corespondența dintre distribuțiile reale $P(Y|X)$ (pe de o parte) și $P(Y)$ și $P(X|Y)$ (pe de altă parte) este dată de formula lui Bayes. Prin urmare, în aceste condiții cei doi clasificatori vor produce rezultate echivalente.

b. Regresia logistică va produce rezultate mai bune, fiindcă ea nu lucrează cu presupoziția de independență condițională. (Vedeți pe de o parte problema 32 de la capitolul *Metode de regresie*, iar pe de altă parte problema 10 de la prezentul capitol.)

c. Da, algoritmul Bayes Naiv este un clasificator de tip *generativ* (engl., *generative classifier*). Putem calcula $P(X)$ prin „marginalizarea” distribuției condiționate $P(X|Y)$ în raport cu eticheta / clasa Y , și anume: $P(X) = \sum_y P(X|Y = y) \cdot P(Y = y)$.

d. Nu, nu este posibil. Așa cum am precizat la punctul a, regresia logistică estimează $P(Y|X)$ (nu $P(X|Y)$ și $P(Y)$, cum calculează clasificatorii de tip Bayes Naiv).

condițional două câte două, sau cel al problemei 17, în care toate variabilele condiționate, $X_i|Y$ pentru $i = 1, \dots, n$ urmează distribuții gaussiene — având varianțele $\sigma_{i0} = \sigma_{i1}$ — și sunt, de asemenea, independente condițional două câte două. (Similar este și cazul problemei 44, care constituie o combinație a precedentelor două tipuri.)

23. (Clasificarea bayesiană gaussiană vs. regresia logistică:
Adevărat sau Fals?)

□ ◦ CMU, 2010 fall, Aarti Singh, midterm, pr. 1.2
CMU, (?) 15-781, midterm example questions, pr. 1.d

a. Corespondența dintre regresia logistică și clasificatorul Bayes Naiv de tip gaussian înseamnă — în cazul în care matricele de covarianță corepunzătoare claselor sunt toate egale cu matricea-identitate³⁴¹ — că există o corespondență 1-la-1 între parametrii celor doi clasificatori.

b. Presupunând că lucrăm cu un număr fix de attribute, putem învăța un clasificator Bayes Optimal de tip gaussian în timp liniar în raport cu numărul de instanțe din setul de date de antrenament.

Răspuns:

a. Fals. Se poate preciza de la început că deși cei doi clasificatori învață separatori decizionali care au aceeași formă (și anume, o formă liniară, vedeți problema 17) nu rezultă în mod neapărat că pe un același set de date cei doi separatori decizionali învățați coincid. Faptul că matricele de covarianță sunt, toate, matrice identitate / diagonale înseamnă că presupuziția de independență condițională este satisfăcută.³⁴² Suntem, așadar, în condiții similare cu cele de la problema 22.a, însă aici nu avem neapărat satisfăcută ipoteza că numărul de instanțe de antrenament tinde la infinit, caz în care rezultatele de clasificare furnizate de LR și NB ar fi echivalente. Așadar, nu există (în general) o corespondență de tip 1-la-1 între parametrii w_{LR} calculați de regresia logistică³⁴³ și parametrii w_{NB} corespunzători clasificatorului Bayes Naiv gaussian.³⁴⁴

Observație: Rezultatul acesta este valabil și în cazul algoritmului Bayes Naiv cu variabile categoricale.

b. Adevărat. În cazul clasificatorului Bayes Optimal de tip gaussian, regula de calcul pentru clasificarea unei instanțe noi $x \in \mathbb{R}^d$ este următoarea:

$$\begin{aligned} y_{GNB} &\stackrel{\text{def.}}{=} \operatorname{argmax}_{y \in \text{Val}(Y)} P(Y = y | X = x) = \operatorname{argmax}_{y \in \text{Val}(Y)} \frac{P(X = x | Y = y)P(Y = y)}{P(X = x)} \\ &= \operatorname{argmax}_{y \in \text{Val}(Y)} P(X = x | Y = y)P(Y = y) = \operatorname{argmax}_{y \in \text{Val}(Y)} \mathcal{N}(x; \mu_y, \Sigma_y)P(Y = y). \end{aligned}$$

Așadar, pentru fiecare $y \in \text{Val}(Y)$ vom avea de estimat câte o pereche de parametri, care caracterizează o distribuție gaussiană multidimensională: vectorul de medii μ_y și matricea de covarianță Σ_y . Conform problemei 48 de la capitolul de *Fundamente*, estimările celor doi parametri sunt media la eșantionare și respectiv matricea de covarianță la eșantionare, deci se calculează în timp liniar în raport cu numărul de instanțe de antrenament.

³⁴¹LC: Chiar mai general, putem considera că aceste matrice sunt diagonale.

³⁴²Vedeți problema 31 de la capitolul de *Fundamente*.

³⁴³Vedeți problema 12 de la capitolul *Metode de regresie*.

³⁴⁴În primul rând, din punctul de vedere al terminologiei, dacă ne referim la parametrii calculați de NB ca fiind parametrii distribuțiilor $P(X|Y = y)$ și $P(Y)$, este imediat că nu există o corespondență 1-la-1 între aceștia și parametrii w_{LR} calculați de LR, care servesc la estimarea / aproximarea distribuției $P(Y|X)$.

2.2 Clasificare bayesiană — Probleme propuse

2.2.1 Ipoteze de probabilitate maximă a posteriori (MAP)

24. (Formula lui Bayes; inferențe statistice; ilustrarea noțiunii de ipoteze MAP)

* CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW1, pr. 5

Imaginează-ți că te afli în fața a trei cutii, care sunt etichetate cu literele A, B, C . Două dintre ele sunt goale, iar una conține un premiu. Tu nu știi în care dintre ele se află premiul; trebuie să ghicești. Procedezi în felul următor:

Mai întâi alegi la întâmplare o cutie X (să zicem că $X = A$). Totuși, chiar înainte de a deschide cutia X observi că altcineva, înaintea ta, a deschis cutia Y (să zicem că $Y = B$). Ai dreptul să privești în cutia Y , ca să vezi dacă ea conține sau nu premiul. Dacă ea nu conține premiul, vei avea dreptul să schimbi alegerea pe care ai făcut-o inițial.

În vederea luării unei decizii cât mai bune, ți se comunică *strategia* după care a fost aleasă cutia Y , și anume, folosind una din următoarele trei *variante*:

a. Dacă cutia pe care ai ales-o inițial conține premiul, atunci cutia Y este aleasă cu probabilitate de $1/2$ una din cele două cutii goale (diferite de cutia X). Dacă X este vidă, atunci Y se alege ca fiind cutia goală diferită de X .

b. Se alege aleatoriu cu probabilitate de $1/2$ una din cutiile diferite de cea pe care ai ales-o tu inițial, X . (În consecință, cutia Y poate sau nu să conțină premiul. Dacă Y conține premiul, ai pierdut jocul.)

c. Se alege în mod aleatoriu cu probabilitate de $1/2$ una din cutiile goale. (Deci este posibil ca $Y = X$. Așadar, în cazul $Y = X$ observi că a fost deschisă anterior chiar cutia X . În continuare vei putea alege una din celelalte două cutii.)

Considerând (pentru simplitate) că $X = A$, $Y = B$, iar cutia B este vidă, pentru fiecare din cele trei *variante* de mai sus decide ce cutie ar trebui să alegi în final pentru a-ți maximiza șansele de a obține premiul. Justifică-ți decizia, elaborând calculul probabilistic aferent.

25. (Adevărat sau Fals?)

CMU, 2002 fall, Andrew Moore, final exam, pr. 11.a

Fie D un set de exemple (date de antrenament), iar H o mulțime de ipoteze pentru (un algoritm de) învățare automată pe datele D . Precizați care este *valoarea de adevăr* a următoarelor afirmații:

$\operatorname{argmax}_{h \in H} P(D|h)$ este ipoteză de probabilitate maximă a posteriori, [iar]
 $\operatorname{argmax}_{h \in H} P(h|D)$ este ipoteză de verosimilitate maximă.

2.2.2 Algoritmii Bayes Naiv și Bayes Optimal

26. (Algoritmii Bayes Naiv și Bayes Optimal; aplicare; numărul minimal de parametri de estimat)

• ◦ *Liviu Ciortuz, 2017, pornind de la setul de date din Machine Learning, Tom Mitchell, 1997, ch. Decision Trees, page 59*

Considerăm următorul set de date de antrenament, în care variabila de ieșire este *EnjoyTennis*:

Day	Outlook	Temperature	Humidity	Wind	EnjoyTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

- a. Determinați decizia luată de către algoritmul Bayes Naiv pentru instanța de test

$$X = \langle \text{Outlook} = \text{sunny}, \text{Temp} = \text{cool}, \text{Humidity} = \text{high}, \text{Wind} = \text{strong} \rangle,$$

precum și probabilitatea cu care este luată această decizie.

- b. Care este numărul *minim* de parametri pe care trebuie să-l estimeze algoritmul Bayes Naiv pe aceste date [pentru a face apoi predicții pe un set oarecare de instanțe de test]? Dar în cazul clasificatorului Bayes Optimal?

- c. Implementați algoritmul Bayes Naiv, iar apoi cu ajutorul acestei implementări calculați eroarea la antrenare și eroarea la CVLOO pe acest set de date.

27. (Algoritmul Bayes Naiv: aplicare)

* *CMU, 2004 fall, T. Mitchell Z. Bar-Joseph, midterm, pr. 6.a*

Se dă setul de date din tabelul alăturat, în care A, B, C sunt attribute (de intrare) binare, iar Y este atribut de ieșire. Care va fi răspunsul algoritmului de clasificare Bayes Naiv pentru intrarea $A = 0, B = 0, C = 1$?

A	B	C	Y
0	0	1	0
0	1	0	0
1	1	0	0
0	0	1	1
1	1	1	1
1	0	0	1
1	1	0	1

28. (Algoritmul Bayes Naiv și algoritmul Bayes Optimal: aplicare)

*prelucrare de Liviu Ciortuz, după**• * CMU, 2002 fall, Andrew Moore, final exam, pr. 4.b-e*

Se dă setul de date alăturat, cu A și B variabile de intrare, iar C variabilă de ieșire.

a. Care este numărul minim de probabilități ce trebuie estimate pentru a putea construi după aceea (pe acest set de date) un clasificator de tip Bayes Naiv? Justificați.

b. Similar, pentru clasificatorul Bayes Optimal. Justificați.

c. Care este decizia clasificatorului Bayes Naiv pentru $A = 0, B = 1$? Precizați cu ce probabilitate este luată această decizie.

d. Care este decizia clasificatorului Bayes Optimal pentru $A = 0, B = 1$? Precizați cu ce probabilitate este luată această decizie.

e. Dacă rezultatele obținute la punctele c și d diferă (fie și numai în privința probabilităților cu care sunt luate deciziile), care este explicația? Justificați în mod riguros.

A	B	C	nr. apariții
0	0	1	3
0	1	0	1
0	1	1	4
1	0	0	5
1	1	0	2
1	1	1	1

29. (Aplicarea algoritmului Bayes Naiv la clasificarea de texte)

*• * Edinburgh, 2009 fall, C. Williams, V. Lavrenko, tutorial 2, pr. 2*

Firma Whizzco decide să implementeze un clasificator de texte. Pentru început, ei vor să clasifice documente aparținând fie clasei *sport* fie clasei *politică*. Ei decid să reprezinte fiecare document ca un vector de attribute descriind prezența ori absența unor cuvinte-cheie:

goal, football, golf, defence, offence, wicket, office, strategy.

Datele de antrenament sunt reprezentate folosind o matrice în care fiecare linie este un vector de valori (0 sau 1) pentru cele 8 attribute.

```

xP=[1 0 1 1 1 0 1 1; % Politica
     0 0 0 1 0 0 1 1;
     1 0 0 1 1 0 1 0;
     0 1 0 0 1 1 0 1;
     0 0 0 1 1 0 1 1;
     0 0 0 1 1 0 0 1]

xS=[1 1 0 0 0 0 0 0; % Sport
     0 0 1 0 0 0 0 0;
     1 1 0 1 0 0 0 0;
     1 1 0 1 0 0 0 1;
     1 1 0 1 1 0 0 0;
     0 0 0 1 0 1 0 0;
     1 1 1 1 1 0 1 0]

```

Folosind algoritmul Bayes Naiv, care este probabilitatea cu care documentul $x = (1, 0, 0, 1, 1, 1, 0)$ va fi clasificat ca aparținând clasei *politică*?

30.

(Aplicarea algoritmului Bayes Naiv:
chestiunea valorilor lipsă (engl., missing values)
în datele de antrenament)

□ • CMU, 2013 fall, A. Smola, G. Gordon, midterm practice, pr. 9

Fie setul de date de antrenament (x, y) și datele de test z :

$$\begin{aligned}x_1 &= (0, 0, 0, 1, 0, 0, 1) & y_1 &= 1 \\x_2 &= (0, 0, 1, 1, 0, 0, 0) & y_2 &= 1 \\x_3 &= (1, 1, 0, 0, 0, 1, 0) & y_3 &= -1 \\x_4 &= (1, 0, 0, 0, 1, 1, 0) & y_4 &= -1 \\z_1 &= (1, 0, 0, 0, 0, 1, 0) \\z_2 &= (0, 1, 1, 0, 0, 1, 1)\end{aligned}$$

Ce problemă va întâmpina clasificatorul Bayes Naiv pe aceste date?

(Indicație: Pentru ca răspunsul dumneavoastră să fie cât mai bine justificat, veți estima toți parametrii necesari și veți aplica algoritmul pe cele două instanțe de test. Veți nota atributele cu $A1, A2, \dots$.)

La curs am prezentat un „remediu” standard pentru o astfel de problemă. Precizați cum se numește „tehnica” respectivă și aplicați-o pe aceste date. După aceea, veți aplica algoritmul Bayes Naiv pentru a clasifica instanțele de test z_1 și z_2 .

31.

(Algoritmului Bayes Naiv:
calculul ratei medii a erorilor)

■ • ◦ CMU, 2010 fall, Aarti Singh, HW1, pr. 4.2

Considerăm următoarea problemă de clasificare:

Fie variabila aleatoare $Y: Hike \in \{T, F\}$ care denotă faptul că Alice și Bob merg sau nu în drumeție în funcție de condițiile vremii: $X_1: Sunny \in \{T, F\}$ și $X_2: Windy \in \{T, F\}$.

Se presupune că au fost estimați următorii parametri:

$$\begin{aligned}P(Hike) &= 0.5 \\P(Sunny | Hike) &= 0.8, & P(Sunny | \neg Hike) &= 0.7 \\P(Windy | Hike) &= 0.4, & P(Windy | \neg Hike) &= 0.5\end{aligned}$$

De asemenea, se consideră că este satisfăcută presupuziția de independență condițională a algoritmului Bayes Naiv.

a. Care este probabilitatea (comună) ca Alice și Bob să meargă în drumeție atunci când vremea este însorită și bate vântul, adică

$$P(Sunny = T, Windy = T, Hike = T) = ?$$

Care este decizia luată de algoritmul Bayes Naiv în acest caz?

b. Completați tabelul următor:

X_1	X_2	Y	$P(X_1, X_2, Y)$	$P_{NB}(Y X_1, X_2)$	decizia algoritmului Bayes Naiv
F	F	F			
F	F	T			
F	T	F			
F	T	T			
T	F	F			
T	F	T			
T	T	F			
T	T	T			

Observație: Calculele de la punctul *a* corespund ultimei linii din tabelul de mai sus.

c. Care este rata medie a erorilor (engl., expected error rate) produse de algoritmul Bayes Naiv? Vă reamintim că această (rată) medie este definită ca fiind suma probabilităților $P(X_1, X_2, Y)$ pentru acele (triplete de) valori ale variabilelor X_1, X_2, Y pentru care decizia luată de algoritmul Bayes Naiv diferă de valoarea variabilei Y .

În cele ce urmează se presupune că se obțin mai multe informații despre vreme. Se introduce o nouă trăsătură $X_3: \text{Rainy} \in \{T, F\}$. Se presupune că în fiecare zi vremea poate fi fie *Rainy* fie *Sunny*, dar nu și *Rainy* și *Sunny*. Similar, se presupune că vremea nu poate fi într-o zi $\neg \text{Rainy}$ și $\neg \text{Sunny}$.

d. În noile condiții, presupoziția de independență condițională rămâne oare adevărată? Justificați.

e. Calculați $P(\text{Sunny} = T, \text{Windy} = T, \text{Rainy} = F, \text{Hike} = T)$.

f. Care este rata medie a erorilor produse de clasificatorul Bayes Naiv când se folosesc toate cele 3 atribute de intrare?

g. S-a îmbunătățit performanța algoritmului Bayes Naiv prin adăugarea atributului *Rainy*? Explicați de ce.

32.

(Algoritmul Bayes Naiv:
calculul ratei medii a erorii – exemplificare;
comparație cu regresia logistică)

prelucrare de Livi Ciortuz, după

• * CMU, 2009 fall, Carlos Guestrin, HW1, pr. 4.1.4

Considerăm o problemă de clasificare binară în care fiecare exemplu de antrenament are două atribute binare $X_1, X_2 \in \{T, F\}$ și eticheta / clasa $Y \in \{T, F\}$. Presupunem că $P(Y = T) = 0.5$, iar $P(X_1 = T|Y = T) = 0.8$, $P(X_1 = F|Y = F) = 0.7$, $P(X_2 = T|Y = T) = 0.5$ și $P(X_2 = F|Y = F) = 0.9$. (Se poate observa că atributul X_1 furnizează / constituie un indiciu întrucâtva mai puternic decât atributul X_2 în ce privește determinarea clasei unei instanțe oarecare.)

În cele ce urmează vom presupune că X_1 și X_2 sunt independente în raport cu Y .

a. Calculați probabilitățile $P(X_1 = F|Y = T)$, $P(X_1 = T|Y = F)$, $P(X_2 = F|Y = T)$ și $P(X_2 = T|Y = F)$. Asociați răspunsului dumneavoastră o justificare generală, sub forma unei formule din teoria probabilităților:

$$P(\neg A|B) = \dots, \text{ unde } A \text{ și } B \text{ sunt evenimente aleatoare oarecare.}$$

b. Scrieți regula de decizie a algoritmului Bayes Naiv pentru $X_1 = x_1$ și $X_2 = x_2$, justificând în mod succint obținerea ei.

c. Calculați rata medie a erorii produse de algoritmul Bayes Naiv, atunci când se folosesc ambele atribute, X_1 și X_2 . (Veți da în prealabil definiția ratei medii a erorii.) Este oare această rată mai bună decât în cazul în care se folosește un singur atribut (X_1 sau X_2)? De ce?

d. Să presupunem acum că se crează un nou atribut, X_3 , care este o copie exactă a lui X_2 . Așadar, pentru fiecare exemplu de antrenament, attributele X_2 și X_3 au aceeași valoare, $X_2 = X_3$. Răspundeți la următoarele întrebări:

- Sunt X_2 și X_3 independente condițional în raport cu Y ?
- Cât este rata medie a erorii pentru Bayes Naiv acum? (Atenție! Distribuția „adevărată” a datelor nu s-a modificat.)
- Explicați ce se întâmplă cu algoritmul Bayes Naiv. Oare *regresia logistică* are aceeași problemă? Explicați de ce.

33. (Clasificare bayesiană: calculul ratei medii a erorilor pentru diverși clasificatori bayesieni)

prelucrare de L. Ciortuz, după

■ ● ○ CMU, 2004 fall, T. Mitchell Z. Bar-Joseph, HW3, pr. 1.2

Fie funcția $Y = (A \wedge B) \vee \neg(B \vee C)$, unde A, B și C sunt variabile aleatoare binare independente, fiecare dintre ele având posibilitatea să ia valoarea 0 cu probabilitate de 50%.

a. Câți parametri trebuie să estimeze clasificatorul Bayes Naiv pentru a învăța funcția Y ? Enumerați acești parametri. *Atenție:* $P(\neg x)$ nu va fi socotit ca parametru dacă $P(x)$ a fost deja estimat ca parametru.

b. Care este rata medie a erorii la antrenare pentru clasificatorul Bayes Naiv la învățarea conceptului Y , presupunând că avem o infinitate de date de antrenament?

Indicație: Scrieți mai întâi tabela de adevăr a funcției Y , apoi estimați valorile parametrilor (în sensul verosimilității maxime, MLE). Pentru conveniență, centralizați toate calculele făcute de către algoritmul Bayes Naiv într-un tabel. (Sau, altfel, puteți adăuga niște coloane suplimentare la tabela de adevăr a funcției Y .)

Convenție: În cazul în care, pentru o setare oarecare a variabilelor A, B și C , cele două probabilități calculate de către algoritmul Bayes Naiv în vederea determinării valorii y_{NB} sunt egale, convenim că algoritmul va lua decizia $y_{NB} = 1$.

c. Câți parametri trebuie să estimeze clasificatorul Bayes Optimal pentru a „învăța” funcția Y ? Justificați în detaliu.

d. Care este rata medie a erorii la antrenare pentru clasificatorul Bayes Optimal la învățarea conceptului Y , presupunând același lucru ca mai sus? *Atenție:* Nu este nevoie să calculați efectiv această rată; este suficient să indicați valoarea ei și să o justificați printr-un *raționament calitativ*.³⁴⁵

e. Considerăm un alt clasificator de tip Bayes, care presupune că A este independent de C , condiționat de B și Y — în contrast cu clasificatorul Bayes Naiv, care presupune că variabilele A, B și C sunt independente două câte două în raport cu Y .

Arătați că acest clasificator Bayes va avea nevoie să estimeze mai puțini parametri decât clasificatorul Bayes Optimal la învățarea conceptului Y , și totuși va obține aceeași rată medie a erorii la antrenare (considerând că este valabilă aceeași presuposiție în legătură cu datele de antrenament).

Indicații:

- Scrieți tabela de adevăr a funcției Y .
- Folosind această tabelă de adevăr, dovediți că într-adevăr proprietatea de independență condițională formulată mai sus este satisfăcută. Și anume: pentru fiecare pereche de valori b și y ale variabilelor aleatoare B și respectiv Y , arătați că

$$P(A = a, C = c | B = b, Y = y) = P(A = a | B = b, Y = y) \cdot P(C = c | B = b, Y = y),$$

pentru $\forall a \in \text{Val}(A)$ și $\forall c \in \text{Val}(C)$.

- Scrieți regula de decizie pentru acest nou clasificator de tip Bayes și arătați că ea este echivalentă cu regula de decizie a clasificatorului Bayes Optimal.
- Calculați numărul minim de parametri de estimat de către noul nostru clasificator Bayes și enumerați-i.

34. (Algoritmul Bayes Naiv: independența [condițională a] atributelor de intrare; calculul ratei medii a erorii)

□ • CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, HW3, pr. 1.1

Presupunem că A și B sunt variabile aleatoare binare independente, fiecare dintre ele având posibilitatea de a lua valoarea 0 cu o probabilitate de 50%.

Definiți o funcție booleană $y = f(A, B)$ în așa fel încât variabila A să *nu* fie independentă de variabila B în raport cu y (văzut și el ca variabilă aleatoare), însă clasificatorul Bayes Naiv să producă o rată medie a erorii de 0% (presupunând că datele de antrenament sunt în număr infinit).

Demonstrați că acest clasificator are într-adevăr rata erorii de 0%.

³⁴⁵LC: Totuși, este recomandabil să procedați așa după ce în prealabil ați văzut ce valoare produce algoritmul Bayes Optimal pentru [măcar] una dintre combinațiile de valori ale variabilelor, de exemplu, $A = B = C = 0$.

35.

(Algoritmul Bayes Naiv:
comparație cu alți clasificatori)*prelucrare de L. Ciortuz, după*** CMU, 2004 fall, T. Mitchell, Z. Bar-Joseph, midterm, pr. 6.b*

Considerăm un clasificator Bayes Naiv care lucrează pe un set de date descrise de atributele de intrare A și B și atributul de ieșire Y . (Exemplu: $Y = A \text{ XOR } B$). A și B sunt variabile aleatoare independente între ele.

a. Este posibil ca, în situația aceasta, vreun alt clasificator — de exemplu, ID3, regresia logistică (eventual kernel-izată) sau SVM — să lucreze mai bine decât clasificatorul Bayes Naiv?

b. Care este motivul?

36.

(Comparație între clasificatorul Bayes Naiv
și algoritmul ID3) $\square \bullet \circ$ *CMU, 2010 fall, Ziv Bar-Joseph, midterm, pr. 5.b*

Care dintre afirmațiile de mai jos sunt adevărate atât pentru clasificatorul Bayes Naiv cât și pentru algoritmul ID3 pentru învățarea de arbori de decizie? (Veți putea alege nu neapărat una singură dintre aceste afirmații.)

1. În cazul ambilor clasificatori se presupune că orice pereche de atribute X_i și X_j cu $i \neq j$ — văzute ca variabile aleatoare — sunt independente.
2. În cazul ambilor clasificatori se presupune că orice pereche de atribute X_i și X_j cu $i \neq j$ sunt dependente.
3. În cazul ambilor clasificatori se presupune că orice pereche de atribute sunt independente în raport cu eticheta (adică variabila care reprezintă clasa).
4. În cazul ambilor clasificatori se presupune că orice pereche de atribute sunt dependente în raport cu eticheta.

37.

(Algoritmul Bayes Naiv – clasificator de tip MAP;
o condiție [suficientă] pentru echivalența cu clasificarea de tip ML) $\square \bullet \circ$ *CMU, 2012 fall, T. Mitchell, Z. Bar-Joseph, HW3, pr. 4.2*

Algoritmul Bayes Naiv asociază unui exemplu x clasa c dacă aceasta maximizează probabilitatea $P(c|x)$.

Când este această condiție [din formularea algoritmului Bayes Naiv] echivalentă cu selecționarea acelei clase c care maximizează probabilitatea $P(x|c)$?

38.

(Algoritmii Bayes Naiv și Bayes Optimal: Adevărat sau Fals?)

 $\bullet \circ$ *CMU, 2005 spring, C. Guestrin, T. Mitchell, midterm, pr. 2.b.5*
CMU, 2011 spring, Tom Mitchell, midterm, pr. 1.1.ab

a. Clasificatorul Bayes Optimal poate să obțină rata de eroare 0 [la antrenare] pentru orice set de date. Justificați.

- b. Dacă antrenăm un clasificator Bayes Naiv folosind un număr infinit de date de antrenament care satisfac toate presuposițiile luate în calcul de acest tip de modelare (de exemplu, independența condițională), atunci acest clasificator va produce *eroare* zero pe setul de exemple de antrenament considerat.
- c. Dacă antrenăm un clasificator Bayes Naiv folosind un număr infinit de date de antrenament care satisfac toate presuposițiile luate în calcul de acest tip de modelare (de exemplu, independența condițională), atunci acest clasificator va produce *eroare* „adevărată” zero pentru exemple de test generate conform aceleiași distribuții.

2.2.3 Clasificare bayesiană [cu atribute de intrare] de tip gaussian

39. (Distribuția gaussiană unidimensională: exemplificare pentru estimarea mediei în sens MLE și calcularea unei probabilități a posteriori)

• ◦ * *U. Edinburgh, 2009 fall, C. Williams, V. Lavrenko, HW2, pr. 1*

Considerăm un set de date [de antrenament] care constă din *exemple* pentru două *clase*. Exemplele pentru clasa 1 sunt 0.5, 0.1, 0.2, 0.4, 0.3, 0.2, 0.2, 0.1, 0.35, 0.25, iar exemplele pentru clasa 2 sunt 0.9, 0.8, 0.75, 1.0.

- a. Vă cerem să antrenați (engl., fit) câte o distribuție gaussiană unidimensională pentru fiecare din cele două clase, folosind metoda estimării în sensul verosimilității maxime (engl., Maximum Likelihood estimation, MLE). Veți presupune că varianța distribuției corespunzătoare clasei 1 este 0.0149, iar varianța distribuției pentru clasa 2 este 0.0092.
- b. Estimați de asemenea în sens MLE probabilitățile *a priori* de selecție pentru [generarea de exemple din] cele două clase.
- c. Care este probabilitatea *a posteriori* ca punctul $x = 0.6$ să aparțină clasei 1?

40. (Algoritmul Bayes [Naiv] gaussian: exemplificare pe date din \mathbb{R} , granițe de decizie)

◻ • ◦ *CMU, 2009 spring, Tom Mitchell, midterm, pr. 5*

În acest exercițiu vom considera mai mulți clasificatori de tip Bayes Naiv gaussian (GNB) pentru un set de date având un singur atribut x și două clase, 0 și 1.³⁴⁶ Ca de obicei pentru clasificatori bayesieni, vom clasifica o instanță x ca aparținând clasei 1 dacă

$$P(y = 1|x) \geq P(y = 0|x) \Leftrightarrow \ln \frac{P(y = 1|x)}{P(y = 0|x)} \geq 0.$$

³⁴⁶Fiind dat un singur atribut de intrare, putem renunța la termenul „Naiv” din expresia care desemnează tipul clasificatorului.

3 Învățare bazată pe memorare

Sumar

Noțiuni preliminare

- măsuri de distanță, măsuri de similaritate: ex. 2;
- normă într-un spațiu vectorial; [măsura de] distanță indusă de către o normă: ex. 7;
- k -NN vecinătate a unui punct din \mathbb{R}^d .

Algoritmul k -NN

- pseudo-cod: cartea ML, pag. 232;
- bias-ul inductiv: „Cine se aseamănă se adună” (sau: „Spune-mi cu cine te împrietenеști, ca să-ți spun cine ești”): ex. 15.a;
- exemple (simple) de aplicare: ex. 1, ex. 2;
- complexitate de spațiu: $\mathcal{O}(dn)$
complexitate de timp:
 - la antrenare: $\mathcal{O}(dn)$
 - la testare: $\mathcal{O}(dn \log n)$
[LC: $\mathcal{O}(dnk \log k)$ pt. $k > 1$ (worst case) și $\mathcal{O}(dn)$ pt. $k = 1$],
- unde d este numărul de attribute, iar n este numărul de exemple;
- arbori kd (engl., kd -trees): *Statistical Pattern Recognition*, Andrew R. Webb, 3rd ed., 2011, Willey, pag. 163-173;
- k -NN ca algoritm ML “lazy” (vs. “eager”):
suprafețe de decizie și granițe de decizie:
diagrame Voronoi pentru 1-NN: ex. 4, ex. 11.a, ex. 18, ex. 19, ex. 20.a;
- analiza erorilor:
 - 1-NN pe date consistente: eroarea la antrenare este 0: ex. 2, ex. 12.a;
 - variația numărului de erori (la antrenare și respectiv testare) în funcție de valorile lui k : ex. 22, ex. 23.ab;
 k -NN ca metodă neparametrică; alegerea lui k : CV: ex. 23.c;
 - CVLOO: ex. 3, ex. 12.b, ex. 15.bc, ex. 16, ex. 24.a, ex. 20.b;
 - sensibilitatea / robustețea la „zgomote”: ex. 5, ex. 15;
 - eroarea asimptotică: ex. 10, ex. 25.
- efectul trăsăturilor redundante sau irelevante;
- alegerea valorii convenabile pentru k : ex. 21.

Proprietăți ale algoritmului k -NN

- (P0) output-ul algoritmului k -NN pentru o instanță oarecare de test x_q depinde de valoarea lui k : ex. 1;
- (P1) pe seturi de date de antrenament *consistente*, eroarea la antrenare produsă de algoritmul 1-NN este 0: ex. 2, ex. 12.a;
- (P2) output-ul algoritmului k -NN, precum și suprafețele de decizie și separatorii decizionali depind de *măsura de distanță* folosită: ex. 7;
- (P3) „blestemul marilor dimensiuni” (engl., the curse of dimensionality): în anumite condiții, numărul de instanțe de antrenament necesare pentru a avea un *cel mai apropiat vecin* situat la distanță *rezonabilă* față de instanța de test x_q crește exponențial în funcție de numărul de atribute folosite: ex. 9;
- (P4) în anumite condiții, rata medie a *erorii asimptotice* a algoritmului 1-NN este mărginită superior de dublul ratei medii a erorii algoritmului Bayes Optimal: ex. 10, ex. 25.

Comparații cu alți algoritmi de clasificare automată

- ID3: ex.11.b, ex. 13.ab;
- SVM: ex.12, ex. 13.c, ex. 24.b;
- regresia logistică: ex. 24.b;
- 1-NN cu mapare cu RBF: ex. 14.

Variante ale algoritmului k -NN

- k -NN folosind alte măsuri de distanță (decât distanța euclidiană): ex. 7;
- k -NN cu *ponderarea distanțelor* (engl., distance-weighted k -NN): cartea ML, pag. 236-238 (formulele 8.2, 8.3, 8.4);³⁵¹
- algoritmul lui Shepard: ex.8.

Alte metode de tip IBL

- rețele RBF: cartea ML, pag. 238-240;
- raționare bazată pe cazuri (engl., case-based reasoning): cartea ML, pag. 240-244.

³⁵¹Sectiunea 8.3 din cartea ML (pag. 236-238) se referă la regresia [liniară] local-ponderată ca o formă mai generală de aproximare a [valorilor] funcțiilor, în raport cu cele calculate de către algoritmul k -NN atunci când se folosește ponderarea distanțelor.

3.1 Învățare bazată pe memorare — Probleme rezolvate

1. (Algoritmul k -NN: aplicare în \mathbb{R}^2 pentru diferite valori ale lui k)

■ • CMU, 2006 fall, E. Xing, T. Mitchell, final exam, pr. 2

Fie setul de instanțe de antrenament din tabelul alăturat:

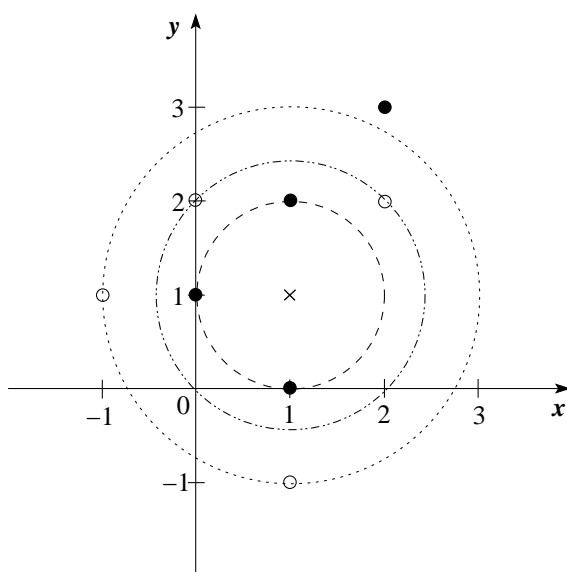
a. Vizualizați datele într-un sistem de axe din \mathbb{R}^2 .

b. Presupunând că se folosește distanța euclidiană, care va fi predicția făcută pentru punctul (1,1) de către

- clasificatorul 3-NN?
- clasificatorul 5-NN?
- clasificatorul 7-NN?

x	y	
-1	1	-
0	1	+
0	2	-
1	-1	-
1	0	+
1	2	+
2	2	-
2	3	+

Răspuns:



În figura alăturată am reprezentat:

- datele de antrenament, folosind cerculețe albe pentru cele clasificate cu - și cerculețe negre pentru cele clasificate cu +;
- punctul care trebuie clasificat, marcat cu \times ;
- vecinătățile luate în considerare de către cei trei clasificatori precizați în enunț; aceste vecinătăți sunt reprezentate prin cele trei cercuri concentrice având centrul în punctul (1,1).

Analizând pe rând etichetele instanțelor din aceste trei vecinătăți, ajungem la concluzia că rezultatele obținute de cei trei clasificatori vor fi următoarele:

- 3-NN: +
- 5-NN: +
- 7-NN: -

Observație: Acest exercițiu pune în evidență faptul că la clasificarea cu algoritmul k -NN rezultatul poate diferi în funcție de diversele valori ale lui k , adică în funcție de cum se lărgște (sau se restrânge) vecinătatea punctului de test considerat.

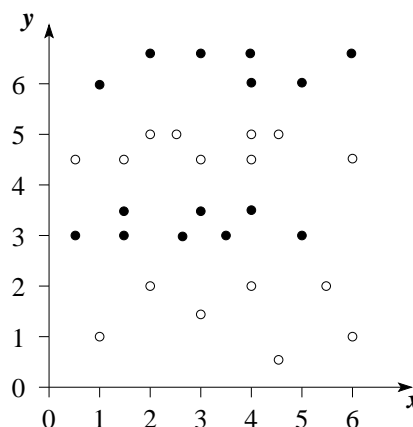
2.

(Algoritmul 1-NN: eroarea la antrenare)

CMU, 2002 fall, Andrew Moore, final exam, pr. 6.e

Figura alăturată prezintă un set de date având două atribute de intrare x și y , și un atribut de ieșire, ale cărui valori sunt reprezentate prin culoarea punctului (alb sau negru).

Putem alege o *metrică* (adică, *măsură* sau *funcție de distanță*) astfel încât, folosind algoritmul de învățare 1-NN (engl., [one] nearest neighbour), să obținem eroare 0 la antrenare pe setul acesta de date?



Răspuns:

Fie $d : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ o metrică oarecare. Conform definiției matematice, d îndeplinește următoarele *condiții*:

$d(x, y) \geq 0, \forall x, y \in \mathbb{R}^2$	(nenegativitatea)
$d(x, y) = 0 \Leftrightarrow x = y$	(identitatea indiscernabililor)
$d(x, y) = d(y, x)$	(simetria)
$d(x, y) \leq d(x, z) + d(z, y), \forall x, y, z \in \mathbb{R}^2$	(inegalitatea triunghiului)

În particular, d poate fi metrica euclidiană, definită astfel:

$$d((x_1, y_1), (x_2, y_2)) = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2},$$

pentru orice $(x_1, y_1), (x_2, y_2) \in \mathbb{R}^2$.

În general, când algoritmul 1-NN primește o instanță de test (x, y) , el caută în setul de antrenament punctul (x', y') cu proprietatea că distanța $d((x, y), (x', y'))$ este minimă.

În cazul particular în care punctul (x, y) însuși aparține datelor de antrenament, urmează că $(x', y') = (x, y)$, datorită primelor două proprietăți ale lui d enunțate mai sus.

Cum setul de date din enunț este *consistent* — adică nu există nicio instanță care să apară de două sau mai multe ori, însă cu etichete diferite —, vom avea de luat în considerare o singură valoare pentru calculul atributului de ieșire pentru punctul de test (x, y) . Evident, algoritmul 1-NN o va folosi pe aceasta ca rezultat al clasificării. Concluzia acestui raționament este că eroarea la antrenare produsă de algoritmul 1-NN pe acest set de date este 0.

Facem *observația* că nu doar pe acest set de date ci pe orice set de date de antrenament fără zgomote / inconsistențe în ce privește etichetarea, algoritmul 1-NN va avea eroarea la antrenare 0, indiferent de metrica folosită (bineînțeles, dacă există o metrică în spațiul instanțelor respective).

3. (Algoritmul k -NN: calculul erorii la CVLOO pentru diferite valori ale lui k)

CMU, 2003 fall, T. Mitchell, A. Moore, final exam, pr. 5.ab

Fie setul de date de antrenament din tabelul alăturat.

Vom folosi algoritmul k -NN cu distanța euclidiană (neponderată) pentru a prezice valorile atributului de ieșire Y (de tip boolean) plecând de la atributul de intrare X , care ia valori reale.

Care este eroarea produsă de algoritmul k -NN la cross-validare cu metoda “Leave-One-Out” în cazurile următoare:

a. $k = 1$.

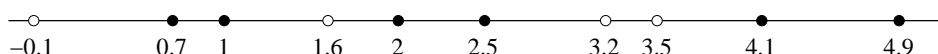
b. $k = 3$.

Exprimați răspunsul sub forma numărului de instanțe clasificate eronat.

X	Y
-0.1	-
0.7	+
1.0	+
1.6	-
2.0	+
2.5	+
3.2	-
3.5	-
4.1	+
4.9	+

Răspuns:

Figura următoare vizualizează pe o axă datele de antrenament, precum și clasificările acestora (○ pentru instanțe negative și ● pentru instanțe pozitive):



Comportamentul celor doi clasificatori la cross-validare prin metoda “Leave-One-Out” este cel explicat mai jos:

• Clasificatorul 1-NN:

Data	Eticheta	Vecinătate	Clasificare la CVLOO	Eroare?
-0.1	-	{0.7}	+	da
0.7	+	{1.0}	+	nu
1.0	+	{0.7}	+	nu
1.6	-	{2.0}	+	da
2.0	+	{1.6}	-	da
2.5	+	{2.0}	+	nu
3.2	-	{3.5}	-	nu
3.5	-	{3.2}	-	nu
4.1	+	{3.5}	-	da
4.9	+	{4.1}	+	nu

• Clasificatorul 3-NN:

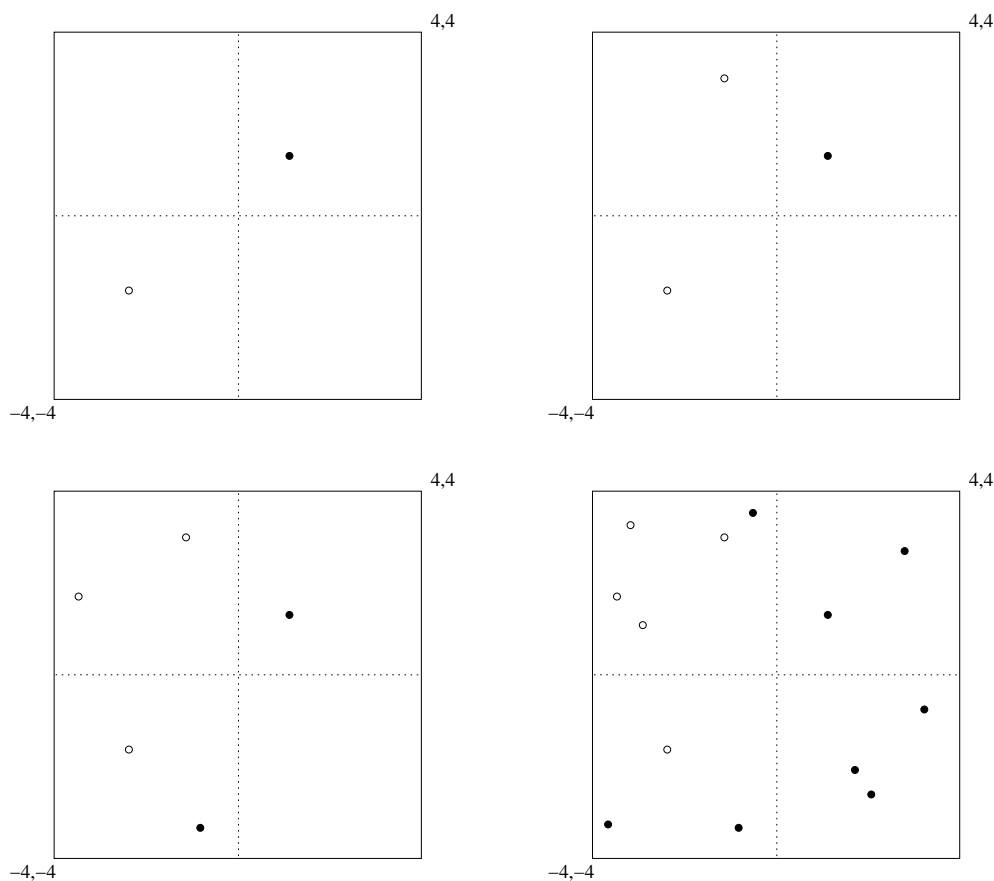
Data	Eticheta	Vecinătate	Clasificare la CVLOO	Eroare?
-0.1	-	{0.7; 1.0; 1.6}	+	da
0.7	+	{-0.1; 1.0; 1.6}	-	da
1.0	+	{0.7; 1.6; 2.0}	+	nu
1.6	-	{1.0; 2.0; 0.7/2.5}	+	da
2.0	+	{1.0; 1.6; 2.5}	+	nu
2.5	+	{1.6; 2.0; 3.2}	-	da
3.2	-	{2.5; 3.5; 4.1}	+	da
3.5	-	{2.5; 3.2; 4.1}	+	da
4.1	+	{3.2; 3.5; 4.9}	-	da
4.9	+	{3.2; 3.5; 4.1}	-	da

În concluzie, la cross-validare cu metoda “Leave-One-Out” avem 4 erori la clasificarea cu algoritmul 1-NN și 8 erori la clasificarea cu algoritmul 3-NN. Putem concluziona că rezultatul algoritmului 3-NN este foarte afectat de către puternica „mixare” (adică, de frecvențele schimbări de clasă, de la o instanță oarecare la vecinii ei) din setul de antrenament.

4. (Algoritmul 1-NN: granițe / suprafețe de decizie; diagrame Voronoi ca modalitate de învățare rapidă / “eager”)

■ • CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, HW1, pr. 3.1-2

În fiecare din figurile următoare se dau câteva puncte în spațiul euclidian bidimensional. Fiecare dintre aceste puncte este etichetat fie pozitiv (cerc plin) fie negativ (cerc simplu).



a. Presupunând că folosim ca metrică distanța euclidiană, desenați suprafețele de decizie corespunzătoare clasificatorului 1-NN, în fiecare din aceste patru cazuri.

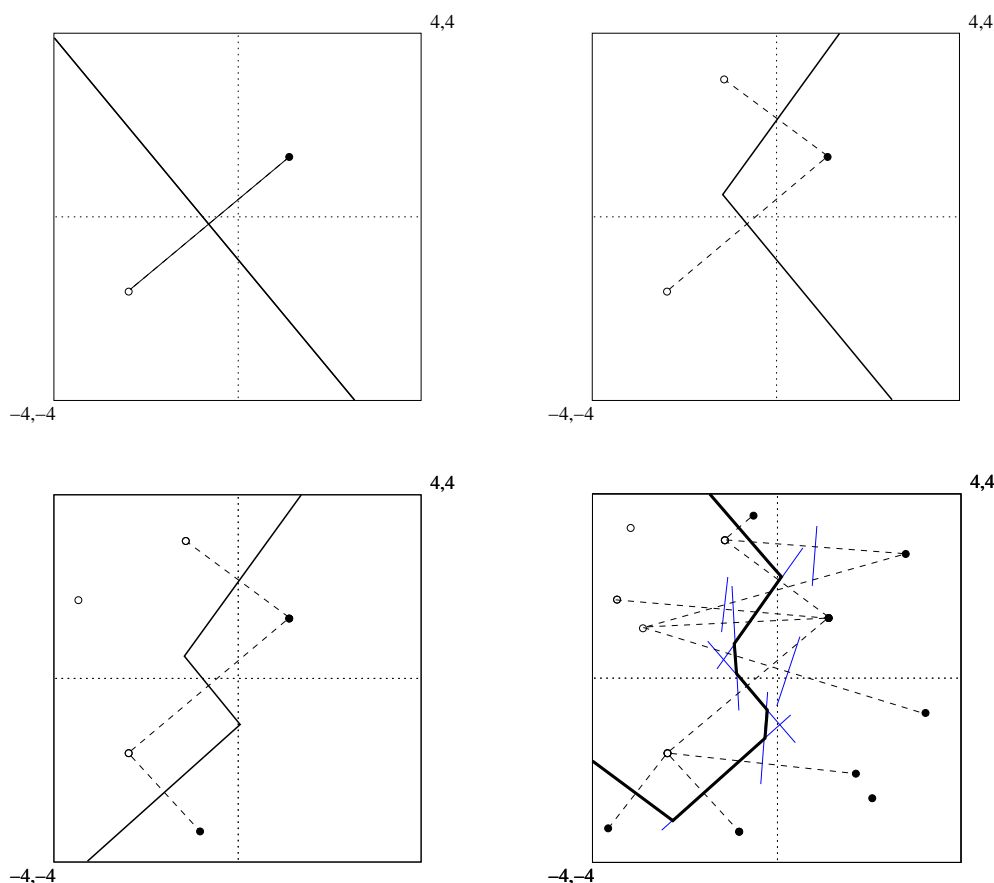
b. La curs am afirmat despre k -NN că este un clasificator lent (engl., lazy), care doar memorează toate instanțele de antrenament până ajunge la faza de test.

Totuși, la punctul precedent am văzut că putem să trasăm suprafețele de decizie pentru clasificatorul 1-NN, înainte de a intra în faza de test / generalizare. Atunci, în această fază, în loc să calculăm diverse distanțe și apoi să determinăm care sunt cei mai apropiați vecini față de punctul de test dat (notat x_q), pur și simplu i se va asigna lui x_q clasa / eticheta corespunzătoare suprafeței de decizie în care se plasează.

Dacă am decide să memorăm toate aceste suprafețe de decizie (ca linii poligonale) în loc să memorăm toate datele de antrenament, am obține *întotdeauna* o îmbunătățire în ceea ce privește consumul de memorie necesar pentru acest clasificator?

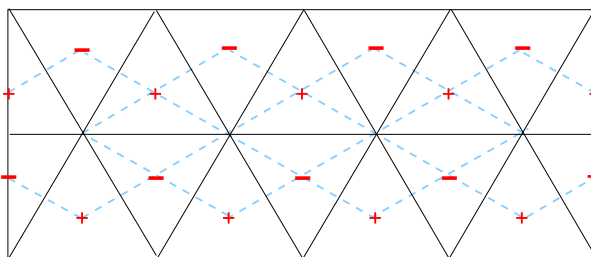
Răspuns:

a. Suprafețele de decizie corespunzătoare clasificatorului 1-NN pentru cele patru cazuri sunt:



b. Răspunsul este negativ, adică: nu întotdeauna memorarea separatorului decizional este mai convenabilă decât memorarea datelor de antrenament. Vom justifica dând un exemplu, care reprezintă o situație-limită, și anume, cazul care apare atunci când se creează tot atâtea suprafețe de decizie câte instanțe sunt în setul de antrenament.

În figura alăturată, pentru cele n instanțe de antrenament avem nevoie să memorăm $3\frac{n}{2}$ puncte care determină triunghiurile – suprafețele de decizie pentru clasa +. (Clasa – va fi determinată prin exclusiune / complementaritate.)



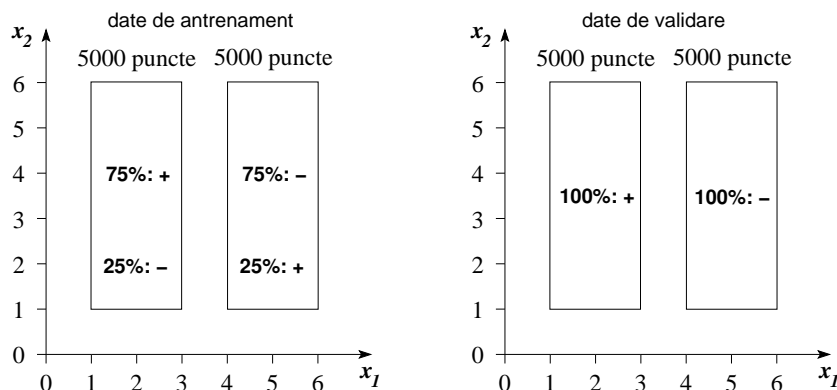
Se poate demonstra ușor următoarea proprietate: Chiar dacă am memora o singură dată cele $3\frac{n}{2}$ puncte — fiecare punct având câte 2 coordonate — care determină contururile suprafețelor de decizie și apoi am folosi indecși pentru a indica ce puncte determină fiecare suprafață de decizie, consumul de memorie ar fi mai mare decât dacă am memora doar instanțele de antrenament.

5.

(Algoritmul k -NN:
rata medie a erorii la antrenare, la CVLOO și la testare
în prezența unor „zgomote” în datele de antrenament)

• ◦ CMU, 2002 fall, Andrew Moore, final exam, pr. 7

Se constituie un set de date de antrenament și un set de date de validare, alegând în mod uniform aleatoriu puncte situate în anumite regiuni dreptunghiulare, conform imaginii următoare:



Se observă că datele de antrenament sunt „bruiate” (engl., noisy): în fiecare regiune, 25% din date provin din clasa adversă. Datele de validare nu sunt bruiate.

Pentru fiecare dintre întrebările următoare, încercuiți fracția care este *cea mai apropiată de media* / rata erorii în cazul respectiv. Justificați în mod riguros alegerea făcută.

a. Care este rata medie a erorii la antrenare pentru clasificatorul 1-NN ?

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

b. Care este rata medie a erorii la cross-validare cu metoda “Leave-One-Out” pentru clasificatorul 1-NN pe setul de antrenare?

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

- c. Care este rata medie a erorii la testare pentru clasificatorul 1-NN pe setul de validare? (Antrenarea se face pe setul de antrenare.)

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

- d. Care este rata medie a erorii la antrenare pentru clasificatorul 21-NN?

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

- e. Care este rata medie a erorii la cross-validare cu metoda “Leave-One-Out” pentru clasificatorul 21-NN pe setul de antrenare?

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

- f. Care este rata medie a erorii la testare pentru clasificatorul 21-NN pe setul de validare? (Antrenarea se face pe setul de antrenare.)

0 1/8 1/4 3/8 1/3 1/2 5/8 2/3 3/4 7/8 1

Răspuns:

a. În acest caz, rata medie a erorii este 0, întrucât orice punct din datele de antrenament este cel mai apropiat vecin în raport cu el însuși, iar probabilitatea ca două puncte care au fost generate în mod aleatoriu (în oricare din cele două dreptunghiuri) să aibă exact aceleași coordonate și aceeași etichetă este foarte mică (este practic 0).

b. 3/8.

Să analizăm cazul primului dreptunghi. Constituie eroare cazul când un exemplu pozitiv (care apare cu probabilitatea „așteptată” de 3/4) este clasificat negativ (iar aceasta se întâmplă cu probabilitatea 1/4) sau, invers, când un exemplu negativ (probabilitate: 1/4) este clasificat pozitiv (probabilitate: 3/4). Așadar, probabilitatea de a clasifica greșit un exemplu din primul dreptunghi o aflăm astfel: $3/4 \cdot 1/4 + 1/4 \cdot 3/4 = 3/8$. Pentru dreptunghiul din dreapta raționamentul este similar, iar selecția exemplurilor din cele două dreptunghiuri se face cu aceeași probabilitate (1/2), deci rezultatul final este 3/8.

c. 1/4.

Probabilitatea de eroare este 1/4 pentru fiecare dreptunghi, fiindcă (în general) pentru un punct (x_1, x_2) selectat în mod aleatoriu din datele de validare din dreptunghiul respectiv aceasta (adică 1/4) este probabilitatea ca vecinul cel mai apropiat de (x_1, x_2) în setul de date de antrenament să aibă semnul opus (față de semnul lui (x_1, x_2)). Așadar, media erorii la testare cu clasificatorul 1-NN pe setul de validare este $1/2 \cdot 1/4 + 1/2 \cdot 1/4 = 1/4$.

d. 1/4.

Folosind algoritmul 21-NN, în dreptunghiul din stânga 3/4 din datele de antrenament sunt (în general) clasificate corect, iar restul de 1/4 (și anume, cele având semnul $-$) sunt clasificate eronat. Într-adevăr, la testarea unui punct oarecare (x_1, x_2) din acel dreptunghi, în 21-NN vecinătatea lui (x_1, x_2) , unul (cel mai apropiat vecin) este însuși (x_1, x_2) , iar dintre ceilalți 20 cei mai apropiați vecini 3/4 sunt (în general) de semn $+$, iar 1/4 de semn $-$. Analog se raționează pentru dreptunghiul din dreapta. Deci rata medie a erorii la antrenare este:

$$\frac{1}{2} \left(\frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 1 \right) + \frac{1}{2} \left(\frac{3}{4} \cdot 0 + \frac{1}{4} \cdot 1 \right) = \frac{1}{4}.$$

e. $1/4$.

În cazul primului dreptunghi, apare eroare atunci când un exemplu pozitiv (probabilitate de selectare: $3/4$) este clasificat negativ (probabilitate: 0, fiind vorba de 21-NN) sau când un exemplu negativ (probabilitate de selectare: $1/4$) este clasificat pozitiv (probabilitate: 1).³⁵² Așadar, probabilitatea de a clasifica greșit un exemplu din primul dreptunghi este $3/4 \cdot 0 + 1/4 \cdot 1 = 1/4$. Cazul dreptunghiului din dreapta este similar, prin urmare rezultatul final este $1/4$.

f. 0.

Datorită distribuției uniforme a datelor de antrenament, fiecare instanță (x_1, x_2) din setul de validare va avea majoritatea vecinilor — dintre cei mai apropiați, selectați de către algoritmul 21-NN din setul de date de antrenament — de același semn cu semnul lui (x_1, x_2) . (Și anume, în medie de 3 ori mai mulți vecini decât cei de semn contrar.) Prin urmare, fiecare punct din setul de validare va fi clasificat corect.

Observație (1): Comparând rezultatele obținute la punctele a și c pe de o parte cu cele de la punctele d și f (sau chiar b și e) pe de altă parte, se observă că are loc o relație exact de același gen cu cea din definiția fenomenului de *overfitting* (sau, *supra-specializare*):³⁵³ erorile produse de algoritmi 1-NN și 21-NN la antrenare sunt în relația $0 < 1/4$ dar în relație inversă ($1/4 > 0$) la testare (respectiv $3/8 > 1/4$ la cross-validare).

Observație (2): Remarcați „sublinerea“ din expresia „fracția cea mai apropiată de media / rata erorii“ din enunț. Generarea aleatorie a datelor poate conduce la rezultate ușor diferite față de cele pe care le-am obținut mai sus. În urma realizării unei implementări,³⁵⁴ au fost obținute următoarele rezultate de tip *eroare medie*, calculată în urma repetării de 100 de ori a generării datelor și aplicării algoritmilor k -NN, conform cerințelor din enunț:

a. 0, b. 0.374022, c. 0.250472, d. 0.249342, e. 0.253088, f. 0.006436.

Constatăm că se verifică „previziunile“ noastre din rezolvarea de mai sus.

³⁵²În 21-NN vecinătatea punctului considerat (pentru testare), spre deosebire de cazul erorii la antrenare, la CVLOO nu se mai include punctul respectiv.

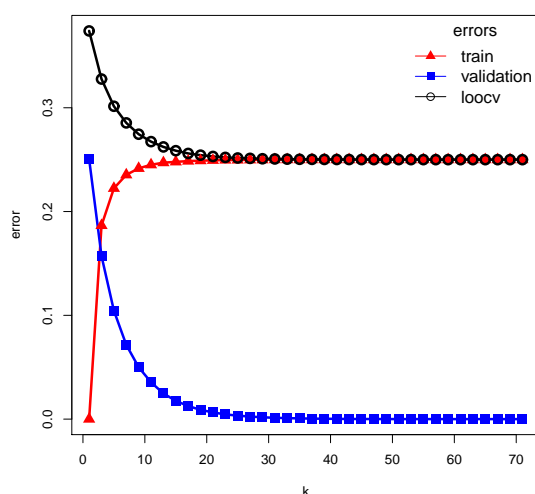
³⁵³Vă readucem aminte definiția lui Tom Mitchell pentru fenomenul de *overfitting* (vedeți cartea *Machine Learning*, pag. 67): două ipoteze h și h' obținute (eventual cu un același algoritm de clasificare automată) pe un set de date sunt în raport de overfitting dacă

$$error_{train}(h) < error_{train}(h'), \quad \text{dar} \quad error_{\mathcal{D}}(h) > error_{\mathcal{D}}(h'),$$

unde \mathcal{D} este distribuția reală a datelor.

³⁵⁴Implementarea a fost făcută de către studentul Sebastian Ciobanu de la Facultatea de Informatică a Universității „Al. I. Cuza“ din Iași în semestrul I al anului universitar 2016-2017.

Evoluția celor trei tipuri de eroare (la antrenare, la validare și cross-validare cu metoda “leave-one-out”), pentru $k = 1, 3, \dots, 69, 71$ este prezentată în figura alăturată. Se observă convergența la aceeași valoare (aprox., 0.25) pentru eroarea la antrenare și eroarea CV-LOO (pe setul de date de antrenare) începând din jurul valorii $k = 31$, precum și convergența la valori foarte apropiate de 0 pentru eroarea la testare pe setul de date de validare, începând din jurul valorii $k = 35$.



6. (O versiune ipotetică pentru algoritmul k -NN: selectarea celor mai apropiați vecini nu pe baza calculării funcției de distanță, ci folosind un oracol / “black box”)

CMU, 2009 spring, Ziv Bar-Joseph, final exam, pr. 2.2-3

Se încearcă clasificarea unor puncte în spațiul euclidian bidimensional. Sunt date n instanțe P_1, P_2, \dots, P_n , precum și etichetările corespunzătoare c_1, c_2, \dots, c_n , unde c_1, c_2, \dots, c_n reprezintă valori dintr-o mulțime C . În schema de clasificare k -NN, fiecare element nou Q este clasificat cu eticheta majoritară obținută în cadrul vecinătății formate din cei mai apropiați k vecini.

Să presupunem că măsura de distanță nu este dată în mod explicit. În locul acesteia, aveți la dispoziție un “black box”. Dacă se introduc instanțele $P_{i_1}, P_{i_2}, \dots, P_{i_l}$ (unde l este un număr natural oarecare) și un punct Q , black box-ul returnează cel mai apropiat vecin al lui Q , adică un element $P_{i_0} \in \{P_{i_1}, P_{i_2}, \dots, P_{i_l}\}$, precum și clasificarea corespunzătoare (c_{i_0}).

a. Este posibil să se construiască un algoritm de tip k -NN bazat doar pe acest black box? Dacă da, explicați cum anume, iar dacă nu, explicați de ce nu este posibil.

b. Dacă, în schimb, acel black box returnează cei mai apropiați j vecini (și etichetele corespunzătoare), iar $j \neq k$, este posibil să se construiască un algoritm de tip k -NN bazat doar pe acest black box? Dacă da, explicați cum anume, iar dacă nu, explicați de ce nu este posibil.

Răspuns:

a. Da. Se folosește black box-ul dându-i ca intrare mai întâi mulțimea de exemple P_1, P_2, \dots, P_n ; se obține cel mai apropiat vecin al lui Q și eticheta sa. Apoi se scoate instanța / punctul returnat din mulțimea de exemple. Se repetă acest procedeu de k ori, iar la final se va alege pentru Q eticheta corespunzătoare majorității din mulțimea de k vecini care au fost identificați de către black box.

b. Dacă $j < k$, atunci se folosește black box-ul de $\lceil k/j \rceil$ ori, obținându-se $j * \lceil k/j \rceil$ cei mai apropiați vecini și clasificările acestora. Notăm cu V_1 mulțimea formată din acești vecini. În cazul în care $k \neq j * \lceil k/j \rceil$, pentru a obține restul de vecini necesari pentru k -NN, adică încă $k - j * \lceil k/j \rceil$ vecini, vom folosi black box-ul încă o dată. Vom nota cu V_2 mulțimea acestor noi j vecini. Dintre aceștia va trebui să alegem doar $k - j * \lceil k/j \rceil$ instanțe. Vom proceda astfel: considerăm o nouă mulțime de instanțe alcătuită din elementele lui V_2 și $j - (k - j * \lceil k/j \rceil)$ dintre toți ceilalți vecini obținuți anterior (V_1). Aplicăm black box-ul pe acest nou set de date și vom obține cei mai apropiați j vecini pentru punctul Q . Printre aceștia se vor afla cei $k - j * \lceil k/j \rceil$ vecini căutați.

Dacă $j > k$ iar black box-ul nu acceptă intrări duplicate, atunci el nu poate fi folosit pentru a determina cei mai apropiați k vecini ai instanței de test.

Dar dacă $j > k$ și black box-ul acceptă intrări duplicate, este posibil să rezolvăm problema, cel puțin în situația în care cei mai apropiați j vecini ai lui Q se află toți la distanțe diferite față de acesta.³⁵⁵ De exemplu, pentru $k = 1$ și $j = 3$ vom proceda astfel:

Pasul 1: Aplicăm black box-ul inputului P_1, \dots, P_n . Vom obține outputul $P_{i_1}, P_{i_2}, P_{i_3}$.

Pasul 2: Corespunzător fiecărui punct $P_{i_1}, P_{i_2}, P_{i_3}$, vom aplica pe rând black box-ul inputului

- $P_{i_1}, P_{i_1}, P_{i_2}, P_{i_3}$,
- $P_{i_1}, P_{i_2}, P_{i_2}, P_{i_3}$ și respectiv
- $P_{i_1}, P_{i_2}, P_{i_3}, P_{i_3}$.

Se poate constata că într-unul singur din aceste cazuri black box-ul returnează outputul $P_{i_1}, P_{i_2}, P_{i_3}$.³⁵⁶ De pildă, dacă outputurile în aceste trei cazuri sunt

- $P_{i_1}, P_{i_1}, P_{i_2}$,
- $P_{i_1}, P_{i_2}, P_{i_2}$,
- $P_{i_1}, P_{i_2}, P_{i_3}$,

rezultă că P_{i_3} este cel mai distant dintre cei trei vecini ai lui Q .

Pasul 3: Considerând că lucrurile stau ca la finalul pasului precedent, pentru fiecare dintre punctele P_{i_1}, P_{i_2} vom aplica black box-ului următorul input:

- $P_{i_1}, P_{i_1}, P_{i_1}, P_{i_2}$ și respectiv
- $P_{i_1}, P_{i_2}, P_{i_2}, P_{i_2}$.

Dacă vom obține outputul $P_{i_1}, P_{i_1}, P_{i_1}$, va rezulta că P_{i_1} este mai apropiat de Q decât punctul P_{i_2} . Invers, dacă obținem outputul $P_{i_2}, P_{i_2}, P_{i_2}$, va rezulta că P_{i_2} este mai apropiat de Q decât punctul P_{i_1} .

³⁵⁵Rămâne de analizat varianta contrară.

³⁵⁶Am presupus că black box-ul returnează exact j instanțe, chiar dacă între P_1, \dots, P_n există mai multe instanțe egal depărtate față de punctul Q , în speță mai multe instanțe situate la maximumul distanțelor dintre fiecare din cele j pe de o parte și punctul Q pe de altă parte.

7. (Algoritmul 1-NN: suprafețele de decizie [și separatorii decizionali] depind de măsurile de distanță folosite)

CMU, 2008 fall, Eric Xing, HW1, pr. 3.1.2

Se dau două puncte din spațiul euclidian bidimensional: punctul $(-1, 0)$ clasificat negativ și punctul $(1, 0)$ clasificat pozitiv.

Clasificatorul 1-NN care folosește distanța euclidiană și are ca set de date de antrenament cele două puncte de mai sus are următoarea *formă analitică* (ușor de dedus):

- dat fiind un punct arbitrar (x, y) ,
în cazul în care $x > 0$, eticheta asignată punctului respectiv este $+$, iar în cazul $x < 0$ eticheta asignată este $-$;
- dreapta $x = 0$ este *granița de decizie* (engl., decision boundary) corespunzătoare acestui clasificator.

a. Care va fi forma analitică a clasificatorului 1-NN dacă în locul distanței euclidiene (indusă de norma L_2) se folosește distanța Manhattan (indusă de norma L_1)? Vă reamintim că distanța Manhattan dintre două puncte (x_1, y_1) și (x_2, y_2) este $|x_1 - x_2| + |y_1 - y_2|$.

b. Dar dacă se folosește distanța [indusă de norma] L_∞ , definită în \mathbb{R}^2 prin

$$d((x_1, y_1), (x_2, y_2)) = \max\{|x_1 - x_2|, |y_1 - y_2|\} ?$$

Răspuns:

a. Putem exprima distanța Manhattan dintre un punct oarecare (x, y) din plan și de fiecare dintre cele două puncte date în enunț — punctul $(1, 0)$ clasificat pozitiv și punctul $(-1, 0)$ clasificat negativ — astfel:

$$\begin{aligned} d_+ &\stackrel{\text{not.}}{=} d((x, y), (1, 0)) = |x - 1| + |y| \\ d_- &\stackrel{\text{not.}}{=} d((x, y), (-1, 0)) = |x + 1| + |y| \end{aligned}$$

În consecință, a stabili care dintre cele două distanțe (d_+ și d_-) este mai mică revine la a compara (doar) expresiile $|x - 1|$ și $|x + 1|$.

Avem următoarele cazuri:

- dacă $x > 1$, atunci

$$\left. \begin{aligned} d_+ &= x - 1 + |y| \\ d_- &= x + 1 + |y| \end{aligned} \right\} \Rightarrow d_+ < d_- \Rightarrow \text{punctul } (x, y) \text{ va fi clasificat } +$$

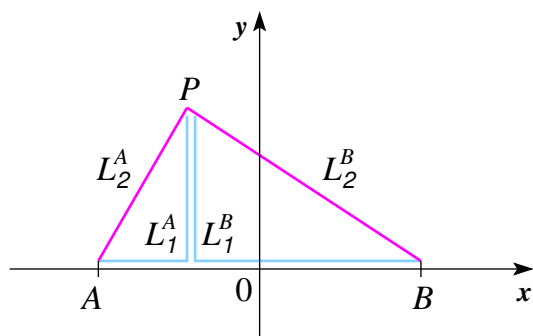
- dacă $x < -1$:

$$\left. \begin{aligned} d_+ &= -x + 1 + |y| \\ d_- &= -x - 1 + |y| \end{aligned} \right\} \Rightarrow d_+ > d_- \Rightarrow \text{punctul } (x, y) \text{ va fi clasificat } -$$

- dacă $-1 \leq x \leq 1$:

$$\left. \begin{aligned} d_+ &= -x + 1 + |y| \\ d_- &= x + 1 + |y| \end{aligned} \right\} \Rightarrow \begin{cases} d_+ < d_- & \text{pentru } x > 0 \text{ deci } (x, y) \text{ va fi clasificat } + \\ d_+ > d_- & \text{pentru } x < 0 \text{ deci } (x, y) \text{ va fi clasificat } - \\ d_+ = d_- & \text{dacă și numai dacă } x = 0. \end{cases}$$

Rezultă că granița de decizie a acestui clasificator este dreapta $x = 0$. Chiar mai mult: forma analitică a clasificatorului 1-NN care folosește distanța Manhattan pe setul de date din enunț este aceeași cu a clasificatorului 1-NN care folosește distanța euclidiană.



Observație: La concluzia de mai sus se putea ajunge mult mai ușor ținând cont de următoarea *proprietate* (demonstrabilă imediat pe cale geometrică): pentru orice două puncte A și B din \mathbb{R}^2 au loc următoarele relații de echivalență:

$$L_1(P, A) < L_1(P, B) \Leftrightarrow L_2(P, A) < L_2(P, B) \text{ pentru } \forall P \in \mathbb{R}^2 \text{ și}$$

$$L_1(P, A) = L_1(P, B) \Leftrightarrow L_2(P, A) = L_2(P, B) \text{ pentru } \forall P \in \mathbb{R}^2.$$

b. Dacă se folosește distanța L_∞ , atunci conform definiției acestei metrici vom avea:

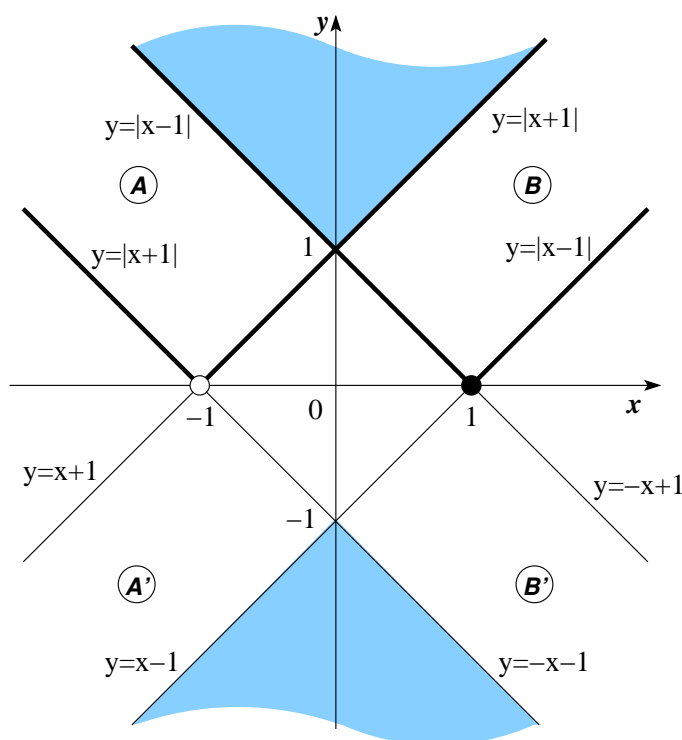
$$d_+ = d((x, y), (1, 0))$$

$$= \max\{|x - 1|, |y|\}$$

$$d_- = d((x, y), (-1, 0))$$

$$= \max\{|x + 1|, |y|\}.$$

Mai întâi vom trasa graficele funcțiilor $y = |x - 1|$ și $y = |x + 1|$, și vom obține rezultatul din figura alăturată.



Cazul i: Se observă că $|y| \geq |x - 1|$ și $|y| \geq |x + 1|$ pentru toate punctele (x, y) situate în zonele unghiulare hașurate. În consecință, d_+ va fi egal cu d_- pentru orice punct din aceste zone. Așadar, zonele hașurate vor aparține graniței / suprafeței de decizie a clasificatorului nostru.

Cazul ii: Considerăm acum zonele identificate prin literele A, A', B și B' în figura de mai sus.³⁵⁷ Se observă ușor că pentru zonele A și A' avem $d_+ =$

³⁵⁷ Analitic, zona A este definită de punctele (x, y) care satisfac inecuațiile $|x + 1| < |y| < |x - 1|$ și $y > 0$; zona A' : $|x + 1| < |y| < |x - 1|$ și $y < 0$; zona B : $|x - 1| < |y| < |x + 1|$ și $y > 0$; zona B' : $|x - 1| < |y| < |x + 1|$ și $y < 0$.

$\max\{|x-1|, |y|\} = |x-1|$ și $d_- = \max\{|x+1|, |y|\} = |y|$, iar pentru zonele B și B' avem $d_+ = \max\{|x-1|, |y|\} = |y|$ și $d_- = \max\{|x+1|, |y|\} = |x+1|$. În consecință, $d_- = |y| < |x-1| = d_+$ pentru zonele A și A' , iar $d_+ = |y| < |x+1| = d_-$ pentru zonele B și B' . Așadar, zonele A și A' vor fi clasificate negativ, iar zonele B și B' vor fi clasificate pozitiv.

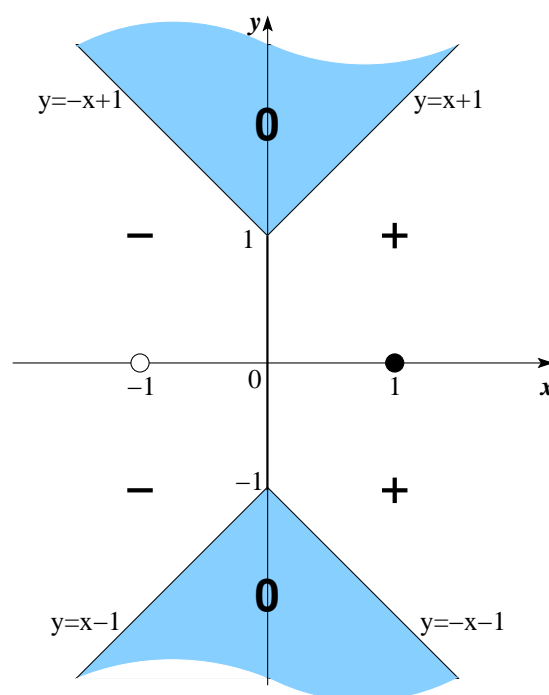
Cazul iii: Pentru toate celelalte zone rămase în discuție — adică pentru orice punct (x, y) situat în afara celor două zone hașurate și în afara zonelor A, A', B, B' —, vom avea $d_+ = \max\{|x-1|, |y|\} = |x-1|$ și $d_- = \max\{|x+1|, |y|\} = |x+1|$. Din grafic se observă că $|x+1| > |x-1|$ pentru $x > 0$ și $|x-1| > |x+1|$ pentru $x < 0$, iar $|x+1| = |x-1|$ pentru $x = 0$. Așadar, sumărind, în zone avem: $d_+ = \min\{d_+, d_-\}$ pentru $x > 0$ și $d_- = \min\{d_+, d_-\}$ pentru $x < 0$, iar $d_+ = d_-$ pentru $x = 0$.

Concluzionând, forma analitică a clasificatorului 1-NN care folosește distanța L_∞ este următoarea:

- (x, y) va fi etichetat cu $+$ dacă $x > 0$ și $-x-1 < y < x+1$;
- (x, y) va fi etichetat cu $-$ dacă $x < 0$ și $x-1 < y < -x+1$;
- în rest este vorba de suprafața de separare, adică locul geometric al punctelor (x, y) pentru care distanța față de cele două puncte din enunț este egală.

Reprezentarea grafică a suprafețelor de decizie este dată în figura alăturată.

Este de remarcat faptul că pentru acest clasificator granița / suprafața de decizie nu este formată doar din drepte, ci este reuniunea unei drepte (axa Oy) cu două intersecții de (câte două) semiplane.



8. (Algoritmul / metoda lui Shepard – aplicare)
CMU, 2002 fall, Andrew Moore, final exam, pr. 6.a-d

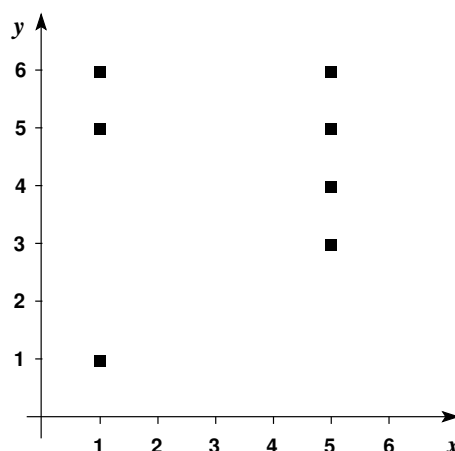
Figura de mai jos prezintă un set de date de antrenament cu un atribut de intrare $x \in \mathbb{R}$ și un atribut de ieșire $y \in \mathbb{R}$.

Vom estima din aceste date câteva valori ale unei funcții continue $f : \mathbb{R} \rightarrow \mathbb{R}$, folosind *metoda lui Shepard*. Aceasta este o variantă (de tip *regresie*) a algoritmului k -NN în care se iau în considerare toate punctele de antrenament, dar se aplică ponderi în funcție de distanță:

$$\hat{f}(x) \leftarrow \frac{\sum_i w(x, x_i) f(x_i)}{\sum_i w(x, x_i)}$$

Se va considera

$$w(x, x_i) = \begin{cases} 1, & \text{dacă } |x - x_i| \leq 3 \\ 0, & \text{în rest.} \end{cases}$$



Care va fi valoarea prezisă pentru funcția f pentru

- | | |
|--------------|--------------|
| a. $x = 1$? | c. $x = 5$? |
| b. $x = 3$? | d. $x = 6$? |

Răspuns:

Din modul cum au fost definite ponderile w ajungem la concluzia că valoarea lui f pentru un punct oarecare x va fi calculată ca medie aritmetică a valorilor / componentelor y din acele date de antrenament pentru care abscisa (x') este situată la distanță de cel mult 3 unități de punctul x care ne interesează.

a. Avem $|1 - 1| = 0 \leq 3$ și $|1 - 5| = 4 > 3$, prin urmare vor fi luate în considerare doar valorile învățate pentru $x = 1$.

$$\hat{f}(1) = \frac{1 + 5 + 6}{3} = \frac{12}{3} = 4.$$

b. Avem $|3 - 1| = 2 \leq 3$ și $|3 - 5| = 2 \leq 3$, prin urmare vor fi luate în considerare și valorile învățate pentru $x = 1$ și cele pentru $x = 5$.

$$\hat{f}(3) = \frac{1 + 5 + 6 + 3 + 4 + 5 + 6}{7} = \frac{30}{7}.$$

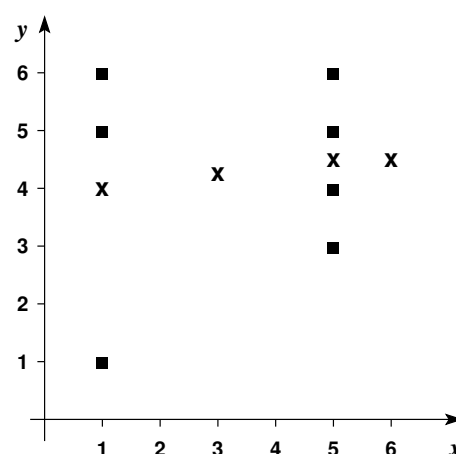
c. Avem $|5 - 1| = 4 > 3$ și $|5 - 5| = 0 \leq 3$, prin urmare vor fi luate în considerare doar valorile învățate pentru $x = 5$.

$$\hat{f}(5) = \frac{3 + 4 + 5 + 6}{4} = \frac{18}{4} = \frac{9}{2}.$$

d. Avem $|6 - 1| = 5 > 3$ și $|6 - 5| = 1 \leq 3$, prin urmare vor fi luate în considerare doar valorile învățate pentru $x = 5$, ca și în cazul precedent.

$$\hat{f}(6) = \hat{f}(5) = 4.5$$

Dacă vom plasa rezultatele de mai sus pe grafic și vom reprezenta punctele $(x, f(x))$ sub forma unor cruciulițe, vom obține figura alăturată.



9.

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

■ □ ● CMU, 2010 fall, Aarti Singh, HW2, pr. 2.2

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 unitatea³⁵⁸ 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 $\{x_1, x_2, \dots, x_n\}$:

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

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.

a. În cazul particular $p = 1$, calculați expresia *funcției de distribuție cumulativă*³⁵⁹ 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 :

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

³⁵⁸Termenul folosit în limba engleză pentru o astfel de sferă este *unit ball*.

³⁵⁹Engl., cumulative distribution function, c.d.f.

unde Γ reprezintă funcția Gamma a lui Euler,³⁶⁰ care are proprietățile:

$$\Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}, \quad \Gamma(1) = 1, \quad \text{iar } \Gamma(x+1) = x\Gamma(x) \text{ pentru } x > 0.^{361}$$

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 ?

Răspuns:

a. Pentru $p = 1$, sfera de rază 1 este intervalul $[-1, 1]$, iar funcția de distribuție cumulativă va avea expresia:

$$F_{n,1}(t) \stackrel{not.}{=} P(d^* \leq t) = 1 - P(d^* > t) = 1 - P(\|x_i\| > t, i = 1, 2, \dots, n)$$

Ținem cont de presupoziția de independență la generarea punctelor x_i , rezultă:

$$F_{n,1}(t) = 1 - \prod_{i=1}^n P(\|x_i\| > t) = 1 - (1-t)^n.$$

³⁶⁰Vedeți problema 28.b de la capitolul de *Fundamente*.

³⁶¹Se verifică ușor că pentru $p = 3$ se obține volumul sferei: $V_3(r) = \frac{(r\sqrt{\pi})^3}{\frac{3}{4}\sqrt{\pi}} = \frac{4\pi r^3}{3}$. Pentru demonstrarea unora dintre proprietățile funcției Γ indicate mai sus, vedeți problema 41 de la capitolul de *Fundamente*.

b. În cazul general, adică pentru p un număr natural oarecare nenul, fixat, vom exprima, mai întâi $P(d^* \leq t)$ exact ca mai înainte:

$$\begin{aligned} F_{n,p}(t) \stackrel{not.}{=} P(d^* \leq t) &= 1 - P(d^* > t) = 1 - P(\|x_i\| > t, i = 1, 2, \dots, n) \\ &\stackrel{indep. cdt.}{=} 1 - \prod_{i=1}^n P(\|x_i\| > t). \end{aligned}$$

Apoi, ținând cont de presupuziția de uniformitate la generarea punctelor x_i și, folosind notația $V_p(t)$ pentru volumul sferei de rază t , obținem:

$$F_{n,p}(t) = 1 - \left(\frac{V_p(1) - V_p(t)}{V_p(1)} \right)^n = 1 - \left(1 - \frac{V_p(t)}{V_p(1)} \right)^n.$$

În sfârșit, folosind formula *sugerată* în enunț pentru V_p , rezultă imediat că $F_{n,p}(t) = 1 - (1 - t^p)^n$.

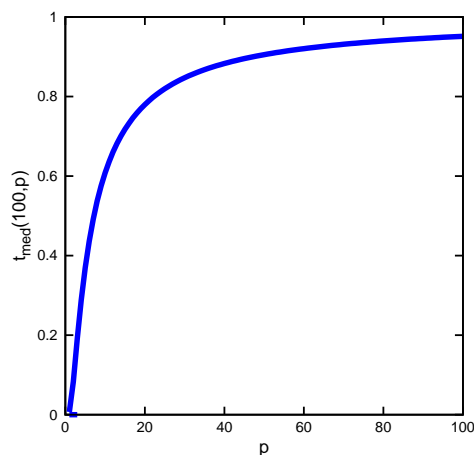
c. Pentru a afla valoarea mediană corespunzătoare variabilei aleatoare d^* , vom rezolva ecuația $P(d^* \leq t) = 1/2$ în funcție de t :

$$\begin{aligned} P(d^* \leq t) = \frac{1}{2} &\Leftrightarrow F_{n,p}(t) = \frac{1}{2} \stackrel{b.}{\Leftrightarrow} 1 - (1 - t^p)^n = \frac{1}{2} \\ &\Leftrightarrow (1 - t^p)^n = \frac{1}{2} \Leftrightarrow 1 - t^p = \frac{1}{2^{1/n}} \\ &\Leftrightarrow t^p = 1 - \frac{1}{2^{1/n}} \end{aligned}$$

Prin urmare,

$$t_{med}(n, p) = \left(1 - \frac{1}{2^{1/n}} \right)^{1/p}.$$

Graficul funcției $t_{med}(100, p)$ pentru $p = 1, 2, \dots, 100$ este cel din figura alăturată. Se observă că sfera minimală care conține cel mai apropiat vecin (un x_i , cu $i \in \{1, 2, \dots, n\}$) al originii O se lărgeste foarte repede pe măsură ce p crește. Pentru valori ale lui p mai mari decât 10, majoritatea dintre cele 100 de instanțe de antrenament sunt mai aproape de conturul sferei de rază 1 decât de originea O .



d. Putem scrie următorul șir de echivalențe:

$$\begin{aligned} P(d^* \leq 0.5) \geq 0.9 &\Leftrightarrow F_{n,p}(0.5) \geq 0.9 \Leftrightarrow \\ &\stackrel{b.}{\Leftrightarrow} 1 - \left(1 - \frac{1}{2^p} \right)^n \geq \frac{9}{10} \Leftrightarrow \left(1 - \frac{1}{2^p} \right)^n \leq \frac{1}{10} \\ &\Leftrightarrow n \cdot \ln \left(1 - \frac{1}{2^p} \right) \leq -\ln 10 \end{aligned}$$

$$\Leftrightarrow n \geq \frac{\ln 10}{-\ln\left(1 - \frac{1}{2^p}\right)}$$

Se poate vedea imediat că membrul din partea dreaptă a inegalității de mai sus tinde la $+\infty$ pentru $p \rightarrow \infty$. Este necesar să vedem *cât de repede* are loc această tendere la infinit. Pentru aceasta, vom folosi descompunerea lui $-\ln(1 - 1/2^p)$ sub forma unei serii Taylor (luând $x = 1/2^p$):³⁶²

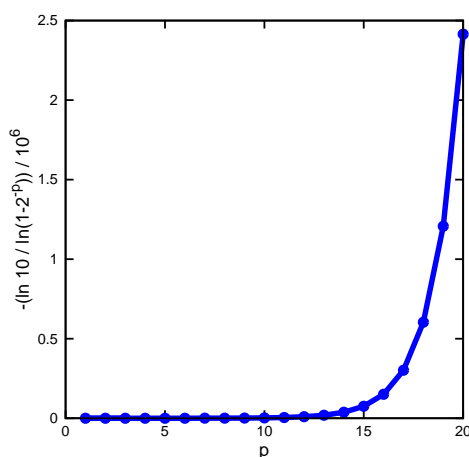
$$\begin{aligned} P(d^* \leq 0.5) \geq 0.9 &\Rightarrow n \geq (\ln 10) 2^p \frac{1}{1 + \frac{1}{2} \cdot \frac{1}{2^p} + \frac{1}{3} \cdot \frac{1}{2^{2p}} + \dots + \frac{1}{n} \frac{1}{2^{(n-1)p}} + \dots} \\ &\Rightarrow n > \frac{4}{3} 2^{p-1} \ln 10. \end{aligned}$$

Pentru obținerea ultimei inegalități de mai sus am ținut cont de faptul că inegalitatea $\frac{1}{n \cdot 2^{(n-1)p}} \leq \frac{1}{2^n} \Leftrightarrow 2^n \leq n \cdot 2^{(n-1)p}$ are loc pentru orice $p \geq 1$, și $n \geq 2$,³⁶³ deci

$$\begin{aligned} &1 + \frac{1}{2} \cdot \frac{1}{2^p} + \frac{1}{3} \cdot \frac{1}{2^{2p}} + \dots + \frac{1}{n} \cdot \frac{1}{2^{(n-1)p}} + \dots \\ &\leq 1 + \frac{1}{2^2} + \frac{1}{2^3} + \dots + \frac{1}{2^n} + \dots \\ &< \left[1 + \frac{1}{2} + \frac{1}{2^2} + \frac{1}{2^3} + \dots + \frac{1}{2^n} + \dots \right] - \frac{1}{2} \\ &\rightarrow \frac{1}{1 - \frac{1}{2}} - \frac{1}{2} = 2 - \frac{1}{2} = \frac{3}{2}. \end{aligned}$$

Observați că limita de mai sus este o limită *superioară*, de aceea inegalitatea se păstrează până la final.

Așadar, rezultă că n crește în mod exponențial în raport cu p .³⁶⁴ Graficul pentru marginea inferioară dedusă mai sus ($-\ln 10 / \ln(1 - 2^{-p})$) este cel din figura alăturată.



Se observă că, într-adevăr, creșterea acestei margini inferioare (deci și a lui n) este exponențială.

³⁶² $-\ln(1 - x) = x + \frac{1}{2}x^2 + \frac{1}{3}x^3 + \dots + \frac{1}{n}x^n + \dots$ pentru orice $x \in (-1, +1)$.

³⁶³ Demonstrația se poate face prin inducție după valorile lui p .

³⁶⁴ Mai detaliat: n , numărul de instanțe de antrenament necesare pentru a ne asigura că d^* (distanța până la cel mai apropiat vecin al originii O) este cu o probabilitate mare (și anume, $9/10$) mai mică decât 0.5 crește în mod exponențial în raport cu p .

e. Conform intuiției, clasificatorul k -NN se comportă bine atunci când instanța de test x_q este situată într-o vecinătate densă de instanțe de antrenament. Totuși, analiza teoretică de mai sus ne arată că pentru a ne asigura că punctul x_q are o vecinătate densă, numărul tuturor instanțelor de antrenament trebuie să crească exponențial în raport cu p , ceea ce nu este fezabil pentru valori mari ale lui p . (Parametrul p este dimensiunea spațiului în care se lucrează, adică numărul de trăsături ale instanțelor de antrenament și, respectiv, de test).

În consecință, pentru aplicațiile practice în care se folosesc date cu multe atribute, este recomandat ca execuția algoritmului k -NN să fie precedată de efectuarea unei „selecții de trăsături” (engl., feature selection).

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

■ □ ● CMU, 2005 spring, C. Guestrin, T. Mitchell, HW3, pr. 1

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

La 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 $\text{Error}_{\text{Bayes}}(x)$.

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 .³⁶⁵ 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:³⁶⁶

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

Răspuns:

a. Conform enunțului, $p_1(x) \stackrel{not.}{=} P(X = x|Y = 1)$, $p_0(x) \stackrel{not.}{=} P(X = x|Y = 0)$, $q(x) \stackrel{not.}{=} P(Y = 1|X = x)$ și $\theta \stackrel{not.}{=} P(Y = 1)$. Putem calcula probabilitatea $q(x)$ în funcție de $p_1(x)$, $p_0(x)$ și θ folosind formula lui Bayes:

$$\begin{aligned} q(x) &\stackrel{Bayes}{=} \frac{P(X = x|Y = 1)P(Y = 1)}{P(X = x)} \\ &= \frac{P(X = x|Y = 1)P(Y = 1)}{P(X = x|Y = 1)P(Y = 1) + P(X = x|Y = 0)P(Y = 0)} \\ &= \frac{p_1(x) \theta}{p_1(x) \theta + p_0(x)(1 - \theta)} \end{aligned}$$

b. Este imediat faptul următor: probabilitatea ca algoritmul Bayes Optimal să greșească este $P(Y = 0|X = x)$ în cazul în care $P(Y = 1|X = x) \geq P(Y = 0|X = x)$, respectiv $P(Y = 1|X = x)$ atunci când $P(Y = 0|X = x) \geq P(Y = 1|X = x)$.³⁶⁷ Altfel spus,

$$\begin{aligned} Error_{Bayes}(x) &= \min\{P(Y = 0|X = x), P(Y = 1|X = x)\} \\ &= \min\{1 - q(x), q(x)\} = \begin{cases} q(x) & \text{în cazul } q(x) \in [0, 1/2] \\ 1 - q(x) & \text{în cazul } q(x) \in (1/2, 1] \end{cases} \end{aligned}$$

c. Algoritmul 1-NN greșește atunci când instanța de antrenament x are eticheta 1 iar x' , cel mai apropiat vecin al lui x , are eticheta 0, sau invers, adică

³⁶⁵ Adică, $P(Y = 1|X = x') \rightarrow P(Y = 1|X = x)$. Aceasta se justifică ținând cont de continuitatea lui $p_y(x) \stackrel{not.}{=} P(X = x|Y = y)$ care a fost asumată în enunț și, de asemenea, de rezultatul obținut la punctul a.

³⁶⁶ Cititorul 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 presupuzițiilor din enunț, rezultă $x \rightarrow x'$ pentru orice x .

³⁶⁷ Am ținut cont de proprietatea (evidentă) $P(X = x, Y = y) = P(X = x|Y = y)P(Y = y) = P(Y = y|X = x)P(X = x)$, care are loc pentru orice x și y .

atunci când x are eticheta 0 iar x' are eticheta 1. În consecință, folosind algoritmul 1-NN, eroarea „așteptată” la clasificarea lui x este:

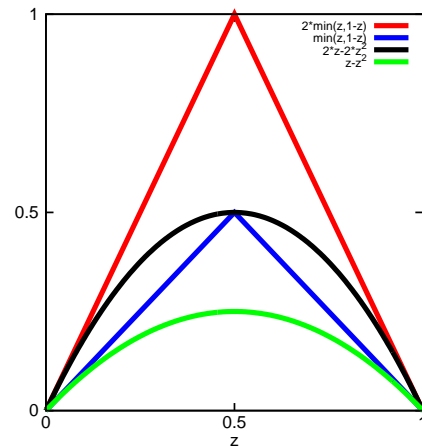
$$\begin{aligned} \text{Error}_{1\text{-NN}}(x) &= P(Y = 1|X = x)P(Y = 0|X = x') + \\ &\quad P(Y = 0|X = x)P(Y = 1|X = x') \\ &= q(x)(1 - q(x')) + (1 - q(x))q(x'). \end{aligned}$$

d. Este imediat că $\lim_{x' \rightarrow x} \text{Error}_{1\text{-NN}}(x) = 2q(x)(1 - q(x))$.

e. Se arată imediat că $z - z^2 \leq z$ pentru $\forall z$, deci și pentru $z \in [0, 1/2]$, iar $z - z^2 \leq 1 - z$ pentru $\forall z$, deci și pentru $z \in [1/2, 1]$. Așadar, pentru orice x , vom avea:

$$q(x)(1 - q(x)) \leq \begin{cases} q(x) & \text{dacă } q(x) \in [0, 1/2] \\ 1 - q(x) & \text{dacă } q(x) \in (1/2, 1]. \end{cases}$$

Coroborând cu rezultatul de la punctul b, obținem: $2q(x)(1 - q(x)) \leq 2\text{Error}_{\text{Bayes}}(x)$ pentru orice x .



Combinând acest rezultat cu egalitatea de la punctul d, rezultă că inegalitatea

$$\lim_{n \rightarrow \infty} \text{Error}_{1\text{-NN}}(x) = \lim_{x' \rightarrow x} \text{Error}_{1\text{-NN}}(x) \leq 2\text{Error}_{\text{Bayes}}(x)$$

este adevărată pentru orice x . Înmulțind ambii membri ai acestei inegalități cu $P(x)$ și însumând după toate valorile lui x — de fapt, integrând în raport cu x —, obținem:

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

Așadar, am demonstrat că media (sau: rata medie a) erorii asimptotice a algoritmului 1-NN este cel mult dublul mediei (sau: ratei medii a) erorii algoritmului Bayes Optimal.

Observația 1: Această margine a erorii asimptotice nu se păstrează și în cazul neasimptotic, unde numărul de exemple de antrenament este finit.

Observația 2: La fel, se poate arăta că $2z - 2z^2 \geq z$ pentru $\forall z \in [0, 1/2]$ și $2z - 2z^2 \geq 1 - z$ pentru $\forall z \in [1/2, 1]$. Luând din nou $z = q(x)$ și ținând cont de rezultatul de la punctul b, obținem că $2q(x)(1 - q(x)) \geq \text{Error}_{\text{Bayes}}(x)$, pentru orice x . Combinând această inegalitate cu egalitatea de la punctul d, rezultă o nouă inegalitate:

$$\lim_{n \rightarrow \infty} \text{Error}_{1\text{-NN}}(x) = \lim_{x' \rightarrow x} \text{Error}_{1\text{-NN}}(x) \geq \text{Error}_{\text{Bayes}}(x) \text{ pentru orice } x.$$

În final, trecând la medii, obținem următoarea inegalitate (care era de altfel de așteptat):

$$E[\lim_{n \rightarrow \infty} \text{Error}_{1\text{-NN}}] \geq E[\text{Error}_{\text{Bayes}}].$$

Observația 3 (preluată din *An Elementary Introduction to Statistical Learning Theory*, de Sanjeev Kulkarni și Gilbert Harman, 2011, pag. 69): În mod intuitiv, dacă mărim valoarea lui k , ar trebui ca eroarea medie a algoritmului k -NN să se reducă. Într-adevăr, în anumite condiții (dar nu în orice condiții!) se poate arăta că are loc următoarea inegalitate dublă:

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

Este de remarcat faptul că există distribuții probabilistice ale datelor pentru care clasificatorul 1-NN se comportă mai bine decât k -NN pentru orice $k \neq 1$.

Observația 4 (preluată din aceeași lucrare, *An Elementary Introduction to Statistical Learning Theory*, citată mai sus): Dacă lucrăm cu k_n -NN, adică îl fixăm pe k în funcție de n (numărul instanțelor de antrenament), se poate demonstra că în cazul în care $\frac{k_n}{n} \rightarrow 0$ pentru $n \rightarrow \infty$ (de exemplu, $k_n = \sqrt{n}$), se obține:

$$E[\lim_{n \rightarrow \infty} Error_{k_n-NN}] = E[Error_{Bayes}].$$

Aceasta înseamnă că, la limită, algoritmul k_n -NN se comportă la fel de bine ca algoritmul Bayes Optimal!

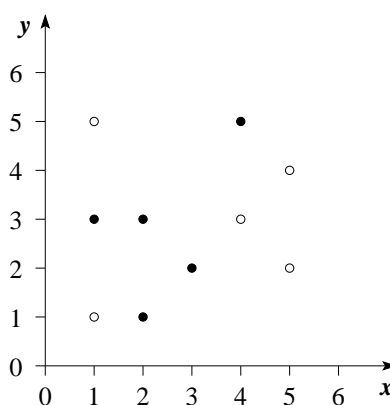
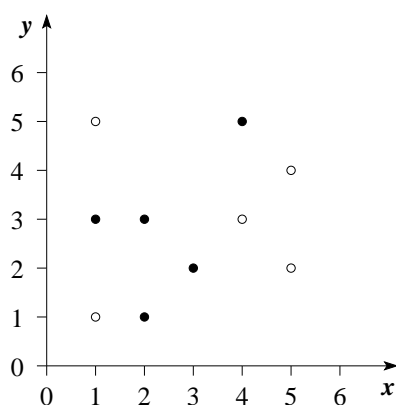
11. (Comparație între algoritmi 1-NN și ID3: zone și suprafețe de decizie)

prelucrare de Liviu Ciortuz, după

■ • CMU, 2007 fall, Carlos Guestrin, HW2, pr 1.4

Pe setul de date de mai jos desenați *granițele de decizie* și apoi hașurați *suprafețele de decizie* produse de

- algoritmul 1-NN (veți obține deci diagrama Voronoi);
- algoritmul ID3 extins cu capacitatea de a procesa atribute cu valori continue.



Răspuns:

a. *Algoritmul 1-NN:*

Pentru a defini suprafețele de decizie în acest caz se procedează în felul următor:

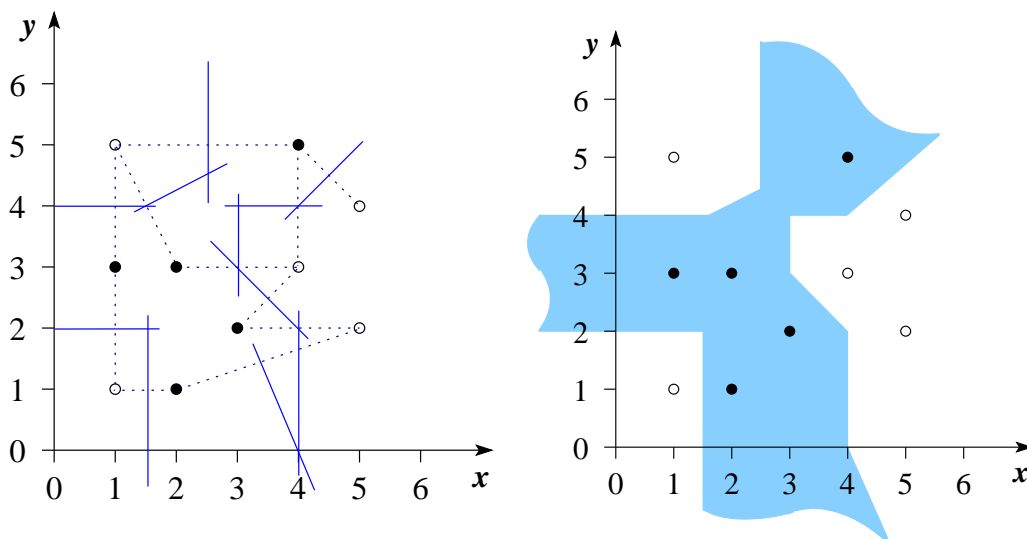
- se trasează mediatoarele segmentelor de dreaptă determinate de perechi de puncte din setul de antrenament care sunt etichetate în mod diferit;
- se stabilesc intersecțiile acestor mediatoare; acest lucru este reprezentat în figura de mai jos, partea stângă;
- apoi se marchează pe aceste mediatoare acele segmente (determinate de intersecții) care determină zonele de decizie corespunzătoare clasificatorului 1-NN.

Observații:

1. De fapt, întrucât nu este necesar să se lucreze cu toate perechile de instanțe cu etichete diferite sunt relevante pentru clasificarea unei instanțe / zone, în timpul „executării” punctelor de mai sus este foarte util să se țină cont de următoarea regulă / euristică de *ghidare*: alegerea perechilor de instanțe și apoi a segmentelor de pe mediatoare se va face urmărind delimitarea zonelor corespunzătoare instanțelor negative (\circ) de zonele corespunzătoare instanțelor pozitive (\bullet).

2. Suplimentar, în jurul fiecărui punct de antrenament A se poate identifica câte o zonă [convexă] care va constitui mulțimea punctelor mai apropiate de A decât de oricare alt punct din setul de date de antrenament. (Toate punctele din această zonă convexă vor avea aceeași clasificare / etichetă ca și punctul A .) Aceasta este ușor de văzut pentru instanțele negative (\circ) din cazul de față.

În figura următoare, în partea dreaptă am hașurat zona corespunzătoare instanțelor pozitive (\bullet).



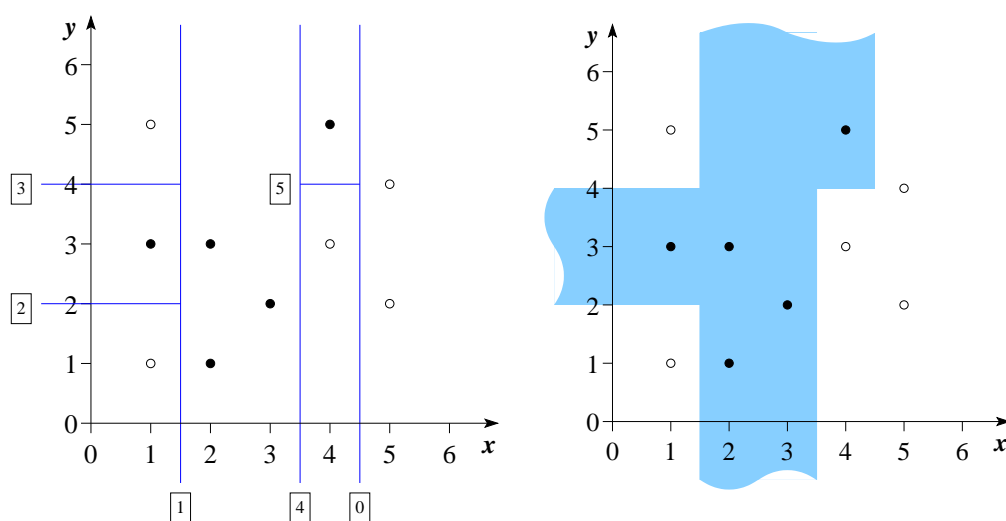
b. Algoritmul ID3:

Pentru construirea suprafețelor de decizie ale algoritmului ID3,

- mai întâi se determină valorile-prag pentru teste (adică, punctele de splitare de pe fiecare axă) și apoi se alege testul corespunzător nodului rădăcină din arborele ID3; se trasează o dreaptă prin punctul respectiv, paralelă cu cealaltă axă. În cazul nostru, testul din nodul rădăcină va fi $x < 4.5$.

– pentru fiecare test / split ulterior se trasează o semidreaptă (sau un segment de dreaptă) mărginit(ă) la un capăt de dreapta corespunzătoare nodului-părinte.

În figura de mai jos, în partea stângă am vizualizat toate testele / split-urile realizate de algoritmul ID3 pentru a învăța complet arborele de decizie, iar în partea dreaptă, ca și mai sus, am păstrat doar frontierele dintre zonele cu clasificări diferite.



Se observă că suprafețele de decizie determinate de cei doi algoritmi nu sunt identice, dar sunt totuși asemănătoare într-o anumită măsură (pentru că ambele sunt consistente cu datele de antrenament).

Observații:

3. Este de reținut faptul că suprafețele de decizie produse de algoritmul ID3 nu sunt neapărat unic determinate, fiindcă sunt situații în care două teste diferite pot conduce la același câștig de informație. De exemplu, dacă în exercițiul nostru am avea de partiționat la un moment dat mulțimea formată din instanțele de antrenament $(1, 1), (1, 3), (2, 1), (2, 3), (3, 2)$, atunci testele $x > 1.5$ și $y > 1.5$ ar produce același câștig de informație, iar suprafețele de decizie rezultate ar fi determinate (în mod diferit!) de ce anume alegem ca prim test.

4. De asemenea, trebuie să scoatem în evidență faptul că pragurile / spliturile care sunt calculate de algoritmul ID3 pentru un același atribut continuu pot diferi de la un nod de test la altul. De exemplu, la nodul rădăcină (nodul 0), atunci când se calculează câștigul de informație maxim, pentru atributul y se iau în calcul pragurile 1.5, 2.5, 3.5, și 4.5, în vreme ce la nodurile 2 și / sau 3 (vedeți figura de mai sus, partea stângă) se analizează pragurile 2 și 4, întrucât seturile / partițiile de instanțe asignate acestor noduri sunt diferite!

12. (Comparație între algoritmi 1-NN, ID3 cu atribute continue și SVM: eroarea la antrenare, eroarea la CVLOO)

prelucrare de Liviu Ciortuz, după

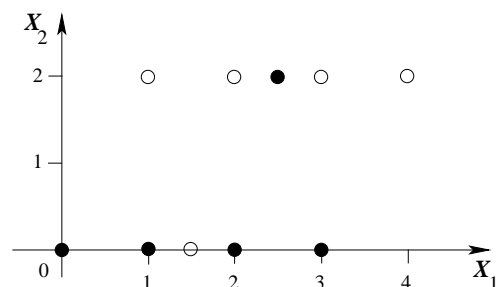
- CMU, 2003 fall, T. Mitchell, A. Moore, midterm exam, pr. 5

Folosind metoda 1-NN cu distanța euclidiană, învățăm un clasificator cu două valori pentru atributul de ieșire, $Y = 0$ și $Y = 1$, pornind de la datele de antrenament din tabelul de mai jos (X_1 și X_2 sunt atribute de intrare).

- | | X_1 | X_2 | Y |
|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-------|-------|-----|
| a. Care este eroarea la antrenare (exprimată ca număr de exemple clasificate eronat)? | 0 | 0 | 1 |
| | 1 | 0 | 1 |
| b. Care este eroarea la cross-validare folosind metoda "Leave-One-Out"? | 2 | 0 | 1 |
| | 2.5 | 2 | 1 |
| c. Răspundeți la întrebările de mai sus, considerând acum arbori de decizie cu atribute numerice continue în locul metodei 1-NN. | 3 | 0 | 1 |
| | 1 | 2 | 0 |
| | 1.5 | 0 | 0 |
| d. În sfârșit, răspundeți la întrebările a și b, considerând mașini cu vectori-suport în locul metodei 1-NN. (Se vor considera doar SVM-uri în cazul liniar cu margine "soft", cu un parametru C suficient de mare pentru a minimiza numărul de instanțe de antrenament clasificate eronat.) | 2 | 2 | 0 |
| | 3 | 2 | 0 |
| | 4 | 2 | 0 |

Răspuns:

Reprezentarea datelor în planul euclidian este cea din figura alăturată.



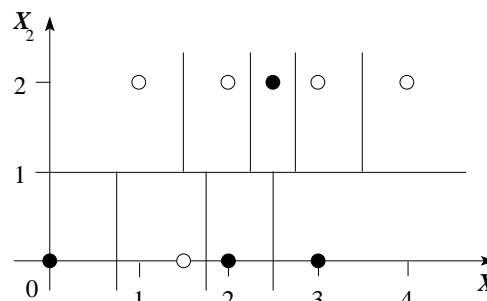
- a. După cum am explicat și la exercițiul 2, întrucât datele de antrenament nu conțin inconsistențe, eroarea la antrenare produsă de algoritmul 1-NN va fi 0.
- b. Comportamentul algoritmului la cross-validare cu metoda "Leave-One-Out" este cel descris în tabelul următor:

Data	Eticheta	Vecinătate	Clasificare la CVLOO	Eroare?
(0; 0)	1	(1; 0)	1	nu
(1; 0)	1	(1.5; 0)	0	da
(2; 0)	1	(1.5; 0)	0	da
(2, 5; 2)	1	(2; 2)/(3; 2)	0	da
(3; 0)	1	(2; 0)	1	nu
(1; 2)	0	(2; 2)	0	nu
(1, 5; 0)	0	(1; 0)/(2; 0)	1	da
(2; 2)	0	(2.5; 2)	1	da
(3; 2)	0	(2.5; 2)	1	da
(4; 2)	0	(3; 2)	0	nu

Deci în total avem 6 erori (din totalul de 10 instanțe), ceea ce indică faptul ca algoritmul 1-NN este foarte puțin adecvat pentru acest gen de date.

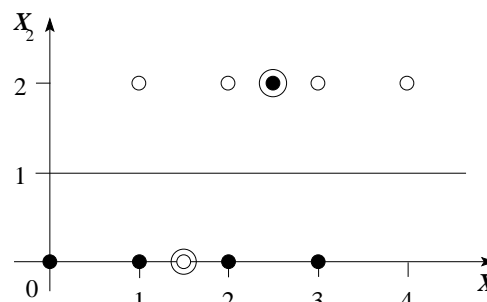
c. Eroarea la antrenare cu algoritmul ID3 pe acest set de date este 0, fiindcă instanțele sunt etichetate în mod consistent.

Eroarea la CVLOO produsă de algoritmul ID3 pe același set de date se poate calcula foarte ușor, ținând cont că este suficient să determinăm granițele de decizie — deci și zonele de decizie — corespunzătoare celor două clase. De exemplu, atunci când se va elimina instanța (1,0), se vor obține granițele de decizie ca în figura alăturată. Este imediat că *modelul* rezultat în acest caz va clasifica instanța (1,0) în mod eronat.



Se constată că eroarea de tip CVLOO produsă de ID3 cu attribute continue este de 6 (din totalul de 10) instanțe. Ba chiar, — este încă o coincidență! — ID3 produce la CVLOO exact aceleași erori (mai precis, el clasifică în mod eronat exact aceleași instanțe) ca și algoritmul 1-NN.

d. Se observă ușor că separatorul liniar care minimizează eroarea la antrenare este cel reprezentat în figura alăturată. Datele clasificate eronat de către acest separator sunt cele două puncte încercuite din figură; ele sunt de asemenea singurele puncte care generează eroare și la testare cu metoda CVLOO.



Mai concret, în cazul CVLOO folosind SVM, deoarece separatorul optimal este „susținut“ de mai mulți vectori-suport pe fiecare parte, lipsa unuia singur dintre ei nu influențează cu nimic construirea separatorului. În fiecare caz în parte se va învăța ca separator dreapta paralelă cu Ox_1 care trece prin punctul (0, 1).

Avem deci în ambele situații, adică atât la antrenare cât și la cross-validare cu metoda “Leave-One-Out”, (doar) două puncte clasificate eronat de către SVM.

Comparând cei trei clasificatori, 1-NN, ID3 cu attribute continue și SVM, rezultă în mod clar că SVM este cel mai convenabil pe acest set de date (deși la antrenare SVM produce două erori, iar 1-NN și ID3 nicio eroare), fiindcă nu produce *overfitting* (aici!).

13. (Comparații între algoritmi 1-NN și ID3: Da sau nu?)

CMU, 2009 spring, Ziv Bar-Joseph, final exam, pr. 2.1

a. Este posibil să se construiască un arbore de decizie (având în fiecare nod intern teste de forma $x > a$, $x \leq b$, $y > c$, sau $y \leq d$, unde a, b, c, d sunt numere reale oarecare) care să producă la clasificare aceleași rezultate ca și algoritmul 1-NN folosind distanța euclidiană? Justificați răspunsul.

(Învățare rapidă / “eager” vs. învățare lentă / “lazy”;
 k -NN vs. ID3)

CMU, 2010 spring, HW1, pr. 3.3

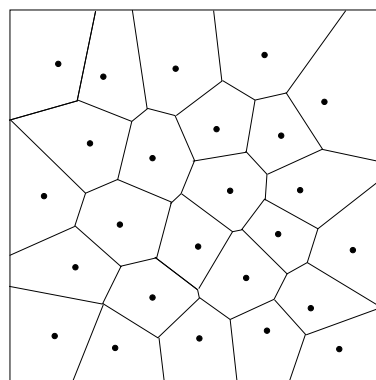
b. Algoritmul ID3 este o metodă de învățare de tip “batch”, care solicită ca toate datele de antrenament să-i fie puse la dispoziție pentru a putea elabora arborele de decizie. Așadar, în situația în care date de antrenament suplimentare ne sunt puse ulterior la dispoziție, acestea trebuie tratate cu atenție fiindcă ele pot modifica arborele de decizie rezultat în urma învățării.³⁶⁸ Algoritmul k -NN suferă și el de această problemă? Justificați.

Răspuns:

a. Nu.

Suprafețele de decizie pentru algoritmul 1-NN corespund diagramei Voronoi și nu sunt neapărat paralele cu axele de coordonate, după cum se observă în figura alăturată.

Suprafețele de decizie pentru un arbore de decizie cu atribute cu valori continue sunt întotdeauna paralele cu axele de coordonate, deoarece deciziile din fiecare nod sunt de forma $x > a$, $x \leq b$, $y > c$, sau $y \leq d$, $\forall a, b, c, d \in \mathbb{R}$.



Așadar, răspunsul este negativ pentru cazul general, deși există situații în care cei doi algoritmi produc exact aceeași clasificare.

b. Nu.

Răspunsul decurge din modul de lucru a algoritmului k -NN, care este un algoritm de învățare de tip “lazy”: k -NN estimează valoarea locală a unei funcții-target \hat{f} pentru una sau mai multe instanțe de test x_q . Valoarea $\hat{f}(x_q)$ nu depinde decât de cei mai apropiați vecini ai lui x_q ; celelalte instanțe de antrenament nu sunt necesare pentru calculul lui $\hat{f}(x_q)$. Evident, dacă pe măsură ce se acumulează noi date de antrenament se modifică și vecinătatea lui x_q , atunci se prea poate să se modifice și $\hat{f}(x_q)$. Însă în sine, algoritmul procedează exact la fel ca mai înainte. Se poate modifica doar outputul lui. Spre deosebire de algoritmul k -NN, la învățarea arborilor de decizie adăugarea de noi instanțe modifică în general și modelul / arborele rezultat, nu doar decizia pentru o instanță de test particulară.

14. (1-NN cu mapare cu RBF: Adevărat sau Fals?)

■ ● ○ CMU, 2003 fall, T. Mitchell, A. Moore, final exam, pr. 7.f

Algoritmul 1-NN folosind distanța euclidiană neponderată este capabil să obțină rezultate mai bune dacă în prealabil intrările sale sunt mapate într-un „spațiu de trăsături” folosind o funcție-nucleu cu baza radială (RBF).

Răspuns:

³⁶⁸Vedeți problema 16 de la capitolul *Arbori de decizie*.

Fals.

Fie $\phi : \mathbb{R}^d \rightarrow \mathbb{R}^n$ funcția de mapare în spațiul de trăsături, astfel încât să avem $K(x, y) \stackrel{\text{not.}}{=} e^{-\frac{\|x-y\|^2}{2\sigma^2}} = \phi(x) \cdot \phi(y), \forall x, y \in \mathbb{R}^d$. (\mathbb{R}^d reprezintă spațiul inițial, \mathbb{R}^n spațiul de trăsături în care se face maparea, iar $e^{-\frac{\|x-y\|^2}{2\sigma^2}}$ este funcția-nucleu cu baza radială.) Avem:

$$\begin{aligned} \|\phi(x) - \phi(y)\|^2 &= (\phi(x) - \phi(y)) \cdot (\phi(x) - \phi(y)) \\ &= \phi(x) \cdot \phi(x) + \phi(y) \cdot \phi(y) - 2 \cdot \phi(x) \cdot \phi(y) = e^{-\frac{\|x-x\|^2}{2\sigma^2}} + e^{-\frac{\|y-y\|^2}{2\sigma^2}} - 2 \cdot e^{-\frac{\|x-y\|^2}{2\sigma^2}} \\ &= e^0 + e^0 - 2 \cdot e^{-\frac{\|x-y\|^2}{2\sigma^2}} = 2 - 2 \cdot e^{-\frac{\|x-y\|^2}{2\sigma^2}} = 2 - K(x, y). \end{aligned}$$

Prin urmare, pentru orice $x, x_i, x_j \in \mathbb{R}^d$ vom avea:

$$\begin{aligned} \|\phi(x) - \phi(x_i)\|^2 \leq \|\phi(x) - \phi(x_j)\|^2 &\Leftrightarrow 2 - K(x, x_i) \leq 2 - K(x, x_j) \Leftrightarrow K(x, x_i) \geq K(x, x_j) \\ \Leftrightarrow e^{-\frac{\|x-x_i\|^2}{2\sigma^2}} &\geq e^{-\frac{\|x-x_j\|^2}{2\sigma^2}} \Leftrightarrow -\frac{\|x-x_i\|^2}{2\sigma^2} \geq -\frac{\|x-x_j\|^2}{2\sigma^2} \Leftrightarrow \|x-x_i\|^2 \leq \|x-x_j\|^2. \end{aligned}$$

Cu alte cuvinte, dacă o instanță de test x are drept cel mai apropiat vecin punctul x_i în spațiul inițial, acest lucru rămâne valabil și în spațiul de trăsături. Așadar, decizia algoritmului 1-NN este identică în ambele spații. Aceeași concluzie este valabilă și pentru k -NN.

Observație: Este însă posibil ca folosind ponderarea sau alte măsuri de distanță (decât cea euclidiană) kernel-izarea să funcționeze cu succes pentru k -NN.

3.2 Învățare bazată pe memorare — Probleme propuse

15. (Algoritmul k -NN: acuratețe; comparație cu un simplu clasificator aleator)

prelucrare de Liviu Ciortuz, după CMU, 2014 fall, W. Cohen, Z. Bar-Joseph, HW1, pr. 5.bc

a. Enunțați [succint] *regula de decizie* a algoritmului k -NN pentru o instanță de test x_q .

Care este *bias*-ul inductiv al algoritmului k -NN?

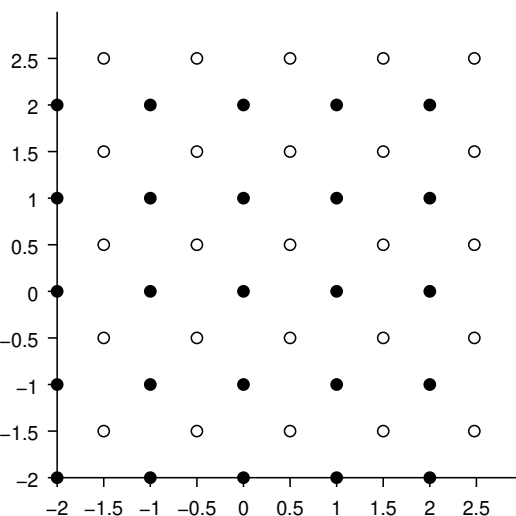
În continuare veți lucra pe datele din figura de mai jos. Veți aplica algoritmul k -NN, folosind distanța euclidiană.

k -NN funcționează bine atunci când instanțele dintr-o aceeași clasă sunt plasate într-una sau mai multe zone din spațiu relativ bine delimitate, fără întrupări puternice.

Obiectivul nostru acum este să analizăm ce se întâmplă atunci când datele sunt puternic mixate. Rezultatele pe care le veți obține la calculul erorilor vor fi exprimate sub formă de numere [fracționare] din intervalul $[0, 1]$.

Observație importantă:

În cazul în care există două sau mai multe instanțe situate exact pe „marginea” [adică, pe conturul circular al] k -NN-vecinătății asociate instanței de clasificat, se va considera că toate aceste instanțe aparțin respectivei vecinătăți, iar fiecare dintre ele dispune de un vot întreg.



b. Pentru $k = 1$, calculați eroarea la antrenare și eroarea la cross-validare cu metoda “leave-one-out” (CVLOO).

Ce puteți spune comparând cele două rezultate? (Care este legătura între *bias*-ul inductiv al lui k -NN și puterea de generalizare a lui 1-NN pe astfel de date? Se produce oare aici un *anumit* fenomen, specific multor situații din clasificarea automată?)

c. Pentru $k = 2$, calculați eroarea la cross-validare cu metoda “leave-one-out” (CVLOO).

d. Considerăm $k = 50$. (Remarcați faptul că în total în setul nostru de date sunt 50 de instanțe.) De această dată, vom impune ca algoritmul k -NN să ia decizia în mod *probabilist*. Aceasta înseamnă că dacă în vecinătatea k -NN a unei instanțe de test există n vecini pozitivi și m vecini negativi,

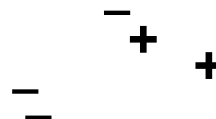
atunci algoritmul k -NN va returna (pentru instanța respectivă) decizia + cu probabilitatea $n/(n+m)$ și decizia – cu probabilitatea $m/(n+m)$. În consecință, pentru întreg setul de date vom putea calcula o *eroare medie*.

Calculați *eroarea medie* la antrenare pentru algoritmul 50-NN pe datele de mai sus. Cunoașteți o metodă de clasificare foarte simplă care obține pe aceste date rezultate la fel de bune / proaste precum 50-NN?

16. (Algoritmul 1-NN: calculul erorii la CVLOO)

* CMU, 2010 spring, E. Xing, T. Mitchell, A. Singh, midterm exam, pr. 1.7

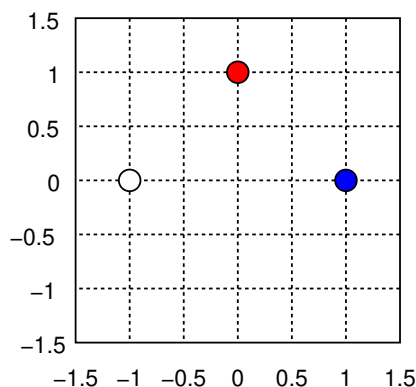
Care este eroarea clasificatorului 1-NN la cross-validare de tip “Leave-One-Out” pe setul de date alăturat?



17. (Algoritmul 1-NN: granițe / suprafețe de decizie)

• CMU, 2013 fall, W. Cohen, E. Xing, final exam, pr. 3.6

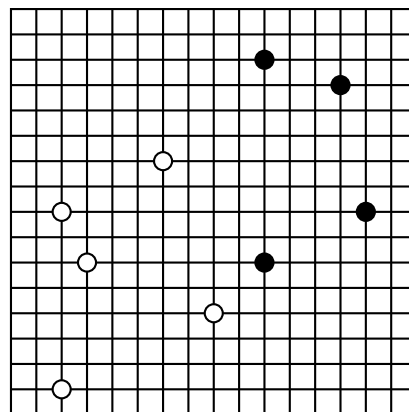
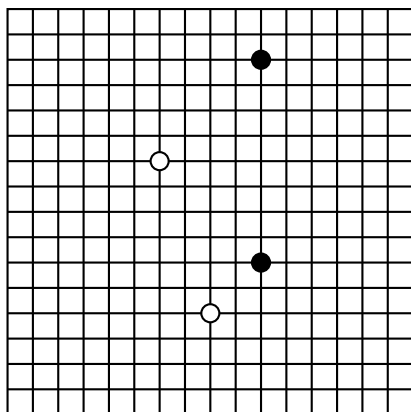
Trasați granițele de decizie (engl., *decision boundaries*) produse de clasificatorul 1-NN la aplicarea pe datele din figura alăturată. Diversele forme ale punctelor reprezintă clase diferite.



18. (Algoritmul 1-NN: granițe / suprafețe de decizie)

* CMU, 2008 fall, Eric Xing, HW1, pr. 3

În fiecare din figurile următoare se dau câteva puncte în spațiul bidimensional, etichetate cu + sau –. Indicați în fiecare caz granițele / suprafețele de decizie pentru algoritmul 1-NN presupunând că se folosește distanța euclidiană.

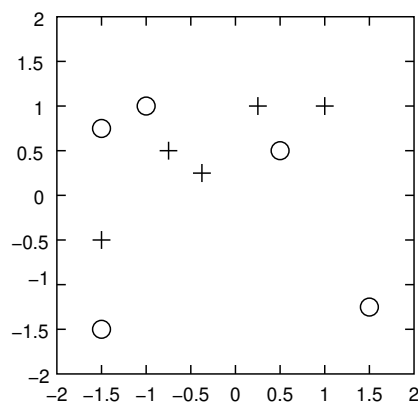
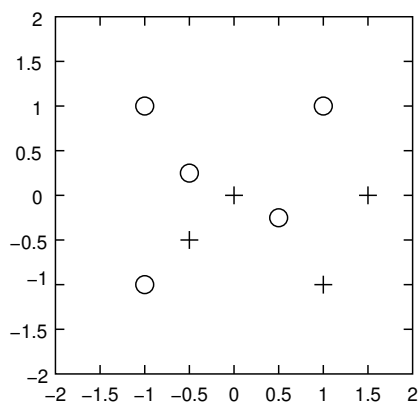
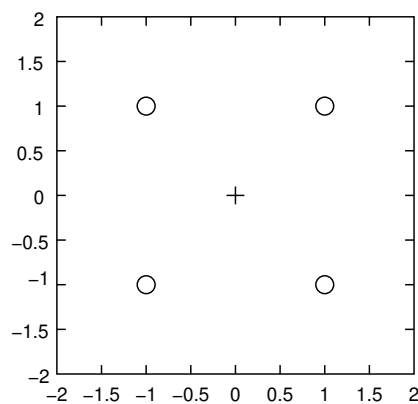
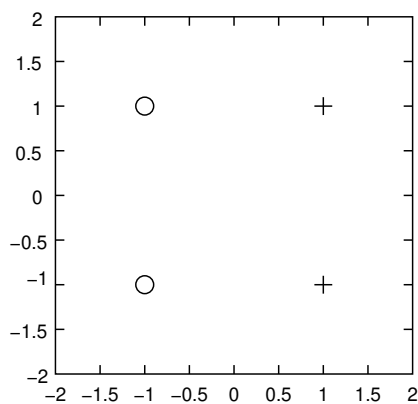


19.

(Algoritmul 1-NN: diagrame Voronoi)

• * CMU, 2010 fall, Ziv Bar-Joseph, HW1, pr. 3.1

Desenați suprafețele de decizie pentru clasificatorul 1-NN pentru fiecare dintre seturile de date din figurile de mai jos. Folosiți distanța euclidiană. Hașurați fin zonele corespunzătoare clasei +.



20. (Algoritmul k -NN: diagrama Voronoi, eroarea la CVLOO; comparație pentru diferite valori ale lui k)

• prelucrare de L. Ciortuz, după

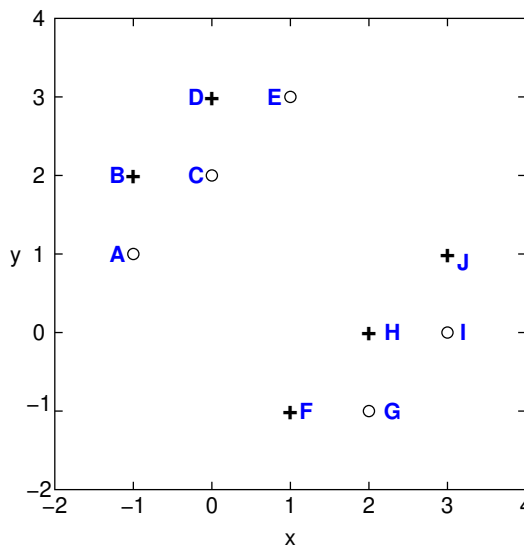
CMU, 2012 spring, Ziv Bar-Joseph, midterm exam, pr. 2

La acest exercițiu veți aplica algoritmul k -NN folosind distanța euclidiană pe setul de date din figura de mai jos. Fiecare punct aparține la una din două clase, desemnate cu $+$ și respectiv \circ .

a. Trasați diagrama Voronoi și hașurați zona / zonele de decizie corespunzătoare etichetei $+$.

b. Care este eroarea la cross-validare cu metoda “Leave-One-Out” (CVLOO) dacă se folosește algoritmul 1-NN?

c. Care dintre următoarele valori ale lui k va conduce la o valoare minimă a erorii de tip CVLOO: 3, 5, 7 sau 9? Comentați succint rezultatul.



Indicații:

1. k -NN-vecinătățile vor fi construite în manieră *inclusivă*.³⁶⁹ Vă *cerem(!)* să puneți în evidență toate cazurile de acest tip. Vedeți, spre exemplu, liniile 2 și 3 din tabelul de mai jos, coloana 1-NN vecinătăților.
2. În caz de *paritate la voturi* (dar doar în acest caz!), se va considera că se aplică (în mod intuitiv) ponderarea distanțelor în sensul prezentat la curs.
3. Dacă veți ști să exploatați *simetriile*, veți avea mult mai puțin de elaborat la nivel de detaliu!

21. (Algoritmul k -NN: alegerea valorii convenabile pentru k)

prelucrare de Liviu Ciortuz, după

CMU, (?) spring, ML course 10-701, HW1, pr. 5

Pe un anumit set de date format din date de antrenament, date de validare și date de test, după ce fost antrenat algoritmul k -NN pentru diferite valori

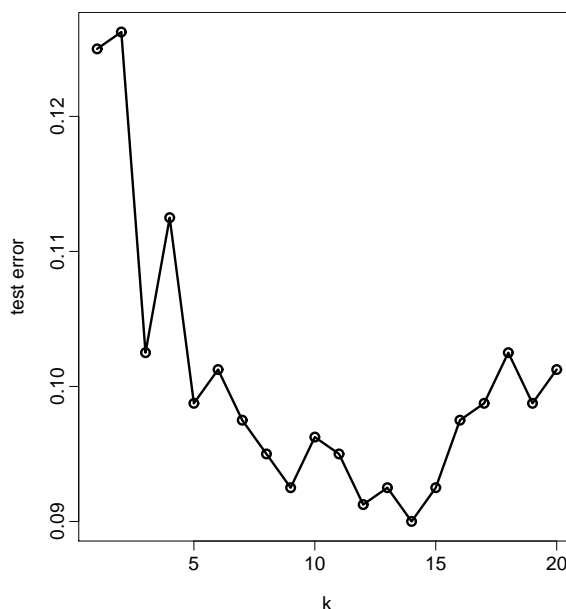
³⁶⁹ Adică, dacă notăm cu

- x_1, x_2, \dots, x_n instanțele de antrenament,
- x o instanță oarecare căreia i se aplică la un moment dat procedura de cross-validare LOO cu algoritmul k -NN (unde k este fixat),
- $d(x, x_{i_1}) \leq d(x, x_{i_2}) \leq \dots \leq d(x, x_{i_n})$ secvența ordonată a distanțelor de la x la fiecare din instanțele de antrenament,

și există $l > 0$ astfel încât $x_{i_k} = x_{i_{k+1}} = \dots = x_{i_{k+l}} < x_{i_{k+l+1}}$ sau $k + l = n$, atunci în k -NN vecinătatea lui x vor fi incluse toate instanțele $x_{i_{k+1}}, \dots, x_{i_{k+l}}$.

ale lui k , rezultatele obținute la validare au fost reprezentate în graficul care urmează.

Care este — în conformitate cu aceste rezultate — valoarea optimă care trebuie aleasă pentru k , în vederea folosirii ulterioare pe datele de test? Justificați alegerea făcută.



22.

(Algoritmul k -NN: acuratețe; comparații pentru diferite valori ale lui k)

• ◦ CMU, 2014 fall, W. Cohen, Z. Bar-Joseph, HW1, pr. 5.a

Considerăm două clase, notate cu C_1 și C_2 , în spațiul euclidian bidimensional. Datele din clasa C_1 sunt distribuite în mod uniform într-un cerc de rază r . Datele din clasa C_2 sunt distribuite în mod uniform într-un alt cerc de rază r . (*Observație:* Numărul de date din cele două clase nu este neapărat același.) Centrele celor două cercuri sunt situate la o distanță strict mai mare decât $4r$.

Arătați că este posibil ca [la antrenare] acuratețea algoritmului 1-NN aplicat pe aceste date să fie *strict* mai mare decât acuratețea algoritmului k -NN, pentru un anumit număr întreg $k \geq 3$, impar, ales în mod convenabil.

23.

(Algoritmul k -NN: întrebări de ordin calitativ)

◻ • ◦ CMU, 2012 spring, Ziv Bar-Joseph, HW1, pr. 3

Știm că la aplicarea algoritmului k -NN, clasificarea unei instanțe date se face pe baza votului majoritar obținut în „vecinătatea“ instanței respective. Presupunem că se dau două clase de instanțe, fiecare clasă având $n/2$ puncte, întrepătrunse într-o anumită măsură, într-un spațiu bidimensional.

a. Descrieți ce se întâmplă cu *eroarea la antrenare* (folosind toate datele disponibile) când numărul k al vecinilor considerați variază de la n la 1.

- b. Schițați grafic cum anume ar evolua *eroarea de generalizare* (de exemplu, reținând o parte din date pentru testare) atunci când k variază. Explicați modul în care ați raționat.
- c. Propuneți o metodă de determinare a unei valori adecvate pentru k .
- d. La folosirea algoritmului k -NN, odată ce s-a stabilit valoarea lui k , toți cei mai apropiați k vecini ai punctului de clasificat au ponderi egale (adică, aceeași importanță) la stabilirea etichetei respectivului punct. Sugerați o modificare a algoritmului k -NN care elimină această limitare.
- e. Dați două motive pentru care este de preferat să nu folosim algoritmul k -NN atunci când dimensiunea spațiului datelor de intrare este mare.

24.

(Algoritmul k -NN:CVLOO: comparație pentru diferite valori ale lui k ;
eroarea la antrenare: comparație cu alți clasificatori)

• ◦ CMU, 2010 fall, Aarti Singh, midterm exam, pr. 2

- a. Care dintre clasificatorii de mai jos realizează o eroare de tip CVLOO (Leave-One-Out Cross-validation) mai mare pe setul de date alăturat?

☐ 1-NN☐ 3-NN

+	+	—	—
	—		—
+	+	—	—

Recomandare: Formulați explicit euristica pe care o veți folosi pentru a trata cazurile în care apar tot atâtea voturi pozitive cât cele negative.

- b. Considerăm setul de date din figura alăturată. Care dintre clasificatorii de mai jos obține / obțin eroare nulă la antrenare pe acest set de date?

○	+
+	○

- ☐ arborii de decizie ID3 de adâncime 2
- ☐ clasificatorul 3-NN
- ☐ regresia logistică
- ☐ SVM (cu nucleu pătratic)

25.

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

* Liviu Ciortuz, 2014, bazat pe un rezultat din

■ “An Elementary Introduction to Statistical Learning Theory”,
S. Kulkarni, G. Harman, 2011, pag. 68-69

La problema 10 am demonstrat că *rata medie a erorii* clasificatorului 1-NN este mărginită asimptotic³⁷⁰ de dublul ratei medii a erorii clasificatorului Bayes Optimal.

³⁷⁰ Adică, atunci când $n \rightarrow \infty$, unde n este numărul de instanțe de antrenament.

Arătați că — în aceleași condiții ca la problema 10 — se poate obține o margine chiar mai bună:

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

26. (Adevărat sau Fals?)

◦ CMU, 2010 fall, Ziv Bar-Joseph, midterm, pr. 1.bc

Care dintre următoarele afirmații sunt adevărate pentru clasificatorii k -NN? (Justificați pe scurt răspunsul, în dreptul fiecărui punct.)

- a. Acuratețea la antrenare crește pe măsură ce crește valoarea lui k .
- b. Suprafața de decizie este mai netedă (engl., smoother) pe măsură ce valoarea lui k scade.
- c. k -NN nu necesită o procedură explicită de antrenare.
- d. Suprafața / granița de decizie este liniară.
- e. Este posibil ca un clasificator binar 1-NN să clasifice întotdeauna orice instanță de test ca fiind pozitivă, chiar dacă în setul de date de antrenament există instanțe negative.

27. (Întrebări calitative despre design-ul unor experimente din Învățarea Automată: OK ori ...problematic?)

• ◦ CMU, 2009, Geoff Gordon, midterm exam, pr. 3

Fiecare din punctele de mai jos prezintă pe scurt design-ul unui experiment practic de învățare automată. Analizați fiecare din aceste cazuri, indicând apoi dacă respectivul experiment este *ok* ori *problematic* (încercuiți varianta pe care o alegeți). Dacă este *problematic*, identificați TOATE defectele [de concepție ale] design-ului respectiv.

- a. O echipă de proiectare raportează o eroare mică la antrenare și susține că metoda folosită este bună.

Ok

Problematic

- b. O echipă de proiectare susține că este un mare succes faptul că a obținut 98% acuratețe la antrenare pentru un task de clasificare binară care are următorul specific: unul din cele două cazuri se întâlnește foarte rar comparativ cu celălalt caz. (O astfel de problemă o constituie, de exemplu, identificarea tranzacțiilor bancare frauduloase.) Datele lor au constatat din 50 de exemple pozitive și 4950 de exemple negative.

Ok

Problematic

- c. O echipă de proiectare și-a împărțit datele de care dispune în date de antrenament și date de test. Folosind datele de antrenament, ei au construit

un *model* de clasificare caracterizat de anumiți *parametri*. Apoi, făcând *cross-validare*, au ales cea mai bună setare a parametrilor. La final, au raportat *eroarea* obținută pe *datele de test*.

Ok

Problematic

d. O echipă de proiectare a efectuat o procedură de *selecție a atributelor* (engl., features) pe toate datele și apoi a redus setul mare de attribute la un set mai mic. După aceea, membrii echipei au împărțit datele în date de test și date de antrenament. Au construit *modelul* de clasificare pe datele de antrenament folosind mai multe setări ale parametrilor modelului, și au raportat cea mai bună *eroare la testare* pe care au obținut-o.

Ok

Problematic