

Învățare automată

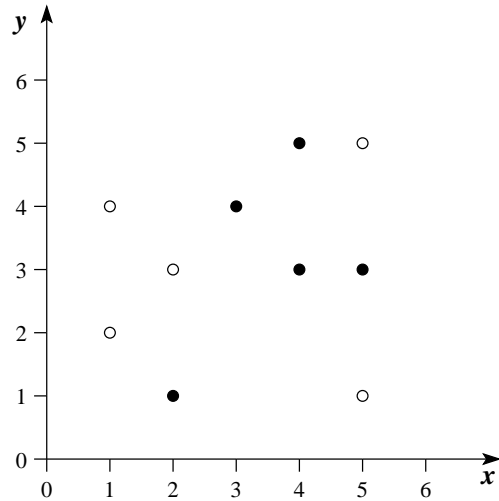
— Licență, anul III, 2018-2019, examenul parțial II —

Nume student:

Grupa:

1. (Comparație între algoritmi 1-NN: zone și suprafețe de decizie; calculul erorii la CV cu metoda “leave-one-out”)

a. Pe setul de date alăturat desenați *granițele de decizie* produse de către algoritmul 1-NN (veți obține deci *diagrama Voronoi*). Apoi hașurați *suprafața de decizie* corespunzătoare clasei +, marcată prin eticheta / simbolul •.



b. Pe același set de date, calculați eroarea produsă la cross-validare cu metoda “leave-one-out” (CVLOO) de către algoritmul 1-NN. Veți exprima această eroare sub forma unei *fracții*.

Atenție: În cazul în care veți obține „paritate” de voturi, veți calcula eroarea CVLOO folosind algoritmul 3-NN. Dacă nici atunci nu reușiți să eliminați paritatea, aplicați algoritmul 5-NN ș.a.m.d.

Răspuns:

Data	Eticheta	Vecinătate	Clasificare la CVLOO	Eroare?
A(1,2)		{ ... }		
B(1,4)				
C(2,1)				
D(2,3)				
E(3,4)				
F(4,3)				
G(4,5)				
H(5,1)				
I(5,3)				
J(5,5)				

Observație: La completarea coloanei *Vecinătate*, veți folosi litere în loc de coordonatele punctelor. (Pentru aceasta, este recomandabil ca, în prealabil, pe desenul de mai sus să puneți literele corespunzătoare punctelor, conform primei coloane din tabel.)

$$err_{CVLOO}(1-NN) = \dots$$

2.

(Distribuția gaussiană uni-variată: estimarea parametrilor;
calcularea unor probabilități a posteriori)

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 și presupunând că instanțele din clasa A au fost generate de o distribuție gaussiană, iar instanțele din clasa B au fost generate de o altă gaussiană, estimați *parametrii* acestor gaussiene, prin metoda verosimilității maxime (MLE).

Atenție!

Veți enunța mai întâi formulele corespunzătoare din capitolul *Estimarea parametrilor; metode de regresie*.

Centralizați rezultatele, completând tabelul de mai jos.

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

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

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

c. [Bonus]

Cum va clasifica varianta gaussiană a algoritmului Bayes Naiv (pe care am prezentat-o în mod succint la curs) punctul $X = 2$?

Veți exprima răspunsul mai întâi în funcție de α și β . Apoi veți face calculele folosind valorile lui α și β , determinate în funcție de *parametrii* calculați la punctul precedent.

Răspuns:

3. (Algoritmul K -means – varianta care folosește variabile-indicator γ_{ij} :
 aplicare pe date din \mathbb{R} ;
 Algoritmul EM/GMM, cazul $\sigma_1^2 = \sigma_2^2 = 1$, $\pi_1 = \pi_2$:
 executarea unei iterații, pe date din \mathbb{R})

A. Fie $X = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n\}$ o mulțime de instanțe de clusterizat, iar K numărul de clustere cu care vom lucra.

Veți folosi următoarea variantă a algoritmului K -means:

- Se inițializează în mod arbitrar centroizii $\mu_1, \mu_2, \dots, \mu_K$ și se ia $C = \{1, \dots, K\}$.
- Repetă:

Pasul 1:

Calculează matricea γ (de dimensiune $n \times K$ și având elemente din mulțimea $\{0, 1\}$) astfel:

$$\gamma_{ij} \leftarrow \begin{cases} 1, & \text{dacă } \|\mathbf{x}_i - \mu_j\|^2 \leq \|\mathbf{x}_i - \mu_{j'}\|^2, \forall j' \in C, \\ 0, & \text{în caz contrar.} \end{cases}$$

În caz de egalitate, alege în mod arbitrar cărui cluster (dintre cele eligibile) să-i aparțină instanța \mathbf{x}_i .

Pasul 2:

Recalculează μ_j folosind matricea γ actualizată:

Pentru fiecare $j \in C$, dacă $\sum_{i=1}^n \gamma_{ij} > 0$, asignează

$$\mu_j \leftarrow \frac{\sum_{i=1}^n \gamma_{ij} \mathbf{x}_i}{\sum_{i=1}^n \gamma_{ij}}.$$

Altfel, menține neschimbat centroidul μ_j .

până când matricea γ nu se mai schimbă de la o iterație la alta.

În continuare se va considera că $n = 2$, $\mathbf{x}_1 = 0.5$ și $\mathbf{x}_2 = 2$, iar valorile inițiale pentru *centroizii* μ_1 și μ_2 sunt 1 și respectiv 2. (Notăție: $\mu_1^{(0)} = 1$, $\mu_2^{(0)} = 2$.)

Aplicați algoritmul K -means (în varianta de mai sus!) pe aceste date.

Răspuns:

Inițializare:

$$\mu_1^{(0)} = 1, \mu_2^{(0)} = 2$$

Iterația 1:

Iterația 2:

Iterația 3:

...

B. Fie un model de mixtură gaussiană (engl., Gaussian mixture model, GMM) cu două componente având varianțe cunoscute și probabilități *a priori* egale pentru selecția celor două distribuții:

$$\frac{1}{2}\mathcal{N}(x; \mu_1, 1) + \frac{1}{2}\mathcal{N}(x; \mu_2, 1), \quad x \in \mathbb{R}.$$

În continuare se va considera (din nou) că $n = 2$, $x_1 = 0.5$ și $x_2 = 2$, iar valorile inițiale pentru *mediile* μ_1 și μ_2 sunt 1 și respectiv 2. (Notăție: $\mu_1^{(0)} = 1$, $\mu_2^{(0)} = 2$.)

Executați în mod manual o iterație a algoritmului EM, versiunea prezentată la curs (preluată din cartea *Machine Learning* a lui Tom Mitchell, pag. 193) pe aceste date, astfel:

a. Pasul E — estimarea probabilităților pentru variabilele „neobservabile“:

Pentru $i \in \{1, 2\}$ și $j \in \{1, 2\}$, calculați $P(Z_{ij} = 1 | X = x_i; \mu_1^{(0)}, \mu_2^{(0)})$, probabilitățile a posteriori de apartenență a datelor observate (x_1 și x_2) la cele două componente ale mixturii. (Vă readucem aminte că $Z_{ij} = 1$ dacă instanța x_i a fost generată de către gaussiană cu media μ_j , iar $Z_{ij} = 0$ în cazul contrar.) Justificați în mod detaliat!

Indicație: În vederea efectuării calculelor, pentru conveniență puteți considera valorile distribuției normale / gaussiene standard $\mathcal{N}(x; \mu = 0, \sigma^2 = 1)$ în punctele 0, 0.5, 1, 1.5 și 2 ca fiind respectiv 0.4, 0.35, 0.24, 0.13 și 0.05.

b. Pasul E (continuare) — calcularea funcției „auxiliare“:

Definiția funcției „auxiliare“ la iterația 1 este următoarea:

$$Q(\mu_1, \mu_2 | \mu_1^{(0)}, \mu_2^{(0)}) = E[\ln P(Y | \mu_1, \mu_2)],$$

unde

$Y \stackrel{\text{not.}}{=} \{y_1, y_2\}$, cu $y_i \stackrel{\text{not.}}{=} (x_i, Z_{i1}, Z_{i2})$, pentru $i \in \{1, 2\}$,

μ_1 și μ_2 sunt din \mathbb{R} și sunt considerați parametri liberi, iar

media E se calculează în funcție de mediile variabilelor „neobservabile“ Z_{ij} calculate [folosind probabilitățile a posteriori deduse] la punctul precedent.

Determinați formula (formulele) de calcul pentru

$$\mu^{(1)} \stackrel{\text{not.}}{=} \arg \max_{\mu \in \mathbb{R}^2} Q(\mu | \mu^{(0)}).$$

Am folosit notațiile $\mu = (\mu_1, \mu_2)$ și $\mu^{(t)} = (\mu_1^{(t)}, \mu_2^{(t)})$.

c. Pasul M — maximizarea funcției „auxiliare“ Q :

Calculați valorile parametrilor μ_1 și μ_2 la iterația 1 (adică $\mu_1^{(1)}$ și $\mu_2^{(1)}$), în funcție de probabilitățile calculate la primul punct. (Justificați în mod detaliat!) Care credeți că va fi tendința de mișcare a mediilor la următoarele iterații?

C. La curs am enunțat un *rezultat teoretic* (demonstrat în carte), care afirmă că în anumite condiții, algoritmul EM pentru mixturi de distribuții gaussiene se comportă *la limită* asemenea algoritmului K -means. Concret, la ce anume se referă această trecere la limită?

Folosind acest rezultat teoretic, propuneți o *schimbare* [minimală!] relativă la *setarea inițială* a parametrilor distribuțiilor gaussiene astfel încât, pe datele $X = \{x_1, x_2\}$ de la punctele A și B de mai sus, *la limită* pozițiile finale ale mediilor obținute de către EM să coincidă cu centroizii obținuți de către algoritmul K -means la convergență.

Răspuns:

a.

$$P(Z_{11} = 1 | x_1; \mu_1^{(0)}, \mu_2^{(0)}) = \dots$$

$$P(Z_{12} = 1 | x_1; \mu_1^{(0)}, \mu_2^{(0)}) = \dots$$

$$P(Z_{21} = 1|x_2; \mu_1^{(0)}, \mu_2^{(0)}) = \dots$$

$$P(Z_{22} = 1|x_2; \mu_1^{(0)}, \mu_2^{(0)}) = \dots$$

b.

$$y_1 \stackrel{not.}{=} (x_1, Z_{11}, Z_{12}) \Rightarrow \ln P(y_1|\mu_1, \mu_2) = \dots$$

$$y_2 \stackrel{not.}{=} (x_2, Z_{21}, Z_{22}) \Rightarrow \ln P(y_2|\mu_1, \mu_2) = \dots$$

$$Y \stackrel{not.}{=} \{y_1, y_2\} \Rightarrow \ln P(Y|\mu_1, \mu_2) = \dots$$

$$\Rightarrow Q(\mu_1, \mu_2|\mu_1^{(0)}, \mu_2^{(0)}) \stackrel{def.}{=} E[\ln P(Y|\mu_1, \mu_2)] = \dots$$

c.

Formula de calcul pentru $\mu_1^{(1)}$:

Justificare (demonstrație!):

Așadar, valoarea lui $\mu_1^{(1)}$ este:

Formula de calcul pentru $\mu_2^{(1)}$:

Justificare (demonstrație!):

Așadar, valoarea lui $\mu_2^{(1)}$ este:

C.