

2. (K-means vs. EM/GMM, cazul uni-variat ( $\pi, \mu, \sigma$ ))

Redăm mai jos pseudo-codul algoritmului EM — de fapt, doar partea sa iterativă, conținând cei doi pași, E și M — pentru clusterizare prin modelare de mixturi de gaussiene (GMM), cazul uni-variat, cu toți parametrii lăsați liberi. Îți cerem să faci schimbările minimale(!) necesare pentru a-l transforma într-un pseudo-cod corespunzător buclei principale a algoritmului de clusterizare K-means. Rescrie pașii care trebuie modificați, folosind spațiul lăsat disponibil sub fiecare pas. Dacă un pas nu necesită schimbări, scrie „Nicio modificare” în spațiul disponibil sub pasul respectiv. Dacă un anumit pas nu este necesar, scrie „Elimină acest pas” în spațiul disponibil sub pasul respectiv.

Pasul E:

Calculează probabilitatea  $p_{ij}^{(t)}$  de asignare a instanței  $x_i$  la clusterul [corespunzător gaussienei]  $j$ , ținând cont de valorile actuale ale parametrilor  $\pi^{(t)}, \mu^{(t)}, \sigma^{(t)}$ , unde  $t$  identifică iterația curentă a algoritmului EM.

$$p_{ij}^{(t)} = P(z_i = j | x_i, \pi^{(t)}, \mu^{(t)}, \sigma^{(t)}) \text{ pentru } i \in \{1, \dots, n\} \text{ și } j \in \{1, \dots, K\}$$

Am  $x_i \in G_j$  în  $\{x_i\}$  dacă  $|x_i - \mu_j| \leq |x_i - \mu_{j'}|$ ,  $\forall j' \in \{1, \dots, K\}$   
 altfel (dacă ... , se alege j-ul cel mai apropiat)

Pasul M:

A. Re-calculează probabilitățile a priori pentru fiecare cluster / componentă a mixturii:

$$\pi_j^{(t+1)} = \sum_{i=1}^n \frac{p_{ij}^{(t)}}{n} \text{ pentru } j \in \{1, \dots, K\}$$

Elimină ac pas

B. Re-calculează [centrozii clusterelor, reprezentați de] mediile distribuțiilor gaussiene care „modelează” clusterule:

$$\mu_j^{(t+1)} = \frac{\sum_{i=1}^n p_{ij}^{(t)} x_i}{\sum_{i=1}^n p_{ij}^{(t)}} \text{ pentru } j \in \{1, \dots, K\}$$

dacă reținem  $p_{ij}^{(t)}$  ca  $\frac{1}{n} \cdot \frac{1}{|G_j|}$  atunci expr de la B devine  $\mu_j = \frac{\sum_{x_i \in G_j} x_i}{|G_j|}$

C. Re-calculează varianțele distribuțiilor gaussiene care „modelează” clusterule:

$$(\sigma_j^2)^{(t+1)} = \frac{\sum_{i=1}^n p_{ij}^{(t)} (x_i - \mu_j^{(t+1)})^2}{\sum_{i=1}^n p_{ij}^{(t)}} \text{ pentru } j \in \{1, \dots, K\}$$

Elimină ac pas

3. (EM/GMM, cazul uni-variat: aplicarea manuală a unei iterații, pentru o mixtură de tipul  $\mu_1 = \mu_2$  (liber),  $\pi_1 = \pi_2, \sigma_1 = 1, \sigma_2 = 2$ )

Considerăm o mixtură de două gaussiene uni-variate, pentru care funcția de densitate de probabilitate (p.d.f.) are pentru o „observație” oarecare  $x$  expresia

$$\frac{1}{2} \mathcal{N}(x | \mu, 1) + \frac{1}{2} \mathcal{N}(x | \mu, 2^2).$$

În această formulă,  $\mathcal{N}(x | \mu, \sigma^2)$  desemnează, pentru  $x$ , valoarea densității distribuției normale uni-variate de medie  $\mu$  și varianță  $\sigma^2$ . Remarcați faptul că în acest model de mixtură probabilitățile de mixare / selecție sunt egale, apoi că mediile celor două componente sunt egale și, de asemenea, că deviațiile standard ale celor două componente au valorile fixate 1 și respectiv 2. Modelul acesta de mixtură are un singur parametru,  $\mu$ .

Presupunem că dorim să estimăm valoarea parametrului  $\mu$  prin metoda maximizării verosimilității, folosind algoritmul EM. Răspundeți la următoarele întrebări privitoare la modul în care operează pașii E și M ai acestui algoritm, atunci când considerăm cele trei instanțe / „observații” de mai jos:

4.0, 4.6, 2.0.

Vă punem la dispoziție un tabel cu anumite valori ale funcției de densitate pentru distribuția normală standard, de care veți avea probabil nevoie în rezolvarea acestui exercițiu:

$x$	0.0	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9	1.0
$\mathcal{N}(x   0, 1)$	.40	.40	.39	.38	.37	.35	.33	.31	.29	.27	.24
$x$	1.1	1.2	1.3	1.4	1.5	1.6	1.7	1.8	1.9	2.0	
$\mathcal{N}(x   0, 1)$	.22	.19	.17	.15	.13	.11	.09	.08	.07	.05	

a. Găsiți valoarea probabilităților condiționale calculate la pasul E, presupunând că la execuția precedentului pas M estimarea obținută pentru parametrul modelului a fost  $\mu = 4$  (iar  $\sigma_1 = 1$  și  $\sigma_2 = 2$ , totdeauna). Remarcați faptul că, întrucât probabilitățile condiționale pentru cele două componente ale mixturii trebuie să se sumeze la valoarea 1, este suficient să se calculeze  $p_{i1}^{(t)} = P(\text{componenta}_1 | x_i)$  pentru  $i = 1, 2, 3$ .

b. Folosind probabilitățile pe care le-ați calculat la punctul a, găsiți estimarea pentru parametrul  $\mu$  care va fi obținută la următoarea execuție a pasului M. Vă reamintim că pasul M maximizează media (engl., expected value) log-probabilității corelate pentru instanțele  $x_1, x_2$  și  $x_3$  și variabilele-indicator [pentru componentele mixturii], care sunt latente. Media aceasta se calculează în raport cu distribuția de probabilitate condițională a variabilelor-indicator [pentru componentele mixturii], care a fost determinată / calculată la precedentul pas E.

Sugestie: Nu este necesar ca în prealabil să faceți elaborarea formulelor algoritmului EM pentru acest caz specific de mixtură de distribuții gaussiene uni-variate. Este suficient să lucrați cu funcția „auxiliară”  $Q$  corespunzătoare iterației respective. Sau, mult mai simplu, puteți folosi formula de actualizare a mediilor de la problema precedentă (fiindcă forma ei se menține și pentru acest caz).

Andașe RE - parțial 2, 12.2018

$$3 \times 0.8 = 2.4$$

fundament calcul

$$\begin{array}{|c|c|c|c|} \hline 8.0, 2 & 8.0, 2 & 8.0, 2 & 4 \times 0.5 \\ \hline 1.6 & 1.6 & 3.2 & 2 \\ \hline 1.6 & 1.6 & 1.6 & 2 \\ \hline \end{array} \quad \begin{array}{|c|c|c|} \hline 3 \times 0.8 & \text{calcul} & \\ \hline 2.4 & 0.6 & 0.6 \\ \hline 3.2 & 3.6 & \\ \hline \end{array}$$

Bonus  
calcul Q

$$\text{penaliz } \frac{1}{3} \text{ pt } \frac{1}{3}$$