

Classification non supervisée mélanges Gaussiens & ALGO EM

I. Algorithme k-means.

Notations :

- Observations $x_1, \dots, x_N \in \mathbb{R}^D \mapsto c_1, \dots, c_N \in \llbracket 1, K \rrbracket$
(N obs, K classes, D features)
- $c_n(t)$: classe de l'individu n après t itérations.
- $r_{nk}(t)$: assignation $\begin{cases} 1 & \text{si } c_n(t) = k \\ 0 & \text{sinon} \end{cases} \quad \left(\begin{array}{l} n \in \llbracket 1, N \rrbracket \\ k \in \llbracket 1, K \rrbracket \end{array} \right)$
- $\mu_k(t)$: centroid / prototype de classe k

$$J(t) = \sum_{n=1}^N \|x_n - \mu_{c_n}(t)\|^2$$
$$= \sum_{n=1}^N \sum_{k=1}^K r_{nk}(t) \cdot \|x_n - \mu_k(t)\|^2$$

$$R = \begin{matrix} & \begin{matrix} c_1 & c_2 & \dots & c_K \end{matrix} \\ \begin{matrix} \text{ind}_1 \\ \vdots \\ \text{ind}_N \end{matrix} & \begin{bmatrix} r_{11} & r_{12} & \dots & r_{1K} \\ \vdots & \vdots & \ddots & \vdots \\ r_{N1} & r_{N2} & \dots & r_{NK} \end{bmatrix} \end{matrix} \quad \begin{matrix} Z=1 \\ \vdots \\ Z=K \end{matrix}$$

Algorithme K-means :

① Initialiser les centroids

② Répéter :

a. minimiser $J(t)$ par rapport aux $(r_{nk})_{n,k}$

$$c_n(t+1) = \underset{1 \leq j \leq K}{\operatorname{argmin}} \|x_n - \mu_j(t)\|$$

b. minimiser $J(t)$ par rapport aux $(\mu_k)_k$

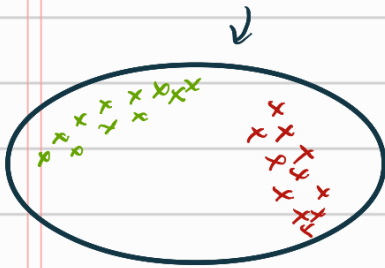
$$\frac{\partial J(t)}{\partial \mu_k(t)} = -2 \sum_n r_{nk}(t) (x_n - \mu_k(t))$$

$$= -2 \cdot \sum_{n=1}^N (x_n - \mu_k(t)) \cdot r_{nk}(t)$$

$$\frac{\partial J(t)}{\partial \mu_k} = 0 \iff \mu_k(t) = \frac{\sum_n r_{nk}(t) \cdot x_n}{\sum_n r_{nk}(t)} = \frac{1}{N_k} \sum_n x_n$$

$$\rightarrow \mu_k(t+1) = \frac{1}{N_k} \cdot \sum_{\substack{n=1 \\ c_n=k}}^N x_n$$

Rmq : . l'algorithme converge rapidement
 . Il fonctionne mal quand les clusters sont irréguliers, inhomogènes (variances différentes), ou anisotropes.



. Il n'est pas normalisé, donc on ne peut évaluer la performance de l'algo de façon standardisée pour toutes les applications (on sait qu'on veut J petit, mais pas à quel point)

. En particulier en grande dimension, J peut être très grand \rightarrow on peut utiliser une ACP en amont pour réduire la dim.

. K-means converge vers un min local qui dépend de l'initialisation : essayer plusieurs fois avec plusieurs initialisations différentes & sélectionner le meilleur J .

. kmeans ++ : initialisation qui espace les centroides les uns des autres le + possible.

\hookrightarrow accélérer & augmenter la perf.

. très sensible aux outliers.
 . ne prend pas en compte l'incertitude proche des frontières entre les clusters.

