

CLASSIFICATION NON SUPERVISEE, MELANGES

GAUSSIENS ET ALGORITHME 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 observations, K classes)
 $c_n(t)$: classe de l'individu x_n après t itérations
 $z_{nk}(t)$: assignation $\begin{cases} 1 & \text{si } c_n(t) = k \\ 0 & \text{sinon} \end{cases}$ ($n \in \llbracket 1, N \rrbracket, k \in \llbracket 1, K \rrbracket$)
 $\mu_k(t)$: centroïde / prototype de la classe k

$$\begin{aligned} J(t) &= \sum_{n=1}^N \|x_n - \mu_{c_n(t)}\|^2 \\ &= \sum_{n=1}^N \sum_{k=1}^K z_{nk}(t) \|x_n - \mu_k(t)\|^2 \end{aligned}$$

Algorithme K-means :

1. Initialiser les centroïdes.

2. Répéter :

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

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

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

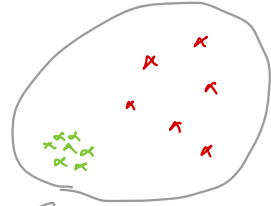
$$\frac{\partial J(t)}{\partial \mu_k(t)} = -2 \sum_{n=1}^N z_{nk}(t) (x_n - \mu_k(t))$$

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

$$\frac{\partial \mathcal{J}(t)}{\partial \mu_k(t)} = 0 \quad (\Rightarrow) \quad \mu_k = \frac{\sum_{n=1}^N r_{nk}(t) x_n}{\sum_{n=1}^N r_{nk}(t)}$$

$$= \frac{1}{N_k} \sum_{\substack{n=1 \\ c_n=k}}^N x_n$$

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



Rmq: ① l'algorithme converge rapidement

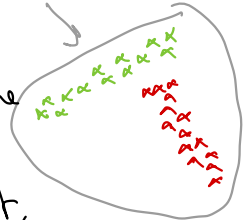
② Il fonctionne mal si les clusters sont irréguliers, inhomogènes (variances différentes), ou anisotropes

③ $\mathcal{J}(t)$ n'est pas normalisée, donc on ne peut évaluer la performance de l'algorithme de façon standardisée pour toutes les applications (on sait qu'on veut \mathcal{J} petit, mais pas à quel point).

④ En particulier en grande dimension, \mathcal{J} peut être très grand \rightarrow on peut utiliser une ACP en amont pour réduire la dimension.

⑤ K-means converge ^{vers} un minimum local qui dépend de l'initialisation : essayer plusieurs fois avec plusieurs initialisations différentes et sélectionner le meilleur \mathcal{J} .

⑥ kmeans++ : initialisation qui espace les centroïdes les uns des autres le plus possible.
↳ accélérer et augmenter la perf.



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