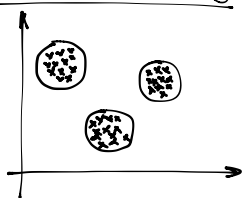


K-means clustering



Assumption: k clusters centred at unknown points (centroids)
 μ_1, \dots, μ_k

- find disjoint sets G_1, \dots, G_k of observations to minimize

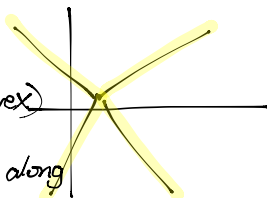
$$\sum_{j=1}^k \sum_{i \in G_j} \|x_i - \mu_j\|^2$$

$= \sum_{i=1}^n \min(\|x_i - \mu_1\|^2, \dots, \|x_i - \mu_k\|^2)$ total ^{within group} sum of squares

R function: `kmeans` ← several different options for algorithm
 ← can allow the algorithm to use more than 1 starting value.

- works best if have spherical clusters.
- transform variables?

Alternative approach:



- divide \mathbb{R}^p into k disjoint (convex) regions.
- Assume density of observations along boundaries is low.
- how to express this mathematically?

Model-based clustering

Model: Mixture model

X has $f(x) = \lambda_1 f_1(x, \theta_1) + \lambda_2 f_2(x, \theta_2) + \dots + \lambda_k f_k(x, \theta_k)$

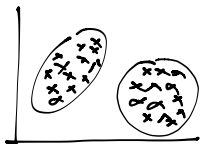
- $\lambda_1, \lambda_2, \dots, \lambda_k \geq 0$ (proportions) with $\lambda_1 + \dots + \lambda_k = 1$
 unknown

- $\theta_1, \dots, \theta_k$ are unknown parameter vectors

- Assumption: For each x , $f_i(x, \theta_i) f_j(x, \theta_j) = 0$ if $i \neq j$

Example: $f_j(x, \theta_j) = N_p(\mu_j, C_j)$

↑ cluster centre
 ↓ describe shape of cluster



Problem: Given data x_1, \dots, x_n , estimate $\lambda_1, \dots, \lambda_k$ and $\theta_1, \dots, \theta_k$

Maximum Likelihood estimation

Maximize $\ln L(\lambda_1, \dots, \lambda_k, \theta_1, \dots, \theta_k) = \sum_{i=1}^n \ln(\lambda_1 f_1(x_i, \theta_1) + \dots + \lambda_k f_k(x_i, \theta_k))$

EM algorithm

E: expectation M: maximization

Idea: Assume first that we know which cluster x_i belongs to.

i.e. we observe $(C_1, X_1), \dots, (C_n, X_n)$
 $\uparrow \qquad \qquad \uparrow$
cluster indicator where C_i takes values in $\{1, \dots, k\}$

Then the likelihood function becomes

$$L_c(\lambda_1, \dots, \lambda_k, \theta_1, \dots, \theta_k) = \prod_{i=1}^n \left\{ \prod_{j=1}^k [\lambda_j f_j(X_i, \theta_j)]^{I(C_i=j)} \right\} \quad \nearrow = \begin{cases} 1 & \text{if } C_i=j \\ 0 & \text{otherwise} \end{cases}$$

$$\ln L_c(\lambda_1, \dots, \lambda_k, \theta_1, \dots, \theta_k) = \sum_{i=1}^n \sum_{j=1}^k \{ I(C_i=j) \ln(\lambda_j) + I(C_i=j) \ln f_j(X_i, \theta_j) \}$$

$$= \ln L_c^{(1)}(\lambda_1, \dots, \lambda_k) + \ln L_c^{(2)}(\theta_1, \dots, \theta_k)$$

MLEs of $\lambda_1, \dots, \lambda_k$:

$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n I(C_i=j) = \text{proportion of cluster } j \text{ in sample}$$

But, we don't observe C_1, \dots, C_n !

EM algorithm: estimate $\Delta_{ij} = I(C_i=j)$. (for $i=1, \dots, n$

E-step: Estimate Δ_{ij} by $E[\Delta_{ij} | X_i; \underbrace{\hat{\lambda}_1, \dots, \hat{\lambda}_k, \hat{\theta}_1, \dots, \hat{\theta}_k}_{\text{correct estimates}}]$

$$= \frac{\hat{\lambda}_j f_j(X_i, \hat{\theta}_j)}{\sum_{l=1}^k \hat{\lambda}_l f_l(X_i, \hat{\theta}_l)} \quad (\text{Bayes rule})$$

M-step: Update estimates of $\hat{\lambda}_1, \dots, \hat{\lambda}_k$ and $\hat{\theta}_1, \dots, \hat{\theta}_k$

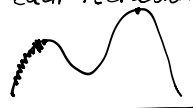
$$\hat{\lambda}_j = \frac{1}{n} \sum_{i=1}^n \hat{\Delta}_{ij}$$

and $\hat{\theta}_1, \dots, \hat{\theta}_k$ maximize $\ln L_c^{(2)}(\theta_1, \dots, \theta_k) = \sum_{i=1}^n \sum_{j=1}^k \hat{\Delta}_{ij} \ln f_j(X_i, \theta_j)$

Now, iterate E- and M-step until convergence!

Notes:

- ① Don't necessarily have convergence to MLEs i.e. maximizers of $\ln L(\lambda_1, \dots, \lambda_k, \theta_1, \dots, \theta_k)$
But at each iteration of E- and M-steps maximized $\ln L$ will increase.



- ② Performance of EM algorithm depends very strongly on initial estimates of $\theta_1, \dots, \theta_k$
- can usually take initial estimates of $\lambda_1, \dots, \lambda_k$ equal to $\frac{1}{k}$.