

L12

✓ Recap : PCA

□ Clustering

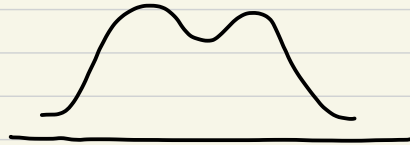
□ Expectation Maximization

□ k-Means

probability distribution



Mixture Model



$$p(x) = \sum_{k=1}^K \underbrace{p(X|c=k)}_{\text{conditional prob of } k\text{-th distribution}} \underbrace{p(c=k)}_{\text{probability of } k\text{-th distribution}}$$

$c = k$  的情况  $X$  的概率

Gaussian mixture model (GMM)

$$p(x) = \sum_{k=1}^K \frac{1}{\sqrt{2\pi}\sigma_k} e^{-\frac{(x-\mu_k)^2}{2\sigma_k^2}} p(c=k)$$

---

$$= \sum_{k=1}^K N(\mu_k, \sigma_k^2) p(c=k)$$

( we know  $k$  in advance )

Goal: Given samples from  $P(x)$ , estimate the parameters,  $\mu_k, \sigma_k$

Assume  $P(C=k) = \frac{1}{K}$

Approach: Maximum likelihood

Input: I.I.D. samples  $\{x_1, x_2, \dots, x_n\}$

Output:  $\{(\mu_k, \sigma_k^2)\}_{k=1}^K$

Optimize: log-likelihood

$$P = \prod_{i=1}^n \sum_{k=1}^K N(x_i, \mu_k, \sigma_k^2) \underbrace{\frac{1}{K}}_{\text{ignore}}$$

log-likelihood

$$L = \sum_{i=1}^n \log \left( \sum_{k=1}^K N(x_i, \mu_k, \sigma_k^2) \right)$$

Simpler case:  $K=2$

$$L = \sum_{i=1}^n \log \left[ \frac{1}{\sqrt{2\pi}\sigma^2} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right]$$

$$L \approx \sum_{i=1}^n -\frac{1}{2} \log(2\pi\sigma^2) - \left[ \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

Maximize wrt  $\mu, \sigma$

$$\frac{\partial L}{\partial \mu} = 0$$

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial L}{\partial \sigma^2} = 0$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mu})^2$$

General Case :

$$\frac{\partial L}{\partial \mu_k} = 0 \quad \downarrow \frac{\partial}{\partial \mu_k} :$$

$$\sum_{i=1}^n \left[ \frac{1}{\sum_{i=1}^n N(X_i; \mu_k, \sigma_k^2)} N(X_i; \mu_k, \sigma_k^2) \cdot \frac{(X_i - \mu_k)}{\sigma_k^2} \right] = 0$$

$\rightarrow \gamma_i(k) : \text{itself a probability distribution}$

Posterior distribution:  $P(c=k | X=X_i)$

$$\sum_{i=1}^n \gamma_i(k) \cdot \left( \frac{X_i - \mu_k}{\sigma_k^2} \right) = 0$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \gamma_i(k) X_i}{\sum_{i=1}^n \gamma_i(k)}$$

$$\hat{\sigma}_k^2 = \frac{\sum_{i=1}^n \gamma_i(k) (X_i - \hat{\mu}_k)^2}{\sum_{i=1}^n \gamma_i(k)}$$

if Oracle tells us posterior distribution  $(\gamma_i(k))$   
then we can compute  $\mu_k, \sigma_k^2$

[ Expectation E-step ]

if oracle tells us means  $\mu_k$ , variances  $\sigma_k^2$ ,  
then we can compute  $\gamma(k)$

[Maximization, M-step]

Alternatively iterate  $\rightarrow$  EM

Once  $\mu_k, \sigma_k^2$  are computed, then  $\gamma(k)$  gives  
 $P(c=k | X=x_i)$



k-means

$\hookrightarrow$  Hard clustering

EM soft-clustering

知道几类, 但没label.

初始时  $k$  个簇中心

input :  $\{X_1, X_2, \dots, X_n\}$

Output :  $S_1 \cup S_2 \cup \dots \cup S_k = X$       subsets

partition of data

$\{\mu_1, \mu_2, \dots, \mu_k\} \rightarrow$  "cluster centers"

Loss function : "k-means objective function"

$$L = ([S_1, S_2, \dots, S_k], \{\mu_1, \dots, \mu_k\})$$

$$= \sum_{k=1}^K \sum_{x_i \in S_k} \|x_i - \mu_k\|_2^2$$

Simpler Case :  $k=1$

Optimize :  $\sum_{x_i \in S} \|x_i - \mu\|_2^2$

write  $\mu$

Minimize Sample mean

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

General Case

?

Assume that cluster identities of each data points are known

Optimize over  $\sum_{x_i \in S_k} \|x_i - \mu_k\|_2^2 \leftarrow O(nd)$

$$\hat{\mu}_k = \frac{1}{|S_k|} \sum_{x_i \in S_k} x_i$$

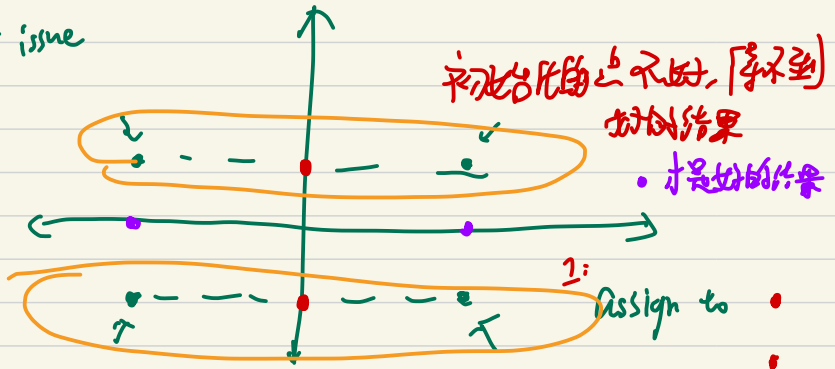
— Assume clusters centers are known. Assign each data points to the nearest  $\mu_k$ .

For each  $i$ , assign  $\hat{c}_i = \arg \min_k \|x_i - \mu_k\|_2 \leftarrow O(ndk)$

— alternative between these steps  
→ k-Means

Worst case:  $k^n$  possible possible cluster choices

Other issue



Caveats : 1) sensitive to initialization

2) exponential running time

3) generally circular/clusters  
spherize

2: compute mean of



k-Means ++

↳ smart initialization

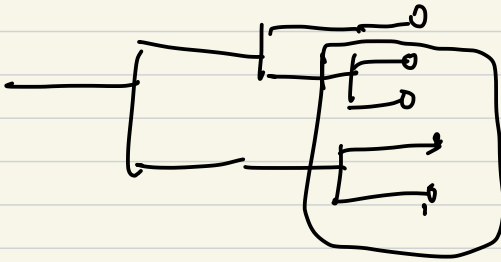
Approaches requires no knowledge of  $k$ .

— Hierarchical clustering

→ Bioinformatics

→ Genomics

→ Linguistics



— Bayesian methods

---