

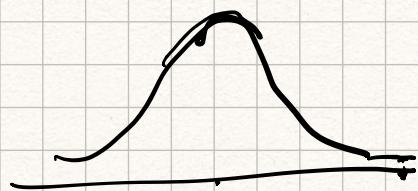
ECE-GY 6143  
Intro to ML

✓ Recap: PCA

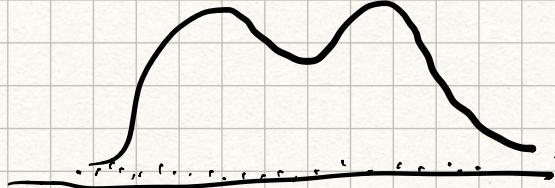
✓ Clustering

- ✓ Expectation - Maximization
- ✓ K-Means

Probability distribution:



Mixture model:



$$P(x) = \sum_{k=1}^K P(X | c=k) P(c=k)$$

Conditional prob.  
of  $k^{th}$  distribution

Probability  
of  $k^{th}$  distrib.

Gaussian mixture model (GMM):

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

$$:= \frac{\sum_{k=1}^K N(\mu_k, \sigma_k^2)}{P(C=k)}$$

Goal: Given samples from  $P(x)$ ,

estimate the parameters  $\mu_k, \sigma_k$ .  $\equiv \begin{cases} \text{Assume } P(C=k) \\ = 1/K \end{cases}$

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.

$$\text{(Likelihood)} P = \prod_{i=1}^n \sum_{k=1}^K N(x_i; \mu_k, \sigma_k^2) \cdot \left(\frac{1}{K}\right)$$

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:  $K=1$ .

$$(x - \mu)^2$$

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

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

Maximize w.r.t.  $\mu, \sigma$ .

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

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

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

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

General case:

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

$$\sum_{i=1}^n \frac{N(x_i; \mu_u, \sigma_u^2)}{\sum_{i=1}^k N(x_i; \mu_u, \sigma_u^2)}$$

$$\left( \frac{x_i - \mu_u}{\sigma_u^2} \right)^2 = 0$$

$y_i(k)$ : itself a probability distribution.

$$\text{Posterior distribution} = p(C=k | X=x_i)$$

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

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

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

- If oracle tells us posterior distribution ( $r_i(k)$ ) 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  $r_i(k)$ .

[ Maximization, M-step ]

Alternatively iterate  $\rightarrow$  [EM]

Once  $\mu_u, \tau_u$  are computed, then  
 $\gamma_i(u)$  gives  $P(C=k \mid X=x_i)$ .

K-means

EM: soft-clustering

$\hookrightarrow$  Hard clustering

Input:  $\{x_1, x_2, \dots, x_n\} = X$

Output:  $S_1 \cup S_2 \cup \dots \cup S_K = X$   
 partition of  $X$ .

$\{\mu_1, \mu_2, \dots, \mu_K\} \rightarrow$  "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  
w.r.t  $\mu$

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

Minimizer : Sample mean

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

General case

- Assume that cluster identities of each data point are known.

Optimize over  
 $\mu_k$

$$\sum_{x_i \in S_k} \|x_i - \mu_k\|_2^2$$

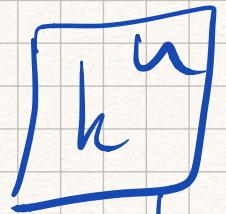
$$\hat{\mu}_k = \frac{1}{|S_k|} \sum_{x_i \in S_k} x_i$$

$O(nd)$

- Assume cluster centers are known.

- Assign each data point to the nearest  $\mu_k$ .
  - For each  $i$ , assign  $\hat{k}_i = \arg\min_k \|x_i - \mu_k\|_2^2$
  - Alternating between these steps  $O(ndk)$
- K-Means

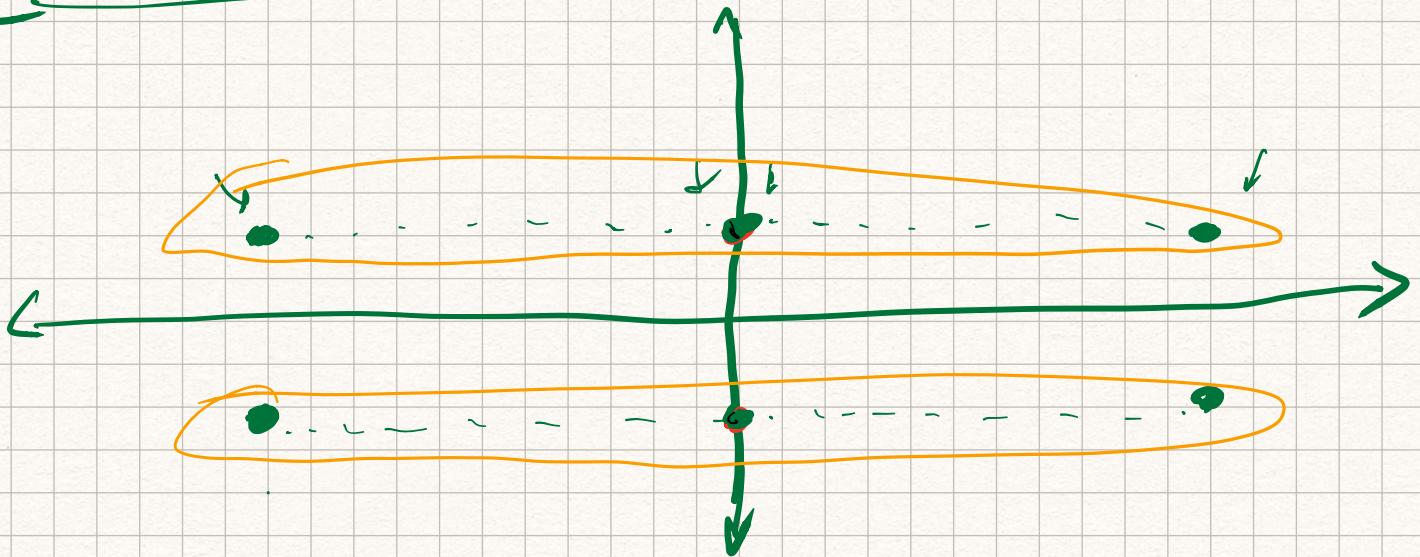
Worst Case:



possible

cluster choices.

Other issue



- Caveats:
- 1) Sensitive to initialization
  - 2) Exponential running time

3) Generally circular/ clusters.  
Spherical.

K-Means ++

→ Smart initialization

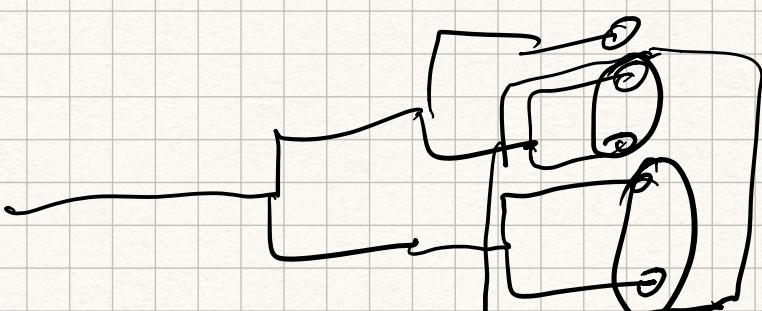
Approaches that do not require  
knowledge of  $k$

- Hierarchical clustering -

→ Bioinformatics

↳ Genomics

↳ Linguistics-



- Bayesian methods

