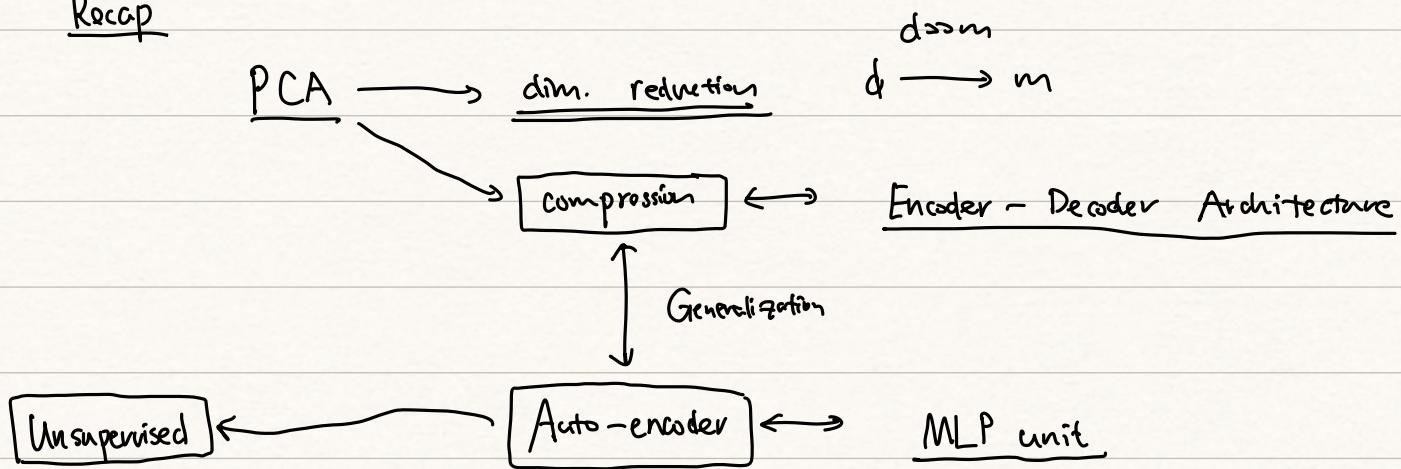


Recap



Unsupervised { cluster: K-means
DE cluster } → GMM

① cluster

$$\text{a) Defn} \quad \rightarrow \textcircled{1} \quad \bigcup_{i=1}^k D_i$$

② within same group, similar
different groups, dissimilar → similarity metric

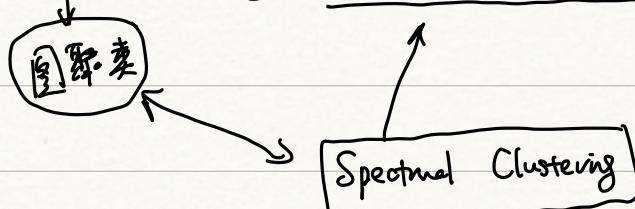
b) Why important?

① information retrieve ← search → match cluster

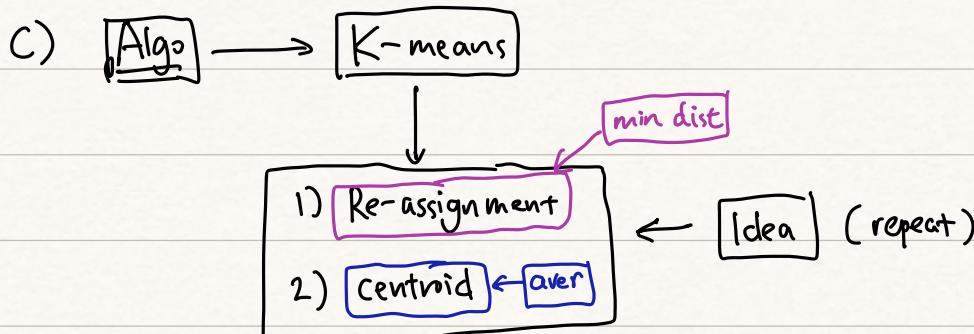
② Genetics

③ Market Research → \hat{P} → partition ← target doctor

④ Social Network Analysis



⑤ Anomaly Detection → focus on outliers



→ Mathematically formulate

① Assignment Matrix $R \in \mathbb{R}^{N \times K}$

$$r_{n,k} = \begin{cases} 1, & x_n \text{ is assigned to cluster } k \\ 0, & \text{otherwise} \end{cases}$$

property →

$$\left\{ \begin{array}{l} R \cdot \mathbf{1}_K = \mathbf{1}_N \leftrightarrow \text{行和为 1} \\ R^T \cdot \mathbf{1}_N = (\# \text{ of points in each cluster}) \end{array} \right.$$

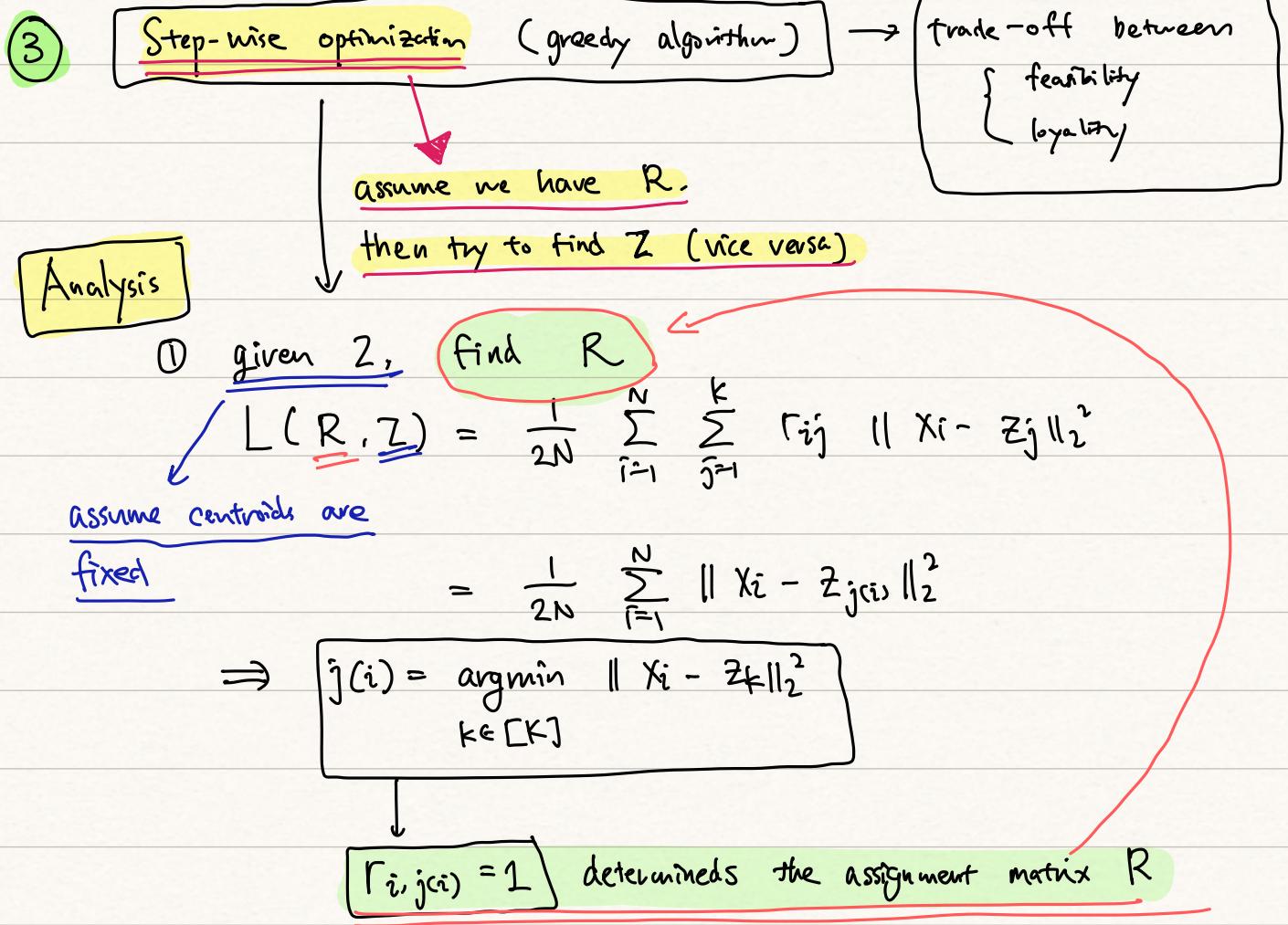
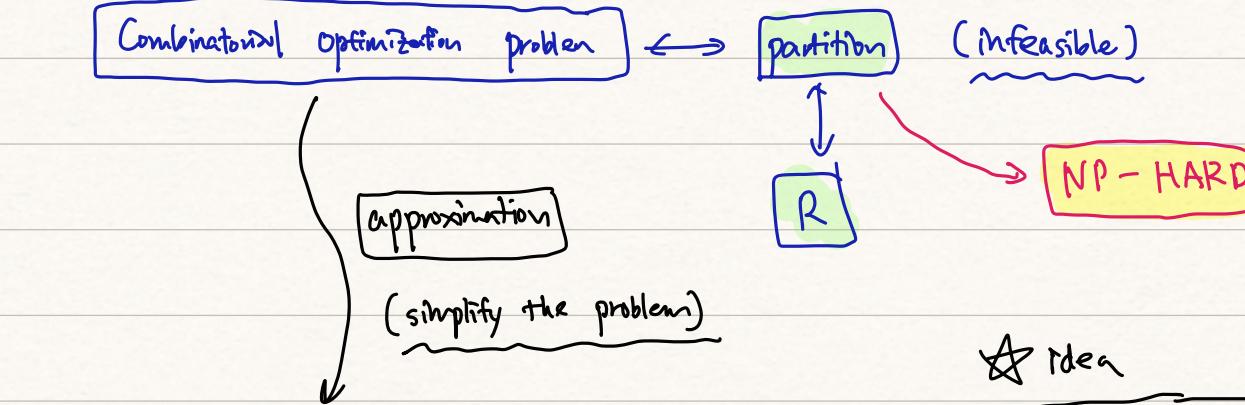
② Minimize some loss

$$J(R, Z) = \sum_{i=1}^N \sum_{j=1}^K r_{ij} \|x_i - z_j\|_2^2$$

assignment Centroid

Each row (i -th), there is only
 One non-zero entry

$$= \sum_{i=1}^N \cdot r_{i,r(i)} \|x_i - z_{r(i)}\|_2^2$$



② given R , find Z . → intuitively, we should choose average within each cluster

assume that assignment is fixed

$$L(R, Z) = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^K r_{ij} \|x_i - z_j\|_2^2$$

PFI

$$= \frac{1}{2N} \sum_{j=1}^K \left(\sum_{i=1}^N r_{ij} \|x_i - \bar{z}_j\|_2^2 \right)$$

$$\hat{z}_j = \frac{\sum_{i=1}^N r_{ij} x_i}{\sum_{i=1}^N r_{ij}}$$

average within cluster

Pf 2

calculate $\nabla_z L(R, z)$ directly

$$\frac{\partial L(R, z)}{\partial z_p} = \sum_{i=1}^N r_{ip} (x_i - \hat{z}_p) = 0$$

$$\Rightarrow \hat{z}_p = \frac{\sum_{i=1}^N r_{ip} x_i}{\sum_{i=1}^N r_{ip}} := N_p$$

④

Initialization → generate Z^0 randomly

have some Guideline to modify the algorithm

k-means ++

⑤

Guarantee to converge with finite steps!

⇒ although we may not achieve the optimal

Solution

highly likely not to

depends on the initialization

⑥

[Limitation]

a) Boundary points → lack of uncertainty

One data can only belong

to one cluster (hard assignment)

Motivation for GMM

soft assignment

② GMM

$$\sum_{k=1}^K p(x, z=z_k) = \sum_{k=1}^K p(z=z_k) \cdot p(x|z=z_k)$$

Model:

$$p(x) = \sum_{k=1}^K \pi_k \cdot P(x|C_k)$$

P(x | C_k)

Gaussian Model
P(x | μ_k, Σ_k)

is actually a distribution

Sampling (Monte Carlo)

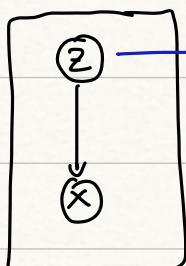
- ① from $\pi = (\pi_1, \dots, \pi_K) \rightarrow c$
- ② from $X|C_k=c \rightarrow \text{Gaussian}$

Graphic Model

(latent)

hidden variable

can apply EM to make approximation



$$\ln = \sum_{i=1}^n \log P(x_i | \theta) \quad \underline{\theta = \{\{\pi_k\}_{k=1}^K, \{\mu_k, \Sigma_k\}_{k=1}^K\}}$$

$$= \sum_{i=1}^n \log \sum_{k=1}^K \pi_k \cdot P(x_i | C_k) \rightarrow \text{difficult to optimize}$$

$$= \sum_{i=1}^n \log \mathbb{E}_{z \sim \pi} P(x_i | z)$$

optimize

Cross-term

Preparation for EM

Posterior Dist for π_k

$$\Rightarrow P(z=z_k | x) = \frac{P(x|C_k) \cdot P(C_k)}{\sum_k P(x|C_k) \cdot P(C_k)}$$



improvement compared with hard assignment K-means

⇒ we are able to say:

for one data point x , its component in each class → means a lot in application of cluster

* we can say, for one user:

{ 50% read books
{ 50% watch TV something like that!



补充完3

Optimization

$$\Theta = \{\mu_k, \Sigma_k, \pi_k\}_{k=1}^K$$

$$P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)$$

→ Directly solve MLE! (not EM Framework)

$$= \frac{1}{(\sqrt{2\pi})^d} \frac{1}{|\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right)$$

$$\underbrace{LL := \ell_n(\Theta)}_{\text{log-likelihood}} = \sum_{i=1}^N \log \left\{ \sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\}$$

log-likelihood

$$\textcircled{1} \quad \frac{\partial \ell_k}{\partial \mu_k} = \sum_{i=1}^N \frac{\pi_k \cdot \frac{\partial}{\partial \mu_k} P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)}{\left\{ \sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\}}$$

$$= \sum_{i=1}^N \frac{\pi_k \cdot (-\Sigma_k^{-1} \cdot (x_i - \mu_k)) \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)}{\sum_{k=1}^K \pi_k P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)}$$

Not closed-form !!

$$= (-\Sigma_k^{-1}) \cdot \sum_{i=1}^N r_{ik} \cdot (x_i - \mu_k)$$

$$r_{ik} = P(C_k | x_i)$$

$$\uparrow \text{posterior dist.}$$

$$\Rightarrow \mu_k = \frac{\sum_{i=1}^N r_{ik} \cdot x_i}{\sum_{i=1}^N r_{ik}}$$

Soft-assignment to cluster k

$$\textcircled{2} \quad S_k := \Sigma_k^{-1} \rightarrow \frac{\partial \ell_k}{\partial S_k} \rightarrow \text{need to check}$$

$$\hat{L} = \sum_{i=1}^N \log \left\{ \sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\}$$

$$\frac{\partial \hat{L}}{\partial \Sigma_p} = \sum_{i=1}^n \frac{\partial}{\partial \Sigma_p} \left[\log \left\{ \sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\} \right]$$

$$= \sum_{i=1}^n \frac{\frac{\pi_p}{\sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)}}{\cdot \frac{\partial}{\partial \Sigma_p} P_{\text{Gauss}}(x_i | \mu_p, \Sigma_p)}$$

$$P_{\text{Gauss}}(x_i | \mu_p, \Sigma_p) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_p|^{\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (x_i - \mu_p)^T \Sigma_p^{-1} (x_i - \mu_p) \right\}$$

$$\frac{\partial P}{\partial \Sigma_p} = (2\pi)^{-\frac{d}{2}} \left[\frac{\partial \det(\Sigma_p)}{\partial \Sigma_p} \cdot \exp \{ * \} + \frac{\partial \exp \left\{ -\frac{1}{2} (x_i - \mu_p)^T \Sigma_p^{-1} (x_i - \mu_p) \right\}}{\partial \Sigma_p} \cdot \det(\Sigma_p)^{-\frac{1}{2}} \right]$$

$$\begin{aligned} &= (2\pi)^{-\frac{d}{2}} \left[-\frac{1}{2} \det(\Sigma_p)^{-\frac{3}{2}} \cdot \det(\Sigma_p) \cdot \Sigma_p^{-1} \cdot \exp \{ * \} \right. \\ &\quad \left. + \exp \{ * \} \cdot \det(\Sigma_p)^{-\frac{1}{2}} \cdot \left(-\frac{1}{2} \right) \cdot \left(-(x_i - \mu_p)(x_i - \mu_p)^T \cdot \Sigma_p^{-2} \right) \right] \\ &= -\frac{1}{2} (2\pi)^{-\frac{d}{2}} \det(\Sigma_p)^{-\frac{1}{2}} \exp \{ * \} \left[\Sigma_p^{-1} - (x_i - \mu_p)(x_i - \mu_p)^T \cdot \Sigma_p^{-2} \right] \end{aligned}$$

$$= \sum_{i=1}^n \frac{\pi_p}{\sum_{k=1}^K \pi_k \cdot P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k)} \left[-\frac{1}{2} P_{\text{Gauss}}(x_i | \mu_p, \Sigma_p) \cdot \Sigma_p^{-2} (\Sigma_p - (x_i - \mu_p)(x_i - \mu_p)^T) \right]$$

$$= -\frac{1}{2} \Sigma_p^{-2} \sum_{i=1}^n \gamma_{ip} (\Sigma_p - (x_i - \mu_p)(x_i - \mu_p)^T)$$

$$\Rightarrow \Sigma_p = \frac{\sum_{i=1}^n \gamma_{ip} (x_i - \mu_p)(x_i - \mu_p)^T}{\sum_{i=1}^n \gamma_{ip}}$$

Not closed form

③ \rightarrow $\text{Lagrange Multiplier}$

$$L = L_L + \lambda \left(1 - \sum_{k=1}^K \pi_R \right)$$

$$\frac{\partial \mathcal{L}}{\partial \pi_e} = \sum_{i=1}^N \frac{1}{\pi_e} \frac{\pi_e P_{\text{Gauss}}(x_i; \mu_e, \Sigma_e)}{\sum \dots - \dots} - \lambda = 0$$

$$\Rightarrow \sum_{i=1}^N x_{ip} - \pi_p \cdot \lambda = 0 \quad \left\{ \begin{array}{l} \sum_{p=1}^K \pi_p = 1 \\ \sum_{p=1}^K \sum_{i=1}^N x_{ip} = \sum_{i=1}^N 1 = N \end{array} \right.$$

$$\Rightarrow N - \lambda = 0 \Rightarrow \boxed{\lambda = N}$$

$$\Rightarrow \Pi_p = \frac{\sum_{i=1}^N \delta_{ie}}{N}$$

- not closed form
- only solve iteratively

To conclude: after solving MLE, we achieve:

$$\textcircled{1} \quad M_p = \frac{1}{\sum_{i=1}^n r_{ip}} \cdot \sum_{i=1}^n r_{ip} x_i$$

$$\textcircled{2} \quad \Sigma_e = \frac{1}{\sum_{i=1}^n f_{ie}} \cdot \sum_{i=1}^n f_{ie} (x_i - \bar{x}_e)(x_i - \bar{x}_e)^T$$

$$\textcircled{3} \quad \bar{x}_{\text{P}} = \frac{\sum_{i=1}^n x_{i\text{P}}}{N}$$

$$\text{where } \pi_{iR} = \frac{\pi R \text{Gauss}(x_i | \mu_e, \Sigma_e)}{\sum_{k=1}^K \pi_R \text{Gauss}(x_i | \mu_k, \Sigma_k)}$$

Self - Recall

EM \rightarrow GMM

Notation

$z_{ik} \rightarrow i\text{-th } x \Rightarrow k\text{-th class}$

Primal formulation:

$$\begin{aligned} \mathcal{L}_1 &= \log L_n(X|\theta) = \log \left\{ \prod_{i=1}^n \sum_{k=1}^K p(z_{ik}) \cdot P(x_i|z_{ik}) \right\} \\ &= \sum_{i=1}^n \log \left\{ \sum_{k=1}^K p(z_{ik}) \cdot P(x_i|z_{ik}) \right\} \end{aligned}$$

EM Approximation

Effectiveness of EM Algo can be proved by MM

$$\mathcal{L}_2 = \mathbb{E}_{Z \sim p_C(X, \theta^e)} [\log L_n(X, Z|\theta)] \quad \text{Framework}$$

$$\theta^{t+1} = \arg \max_{\theta \in \Theta} \mathcal{L}_2 = \arg \max_{\theta \in \Theta} \mathbb{E}_{Z \sim p_C(X, \theta^e)} [l_n(X, Z|\theta)]$$

$$l_n(X, Z|\theta) = \prod_{i=1}^n P(x_i | z_{i,k(i)})$$

$$= \prod_{i=1}^n \prod_{k=1}^K p(x_i | z_{i,k})$$

It can be proved:
 $l_n(\theta^{t+1}) \geq l_n(\theta^e)$

$$\Rightarrow \log l_n(X, Z|\theta) = \sum_{i=1}^n \sum_{k=1}^K \delta_{i,k} \log P(x_i | z_{i,k})$$

$$= \sum_{k=1}^K \left\{ \sum_{i=1}^n \delta_{i,k} \log P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\}$$

$$\Rightarrow \mathbb{E}_{Z \sim p_C(X, \theta^e)} [l_n(X, Z|\theta)] = \sum_{k=1}^K \left\{ \sum_{i=1}^n \mathbb{E} [\delta_{i,k} | X, \theta^e] \log P_{\text{Gauss}}(x_i | \mu_k, \Sigma_k) \right\}$$

$$\mathbb{E} [\delta_{i,k} | X, \theta^e] = P(x_i \rightarrow \text{class } k | X, \theta^e)$$

$$= \frac{\pi_k^e P_G(x_i | \mu_k^e, \Sigma_k^e)}{\sum_{k=1}^K \pi_k^e \cdot P_G(x_i | \mu_k^e, \Sigma_k^e)}$$