

Unsupervised Learning

Task

- a) Dimensionality Reduction
- b) Clustering
- c) Density Estimation \leftrightarrow Generative Model

① PCA Principal Component Analysis

Big Idea: Reduced Representation via linear projection (transformation)

Dataset $\mathcal{D} = \{(x_i)\}_{i=1}^N \quad x_i \in \mathbb{R}^d$

1. Maximize Variance Formulation

$z_i = u^\top x_i \longrightarrow$ Interpretation: the coordinate after projecting to u.



Variance (estimator) of $\{z_1, \dots, z_n\}$ ($\text{Var}[z] = \mathbb{E}[(z - \mathbb{E}[z])^2]$)

$$= \frac{1}{N} \sum_{i=1}^N (z_i - \bar{z})^2, \quad \bar{z} = \frac{1}{N} \sum_{i=1}^N u^\top x_i = u^\top \bar{x}$$

$$= \frac{1}{N} \sum_{i=1}^N u^\top (x_i - \bar{x})(x_i - \bar{x})^\top u$$

Therefore, one realization to calculate variance is:

$$\mathcal{D} = \{x_i\}_{i=1}^N \xrightarrow{x_i - \bar{x}} \widetilde{\mathcal{D}} = \{\tilde{x}_i\}_{i=1}^N \longrightarrow S = \frac{1}{N} \sum_{i=1}^N \tilde{x}_i \tilde{x}_i^\top$$

\Rightarrow variance of $\{z_1, z_2, \dots, z_n\}$ $z_i = u^\top x_i$

$$= u^\top S u$$

Maximize variance \Leftrightarrow

$$\max_u u^T S u \Leftrightarrow$$

s.t. $\|u\|_2^2 = 1$

1-st principal component

u : eigenvector corresponds to the largest eigenvalue

to calculate the 2-nd principal component, the formulation is:

$$\begin{bmatrix} \max_u u^T S u \\ \text{s.t. } \|u\|_2^2 = 1 \text{ & } \underline{u^T u_1 = 0} \end{bmatrix}$$

\rightarrow u : eigenvector corresponds to the second largest eigenvalue

Thus, given m (# of Principal Components we need), we can achieve

$$\rightarrow ① \underbrace{Z_m = X U_m}_{m \text{ principal component scores}}$$

$$\begin{bmatrix} U_m = (u_1, \dots, u_m) \in \mathbb{R}^{d \times m} \\ X = \begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix} \in \mathbb{R}^{N \times d} \end{bmatrix}$$

本质: Basis transformation

Z_m

Interpretation (the projection coordinates for x_1, \dots, x_N towards direction u_1, \dots, u_m)

② U_m \rightarrow projection direction

Interpretation: for the given principal component, the contribution for each dimension (row).

2. Minimize Re-construct Error Formulation

given a set of orthonormal basis $\{u_1, \dots, u_d\}$

$$\text{then } x_i = \sum_{j=1}^d (u_j^\top x_i) \cdot u_j \quad \hat{x}_i = \sum_{j=1}^m \beta_{ij} u_j$$

Our aim is: choose $\{\beta_{ij}\}$ & $\{u_1, \dots, u_d\}$

so that we can minimize $\frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|_2^2$

$$\text{Calculation: } \frac{1}{N} \sum_{i=1}^N \|x_i - \hat{x}_i\|_2^2$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\sum_{j=1}^m (\alpha_{ij} - \beta_{ij})^2 + \sum_{j=m+1}^d \alpha_{ij}^2 \right]$$

$$= \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^m (\alpha_{ij} - \beta_{ij})^2 + \frac{1}{N} \sum_{i=1}^N \sum_{j=m+1}^d \alpha_{ij}^2$$

$$\boxed{\text{Here, } \alpha_{ij} = u_j^\top x_i}$$

$$\text{Analysis: } ① \beta_{ij} = u_j^\top x_i \quad i=1, 2, \dots, N \quad j=1, 2, \dots, m$$

when u_j is determined

$$② \{ \hat{u}_1, \dots, \hat{u}_d \} = \underset{\{u_1, \dots, u_d\}}{\operatorname{argmin}} \frac{1}{N} \sum_{i=1}^N \sum_{j=m+1}^d u_j^\top x_i x_i^\top u_j$$

$$= \underset{\{u_1, \dots, u_d\}}{\operatorname{argmin}} \sum_{j=m+1}^d u_j^\top S u_j$$

$\{u_1, \dots, u_d\}$ is orthonormal basis



$$\begin{cases} u_i^\top u_j = 0 & i \neq j \\ u_i^\top u_i = 1 & i=1, 2, \dots, d \end{cases}$$

\Rightarrow one choice of $\{\hat{u}_1, \dots, \hat{u}_d\}$ is:

the set of eigenvector of S - ordered in decreasing eigenvalue

$$\text{Re-construction : } \mathbf{x} \mapsto \sum_{j=1}^m (\mathbf{u}_j^\top \mathbf{x}) \cdot \mathbf{u}_j = \sum_{j=1}^m z_j \mathbf{u}_j = \mathbf{u}_m \cdot \mathbf{z}$$

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{pmatrix} \mapsto \mathbf{X}(\mathbf{u}_1, \dots, \mathbf{u}_m) \begin{pmatrix} \mathbf{u}_1^\top \\ \vdots \\ \mathbf{u}_m^\top \end{pmatrix}$$

$$= \underbrace{\mathbf{X} \mathbf{u}_m \cdot \mathbf{u}_m^\top}_{\downarrow}$$

$$\mathbf{z}_m = \begin{pmatrix} z_1^\top \\ \vdots \\ z_N^\top \end{pmatrix}$$

(coordinates)
Principal component scores

② Kernel PCA $\xrightarrow{\star}$ not so convenient

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{pmatrix}$$

$$\text{For PCA, } \underline{\text{our interest}} \text{ is : } S = \underbrace{\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \cdot \mathbf{x}_i^\top}_{\sim} = \frac{1}{N} \mathbf{X}^\top \mathbf{X}$$

In Kernel PCA Framework, $\mathbf{x} \in \mathbb{R}^d \longleftrightarrow \phi(\mathbf{x}) \in \mathbb{R}^D$

Then. Algo becomes:

a) $\Phi_{ij} := \phi_j(\mathbf{x}_i) \in \mathbb{R}^{M \times D}$

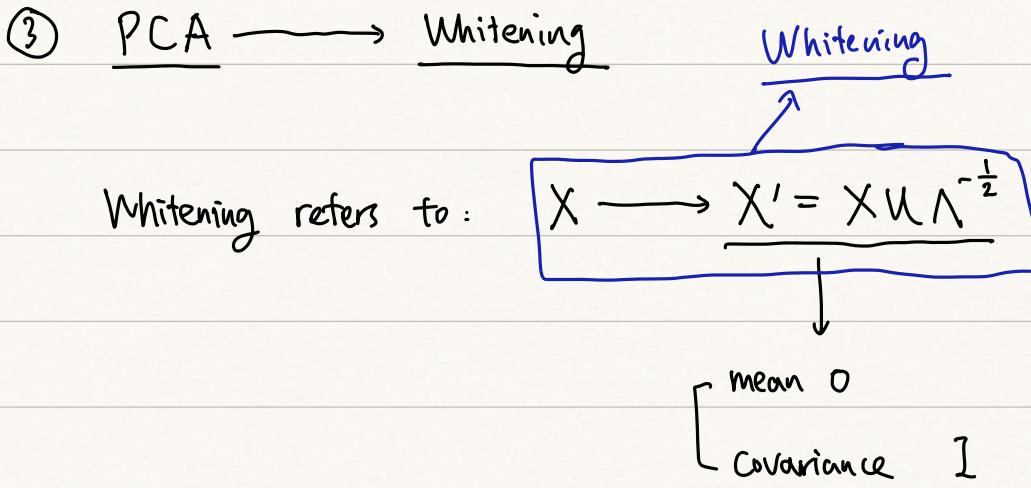
b) Centralized $\bar{\Phi} = \Phi - \frac{1}{N} \mathbf{1}_N \mathbf{1}_N^\top \Phi$

c) $S_\Phi = \frac{1}{N} \bar{\Phi}^\top \bar{\Phi}$

d) Eigenvalue / vector Decomposition

$$\begin{cases} \mathbf{U}_m = (\mathbf{u}_1, \dots, \mathbf{u}_m) & \mathbf{u}_i \in \mathbb{R}^D \\ \Lambda_m = (\lambda_1, \dots, \lambda_m) \end{cases}$$

e) $\mathbf{z}_m = \bar{\Phi} \mathbf{u}_m$

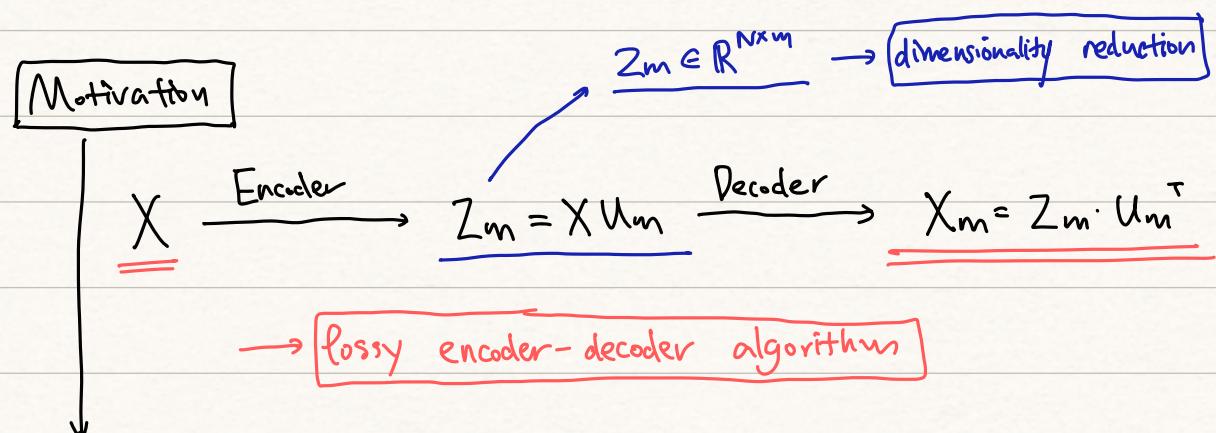


Check: Covariance of X'

$$\begin{aligned}
 &= \frac{1}{N} X'^T X' \\
 &= \frac{1}{N} \Lambda^{-\frac{1}{2}} U^T X^T X U \Lambda^{-\frac{1}{2}} \\
 &= \Lambda^{-\frac{1}{2}} U^T S U \Lambda^{-\frac{1}{2}} \\
 &= \Lambda^{-\frac{1}{2}} \Lambda \Lambda^{-\frac{1}{2}} \\
 &= I
 \end{aligned}$$

Note:
 $S \cdot U = U \cdot \Lambda$
 $\Rightarrow U^T S U = \Lambda$

④ Auto-Encoder —→ Encoder-Decoder Architecture



Generalization

$$\begin{array}{c}
 x \in \mathbb{R}^d \xrightarrow{\text{encoder}} z_m \in \mathbb{R}^m \xrightarrow{\text{decoder}} x_m \in \mathbb{R}^d \\
 | \qquad \qquad \qquad T_{\text{enc}}: \mathbb{R}^d \rightarrow \mathbb{R}^m \qquad \qquad T_{\text{dec}}: \mathbb{R}^m \rightarrow \mathbb{R}^d
 \end{array}$$

↓
q-dim ↓
q-dim

optimization

$$\min_{\theta, \phi} \frac{1}{N} \sum_{i=1}^N \|x_i - T_{dec}^\phi(T_{enc}(x_i))\|_2^2$$

$$\begin{cases} T_{enc}^\theta(x) = A \cdot g(Wx + b) & A \in \mathbb{R}^{m \times q} \\ T_{dec}^\phi(x) = B \cdot g(Vx + c) & B \in \mathbb{R}^{d \times q} \quad V \in \mathbb{R}^{q \times m} \end{cases}$$

⑤ K-means Clustering

$$\mathcal{D} = \{x_i\}_{i=1}^N \quad x_i \in \mathbb{R}^d$$

{ similar within groups
different between groups

Aim: Partition into K groups (K clusters)

(assignment)

$$\rightarrow \begin{cases} 1. \text{ Partition } \{a_{ij}\}_{i=1, j=1}^{N, K} & a_{ij} = \begin{cases} 1, & x_i \rightarrow \text{cluster } j \\ 0, & \text{o/w} \end{cases} \\ 2. \text{ Centroid } \{z_1, \dots, z_K\} \end{cases}$$

1. Formulation:

$$\min_{a_{ij}, z_j} \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^K a_{ij} \|x_i - z_j\|_2^2$$

{ $A \in \mathbb{R}^{N \times K}$ → assignment matrix

{ $Z = (z_1, \dots, z_K) \in \mathbb{R}^{d \times K}$ → centroid matrix

→ NP-Hard → intractable

2. Trade-off → greedy algorithm

iterative algorithm

→ a) suppose we know Z

then update the assignment matrix A

$$\text{by } a_{ij} = \begin{cases} 1, & j = \underset{j'}{\operatorname{argmin}} \|x_i - z_{j'}\|_2^2 \\ 0, & \text{o/w} \end{cases}$$

→ re-assignment

→ b) suppose we know A

then we update the centroid matrix Z

$$\text{by } z_k = \frac{\sum_{i=1}^n a_{ik} x_i}{\sum_{i=1}^n a_{ik}}$$

→ re-center

3. Remarks:

- ① This algorithm guarantees to converge (in finite steps)
- ② This algorithm cannot guarantee to converge to global optima
- ③ $J_{k+1} \leq J_k$
- ④ K-means → hard-assignment $A \leftrightarrow a_{ij}$

$$a_{ij} = \begin{cases} 1 & x_i \rightarrow \text{class } j \\ 0 & \text{o/w} \end{cases}$$

Soft-assignment

Centroid

GMM

6. GMM → Gaussian Mixture Model



Density Estimation Model, more appropriate for the cases that

Data Lying Roughly In-between 2 classes (difficult cases)

① Model: (Gaussian Mixture)

$$X \sim \text{GMM}(\pi_k, z_k, \Sigma_k)$$

$$\rightarrow P_x(x) = \sum_{k=1}^K \pi_k P_G(x; z_k, \Sigma_k)$$

$$= \sum_{k=1}^K \pi_k \cdot (2\pi)^{-\frac{d}{2}} |\det \Sigma_k|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x - z_k)^\top \Sigma_k^{-1} (x - z_k) \right\}$$



② How to make use of this model?

Assume we already have the GMM density $P_x(x)$

then our interest is: $p(r_k=1 | x)$ for $k=1, 2, \dots, K$

Calculation

$$\begin{aligned} p(r_k=1 | x) &= \frac{p(x, r_k=1)}{p(x)} \\ &= \frac{\pi_k P_G(x; z_k, \Sigma_k)}{\sum_{l=1}^K \pi_l P_G(x; z_l, \Sigma_l)} \end{aligned}$$

when we have a data point $x \in \mathbb{R}^d$

we can calculate $P(\Gamma_1=1|X) \dots P(\Gamma_K=1|X)$

① Soft-assignment

② consider the component of X

$\begin{cases} 20\% \text{ 娱乐} \\ 30\% \text{ 音乐} \\ 50\% \text{ 学习} \end{cases}$

the 'responsibility' of
component 1 for explaining
the outcome of X

may have some interesting interpretation

③ How to learn the parameters?

← MLE

$$\max_{\Theta} \log \prod_{i=1}^N P_{GMM}(x_i; \pi_k, z_k, \Sigma_k) \quad \Theta = \{\{\pi_k\}_{k=1}^K, \{z_k\}_{k=1}^K, \{\Sigma_k\}_{k=1}^K\}$$

$$= \max_{\Theta} \sum_{i=1}^N \log \left\{ \sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j) \right\} := \varphi(\Theta)$$

a) Brute-Force (Optimization)

$$\textcircled{1} \rightarrow \frac{\partial \varphi}{\partial z_p} = \sum_{i=1}^N \frac{\pi_p \frac{\partial P_G}{\partial z_p}(x_i; z_p, \Sigma_p)}{\sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j)}$$

$$= \sum_{i=1}^N \frac{\pi_p \cdot (-\Sigma_p^{-1}(x_i - z_p)) \cdot P_G(x_i; z_p, \Sigma_p)}{\sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j)}$$

$$= (-\Sigma_p^{-1}) \sum_{i=1}^N \alpha_{ip} (x_i - z_p)$$

$$\textcircled{2} \rightarrow \frac{\partial \varphi}{\partial \Sigma_p} = \sum_{i=1}^N \frac{\pi_p \frac{\partial P_G}{\partial \Sigma_p}(x_i; z_p, \Sigma_p)}{\sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j)}$$

$$P_G(x_i; \mu_p, \Sigma_p) = (2\pi)^{-\frac{d}{2}} |\det(\Sigma_p)|^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_p)^T \Sigma_p^{-1} (x_i - \mu_p) \right\}$$

$$= (2\pi)^{-\frac{d}{2}} \det(\Sigma_e)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_e)^\top \Sigma_e^{-1} (x_i - \mu_e) \right\}$$

$$\begin{aligned} \frac{\partial P_G}{\partial \Sigma} (x_i; \mu_e, \Sigma_e) &= (2\pi)^{-\frac{d}{2}} \left[-\frac{1}{2} \det(\Sigma_e)^{-\frac{3}{2}} \frac{\partial \det(\Sigma_e)}{\partial \Sigma_e} \exp \left\{ -\frac{1}{2} (x_i - \mu_e)^\top \Sigma_e^{-1} (x_i - \mu_e) \right\} \right. \\ &\quad \left. + \det(\Sigma_e)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (x_i - \mu_e)^\top \Sigma_e^{-1} (x_i - \mu_e) \right\} \frac{\partial \left(-\frac{1}{2} (x_i - \mu_e)^\top \Sigma_e^{-1} (x_i - \mu_e) \right)}{\partial \Sigma_e^{-1}} \cdot \frac{\partial \Sigma_e^{-1}}{\partial \Sigma_e} \right] \end{aligned}$$

$$\boxed{\frac{d \det(A)}{d A} = \det(A) \cdot (A^{-1})^\top} = (2\pi)^{-\frac{d}{2}} \left[-\frac{1}{2} \det(\Sigma_e)^{-\frac{3}{2}} \det(\Sigma_e) (\Sigma_e^{-1}) \exp(*) \right. \\ \left. + \det(\Sigma_e)^{-\frac{1}{2}} \exp(*) \left(-\frac{1}{2} (x_i - \mu_e) (x_i - \mu_e)^\top \right) (-\Sigma_e^{-2}) \right]$$

$$= (2\pi)^{-\frac{d}{2}} \left(-\frac{1}{2} \right) \underline{\det(\Sigma_e)^{-\frac{1}{2}}} \underline{\exp(*)} \left[\Sigma_e^{-1} - (x_i - \mu_e) (x_i - \mu_e)^\top \Sigma_e^{-2} \right] \\ = -\frac{1}{2} P_G(x_i; z_e, \Sigma_e) \left[\Sigma_e^{-1} - (x_i - z_e) (x_i - z_e)^\top \Sigma_e^{-2} \right]$$

$$\frac{\partial \varphi}{\partial \Sigma_e} = \sum_{f=1}^N \frac{\pi_f \frac{\partial P_G}{\partial \Sigma_e} (x_i; z_f, \Sigma_e)}{\sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j)}$$

$$= -\frac{1}{2} \sum_{f=1}^N \alpha_{if} \left[\Sigma_e^{-1} - (x_i - z_e) (x_i - z_e)^\top \Sigma_e^{-2} \right]$$

$$\textcircled{3} \rightarrow \mathcal{L}(\theta; \lambda) = \varphi(\theta) + \lambda \left(1 - \sum_{f=1}^K \pi_f \right)$$

$$\Rightarrow \frac{\partial \varphi}{\partial \pi_f} = \sum_{i=1}^N \frac{P_G(x_i; z_f, \Sigma_e)}{\sum_{j=1}^K \pi_j P_G(x_i; z_j, \Sigma_j)} - \lambda \\ = \sum_{i=1}^N \alpha_{if} \cdot \frac{1}{\pi_f} - \lambda$$

To conclude:

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial z_e} = 0 \Rightarrow (-\Sigma_e^{-1}) \sum_{f=1}^N \alpha_{if} (x_i - z_e) = 0 \\ \frac{\partial \mathcal{L}}{\partial \Sigma_e} = 0 \Rightarrow \sum_{i=1}^N \alpha_{if} \left[\Sigma_e^{-1} - (x_i - z_e) (x_i - z_e)^\top \Sigma_e^{-2} \right] = 0 \\ \frac{\partial \varphi}{\partial \pi_f} = 0 \Rightarrow \sum_{i=1}^N \alpha_{if} \frac{1}{\pi_f} - \lambda = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} Z_\ell = \frac{\sum_{i=1}^N a_{i\ell} x_i}{\sum_{i=1}^N a_{i\ell}} \\ \Sigma_\ell = \frac{1}{\sum_{i=1}^N a_{i\ell}} \sum_{i=1}^N (x_i - Z_\ell)(x_i - Z_\ell)^\top \\ \pi_\ell = \frac{\sum_{i=1}^N a_{i\ell}}{N} \end{array} \right.$$

b) EM Framework

Recap: Objective is $\rightarrow \max_{\theta} \ell_N(x|\theta)$

$$\begin{aligned} &= \sum_{i=1}^N \log \left\{ \sum_{j=1}^K \pi_j P_G(x_i | \mu_j, \Sigma_j) \right\} \\ &= \sum_{i=1}^N \log \mathbb{E}_{A_i \sim \pi} [L(x_i, A_i | \theta)] \end{aligned}$$

$$= \log \mathbb{E}_{A \sim \pi} [L(x, A | \theta)]$$

EM Framework

$$\Rightarrow \theta^{(t+1)} = \operatorname{argmax}_{\theta} \mathbb{E}_{A \sim p_C(x, \theta^{(t)})} [\log L(x, A | \theta)]$$

$$\rightarrow \log L(x, A | \theta) = \sum_{i=1}^N \log L(x_i, A_i | \theta)$$

$$= \sum_{i=1}^N \sum_{j=1}^K a_{ij} \log P_G(x_i | \mu_j, \Sigma_j)$$

$$\rightarrow \mathbb{E}_{A \sim p_C(x, \theta^{(t)})} [\log L(x, A | \theta)] \xrightarrow{\operatorname{argmax}} \theta^{(t+1)} \rightarrow \underline{\text{M-step}}$$

$$= \sum_{k=1}^K \sum_{j=1}^K \mathbb{E}_{A_{ij} \sim P_C | X_i, \theta^{(e)}} [a_{ij}] P_G(x_i | \mu_j, \Sigma_j)$$

$$\rightarrow \mathbb{E}_{A_{ij} \sim P_C | X_i, \theta^{(e)}} [a_{ij}]$$

$$= P(X_i \in \text{cluster } j | \theta^{(e)}, X_i)$$

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

E-step

11.27

GMM Summary

1. Maximize Log-likelihood:

$$f(\theta) = \sum_{i=1}^n \log \{ f(x_i | \theta) \}$$

Rmk: Actually it is a common way since this can be just viewed as 'parameter estimation'

↳ [e.g.] estimate $N(\mu, \sigma^2)$ from samples

$$\theta := \{(\pi_k, \mu_k, \Sigma_k)\}_{k=1}^K$$

$$= \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \pi_k P_g(x_i; \mu_k, \Sigma_k) \right\}$$

Issue: $f(\theta)$ is highly non-convex w.r.t (π_k, μ_k, Σ_k)



Conventional GD may break down (but still can try)

2. 3 Solutions

- a) GD
- b) Necessary Condition $\nabla_{\theta} f(\hat{\theta}) = 0$
- c) EM Framework

"Augmented Lagrangian Multiplier" → method for constrained optimization problem

2a) ALM method etc.

$$\text{Recap: } \ell(\theta) = \sum_{F=1}^n \log \left\{ \sum_{k=1}^K \pi_k p_g(x_i; \mu_k, \Sigma_k) \right\}$$

From Calculation, we have:

$$\begin{cases} \frac{\partial \ell}{\partial \mu_p} = (-\Sigma_p^{-1}) \sum_{F=1}^n \gamma_{ip} (x_i - \mu_p) \\ \frac{\partial \ell}{\partial \Sigma_p} = -\frac{1}{2} \sum_{F=1}^n \gamma_{ip} [\Sigma_p^{-1} - (x_i - \mu_p)(x_i - \mu_p)^T \Sigma_p^{-2}] \\ \frac{\partial \ell}{\partial \pi_p} = \sum_{F=1}^n \gamma_{ip} \frac{1}{\pi_p} - \lambda \end{cases}$$

Therefore, Toughly speaking, we can apply GD Framework:

$$\mu_p \leftarrow \mu_p - \eta \frac{\partial \ell}{\partial \mu_p}$$

$$\Sigma_p \leftarrow \Sigma_p - \eta \frac{\partial \ell}{\partial \Sigma_p}$$

there are Constraint on (π_1, \dots, π_K)

$$\Rightarrow \boxed{\sum_{F=1}^K \pi_i = 1}$$

2b) Solve $\nabla_{\theta} \ell(\hat{\theta}) = 0$ \longrightarrow Non-linear Equation

$$\begin{cases} \frac{\partial \ell}{\partial \mu_p} = (-\Sigma_p^{-1}) \sum_{F=1}^n \gamma_{ip} (x_i - \mu_p) = 0 \\ \frac{\partial \ell}{\partial \Sigma_p} = -\frac{1}{2} \sum_{F=1}^n \gamma_{ip} [\Sigma_p^{-1} - (x_i - \mu_p)(x_i - \mu_p)^T \Sigma_p^{-2}] = 0 \\ \frac{\partial \ell}{\partial \pi_p} = \sum_{F=1}^n \gamma_{ip} \frac{1}{\pi_p} - \lambda = 0 \Leftrightarrow \pi_p = \frac{1}{n} \sum_{i=1}^n \gamma_{ip} \end{cases}$$

Issue: No closed-form solution

use iterative method to solve this non-linear system:

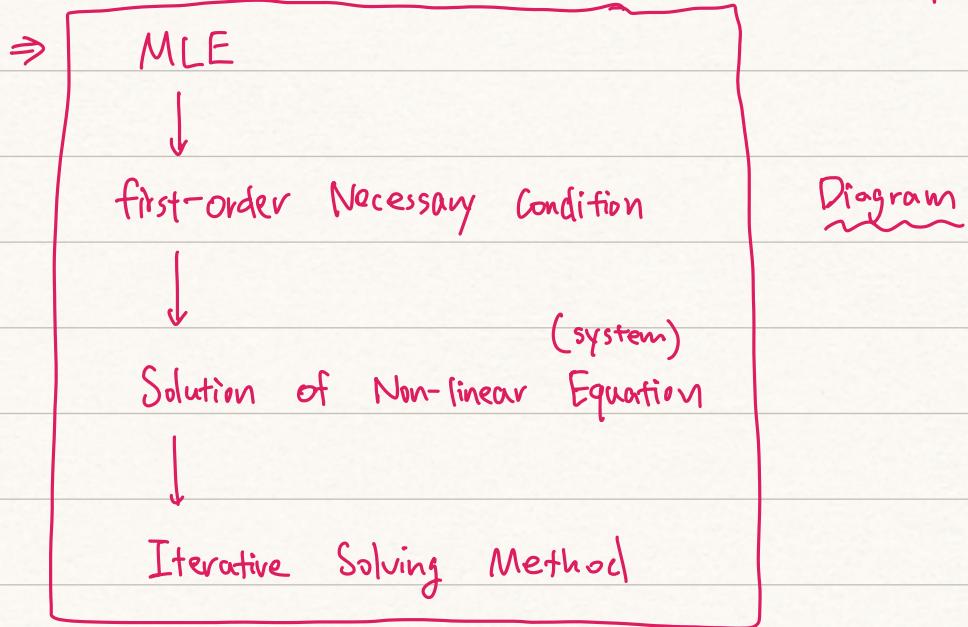
$$\begin{cases} \mu_p = \frac{1}{\sum_{F=1}^n \gamma_{ip}} \sum_{F=1}^n \gamma_{ip} x_i \\ \Sigma_p = \frac{1}{\sum_{F=1}^n \gamma_{ip}} \sum_{F=1}^n \gamma_{ip} (x_i - \mu_p)(x_i - \mu_p)^T \\ \pi_p = \frac{1}{n} \sum_{F=1}^n \gamma_{ip} \end{cases}$$

where $\gamma_{ip} := \frac{\pi_p p_g(x_i; \mu_p, \Sigma_p)}{\sum_{k=1}^K \pi_k p_g(x_i; \mu_k, \Sigma_k)}$

$$\left\{ \begin{array}{l} \mu_p^{(t+1)} = \frac{1}{\sum_{i=1}^n \gamma_{ip}^{(t)}} \sum_{i=1}^n \gamma_{ip}^{(t)} x_i \\ \Sigma_p^{(t+1)} = \frac{1}{\sum_{i=1}^n \gamma_{ip}^{(t)}} \sum_{i=1}^n \gamma_{ip}^{(t)} (x_i - \mu_p^{(t)}) (x_i - \mu_p^{(t)})^\top \\ \pi_p^{(t+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ip}^{(t)} \end{array} \right.$$

Rmk: This actually corresponds to the EM Framework Update.

→ Therefore, this can be viewed as another perspective of EM:



2c) EM Framework → the most successful one

Recap: $\ell(\theta) = \sum_{i=1}^n \log \left\{ \sum_{k=1}^K \pi_k p_g(x_i; \mu_k, \Sigma_k) \right\}$

likelihood of complete observation

$$L(x_i, A_i = p | \theta)$$

$$= p_g(x_i | \mu_p, \Sigma_p)$$

$$= \sum_{i=1}^n \log \{ \mathbb{E}_{A_i \sim \theta} [L(x_i | A_i, \theta)] \}$$

$$= \log \{ \mathbb{E}_{A \sim \theta} [L(X | A, \theta)] \} \rightarrow \text{re-write}$$

A: hidden variable
↓
"assignment"

difficult to optimize directly

EM Framework

Big idea: instead of optimizing log-likelihood directly, we

- 1) Optimize on "another objective" (easy to deal with)
- 2) the solution of "another objective" is still better with respect to the Log-likelihood (discuss later)

primal objective

primal objective

Realization :

- 1) "another objective" is:

$$\mathbb{E}_{A \sim \theta^{(t)}, X} [\log L(X, A | \theta)]$$

Recap: the "primal objective" is: $\log \{ \mathbb{E}_{A \sim \theta} [L(X, A | \theta)] \}$

A: hidden variable
 ↓
"assignment"

in EM Framework, $\theta^{(t+1)} = \operatorname{argmax} \mathbb{E}_{A \sim \theta^{(t)}, X} [\log L(X, A | \theta)]$

- 2) Recap:

we actually want $\hat{\theta}_{MLE} = \operatorname{argmax}_{\theta \in \Theta} f(\theta)$

From the Property of EM Algorithm,

we can show that $\underline{f(\theta^{(t+1)})} \geq \underline{f(\theta^{(t)})}$

* interpretation: the solution from solving "another objective"

is actually better in the sense of "primal objective"

$$f(\theta) = \log \{ \mathbb{E}_{A \sim \theta} [L(X | A, \theta)] \}$$

Q1: What is the solution of $\theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_{A \sim \underline{\theta^{(t)}, X}} [\log L(X, A | \theta)]$



Answer:

$$\underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_{A \sim \underline{\theta^{(t)}, X}} [\log L(X, A | \theta)]$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_{A \sim \underline{\theta^{(t)}, X}} \left[\sum_{i=1}^n \log L(X_i, A_i | \theta) \right]$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} \mathbb{E}_{A \sim \underline{\theta^{(t)}, X}} \left[\sum_{i=1}^n \sum_{k=1}^K a_{ik} \log \{ \pi_k p_g(x_i; \mu_k, \Sigma_k) \} \right]$$

$$\left[\text{here, } a_{ik} := \begin{cases} 1, & x_i \text{ belongs to cluster } k \\ 0, & \text{otherwise} \end{cases} \right]$$

$$= \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \mathbb{E}_{A_i \sim \underline{\theta^{(t)}, X_i}} [a_{ik}] \log \{ \pi_k p_g(x_i; \mu_k, \Sigma_k) \}$$

$$\mathbb{E}_{A_i \sim \underline{\theta^{(t)}, X_i}} [a_{ik}] := \underline{\delta_{ik}^{(t)}} \quad (\text{for the agreement})$$

$$= \mathbb{E}_{A_i \sim \underline{\theta^{(t)}}} [\mathbb{1} \{ x_i \text{ belongs to cluster } k \} \mid X_i = x_i]$$

$$= P_{\theta^{(t)}} (X_i \text{ belongs to cluster } k \mid X_i = x_i)$$

$$= \frac{\pi_k p_g(x_i; \mu_k^{(t)}, \Sigma_k^{(t)})}{\sum_{j=1}^K \pi_j p_g(x_i; \mu_j^{(t)}, \Sigma_j^{(t)})}$$

Q2: what is $\theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \underline{\delta_{ik}^{(t)}} \log \{ \pi_k p_g(x_i; \mu_k, \Sigma_k) \}$?

$$\text{Answer : } \theta^{(t+1)} = \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(t)} \log \left\{ \pi_k (2\pi)^{-\frac{d}{2}} (\det(\Sigma_k))^{-\frac{1}{2}} \right. \\ \left. \exp \left\{ -\frac{1}{2} (x_i - \mu_k)^\top \Sigma_k^{-1} (x_i - \mu_k) \right\} \right\}$$

$$\propto \underset{\theta \in \Theta}{\operatorname{argmax}} \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(t)} \left[\log \pi_k - \frac{1}{2} \log \{ \det(\Sigma_k) \} \right. \\ \left. - \frac{1}{2} (x_i - \mu_k)^\top \Sigma_k^{-1} (x_i - \mu_k) \right]$$

$$:= \underset{\theta \in \Theta}{\operatorname{argmax}} J(\theta)$$

$$\boxed{① \frac{\partial J}{\partial \mu_p} = \sum_{i=1}^n \gamma_{ip}^{(t)} \Sigma^{-1} (x_i - \mu_p)}$$

$$\boxed{② \frac{\partial J}{\partial \Sigma_p} = \sum_{i=1}^n \gamma_{ip}^{(t)} \cdot \left[-\frac{1}{2} \Sigma_p^{-1} - \frac{1}{2} (x_i - \mu_p) (x_i - \mu_p)^\top \cdot (-\Sigma_p^{-2}) \right]} \\ = \sum_{i=1}^n \gamma_{ip}^{(t)} \left[-\frac{1}{2} \Sigma_p^{-1} + \frac{1}{2} (x_i - \mu_p) (x_i - \mu_p)^\top \Sigma_p^{-2} \right]$$

$$\boxed{③ \frac{\partial L}{\partial \pi_p} = \sum_{i=1}^n \gamma_{ip}^{(t)} \cdot \frac{1}{\pi_p} - \lambda}$$

$$\Rightarrow \begin{cases} \frac{\partial J}{\partial \mu_p} = 0 \Rightarrow \mu_p^{(t+1)} = \frac{1}{\sum_{i=1}^n \gamma_{ip}^{(t)}} \sum_{i=1}^n \gamma_{ip}^{(t)} \cdot x_i \\ \frac{\partial J}{\partial \Sigma_p} = 0 \Rightarrow \Sigma_p^{(t+1)} = \frac{1}{\sum_{i=1}^n \gamma_{ip}^{(t)}} \sum_{i=1}^n \gamma_{ip}^{(t)} (x_i - \mu_p^{(t+1)}) (x_i - \mu_p^{(t+1)})^\top \\ \frac{\partial L}{\partial \pi_p} = 0 \Rightarrow \pi_p^{(t+1)} = \frac{\sum_{i=1}^n \gamma_{ip}^{(t)}}{n} \end{cases}$$

$$\underline{\theta^{(t+1)} := \{ \pi_p^{(t+1)}, \mu_p^{(t+1)}, \Sigma_p^{(t+1)} \}_{p=1}^K}$$

#

To conclude, although 2b) and 2c) share the similar formula,
the idea is totally different!

2b) → lecture note

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$$

$(\hat{\theta} \neq \hat{\theta}_{MLE} \text{ may happen})$ ↓ 这一步很粗糙 (necessary condition)

$\hat{\theta}$ may not be
the exact solution
since this is only necessary
condition!

$$\nabla_{\theta} \ell(\hat{\theta}) = 0$$

solution of Non-linear
Equation System

2c) → EM

$$\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$$

↓ difficult to solve directly

$$\text{consider } \theta^{(ea)} = \underset{\theta \in \Theta}{\operatorname{argmax}} J(\theta, \theta^{(e)})$$

terminology: $J(\theta, \theta^{(e)})$ is auxiliary function
of log-likelihood $\ell(\theta)$

⇒ $\theta^{(ea)}$ is exact solution since $J(\theta, \theta^{(e)})$
can be good enough w.r.t θ
($\ell(\theta)$ can be very bad!)

Iterative method for
solving Equation System

↓
indirect method since it will

not solve $\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$ directly

direct method in my personal perspective

* Instead, it tries to find one $\theta^{(L)}$

such that it can maximize $\ell(\theta)$

as much as possible!

directly solve $\hat{\theta}_{MLE} = \underset{\theta \in \Theta}{\operatorname{argmax}} \ell(\theta)$

12.5 More EM Framework

1. Observation

log-likelihood → $\log p(X|\theta)$

$$= \log p(X, Z|\theta) - \log p(Z|X, \theta)$$

$$= \sum_z q(z) \log \frac{p(x, z | \theta)}{q(z)} - \sum_z q(z) \cdot \log \frac{p(z | x, \theta)}{q(z)}$$

$$= \sum_z q(z) \log \frac{p(x, z | \theta)}{q(z)} + \sum_z q(z) \log \frac{q(z)}{p(z | x, \theta)}$$

$$= \underbrace{\mathcal{L}(q, \theta)}_{\geq \mathcal{L}(q, \theta)} + \underbrace{\text{KL}(q || p_z)}$$

$$\geq \mathcal{L}(q, \theta)$$

$$p_z := p(z | x, \theta)$$

To conclude :

$$\begin{aligned} l(\theta) &= \log p(x | \theta) \\ &\geq \sum_z q(z) \log \frac{p(x, z | \theta)}{q(z)} \end{aligned}$$

$$:= \mathcal{L}(q, \theta)$$

$$\text{when } \underbrace{q = p_z = p(z | x, \theta)}_{\text{then}} \text{, then } \underbrace{l(\theta) = \mathcal{L}(q, \theta)}$$

2. Consider the following thing :

$$\textcircled{1} \quad \underset{q}{\operatorname{argmax}} \mathcal{L}(q, \theta) = \underset{q}{\operatorname{argmin}} \text{KL}(q || p_z) = p_z$$

$$\textcircled{2} \quad \underset{\theta}{\operatorname{argmax}} \mathcal{L}(q, \theta) = \underset{\theta}{\operatorname{argmax}} \sum_z q(z) \log \frac{p(x, z | \theta)}{q(z)}$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_z q(z) \log p(x, z | \theta)$$

$$= \underset{\theta}{\operatorname{argmax}} \mathbb{E}_{z \sim q} [\log p(x, z | \theta)]$$

3. 2 Diagrams for illustration:

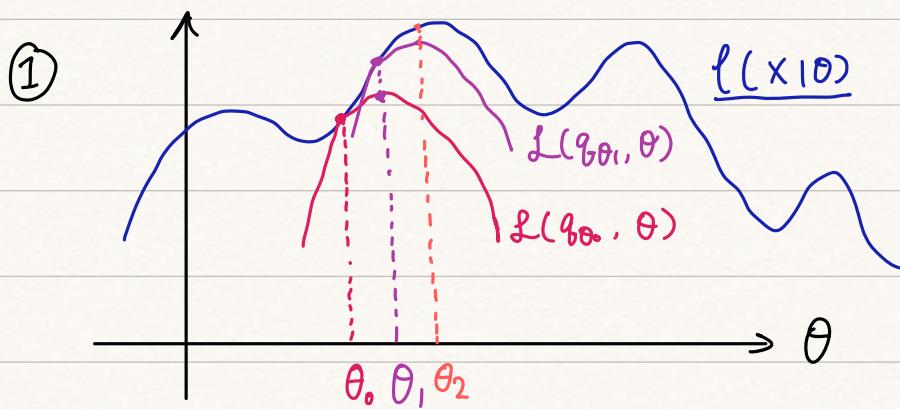
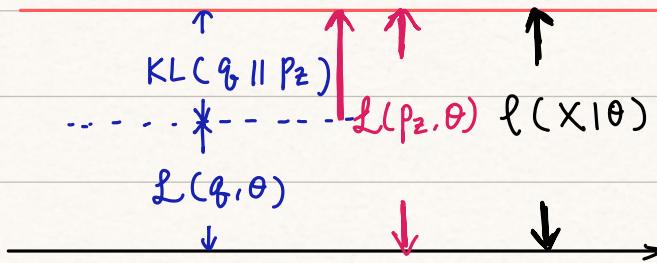


Diagram 1

$$\underline{q_{\theta_0} := P(\cdot | X, \theta_0)} \quad \underline{q_{\theta_1} := P(\cdot | X, \theta_1)}$$

② $\underset{\theta}{\operatorname{argmax}} \quad \underline{\ell(q_\theta, \theta)} = p_z$



$$\underset{\theta}{\operatorname{argmax}} \quad \underline{\ell(q_\theta, \theta)} = \theta^{(\text{new})}$$

