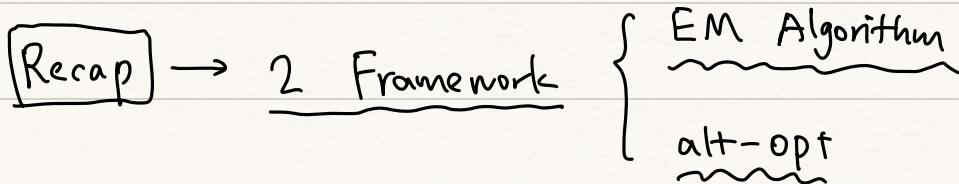
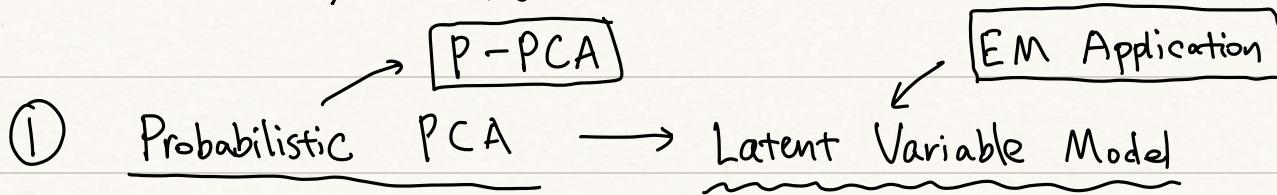


PCA Summary (Family)



Aim: $\hat{\theta}_{MLE} = \underset{\theta}{\operatorname{argmax}} \log p(x|\theta) = \underset{\theta}{\operatorname{argmax}} \log \sum_z p(z|\theta) p(x|z, \theta)$

ILL: Incomplete Log-Likelihood

→ ALT-OPT

1. initialize $\theta = \hat{\theta}$

2. estimate $\hat{z} = \underset{z}{\operatorname{argmax}} \log p(z|x, \hat{\theta})$

3. estimate $\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \log p(x, \hat{z} | \theta)$

slight modified

CLL: Complete Log-Likelihood

→ EM: 1. initialize $\theta = \hat{\theta}$

2. compute posterior distribution $p(z|x, \hat{\theta})$

3. estimate θ by Expected CLL

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

1. Probabilistic PCA → Generative Model



Latent Variable Model for Dimensionality Reduction

Motivation

consider

$$x \in \mathbb{R}^D$$

imagine that x is generated by

scalar form $\leftarrow x \approx \sum_{k=1}^K z_k w_k \quad w_k \in \mathbb{R}^P \quad k = 1, 2, \dots, K$

vector form $\leftarrow \boxed{x \approx Wz}$ Notation

$$\begin{cases} W = (w_1, \dots, w_K) \in \mathbb{R}^{D \times K} \\ z = \begin{pmatrix} z_1 \\ \vdots \\ z_K \end{pmatrix} \in \mathbb{R}^K \end{cases}$$

[Rmk]

Here, z is the low-dimension representation

Question: How to model " \approx "?

- { Sol¹: minimize Re-construct Error \rightarrow PCA
- Sol²: $x|z \sim N(Wz, \sigma^2 I_D) \rightarrow P\text{-PCA}$ (PCA included)

Note: Reconstruct Error = $\|x - zW^\top\|_F^2 \rightsquigarrow$ directly from $x \approx Wz$

2. P-PCA Model

$$\textcircled{1} \quad x|z \sim N(Wz, \sigma^2 I_p)$$

$$\textcircled{2} \quad z \sim N(0, I_K)$$

Q: what is $x \sim ?$

Answer: $p(x) \propto \int_z \exp \left\{ -\frac{1}{2\sigma^2} (x - Wz)^\top (x - Wz) \right\} \exp \left\{ -\frac{1}{2} z^\top z \right\} dz$

$$= \exp \left\{ -\frac{1}{2\sigma^2} x^\top x \right\} \int_z \exp \left\{ -\frac{1}{2\sigma^2} z^\top W^\top W z - \frac{1}{2} z^\top z + \frac{1}{\sigma^2} x^\top W z \right\} dz$$

$$\begin{aligned}
 &= \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{x}^\top \mathbf{x} \right\} \int_{\mathbb{R}^D} \exp \left\{ -\frac{1}{2\sigma^2} \mathbf{z}^\top (\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K) \mathbf{z} + \frac{1}{\sigma^2} \mathbf{z}^\top \mathbf{W}^\top \mathbf{x} \right\} d\mathbf{z} \\
 &\propto \exp \left\{ -\frac{1}{2} \mathbf{x}^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)^{-1} \mathbf{x} \right\}
 \end{aligned}$$

trick 1: $\underbrace{-\frac{1}{2}(\mathbf{z}-\mu)^\top \Sigma^{-1}(\mathbf{z}-\mu)}_{\text{important term}} = -\frac{1}{2} \mathbf{z}^\top \Sigma^{-1} \mathbf{z} + \mathbf{z}^\top \Sigma^{-1} \mu$ $-\frac{1}{2} \mu^\top \Sigma^{-1} \mu$
 trick 2: $\mathbf{W}(\mathbf{W}^\top \mathbf{W} + \sigma^2 \mathbf{I}_K)^{-1} \mathbf{W}^\top$
 $= \mathbf{W}\mathbf{W}^\top (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)^{-1}$
 $= (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D - \sigma^2 \mathbf{I}_D) (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)^{-1}$
 $= \mathbf{I}_D - \sigma^2 (\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)^{-1}$

That is: $\mathbf{x} \sim N(\mathbf{0}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)$ \rightarrow Model Necessary Condition

from assumption

Suppose $\mathbf{x} \sim N(\mathbf{0}, \Sigma)$

then $\mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D \approx \Sigma$

few parameters Many parameters ($\frac{P(P+1)}{2}$)

$(D \cdot K + 1)$ $K \ll D$



Rmk: fewer parameters, less likely to over-fitting!!!

To conclude: in this model

- ① $\mathbf{x} | \mathbf{z} \sim N(\mathbf{W}\mathbf{z}, \sigma^2 \mathbf{I}_D)$
- ② $\mathbf{z} \sim N(\mathbf{0}, \mathbf{I}_K)$
- ③ $\Rightarrow \mathbf{x} \sim N(\mathbf{0}, \mathbf{W}\mathbf{W}^\top + \sigma^2 \mathbf{I}_D)$

3. P-PCA \longrightarrow Learning parameters

$$\|X - ZW^T\|_F^2$$

interpretation for
matrix factorization

a)

Connection with Re-construction Error

(ALT-OPT)

only focus $\log P(X|Z, W, \sigma^2)$

$$\propto -\frac{1}{2\sigma^2} \sum_{i=1}^N \|x_i - Zw_i\|_2^2 - \frac{N\sigma^2}{2} \log(2\pi\sigma^2)$$

trick is that:

$$z \sim N(0, I_k) \quad = -\frac{1}{2\sigma^2} \|X - ZW^T\|_F^2 - \frac{N\sigma^2}{2} \log(\sigma^2)$$

Suppose σ^2 is fixed! (for simplicity)

ALT-OPT Framework is exactly close to minimize
Re-construct Error!

\Rightarrow ① initialize $W = \hat{W}$

② estimate $\hat{Z} = \arg \max_Z P(Z|X, \hat{W}, \sigma^2)$

$$= \arg \max_Z \log \frac{P(X|\hat{Z}, \hat{W}, \sigma^2) P(Z)}{P(X|\hat{W}, \sigma^2)}$$

$$= \arg \max_Z \log P(X|\hat{Z}, \hat{W}, \sigma^2)$$

$$= \arg \min_Z \|X - Z\hat{W}^T\|_F^2$$

③ estimate $\hat{W} = \arg \max_W P(X, \hat{Z}|W, \sigma^2)$

Recap: $Z \sim N(0, I_k)$

$$= \arg \max_W \log P(X|\hat{Z}, W, \sigma^2)$$

$$= \underset{W}{\operatorname{argmin}} \| X - \hat{Z} W^T \|_F^2$$

★

Remark: step ② & ③ are exactly Multi-output Regression

$$\begin{array}{ccc} \boxed{\text{input}} & & \boxed{\text{output}} \\ z \in \mathbb{R}^K & \xrightarrow{W} & \tilde{x} = Wz \in \mathbb{R}^P \quad W \in \mathbb{R}^{D \times K} \end{array}$$

Loss function: $\frac{1}{2} \| \tilde{x} - x \|_2^2$ (L2-loss function)

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} \quad Z = \begin{pmatrix} z_1 \\ \vdots \\ z_N \end{pmatrix} \Rightarrow \boxed{N\text{-observation!}}$$

★

→ All the above discussion is not related to MLE!

b)

MLE & EM Framework to learn P-PCA parameters

① (Brute-Force MLE)

Recap: $X \sim N(0, WW^T + \sigma^2 I_D)$



it is possible to find the log-likelihood

$$\begin{aligned} \Rightarrow \ell(W, \sigma) &= \log P(X | W, \sigma) \\ &= \sum_{i=1}^N \log P(x_i | W, \sigma) \end{aligned}$$

$$\Rightarrow \underline{x_i \sim N(0, WW^T + \sigma^2 I_D)}$$

⇒ exact MLE estimator:

$$\left\{ \begin{array}{l} \hat{W}_{MLE} = U_K (L_K - \hat{\sigma}_{MLE}^2 I)^{\frac{1}{2}} R \\ \hat{\sigma}_{MLE}^2 = \frac{1}{D-K} \sum_{k=K+1}^D \lambda_k \end{array} \right.$$

where

$$\left\{ \begin{array}{l} U_K = (U_1, \dots, U_K) \quad \underbrace{U_i \in \mathbb{R}^D}_{\text{the eigenvector of}} \\ L_K = \text{Diag}(\lambda_1, \dots, \lambda_K) \\ R \rightarrow \text{Rotate Matrix} \end{array} \right.$$

Rmk: ① This solution \hat{W}_{MLE} \rightarrow PCA

$$\text{when } R = I \quad \& \quad \hat{\sigma}_{MLE}^2 \rightarrow 0$$

② It is computationally expensive to achieve
the eigenvalue decomposition for $S = X^T X \in \mathbb{R}^{D \times D}$

② EM Framework

Recap: EM algo

E-step: estimate $p(z_n | x_n, \hat{W}, \hat{\sigma}^2)$



$$M = \hat{W}^T \hat{W} + \hat{\sigma}^2 I_K$$

Conclusion: $z_n | x_n, \hat{W}, \hat{\sigma}^2 \sim N(M^{-1} \hat{W}^T x_n, \hat{\sigma}^2 M^{-1})$

Recap: (from the calculate of $P(x)$)

$$P(z|x) \propto \exp \left\{ -\frac{1}{2\hat{\sigma}^2} z^T (W^T W + \hat{\sigma}^2 I_K) z + \frac{1}{\hat{\sigma}^2} z^T W^T x \right\}$$

$$= \exp \left\{ -\frac{1}{2} z^T \left[\underbrace{\hat{\sigma}^{-2} (W^T W + \hat{\sigma}^2 I_K)}_{\Sigma^{-1}} \right] z + z^T \underbrace{\hat{\sigma}^{-2} W^T x}_{\Sigma^{-1} \mu} \right\}$$

(easy calculation)

$$\text{M-step} : (\hat{W}, \hat{\sigma}^2) = \underset{W, \sigma^2}{\operatorname{argmax}} \mathbb{E}_{Z \sim p_{\cdot \cdot \cdot}(x, z | W, \sigma^2)} [\log P(x, z | W, \sigma^2)]$$

$$\text{Calculation} : \log P(x, z | W, \sigma^2)$$

$$= \sum_{n=1}^N \{ \log P(x_n | z_n, W, \sigma^2) + \log P(z_n) \}$$

$$\propto - \sum_{n=1}^N \left[\frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|x_n\|_2^2 - \frac{1}{\sigma^2} z_n^T W^T x_n \right]$$

omit $\mathbb{E}_{Z \sim p(\cdot | x, \theta)}$ for simplicity

$$+ \frac{1}{2\sigma^2} \operatorname{tr}(z_n z_n^T W^T W) + \frac{1}{2} \operatorname{tr}(z_n z_n^T)$$

$$\Rightarrow \mathbb{E} [\log P(x, z | W, \sigma^2)]$$

$$= - \sum_{n=1}^N \left[\frac{D}{2} \log \sigma^2 + \frac{1}{2\sigma^2} \|x_n\|_2^2 - \frac{1}{\sigma^2} \mathbb{E}[z_n]^T W^T x_n \right]$$

$$+ \frac{1}{2\sigma^2} \operatorname{tr}(\mathbb{E}[z_n z_n^T] W^T W) + \frac{1}{2} \operatorname{tr}(\mathbb{E}[z_n z_n^T]) \]$$

then, since $(\hat{W}, \hat{\sigma}^2) = \underset{W, \sigma^2}{\operatorname{argmax}} \mathbb{E}[\log P(x, z | W, \sigma^2)]$

\Rightarrow (closed-form solution) (Matrix Calculus)

$$\text{then } \begin{cases} \hat{W} = \left[\sum_{n=1}^N x_n \mathbb{E}[z_n]^T \right] \left[\sum_{n=1}^N \mathbb{E}[z_n z_n^T] \right]^{-1} \\ \hat{\sigma}^2 = \frac{1}{ND} \sum_{n=1}^N \{ \|x_n\|_2^2 - 2 \mathbb{E}[z_n]^T \hat{W}^T x_n + \operatorname{tr}(\mathbb{E}[z_n z_n^T] \hat{W}^T \hat{W}) \} \end{cases}$$

$$M = W^T W + \sigma^2 I_K$$

obviously, since $Z_n | x_n, W, \sigma^2 \sim N(M^{-1} W^T x_n, \sigma^2 M^{-1})$

$$\text{then } \begin{cases} \mathbb{E}[z_n] = M^{-1} W^T x_n \\ \mathbb{E}[z_n z_n^T] = \sigma^2 M^{-1} + \mathbb{E}[z_n] \mathbb{E}[z_n]^T \end{cases}$$

Rmk: since when $\hat{\sigma}^2 \rightarrow 0$, MLE \rightarrow PCA

therefore, when we set $\hat{\sigma}^2 = 0$, the EM framework

can be viewed as the efficient algorithm to solve PCA

$\boxed{\text{PCA}} \rightarrow \text{eigenvector decomposition}$

$\boxed{\text{P-PCA}} \rightarrow \text{iterative algorithm to solve eigenvector decomp. issue}$

② Mixture P-PCA Model

1. Motivation

Recap: P-PCA Model



$$\left\{ \begin{array}{l} X_n | Z_n, W, \hat{\sigma}^2 \sim N(WZ_n, \hat{\sigma}^2 I_D) \\ Z_n \sim N(0, I_k) \end{array} \right.$$



Assumption !!!

center at origin



$$X_n \sim N(0, WW^T + \hat{\sigma}^2 I_D)$$

Guideline

⇒ Before doing P-PCA, we should center the data!

→ Another interpretation for centering data

P-PCA with μ
P-PCA without μ

PCA without μ

P-PCA Model (Previous Equivalent form)

$$X_n = WZ_n + \varepsilon_n \quad \leftarrow \left\{ \begin{array}{l} \varepsilon_n \sim N(0, \hat{\sigma}^2 I_D) \\ Z_n \sim N(0, I_k) \end{array} \right.$$

⇒ Center at origin

Modification: $X_n = \mu + W Z_n + \varepsilon_n$

$$\Rightarrow X_n \sim N(\mu, WW^T + \sigma^2 I_D)$$

↓
assumption: data has one cluster

$$\hat{\mu}_{MLE} = \frac{1}{N} \sum_{n=1}^N X_n$$

Interpretation:

Therefore, if our data has one cluster, but not center at the origin, then we can first center the data and apply P-PCA without μ . This is actually equivalent to apply P-PCA with μ .

2. Mixture of P-PCA Model

- ① $C_n \sim \text{multi-noulli}(\pi_1, \dots, \pi_M)$
- ② $Z_n \sim N(0, I_k)$
- ③ $X_n | C_n, Z_n, \Theta \sim N(\mu_{C_n} + W_{C_n} Z_n, \sigma_{C_n}^2 I_D)$

Some induction:

- 1. Parameters: $\{\pi_m, \mu_m, W_m, \sigma_m^2\}_{m=1}^M$
- 2. $X_n | C_n \sim N(\mu_{C_n}, W_{C_n} W_{C_n}^T + \sigma_{C_n}^2 I_D)$
- 3. $X_n \sim \text{GMM (Low-rank GMM)}$

Limitation of P-PCA

3. Advantages & Remarks

1. If the data has several clusters, then we cannot use P-PCA since the assumption is violated!

But Mixture of P-PCA still works!

2. Non-linear Dimensionality Reduction

3. Special Case of Gaussian Mixture Model [GMM]
(Low-rank form)

4. combining { learn the clusters
 learn the dimensionality reduction

③ Standard PCA

- 1. Maximize Variance
- 2. Minimize Re-construction Error
- (3). Changing Basis to de-correlated features

1. Maximize Variance (Projection)

$z_i = u_1^\top x_i \Rightarrow$ project towards u_1 direction

and z_i is the corresponding coordinate

Maximize variance

$$\Leftrightarrow \max_{u_1} \frac{1}{N} \sum_{n=1}^N (z_n - \bar{z})^2$$

$$\Leftrightarrow \max_{u_1} \frac{1}{N} \sum_{n=1}^N u_1^\top (x_n - \bar{x})(x_n - \bar{x})^\top u_1$$

* Center the data (pre-step)

$$\Leftrightarrow \max_{U_1} U_1^T S U_1 \quad S := \frac{1}{N} \sum_{n=1}^N \tilde{x}_n \tilde{x}_n^T$$

$$= \frac{1}{N} X^T X \quad X = \underbrace{\begin{pmatrix} x_1^T \\ \vdots \\ x_N^T \end{pmatrix}}$$

→ first-component \hat{u}_1

→ second-component $\hat{u}_2 = \begin{cases} \arg \max_{U_2} U_2^T S U_2 \\ \text{s.t. } \|U_2\| = 1 \\ U_2^T \hat{u}_1 = 0 \end{cases}$

→ etc...

then $Z = X U_K$ $\begin{cases} U_K = (U_1, \dots, U_K) \in \mathbb{R}^{D \times K} \\ Z \in \mathbb{R}^{N \times K} \Rightarrow \text{embedding matrix} \end{cases}$

principal component score
(terminology)

2. Minimize Re-construction Error (DSAS105)

Trick: ① $\tilde{x}_n = \sum_{d=1}^K (x_n^T u_d) \cdot u_d$

$$= \sum_{d=1}^K u_d \cdot (u_d^T x_n)$$

$$= \sum_{d=1}^K (u_d \cdot u_d^T) x_n$$

$$\textcircled{2} \quad \sum_{n=1}^N \|x_n - \tilde{x}_n\|_2^2$$

$$= \sum_{n=1}^N \|x_n - U_k \cdot z_n\|_2^2$$

$$= \|X - Z \cdot U_k^\top\|_2^2$$

$$= \|X - X U_k U_k^\top\|_2^2$$

$$\textcircled{3} \quad \sum_{n=1}^N \|x_n - \tilde{x}_n\|_2^2$$

$$= \sum_{n=1}^N \|x_n - \sum_{d=1}^k (U_d U_d^\top) x_n\|_2^2$$

$$= \sum_{n=1}^N \|x_n\|_2^2 - 2 \sum_{n=1}^N x_n^\top U x_n + \sum_{n=1}^N x_n^\top U^\top U x_n$$

★

$$\begin{aligned} U^\top U &= \sum_{i=1}^k \sum_{j=1}^k (U_i U_i^\top) (U_j U_j^\top) \\ &= \sum_{i=1}^k \sum_{j=1}^k U_i (U_i^\top U_j) U_j^\top \\ &= \sum_{i=1}^k U_i (U_i^\top U_i) U_i^\top \\ &= \sum_{i=1}^k U_i U_i^\top = U \end{aligned}$$

$$= \sum_{n=1}^N \|x_n\|_2^2 - \sum_{n=1}^N x_n^\top U x_n$$

$$= C - \sum_{n=1}^N \sum_{d=1}^k x_n^\top U_d U_d^\top x_n$$

$$= C - \sum_{d=1}^k \sum_{n=1}^N U_d^\top X_n X_n^\top U_d$$

$$= C - \sum_{d=1}^k U_d^\top S U_d \quad \underbrace{S = X^\top X}_{}$$