

since the solution of A is $[u_1, \dots, u_k]$. where u_i is eigen-vectors.

$$\sum_{i=1}^n \text{trace} (A A^T \tilde{x}_i \tilde{x}_i^T) = \sum_{i=1}^n \text{trace} (A^T \tilde{x}_i \tilde{x}_i^T A)$$

$$= \sum_{i=1}^n \text{trace} \left(\begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} \tilde{x}_i \tilde{x}_i^T [u_1, \dots, u_k] \right)$$

$$= \cancel{n \cdot \text{trace}} \quad n \cdot \frac{1}{n} \sum_{i=1}^n \text{trace} \left(\begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} \tilde{x}_i \tilde{x}_i^T [u_1, \dots, u_k] \right) = n \cdot \text{trace} \left(\begin{bmatrix} u_1 \\ \vdots \\ u_k \end{bmatrix} U \Lambda U^T \right)$$

$$= n \cdot \sum_{j=1}^k \lambda_j$$

$$\text{Thus, } J^* = n \sum_{j=1}^d \lambda_j - n \sum_{j=1}^k \lambda_j = n \sum_{j=k+1}^d \lambda_j = \min_{\vec{\mu}, A, \{\vec{\theta}_i\}} \text{Obj. fun}$$



3) EM Algorithm for Mixed Linear Regression

$$f(y|\vec{x}; \vec{\theta}) \sim \sum_{k=1}^K \varepsilon_k \phi(\vec{y}; \vec{w}_k^T \vec{x} + b_k, \sigma_k^2)$$

$$\text{where } \vec{\theta} = (\varepsilon_1, \dots, \varepsilon_K, \vec{w}_1, \dots, \vec{w}_K, b_1, \dots, b_K, \sigma_1^2, \dots, \sigma_K^2)$$

a. $\vec{x} = (\vec{x}_1, \dots, \vec{x}_n)$ $y = (y_1, \dots, y_n)$

$$\begin{aligned} l(\vec{\theta}; y|\vec{x}) &= \log f(y|\vec{x}; \vec{\theta}) \\ &= \log \prod_{i=1}^n f(y_i|\vec{x}_i; \vec{\theta}) \\ &= \sum_{i=1}^n \log f(y_i|\vec{x}_i; \vec{\theta}) \\ &= \sum_{i=1}^n \log \sum_{k=1}^K \varepsilon_k \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2) \end{aligned}$$

Introduce a hidden variable $S = (S_1, \dots, S_n)$ for each data $\vec{x} = (\vec{x}_1, \dots, \vec{x}_n)$
 Then complete data is $\vec{z} = (\vec{x}, S)$ S_i : describe the component responsible for generating \vec{x}_i
 complete log-likelihood:

$$\begin{aligned} l(\vec{\theta}; y, S|\vec{x}) &= \log L(\vec{\theta}; y, S|\vec{x}) \\ &= \log \prod_{i=1}^n P(y_i, S_i|\vec{x}_i; \vec{\theta}) \\ &= \log \prod_{i=1}^n P(S_i = s_i) \cdot f(y_i|\vec{x}_i, S_i = s_i; \vec{\theta}) \\ &= \sum_{i=1}^n \log P(S_i = s_i; \vec{\theta}) \cdot f(y_i|\vec{x}_i, S_i = s_i; \vec{\theta}) \\ &= \sum_{i=1}^n \log \sum_{k=1}^K \varepsilon_k \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2) \end{aligned}$$

Define $\Delta_{ik} = \begin{cases} 1 & S_i = k \\ 0 & \text{otherwise} \end{cases}$ Then

$$\begin{aligned} l(\vec{\theta}; y|\vec{z}) &= \sum_{i=1}^n \log \left(\sum_{k=1}^K \Delta_{ik} \cdot \varepsilon_k \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2) \right) \\ &= \sum_{i=1}^n \sum_{k=1}^K \Delta_{ik} \log (\varepsilon_k \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2)) \end{aligned}$$



E-step:

$$\begin{aligned} Q(\theta, \theta^{(j)}) &= E_{\sum_i y_i \vec{x}_i} [l(\vec{\theta}; y | \vec{x}, s) | y, \vec{x}; \theta^{(j)}] \\ &= E_{\sum_{i=1}^n \sum_{k=1}^K \Delta_{ik}} [\log \xi_k + \log \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2)] \\ &= \sum_{i=1}^n \sum_{k=1}^K [E_{\sum} [\Delta_{ik}] \log \xi_k \phi(y_i; \vec{w}_k^T \vec{x}_i + b_k, \sigma_k^2)] \end{aligned}$$

Suppose $\gamma_{ik}^{(j)} = E_{\sum} [\Delta_{ik}] = E[\Delta_{ik} | y_i, \vec{x}_i; \theta^{(j)}]$

$$= P(S_i = k | x_i, y_i; \theta^{(j)})$$

Bayes Rule $\rightarrow = \frac{P(S_i = k; \theta^{(j)}) \cdot f(y_i, \vec{x}_i | S_i = k; \theta^{(j)})}{f(y_i, \vec{x}_i; \theta^{(j)})}$

$$= \frac{\xi_k^{(j)} \phi(y_i; \vec{w}_k^{(j)T} \vec{x}_i + b_k^{(j)}, \sigma_k^{(j)2})}{\sum_{l=1}^K \xi_l^{(j)} \phi(y_i; \vec{w}_l^{(j)T} \vec{x}_i + b_l^{(j)}, \sigma_l^{(j)2})}$$

b. M-step:

$$\theta^{(j+1)} = \arg \max_{\theta} Q(\theta, \theta^{(j)})$$

where $Q(\theta, \theta^{(j)}) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(j)} [\log \xi_k + \log ((2\pi)^{-\frac{d}{2}}) / \sigma_k \cdot e^{-\frac{1}{2\sigma_k^2} \|y_i - \vec{w}_k^T \vec{x}_i - b_k\|_{\lambda}^2}]$

$$= \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(j)} [\log \xi_k - \frac{d}{2} \log(2\pi) - \frac{1}{2\sigma_k^2} \|y_i - \vec{w}_k^T \vec{x}_i - b_k\|_{\lambda}^2 - \frac{1}{2} \log \sigma_k^2]$$

First optimize ξ_k : it has constraints $\sum_{k=1}^K \xi_k = 1$. by Lagrange multiplier

theory: $L(\xi_1, \dots, \xi_K) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(j)} \log \xi_k + \lambda (\sum_{k=1}^K \xi_k - 1)$

$$\frac{\partial L}{\partial \xi_k} = \sum_{i=1}^n \gamma_{ik}^{(j)} / \xi_k + \lambda \stackrel{!}{=} 0.$$

$$\Rightarrow \xi_k^* = - \frac{\sum_{i=1}^n \gamma_{ik}^{(j)}}{\lambda}, \text{ since } \sum_{k=1}^K \xi_k^* = 1.$$

$$\Rightarrow - \frac{\sum_{k=1}^K \sum_{i=1}^n \gamma_{ik}^{(j)}}{\lambda} = 1 \Rightarrow -\lambda = \sum_{i=1}^n 1 = n$$

Thus, $\xi_k^* = \frac{1}{n} \sum_{i=1}^n \gamma_{ik}^{(j)}$ i.e. $\xi_k^{(j+1)} = \frac{1}{n} \sum_{i=1}^n \gamma_{ik}^{(j)}$



Second optimize (\vec{w}_k, b_k) , Suppose σ_k^2 is fixed.

$$\text{Then } \max_{(\vec{w}_k, b_k)} Q(\sigma, \sigma^{(j)}) \Leftrightarrow \min_{(\vec{w}_k, b_k)} \frac{1}{2} \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(j)} \|y_i - \vec{w}_k^T \vec{x}_i - b_k\|^2$$

It has the same form as weighted least squared regression.

Thus, the solution:

$$\begin{bmatrix} b_k^{(j+1)} \\ \vec{w}_k^{(j+1)} \end{bmatrix} = (X^T C_k^{(j)} X)^{-1} X^T C_k^{(j)} \vec{y}$$

$$\text{Where } X = \begin{pmatrix} | & \vec{x}_1^T \\ \vdots & \vdots \\ | & \vec{x}_n^T \end{pmatrix}, \quad C_k^{(j)} = \begin{bmatrix} \gamma_{1k}^{(j)} & & 0 \\ & \ddots & \\ 0 & & \gamma_{nk}^{(j)} \end{bmatrix}$$

Last optimize σ_k^2 , plug in results above. use Lagrange multiplier theory

$$L(\sigma_1^2, \dots, \sigma_K^2) = \sum_{i=1}^n \sum_{k=1}^K \gamma_{ik}^{(j)} \left(-\frac{1}{2} \log \sigma_k^2 - \frac{1}{2} \|y_i - \vec{w}_k^T \vec{x}_i - b_k\|^2 / \sigma_k^2 \right)$$

$$\frac{\partial L}{\partial \sigma_k^2} = \sum_{i=1}^n \gamma_{ik}^{(j)} \left(-\frac{1}{2\sigma_k^2} + \frac{\|y_i - \vec{w}_k^T \vec{x}_i - b_k\|^2}{2(\sigma_k^2)^2} \right) \stackrel{!}{=} 0$$

(here we treat σ_k^2 as a variable)

$$\Rightarrow \sigma_k^{(j+1)2} = \frac{\sum_{i=1}^n \gamma_{ik}^{(j)} \|y_i - \vec{w}_k^{(j+1)} \vec{x}_i - b_k^{(j+1)}\|^2}{\sum_{i=1}^n \gamma_{ik}^{(j)}}$$

4) Ncut and Normalized Spectral Clustering

$$K=2. \quad \text{Ncut}(A, \bar{A}) = \frac{1}{2} \left(\frac{C(A, \bar{A})}{\text{Vol}(A)} + \frac{C(A, \bar{A})}{\text{Vol}(\bar{A})} \right) = \frac{1}{2} C(A, \bar{A}) \left(\frac{1}{\text{Vol}(A)} + \frac{1}{\text{Vol}(\bar{A})} \right)$$

$$\text{where } \text{Vol}(A) = \sum_{i \in A} \sum_{j \in V} W_{ij}$$

Given $A \subseteq \{1, 2, \dots, n\}$ define $\vec{f}_A = (f_{A1}, \dots, f_{An})^T \in \mathbb{R}^n$ by

$$f_{Ai} = \begin{cases} +\sqrt{\frac{\text{Vol}(\bar{A})}{\text{Vol}(A)}} & \text{if } i \in A \\ -\sqrt{\frac{\text{Vol}(A)}{\text{Vol}(\bar{A})}} & \text{if } i \notin A. \end{cases}$$



$$\text{Compute: } \vec{f}_A^T \perp \vec{f}_A = \frac{1}{2} \sum_{i,j=1}^n W_{ij} (f_{A_i} - f_{A_j})^2$$

$$= \frac{1}{2} \sum_{i \in A, j \in \bar{A}} W_{ij} \left(\sqrt{\frac{\text{vol}(\bar{A})}{\text{vol}(A)}} + \sqrt{\frac{\text{vol}(A)}{\text{vol}(\bar{A})}} \right)^2$$

$$+ \frac{1}{2} \sum_{i \in \bar{A}, j \in A} W_{ij} \left(-\sqrt{\frac{\text{vol}(\bar{A})}{\text{vol}(A)}} - \sqrt{\frac{\text{vol}(A)}{\text{vol}(\bar{A})}} \right)^2$$

$$= \sum_{i \in A, j \in \bar{A}} W_{ij} \left(\sqrt{\frac{\text{vol}(\bar{A})}{\text{vol}(A)}} + \sqrt{\frac{\text{vol}(A)}{\text{vol}(\bar{A})}} \right)^2$$

$$= \sum_{i \in A, j \in \bar{A}} W_{ij} \left(\frac{\text{vol}(\bar{A})}{\text{vol}(A)} + \frac{\text{vol}(A)}{\text{vol}(\bar{A})} + 2 \right)$$

$$= (\text{vol}(\bar{A}) + \text{vol}(A)) \left(\frac{\sum_{i \in A, j \in \bar{A}} W_{ij}}{\text{vol}(A)} + \frac{\sum_{i \in \bar{A}, j \in A} W_{ij}}{\text{vol}(\bar{A})} \right)$$

$$= (\text{vol}(\bar{A}) + \text{vol}(A)) \left(\frac{C(A, \bar{A})}{\text{vol}(A)} + \frac{C(A, \bar{A})}{\text{vol}(\bar{A})} \right)$$

$$= 2(\text{vol}(\bar{A}) + \text{vol}(A)) \text{Ncut}(A, \bar{A})$$

$$\vec{1}^T D \vec{f}_A = \sum_{i=1}^n d_i f_{A_i} = \sum_{i \in A} d_i \cdot \sqrt{\frac{\text{vol}(\bar{A})}{\text{vol}(A)}} - \sum_{j \in \bar{A}} d_j \sqrt{\frac{\text{vol}(A)}{\text{vol}(\bar{A})}} \quad \text{since } d_i = \sum_{j=1}^n W_{ij}$$

$$= \text{vol}(A) \sqrt{\frac{\text{vol}(\bar{A})}{\text{vol}(A)}} - \text{vol}(\bar{A}) \sqrt{\frac{\text{vol}(A)}{\text{vol}(\bar{A})}} = 0.$$

$$\vec{f}_A^T D \vec{f}_A = \sum_{i=1}^n d_i f_{A_i}^2 = \sum_{i \in A} \frac{d_i \text{vol}(\bar{A})}{\text{vol}(A)} + \sum_{j \in \bar{A}} \frac{d_j \text{vol}(A)}{\text{vol}(\bar{A})}$$

$$= \text{vol}(\bar{A}) + \text{vol}(A)$$

$$\text{Then we claim: } \vec{f}_A^T \perp \vec{f}_A = 2 \vec{f}_A^T D \vec{f}_A \cdot \text{Ncut}(A, \bar{A})$$

There, Ncut can be written as the following optimization problem:

$$\min_{A \in \{1, \dots, n\}} \vec{f}_A^T \perp \vec{f}_A$$

$$\text{s.t. } \vec{1}^T D \vec{f}_A = 0$$

$$\vec{f}_A^T D \vec{f}_A = \text{vol}(A) + \text{vol}(\bar{A})$$

$$= \text{vol}(V)$$

The relaxation:

$$\min_{\vec{f} \in \mathbb{R}^n} \vec{f}^T \perp \vec{f}$$

$$\text{s.t. } \vec{1}^T D \vec{f} = 0$$

$$\vec{f}^T D \vec{f} = \text{vol}(V)$$



Suppose $g = D^{\frac{1}{2}} f$. Then

$$\vec{f}_A^T D \vec{f}_A = (D^{\frac{1}{2}} f)^T (D D^{\frac{1}{2}} f) = g^T g.$$

The optimization problem is:

$$\min_{g \in \mathbb{R}^n} g^T D^{-\frac{1}{2}} L D^{-\frac{1}{2}} g \quad \text{s.t.} \quad D^{\frac{1}{2}} g = 0, \quad g^T g = \text{vol}(V)$$

$$\text{Define: } L_g = D^{-\frac{1}{2}} L D^{-\frac{1}{2}} = D^{-\frac{1}{2}} (D - W) D^{-\frac{1}{2}} = I - D^{-\frac{1}{2}} W D^{-\frac{1}{2}}$$

$$\text{also, } \tilde{L} = D^T L = I - D^T W$$

The optimization problem is:

$$\min_{g \in \mathbb{R}^n} g^T L_g g \quad \text{s.t.} \quad D^{\frac{1}{2}} g = 0, \quad g^T g = \text{vol}(V)$$

Notice L_g is symmetric $n \times n$ matrix, $D^{\frac{1}{2}} \mathbf{1}$ is the first eigenvector $\text{vol}(V)$ is constant, Apply Rayleigh-Ritz Theorem.

The solution of opt problem is the second eigenvector of L_g .

re-substitute $f = D^{-\frac{1}{2}} g$.

suppose u_2 is eigenvector for L_g . Then $L_g u_2 = \lambda_2 u_2$

$$u_2' := D^{-\frac{1}{2}} u_2. \quad \text{Then } \tilde{L} u_2' = D^T L u_2' = D^T L D^{-\frac{1}{2}} u_2 = D^{-\frac{1}{2}} L_g u_2 \\ = D^{-\frac{1}{2}} \lambda_2 u_2 = \lambda_2 u_2'$$

Thus, u_2' is eigenvector of \tilde{L} u_2' corresponds to u_2

Since $f = D^{-\frac{1}{2}} g$, it shows that finding second eigenvector

for L_g is equivalent to finding second eigenvector for \tilde{L}

And the second eigenvector of \tilde{L} is solution for f



b. The condition: $\vec{x}_i \in \mathbb{R}^d$ $\vec{\theta}_i \in \mathbb{R}^d$. Then $\lambda^k \neq \lambda^{k+1}$ is the condition

$$\text{i.e. } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda^k > \lambda^{k+1} \geq \lambda^{k+2} \dots \geq \lambda^d$$

Here λ_i is eigenvalues for sample covariance matrix: $S = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)(x_i - \mu)^T$. And spectral decomposition of S is $S = U \Lambda U^T$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_d)$, $U = [\vec{u}_1, \dots, \vec{u}_d]$, λ_i is eigenvalue, \vec{u}_i is eigenvector.

Suppose $\lambda_k = \lambda_{k+1}$. The corresponding eigenvectors don't have to be the same $\vec{u}_k \neq \vec{u}_{k+1}$. Then, if we choose A to be a $d \times k$ matrix. we can see that $A = [\vec{u}_1, \dots, \vec{u}_k]$ and $A = [\vec{u}_1, \dots, \vec{u}_{k-1}, \vec{u}_{k+1}]$ are both solutions for $\min \sum_{i=1}^n \|\vec{x}_i - \mu - A\vec{\theta}_i\|^2$, however, the subspace $\langle A \rangle$ are different (not unique). vice versa. if A is not unique, only \vec{u}_k is possible to be replaced. it corresponds that $\lambda_k = \lambda_{k+1}$. Thus, our condition is necessary and sufficient



1) PCA

a. obj. fun = $\sum_{i=1}^n \|\tilde{x}_i - \vec{\mu} - A\vec{\theta}_i\|^2 =: J$

Since it's a convex function of $\vec{\theta}_i \in \mathbb{R}^k$.

we compute $\frac{\partial J}{\partial \vec{\theta}_i} = -2A^T(\tilde{x}_i - \vec{\mu} - A\vec{\theta}_i) \stackrel{!}{=} 0$.

$\Rightarrow \vec{\theta}_i = A^T(\tilde{x}_i - \vec{\mu})$. ($\because A^T A = I_{k \times k}$) is optimal $\vec{\theta}_i$

for min J

$\Rightarrow J = \sum_{i=1}^n \|\tilde{x}_i - \vec{\mu} - AA^T(\tilde{x}_i - \vec{\mu})\|^2$, also, J is a convex function of $\vec{\mu} \in \mathbb{R}^d$

$\frac{\partial J}{\partial \vec{\mu}} = -\sum_{i=1}^n 2(\tilde{x}_i - \vec{\mu} - A\vec{\theta}_i) \stackrel{!}{=} 0 \Rightarrow \sum_{i=1}^n (I - AA^T)(\tilde{x}_i - \vec{\mu}) \stackrel{!}{=} 0$.

$\Rightarrow \vec{\mu} = \frac{1}{n} \sum_{i=1}^n \tilde{x}_i = \bar{x}$ is optimal $\vec{\mu}$.

Suppose $\hat{x}_i = \tilde{x}_i - \bar{x}$, then

$$J^* = \sum_{i=1}^n \|\hat{x}_i - AA^T \hat{x}_i\|^2 = \sum_{i=1}^n \text{trace}[(\hat{x}_i - AA^T \hat{x}_i)(\hat{x}_i - AA^T \hat{x}_i)^T]$$

$$= \sum_{i=1}^n \text{trace}[\hat{x}_i \hat{x}_i^T - \hat{x}_i \hat{x}_i^T AA^T - AA^T \hat{x}_i \hat{x}_i^T + AA^T \hat{x}_i \hat{x}_i^T AA^T]$$

$$\downarrow = \sum_{i=1}^n \text{trace}(\hat{x}_i \hat{x}_i^T) - 2 \sum_{i=1}^n \text{trace}(AA^T \hat{x}_i \hat{x}_i^T) + \sum_{i=1}^n \text{trace}(AA^T \hat{x}_i \hat{x}_i^T)$$

(because of the property of trace: $\text{tr}(CA) = \text{tr}(CA^T)$, $\text{tr}(CAB) = \text{tr}(CAB)$.)

$\Rightarrow J^* = \sum_{i=1}^n \text{trace}(\hat{x}_i \hat{x}_i^T) - 2 \sum_{i=1}^n \text{trace}(AA^T \hat{x}_i \hat{x}_i^T)$

since $S = \frac{1}{n} \sum_{i=1}^n (\tilde{x}_i - \bar{x})(\tilde{x}_i - \bar{x})^T = \frac{1}{n} \sum_{i=1}^n \hat{x}_i \hat{x}_i^T$.

$\Rightarrow \text{tr}(S) = \frac{1}{n} \sum_{i=1}^n \text{trace}(\hat{x}_i \hat{x}_i^T) = \sum_{j=1}^d \lambda_j$

$\Rightarrow J^* = n \sum_{j=1}^d \lambda_j - 2 \sum_{i=1}^n \text{trace}(AA^T \hat{x}_i \hat{x}_i^T)$

