

1.  $\theta_1, \dots, \theta_k$  eigen vectors of  $k$ -th largest eigen values of  $C = \frac{1}{n} \phi \phi^T$

$\theta_1$  first principal eigen vector of  $C$ ,  $C\theta_1 = \lambda_1 \theta_1$  and  $\tilde{\phi}_i = (I - \theta_1 \theta_1^T) \phi_i$

$$\tilde{\phi} = (I - \theta_1 \theta_1^T) \phi, \quad \theta_c^T \theta_d = \begin{cases} 0 & c \neq d \\ 1 & c = d \end{cases}$$

$$\begin{aligned} a) \quad \tilde{C} &= \frac{1}{n} \tilde{\phi} \tilde{\phi}^T = \frac{1}{n} (I - \theta_1 \theta_1^T) \phi \phi^T (I - \theta_1 \theta_1^T)^T \\ &= (I - \theta_1 \theta_1^T) C (I - \theta_1 \theta_1^T) \quad [(I - \theta_1 \theta_1^T) \text{ is symmetric}] \\ &= C - \theta_1 \theta_1^T C - C \theta_1 \theta_1^T + \theta_1 \theta_1^T C \theta_1 \theta_1^T \end{aligned}$$

Now,  $C\theta_1 = \lambda_1 \theta_1 \Rightarrow \theta_1^T C^T = \lambda_1 \theta_1^T$  replacing in the eq<sup>n</sup>

$$\tilde{C} = C - \lambda_1 \theta_1 \theta_1^T - C \theta_1 \theta_1^T + \lambda_1 \theta_1 \theta_1^T \theta_1 \theta_1^T \quad (\text{now, } \theta_1^T \theta_1 = I)$$

$$\boxed{\tilde{C} = C - \lambda_1 \theta_1 \theta_1^T}$$

$$b) \quad \text{We have, } \tilde{C} \theta_j = (C - \lambda_1 \theta_1 \theta_1^T) \theta_j = \lambda_j \theta_j - \lambda_1 \theta_1 \theta_1^T \theta_j$$

$$\text{Now if } j=1, \quad \tilde{C} \theta_1 = \lambda_1 \theta_1 - \lambda_1 \theta_1 (\theta_1^T \theta_1) = (\lambda_1 - \lambda_1) \theta_1 = 0$$

$$j \neq 1, \quad \tilde{C} \theta_j = \lambda_j \theta_j \quad [\theta_1^T \theta_j = 0, j \neq 1]$$

Since  $\theta_1, \theta_2, \dots, \theta_k$  are the first  $k$  eigen vectors with largest eigen values in  $C$ , therefore  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$ .

Now for  $\tilde{C}$ ,  $\theta_j$  are the eigenvectors with eigen-values  $(0, \lambda_2, \lambda_3, \dots, \lambda_k)$

Therefore  $\theta_2$  is the eigen vector with the largest eigen value of  $\tilde{C}$ .

c) `findKEigenVectors(C, K, f, lambda-list, eigen-list)`  
 if `lambda-list.length = K`, return;  $\leftarrow$  base case recursive  
`lambda, u = f(C)`  
`C = C - lambda * u * u.transpose`  
`lambda-list.append(lambda)`  
`eigen-list.append(u)`  
`findKEigenVectors(C, K, f, lambda-list, eigen-list)`

## 2. Initialize step:

Initialize  $\mu_c$ ,  $\pi_c$  and  $\Sigma$  matrix.

$\mu_c$  = mean feature vector for class  $c$      $\pi_c$  = Probability of  $y$  taking cluster  $c$   
 $\Sigma$  = Covariance matrix for all classes.

Repeat:

### ① Expectation step:

The superscript  $(m)$  denotes the variable in  $m$ -th iteration.

We calculate the expectations,  $\psi_{ic}^{(m)}$  based on  $\mu_c^{(m-1)}$ ,  $\pi_c^{(m-1)}$  and  $\Sigma^{(m-1)}$

$$\psi_{ic} = E[y_i | x_i] = \frac{\pi_c^{(m-1)} \mathcal{N}(x_i | \mu_c^{(m-1)}, \Sigma^{(m-1)})}{\sum_{c=1}^K \pi_c^{(m-1)} \mathcal{N}(x_i | \mu_c^{(m-1)}, \Sigma^{(m-1)})}$$

### ② Maximization step:

$$\mathcal{L} = E[\log P(x, y | \pi, \mu, \Sigma)] = \sum_{i=1}^n \sum_{c=1}^K \psi_{ic} \{ \log \pi_c + \log \mathcal{N}(x_i | \mu_c, \Sigma) \}$$

Taking the derivative w.r.t  $\pi_c$ ,  $\mu_c$  and  $\Sigma^{-1}$  and setting those to 0, we can get the  $\pi_c^{(m)}$ ,  $\mu_c^{(m)}$  and  $\Sigma^{(m)}$

The second term in the  $\mathcal{L}$  equation would be,

$$\log \mathcal{N}(x_i | \mu_c, \Sigma) = -\frac{1}{2} \log(2\pi) - \frac{1}{2} \log |\Sigma| - \frac{1}{2} (x_i - \mu_c)^T \Sigma^{-1} (x_i - \mu_c)$$

a) Taking derivative w.r.t to  $\mu_c$

$$\frac{\partial \mathcal{L}}{\partial \mu_c} = \sum_{i=1}^n \psi_{ic} \frac{\partial (\log \mathcal{N}(x_i | \mu_c, \Sigma))}{\partial \mu_c}$$

$$0 = \sum_{i=1}^n \psi_{ic} \frac{\partial (-\frac{1}{2} (x_i - \mu_c)^T \Sigma^{-1} (x_i - \mu_c))}{\partial \mu_c}$$

$$\sum_{i=1}^n \psi_{ic} \Sigma^{-1} (x_i - \mu_c) = 0$$

$$\Rightarrow \sum_{i=1}^n \psi_{ic} x_i = \sum_{i=1}^n \psi_{ic} \mu_c^{(m)}$$

$$\therefore \mu_c^{(m)} = \frac{\sum_{i=1}^n \psi_{ic} x_i}{\sum_{i=1}^n \psi_{ic}}$$

b) Taking derivative w.r.t  $\Sigma^{-1}$

$$\frac{\partial \mathcal{L}}{\partial \Sigma^{-1}} = \sum_{i=1}^n \sum_{c=1}^K \eta_{ic} \left[ \left( -\frac{1}{2} \frac{\partial}{\partial \Sigma^{-1}} \log |\Sigma| \right) - \frac{1}{2} \frac{\partial}{\partial \Sigma^{-1}} \left( (x_i - \mu_c)^T \Sigma^{-1} (x_i - \mu_c) \right) \right]$$

① Since  $x^T \Sigma^{-1} x$  is scalar, hence can be rewritten as  $\text{tr}[x x^T \Sigma^{-1}]$

and  $\frac{\partial}{\partial \Sigma^{-1}} \text{tr}[x x^T \Sigma^{-1}] = (x x^T)^T = x x^T$ .

②  $\frac{\partial}{\partial \Sigma^{-1}} (-\log |\Sigma|) = \frac{\partial}{\partial \Sigma^{-1}} \log |\Sigma^{-1}| \times (\Sigma^{-1})^T = \Sigma^{-1}$ .

So,  $0 = \sum_{i=1}^n \sum_{c=1}^K \left[ \eta_{ic} \left( \frac{1}{2} \sum_{i=1}^n \frac{1}{2} (x_i - \mu_c)(x_i - \mu_c)^T \right) \right]$

So,  $\underline{\Sigma}^{(m)} = \frac{\sum_{i=1}^n \sum_{c=1}^K \eta_{ic} (x_i - \mu_c)(x_i - \mu_c)^T}{\sum_{i=1}^n \sum_{c=1}^K \eta_{ic}}$

c) There's a constraint in maximizing  $\pi_c$ ,  $\sum_{c=1}^K \pi_c = 1$

Using Lagrange multiplier,  $\mathcal{L}_\pi = \mathcal{L} + \lambda \left( 1 - \sum_{c=1}^K \pi_c \right)$

$\frac{\partial \mathcal{L}_\pi}{\partial \pi_c} = 0 \Rightarrow \sum_{i=1}^n \eta_{ic} / \pi_c^{(m)} - \lambda = 0$

$\pi_c^{(m)} = \frac{1}{\lambda} \sum_{i=1}^n \eta_{ic}$

now  $\sum_{c=1}^K \frac{1}{\lambda} \sum_{i=1}^n \eta_{ic} = 1 \Rightarrow \lambda = \sum_{c=1}^K \sum_{i=1}^n \eta_{ic} = \sum_{i=1}^n 1 = n$ .

$\left[ \sum_{c=1}^K \eta_{ic} = 1 \right]$ , sum of probabilities of a point being assigned to each the clusters is 1

$\therefore \pi_c^{(m)} = \frac{1}{n} \sum_{i=1}^n \eta_{ic}$

With  $\mu_c^{(m)}$ ,  $\Sigma^{(m)}$  and  $\pi_c^{(m)}$  begin next iteration.