1. Dataset $D = \{(x_i, y_i)\}_{i=1}^{N}$     $f(x) = p(y=c|x) \propto p(x|y=c, \theta)\, p(y=c|\theta)$.

$$\log \; p_{\mathcal{M}}(y=c|x; \pi, \mu, \Sigma)$$

$$= \sum_{i=1}^{N} \log \; p_{\mathcal{M}}(y_i=c|x_i; \pi, \mu, \Sigma)$$

$$= \sum_{i=1}^{N} \log \; p_{\mathcal{M}}(y_i; \pi) + \sum_{i=1}^{N} \log \; p_{\mathcal{M}}(x_i|y_i; \mu, \Sigma)$$

$$= \sum_{i=1}^{N} \log \; \pi_{y_i} + \sum_{i=1}^{N} \log \mathcal{N}(x_i; \mu_{y_i}, \Sigma_{y_i})$$

$$\mathcal{LL} = \sum_{i=1}^{N} \sum_{c=1}^{K} \mathbb{I}\{y_i=c\} \log \pi_c + \sum_{i=1}^{N} \sum_{c=1}^{K} \mathbb{I}\{y_i=c\} \log \mathcal{N}(x_i; \mu_c, \Sigma_c)$$

① Derivating in terms of $\pi_c$ and setting it to 0.

$$\mathcal{LL}_{\pi} = \mathcal{LL} + \lambda\left(1 - \sum_{c=1}^{K} \pi_c\right)$$

$\left[\text{the maximization is subject to } \left[\sum_{c=1}^{N} \pi_c = 1\right] \lambda \text{ lagrange multiplier}\right]$

$$\frac{\partial \mathcal{LL}_{\pi}}{\partial \pi_c} = 0.$$

$$\sum_{i=1}^{N} \frac{\mathbb{I}\{y_i=c\}}{\pi_c} - \lambda = 0 \;\Rightarrow\; \pi_c = \frac{1}{\lambda} N_c \quad \left[N_c = \sum_{i=1}^{N} \mathbb{I}\{y_i=c\}\right]$$

Now, $\sum_{c=1}^{K} \pi_c = 1$

$$\sum_{c=1}^{K} \tfrac{1}{\lambda} N_c = 1.$$

$$\boxed{\pi_c = \frac{N_c}{N}}$$

$$N_1 + N_2 + \dots N_K = \lambda \;\Rightarrow\; \lambda = N.$$

⑪ Derivating in terms of $\mu_c$ and setting it to 0.

$$\frac{\partial \mathcal{LL}}{\partial \mu_c} = 0 \;\Rightarrow\; \sum_{i=1}^{N} \mathbb{I}\{y_i=c\} \frac{\partial}{\partial \mu_c}\left(-\tfrac{1}{2}(x_i - \mu_c)^T \Sigma_c (x_i - \mu_c)\right) = 0$$

$$\sum_{i=1}^{N} \mathbb{I}\{y_i=c\} \Sigma_c (x_i - \mu_c) = 0$$

$$\boxed{\mu_c = \frac{\sum_i x_i}{N_c}}_{\;y_i=c}$$

$$\Rightarrow \sum_{\substack{i \\ y_i=c}} \mu_c = \sum_{\substack{i \\ y_i=c}} x_i \;\Rightarrow\; \mu_c \cdot N_c = \sum_{\substack{i \\ y_i=c}} x_i$$

(III) Deriving in terms of $\Sigma_c$ and setting to 0.

$\dfrac{\partial LL}{\partial \Sigma_c} = 0$

$\sum\limits_{y_i = c} \left[ \left( -\frac{1}{2} \dfrac{\partial}{\partial \Sigma_c^{-1}} \log |\Sigma_c| \right) - \frac{1}{2} \dfrac{\partial}{\partial \Sigma_c^{-1}} \left( (x_i - \mu_c)^T \Sigma_c^{-1} (x_i - \mu_c) \right) \right] = 0$

$\dfrac{\partial}{\partial \Sigma_c^{-1}} \cdot \log |\Sigma_c^{-1}| = \cancel{-\dfrac{1}{\Sigma_c}} \Sigma_c$

$= \Sigma_c^T = \Sigma_c$

scalar $x^T A x = tr[x x^T A]$

$\dfrac{\partial}{\partial A} tr[x x^T A] = x x^T.$

$\sum\limits_{y_i = c} \left[ \frac{1}{2} \Sigma_c - \frac{1}{2} (x_i - \mu_c)(x_i - \mu_c)^T \right] = 0$

$\Rightarrow \quad N_c \Sigma_c = \sum\limits_{y_i = c} (x_i - \mu_c)(x_i - \mu_c)^T \quad \Rightarrow \quad \boxed{\Sigma_c = \dfrac{\sum\limits_{y_i = c} (x_i - \mu_c)(x_i - \mu_c)^T}{N_c}}$

b) $\hat{y}_0 = f(x_0^*) = \underset{y_i}{\arg\max} \; P(y_i)\, P(x_0 | y_i) \qquad \forall \, i = 1 \dots K.$

Let $y_i = c$,

$= \underset{c}{\arg\max} \; P(c)\, P(x_0 | c)$

$= \underset{c}{\arg\max} \; \pi_c \, \mathcal{N}_{x_0}(\mu_c, \Sigma_c)$

$= \underset{c}{\arg\max} \left[ \log \pi_c - \frac{1}{2} (x_0 - \mu_c)^T \Sigma_c^{-1} (x_0 - \mu_c) - \frac{1}{2} \log \Sigma_c \right]$

[taking log and removing constants]

$= \underset{c}{\arg\max} \left[ \log \pi_c + x_0^T \Sigma_c^{-1} \mu_c - \frac{1}{2} \left( x_0^T \Sigma_c^{-1} x_0 + \mu_c^T \Sigma_c^{-1} \mu_c \right) + \frac{1}{2} \log |\Sigma_c| \right]$

This is QDA.

$= \underset{c}{\arg\max} \; x_0^T \Sigma_c^{-1} \mu_c - \frac{1}{2} x_0^T \Sigma_c^{-1} x_0 + b_c$

$\left[ b_c = \log \pi_c - \frac{1}{2} \mu_c^T \Sigma_c^{-1} \mu_c - \frac{1}{2} \log |\Sigma_c| \right]$

②

2. $\min\limits_{\theta_1, \dots, \theta_K} \frac{1}{m} \sum\limits_{i=1}^{n} \left( \log \sum\limits_{c=1}^{K} \exp(\theta_c^T \phi_i) - \theta_{y_i}^T \phi_i \right) + \frac{\lambda}{2} \sum\limits_{c=1}^{K} \| \theta_c \|_1$

Repeat

① Pick a sample $i$ from training data.

② Calculate gradient based on $i$,

For $j = 1, \dots K$

$$\nabla \theta_j L_i = \frac{\exp(\theta_j^T \phi_i) \cdot \phi_i}{\sum\limits_{c=1}^{K} \exp(\theta_c^T \phi_i)} \qquad \left[ \text{where } j \neq y_i \right]$$

$$= \frac{\exp(\theta_j^T \phi_i) \cdot \phi_i}{\sum\limits_{c=1}^{K} \exp(\theta_c^T \phi_i)} - \phi_i \left[ \text{where } j = y_i \right].$$

③ For $j = 1, \dots K$

$$\theta_j = \theta_j - \delta \nabla_{\theta_j} L_i \qquad \text{// Gradient Descent}$$
$$\text{// Step size } = \delta.$$

④ For $j = 1, \dots K$

$$\theta_j = \text{Prox}_{\delta\lambda/2} \theta_j \qquad \text{// Proximal operation.}$$

$$\text{Prox}_{\delta\lambda/2 \|\cdot\|} \theta_j = \begin{cases} \theta_j - \delta\lambda/2 & \theta_j > \delta\lambda/2 \\ 0 & |\theta_j| \leq \delta\lambda/2 \\ \theta_j + \delta\lambda/2 & \theta_j < -\delta\lambda/2 \end{cases}$$

Untill Convergence.

③

3. $\Pr(y_i = c) = \pi_c \quad \Pr(x_i | y_i = c) \sim \text{Multi}(p_c, L_i) = \dfrac{L_i!}{\prod_{j=1}^{d} x_{ij}!} \prod_{j=1}^{d} p_{cj}^{x_{ij}}$

a) $\log \, p_x(x, y; \pi, p)$

$= \sum_{i=1}^{m} \log \, p_x(x_i, y_i; \pi, p)$

$= \sum_{i=1}^{n} \log \, p_x(y_i; \pi) + \sum_{i=1}^{n} \log \, p_x(x_i | y_i; p)$

$= \sum_{i=1}^{n} \log \pi_{y_i} + \sum_{i=1}^{n} \sum_{j=1}^{d} x_{ij} \log p_{y_i j} + \text{Const.}$

$\mathcal{L}\mathcal{L} = \sum_{i=1}^{n} \sum_{c=1}^{k} y_{ic} \log \pi_c + \sum_{i=1}^{n} \sum_{c=1}^{k} \sum_{j=1}^{d} x_{ij} y_{ic} \log p_{cd} + \text{Const.}$

b) E-step: Calculate expectation $\quad y_{ic} = \dfrac{\pi_c \, f(x_i; p_c)}{\sum_{k'=1}^{K} \pi_{k'} \, f(x_i; p_{k'})}$

$f(x_i; p_c) = \dfrac{L_i!}{\prod_{j=1}^{d} x_{ij}!} \prod_{j=1}^{d} p_{cj}^{x_{ij}}$

M-step:

① $\pi$

$\dfrac{\partial \mathcal{L}\mathcal{L}\pi}{\partial \pi_c} = 0$

$\sum_{i=1}^{m} y_{ic}/\pi_c - \lambda_1 = 0$

$\pi_c = \frac{1}{\lambda_1} \sum_{i=1}^{n} y_{ic}$

Constraint $\sum_{c=1}^{k} \pi_c = 1 \Rightarrow \sum_{c=1}^{K} \frac{1}{\lambda_1} \sum_{i=1}^{m} y_{ic} = 1$

$\lambda_1 = \sum_{i=1}^{n} \sum_{c=1}^{K} y_{ic}$

$= \sum_{i=1}^{n} 1 \left[ \sum_{c=1}^{k} y_{ic} = 1 \right]$

$= n$

$\boxed{\pi_c = \frac{1}{n} \sum_{i=1}^{m} y_{ic}}$

$\mathcal{L}\mathcal{L}_\pi = \mathcal{L}\mathcal{L} + \lambda_1 \left(1 - \sum_{c=1}^{k} \pi_c\right)$

$\downarrow \qquad \underbrace{\phantom{xxx}}$

Lagrange    Constraint
Multiplier.

② $p$

$\dfrac{\partial \mathcal{L}\mathcal{L}}{\partial p_{cj}} = 0$

$\sum_{i=1}^{n} \dfrac{x_{ij} \, y_{ic}}{p_{cj}} - \lambda_2 = 0$

$p_{cj} = \frac{1}{\lambda_2} \sum_{i=1}^{n} x_{ij} y_{ic}$

now $\sum_{j=1}^{d} \frac{1}{\lambda_2} \sum_{i=1}^{n} x_{ij} y_{ic} = 1$

$\lambda_2 = \sum_{j=1}^{d} \sum_{i=1}^{n} x_{ij} y_{ic}$

$= \sum_{i=1}^{n} X_i \, y_{ic} \quad \left[ \begin{array}{l} X_i = \text{total length} \\ \text{of words} \end{array} \right]$

⑥ $\boxed{p_{cj} = \dfrac{\sum_{i=1}^{n} x_{ij} y_{ic}}{\sum_{i=1}^{n} X_i \, y_{ic}}}$

$\sum_{j=1}^{d} p_{cj} = 1$

$\mathcal{L}\mathcal{L}_p = \mathcal{L}\mathcal{L} + \lambda_2 \left(1 - \sum_{j=1}^{d} p_{cj}\right)$

④

**4.** Classical scaling:

$$\underset{y_1,\ldots,y_m}{\text{minimize}} \sum_{i=1}^{m}\sum_{j=1}^{m} (\tilde{x}_i^T \tilde{x}_j - y_i^T y_j)^2 \quad \text{where } \tilde{x}_i = x_i - (1/m)\sum_{j=1}^{m} x_j$$

a) The above formulation can be rewritten in matrix form,

$$\underset{Y^T Y}{\min} \; \| \tilde{X}^T \tilde{X} - Y^T Y \|$$

The matrix $\tilde{X}^T \tilde{X}$ is symmetric, therefore can be applied Eigen Value decomposition.

From Eckart Young Theorem,

$$\underset{\tilde{A}}{\min} \| A - \tilde{A} \|^2 \quad \text{s.t} \quad \text{rank}(\hat{A}) = k < n$$

is $\hat{A} = \sum_{i=1}^{k} \sigma_i u_i v_i^T$     $\hat{A} = U_1 \, \Sigma_1 V_1^T$

This is the best approximation of $A$ of rank $k$.

In our problem we use eigen decomposition instead of SVD.

$$\tilde{Y}^T Y = U_1 \Lambda_1 U_1^T = (\Lambda_1^{1/2} U_1^T)^T \Lambda_1^{1/2} U_1^T$$

$$Y = \Lambda_1^{1/2} U_1^T \quad \Big| \quad \begin{array}{l} \Lambda_1 = \text{First } k \text{ eigen values of } \Lambda \\ U_1 = \text{First } k \text{ eigen vectors of } U \end{array}$$

b) $D_{ij} = \| x_i - x_j \|^2 = \| x_i \|^2 + \| x_j \|^2 - 2 x_i^T x_j$

$$= b 1^T + 1 b^T - 2 x^T x$$

Let $b = \begin{bmatrix} \| x_1 \|^2 \\ \| x_2 \|^2 \\ \vdots \\ \| x_m \|^2 \end{bmatrix}$

Now, $x^T x = \frac{1}{2}(b 1^T + 1 b^T - D)$

Now let $H$ be $(I - \frac{1}{m} 1 1^T)$, so, $\tilde{X} = XH$ (Centered data)

Now, $\tilde{X}^T \tilde{X} = H^T X^T X H$

$$= \frac{1}{2}(H b 1^T H + H 1 b^T H - H D H) \qquad \begin{bmatrix} 1^T H = 0 \\ H^T 1 = 0 \end{bmatrix}$$

$$= -\frac{1}{2}(I - \frac{1}{m} 1 1^T) D (I - \frac{1}{m} 1 1^T)$$

$$= -\frac{1}{2} B.$$

⑤