

5.12. From the definitions of S_w and S_b in the textbook we have:

$$S_w = \sum_{i=1}^M P_i E_i[(x - \mu_i)(x - \mu_i)^T] \quad , \quad S_b = \sum_{i=1}^M P_i (\mu_i - \mu_0)(\mu_i - \mu_0)^T$$

$$S_m = E[(x - \mu_0)(x - \mu_0)^T] \quad \text{where } E_i \text{ is } E_{x \in \mathcal{X}_i}$$

$$\rightarrow S_w + S_b = \sum_{i=1}^M P_i \left[E_i[(x - \mu_i)(x - \mu_i)^T] + (\mu_i - \mu_0)(\mu_i - \mu_0)^T \right]$$

if we write $x - \mu_i$ as $x - \mu_0 + \mu_0 - \mu_i$ (so that it connects the expression to both S_b and S_m) we get:

$$\begin{aligned} S_w + S_b &= \sum_{i=1}^M P_i \left[E_i[(x - \mu_0 + \mu_0 - \mu_i)(x - \mu_0 + \mu_0 - \mu_i)^T] + (\mu_i - \mu_0)(\mu_i - \mu_0)^T \right] \\ &= \sum_{i=1}^M P_i \left[E_i[(x - \mu_0)(x - \mu_0)^T + 2(x - \mu_0)(\mu_0 - \mu_i)^T + (\mu_0 - \mu_i)(\mu_0 - \mu_i)^T] + (\mu_i - \mu_0)(\mu_i - \mu_0)^T \right] \end{aligned}$$

Expectation of this term is zero because $E(x - \mu_0) = 0$ and $(x - \mu_0)$ and $(\mu_0 - \mu_i)$ are independent

this term does not depend on x and is therefore a constant w.r.t the expectation

$$\begin{aligned} \Rightarrow S_w + S_b &= \sum_{i=1}^M P_i \left[E_i[(x - \mu_0)(x - \mu_0)^T] + (\mu_0 - \mu_i)(\mu_0 - \mu_i)^T + (\mu_i - \mu_0)(\mu_i - \mu_0)^T \right] \\ &= \sum_{i=1}^M P_i E_i[(x - \mu_0)(x - \mu_0)^T] = E[(x - \mu_0)(x - \mu_0)^T] = \underline{S_m} \end{aligned}$$

(~~sum over all values of $\mu_i - \mu_0$ is 0~~)

5.13. we have $\rho_{ij} = \frac{\sum_{n=1}^N x_{ni} x_{nj}}{\sqrt{\sum_{n=1}^N x_{ni}^2} \sqrt{\sum_{n=1}^N x_{nj}^2}}$, the cauchy-schwarz inequality states that for two vectors x and y we have:

$$|x^T y| \leq \|x\| \|y\|$$

we define $x = \begin{bmatrix} x_{1i} \\ x_{2i} \\ \vdots \\ x_{Ni} \end{bmatrix}$ and $y = \begin{bmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{Nj} \end{bmatrix} \rightarrow |x^T y| = \left| \sum_{n=1}^N x_{ni} x_{nj} \right| \leq \|x\| \|y\|$

dividing both sides by $\|x\| \|y\|$

$$\frac{\left| \sum_{n=1}^N x_{ni} x_{nj} \right|}{\sqrt{\sum_{n=1}^N x_{ni}^2} \sqrt{\sum_{n=1}^N x_{nj}^2}} = |\rho_{ij}| \leq 1$$

5.14. we know that for two same-covariance distributions,

$$d_{ij} = (\mu_i - \mu_j)^T \Sigma^{-1} (\mu_i - \mu_j) \quad \text{since there are only two classes:}$$

$$d_{12} = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2)$$

We will now calculate $S_w^{-1} S_b$:

$$S_w = \sum_{i=1}^2 P_i \Sigma_i = \frac{1}{2} \Sigma + \frac{1}{2} \Sigma = \Sigma \rightarrow S_w^{-1} = \Sigma^{-1}$$

$$S_b = \sum_{i=1}^2 P_i (\mu_i - \mu_0) (\mu_i - \mu_0)^T = \frac{1}{2} (\mu_1 - \mu_0) (\mu_1 - \mu_0)^T + \frac{1}{2} (\mu_2 - \mu_0) (\mu_2 - \mu_0)^T$$

$$\text{using } \mu_0 = \sum_{i=1}^M P_i \mu_i \rightarrow S_b = \frac{1}{2} \left(\frac{1}{2} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T + \frac{1}{2} (\mu_2 - \mu_1) (\mu_2 - \mu_1)^T \right)$$

$$= \frac{1}{4} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T$$

$$\rightarrow S_w^{-1} S_b = \frac{1}{4} \Sigma^{-1} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \rightarrow \text{trace} \{ S_w^{-1} S_b \} = \frac{1}{4} \text{trace} \left\{ \underbrace{\Sigma^{-1}}_A \underbrace{(\mu_1 - \mu_2)}_B \underbrace{(\mu_1 - \mu_2)^T}_C \right\}$$

$$= \frac{1}{4} \text{trace} \left\{ \underbrace{(\mu_1 - \mu_2)^T}_C \underbrace{\Sigma^{-1}}_A \underbrace{(\mu_1 - \mu_2)}_B \right\}$$

trace for the product of three matrices is invariant under cyclic permutations

on the other hand because d_{12} is a scalar we have

$$\text{trace} \{ (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) \} = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) = d_{12}$$

$$\text{trace} \{ S_w^{-1} S_b \} = \frac{1}{4} d_{12}$$

5.17. From the matrix cookbook we know the following:

$$\frac{\partial}{\partial X} \text{trace} \{ X A \} = \frac{\partial}{\partial X} \text{trace} \{ A X \} = A^T \quad (\text{eq 100})$$

$$\frac{\partial}{\partial X} \text{trace} \{ X^T B X \} = (B + B^T) X \quad (\text{eq 108})$$

(where C is symmetric)

$$\frac{\partial}{\partial X} \text{trace} \{ (X^T C X)^{-1} A \} = -(C X (X^T C X)^{-1}) (A + A^T) (X^T C X)^{-1} \quad (\text{eq 125})$$

the task is to prove (eq 126) using these presuppositions.

$$\frac{\partial}{\partial A} \text{trace} \{ (A^T S_1 A)^{-1} (A^T S_2 A) \} = \frac{\partial \text{trace} \{ (A^T S_1 A)^{-1} C_1 \}}{\partial A} + \frac{\partial \text{trace} \{ C_2 (A^T S_2 A) \}}{\partial A}$$

Where C_1 and C_2 are constants w.r.t. their derivatives

$$\text{and } C_1 = A^T S_2 A, \quad C_2 = (A^T S_1 A)^{-1}$$

→ we use eq. 125 and eq. 117 from the matrix cookbook:

$$\frac{\partial}{\partial A} \text{trace} \{ (A^T S_1 A)^{-1} C_1 \} = - (S_1 A (A^T S_1 A)^{-1}) (C_1 + C_1^T) (A^T S_1 A)^{-1}$$

Substituting C_1 and the fact that $C_1 + C_1^T = 2C_1$

$$\rightarrow \frac{\partial}{\partial A} \text{trace} \{ (A^T S_1 A)^{-1} C_1 \} = -2S_1 A (A^T S_1 A)^{-1} (A^T S_2 A) (A^T S_1 A)^{-1}$$

$$\frac{\partial}{\partial A} \text{trace} \{ C_2 (A^T S_2 A) \} = \frac{\partial}{\partial A} \text{trace} \{ (A^T S_2 A) C_2 \} \xrightarrow{\text{eq 117}}$$

$$= S_2 A C_2 + S_2^T A C_2^T \quad \xrightarrow[\substack{S_2 = S_2^T \\ C_2 = C_2^T}]{2S_2 A C_2 = 2S_2 A (A^T S_1 A)^{-1}}$$

$$\rightarrow \frac{\partial}{\partial A} \text{trace} \{ (A^T S_1 A)^{-1} (A^T S_2 A) \} = -2S_1 A (A^T S_1 A)^{-1} (A^T S_2 A) (A^T S_1 A)^{-1} + 2S_2 A (A^T S_1 A)^{-1}$$

5.19. As the hint suggests we will define functions $g_i(x) := f_i(x) - f_{i+1}(x)$ $\forall 1 \leq i \leq M-1$

we prove that these $M-1$ $g(\cdot)$ functions are enough for classification.

In order to make a classification based on the discriminant functions $f_i(\cdot)$ s

we assign x to class (i) if $i = \arg \max_j f_j(x)$ or $f_i(x)$ is the maximum among all possible i 's. This means $M-1$ equalities must hold:

$$\forall j < i \quad f_i(x) > f_j(x)$$

$$\forall i+1 \leq k \leq M \quad f_i(x) > f_k(x)$$

we can construct these inequalities using $g_i(\cdot)$ as follows:

for two arbitrary indices $1 \leq i, j \leq M$, $i \neq j$

we want to prove $f_i(x) > f_j(x)$ or $f_i(x) - f_j(x) > 0$

if $j < i$:

$$f_i(x) - f_j(x) > 0 \rightarrow f_j(x) - f_i(x) < 0$$

$$\begin{aligned} &\rightarrow \cancel{f_j(x)} \\ &(f_j(x) - f_{j+1}(x)) + (f_{j+1}(x) - f_{j+2}(x)) \\ &+ \dots + (f_{i-1}(x) - f_i(x)) \\ &= \sum_{k=j}^{i-1} g_k(x) < 0 \end{aligned}$$

if $j > i$:

$$\begin{aligned} &f_i(x) - f_j(x) > 0 \\ &\rightarrow (f_i(x) - f_{i+1}(x)) + (f_{i+1}(x) - f_{i+2}(x)) \\ &+ \dots + (f_{j-1}(x) - f_j(x)) = \sum_{k=i}^{j-1} g_k(x) > 0 \end{aligned}$$

$\rightarrow M-1$ functions $\{g_i(\cdot)\}$ are enough for classification.

5.20. $S_b = \sum_{i=1}^2 P_i (\mu_i - \mu_0) (\mu_i - \mu_0)^T$. Substituting $\mu_0 = \sum_{i=1}^2 P_i \mu_i$ and expanding we have:

$$S_b = P_1 (P_2 (\mu_1 - \mu_2) P_2 (\mu_1 - \mu_2)^T) + P_2 (P_1 (\mu_1 - \mu_2) P_1 (\mu_1 - \mu_2)^T)$$

$$= P_1 P_2 (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \rightarrow S_b \text{ is a scalar therefore it is of rank } 1$$

$\rightarrow S_w^{-1} S_b$ has exactly one non-zero eigenvalue

\rightarrow Sum of the eigen values of $S_w^{-1} S_b$ is equal to its only non-zero eigenvalue

$$\Rightarrow \lambda = \text{trace} \left\{ \underset{1 \times 1}{P_1} \underset{1 \times 1}{P_2} \underset{1 \times 1}{S_w^{-1}} (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \right\} = P_1 P_2 (\mu_1 - \mu_2)^T S_w^{-1} (\mu_1 - \mu_2)$$

produces a scalar which is equal to the trace

if $S_w^{-1}(\mu_1 - \mu_2)$ is an eigenvector we must have:

$$\underbrace{P_1 P_2 (\mu_1 - \mu_2)^T S_w^{-1} (\mu_1 - \mu_2)}_{\lambda} (S_w^{-1} (\mu_1 - \mu_2)) = (S_w^{-1} S_b) (S_w^{-1} (\mu_1 - \mu_2))$$

$$P_1 P_2 \overset{S_w^{-1}}{(\mu_1 - \mu_2)} (\mu_1 - \mu_2)^T (S_w^{-1} (\mu_1 - \mu_2))$$

$$= P_1 P_2 (\mu_1 - \mu_2)^T S_w^{-1} (\mu_1 - \mu_2) S_w^{-1} (\mu_1 - \mu_2) = \underline{\lambda S_w^{-1} (\mu_1 - \mu_2)}$$

5.21. let A diagonalize both Σ_1 and Σ_2 : $A^T \Sigma_1 A = I \rightarrow \Sigma_1 = A^{-T} A^{-1}$
 $A^T \Sigma_2 A = D \rightarrow \Sigma_2 = A^{-T} D A^{-1}$

$$\rightarrow \Sigma_1^{-1} \Sigma_2 = (A^{-T} A^{-1})^{-1} (A^{-T} D A^{-1}) = (A A^T) (A^{-T} D A^{-1}) = A D A^{-1}$$

let $\lambda_i \neq \lambda_j$ be two distinct eigenvalues of $\Sigma_1^{-1} \Sigma_2$:

$$\begin{aligned} A D A^{-1} v_i &= \lambda_i v_i \rightarrow D \overset{v_i}{(A^{-1} v_i)} = \lambda_i \overset{v_i}{(A^{-1} v_i)} \rightarrow \begin{cases} D v_i' = \lambda_i v_i' \\ D v_j' = \lambda_j v_j' \end{cases} \\ A D A^{-1} v_j &= \lambda_j v_j \rightarrow D \overset{v_j}{(A^{-1} v_j)} = \lambda_j \overset{v_j}{(A^{-1} v_j)} \end{aligned}$$

$\rightarrow v_i'$ and v_j' are eigenvectors of D . Since D is a diagonal matrix it is also symmetrical and the eigenvectors of a symmetrical matrix are orthogonal $\rightarrow v_i'^T v_j' = 0$ (or δ_{ij} if $i=j$ is allowed)

$$\Rightarrow (A^{-1} v_i)^T (A^{-1} v_j) = v_i^T \underset{\Sigma_1}{(A^{-T} A^{-1})} v_j = v_i^T \Sigma_1 v_j = 0 \quad \underline{\text{(or } \delta_{ij})}$$