

1. $X \rightarrow$ Zero mean p -dimensional random vector.

$$E[X] = 0$$

$$\text{Covariance matrix: } R = E[X X^T]$$

$$\text{estimation: } \hat{X} = \sum_{i=1}^M \alpha_i e_i$$

e_i = orthonormal eigenvectors of covariance matrix R

$$\alpha = [\alpha_1 \dots \alpha_p]^T$$

To show: $J = \|X - \hat{X}\|^2 \xrightarrow{\min} \alpha_i = e_i^T X \quad i=1, 2, \dots, M$
as principal component

Proof:

$$\|X - \hat{X}\|^2 = (X - \hat{X})^T (X - \hat{X})$$

$$\text{Substituting } \hat{X} = \sum_{i=1}^M \alpha_i e_i$$

$$J = \|X - \sum_{i=1}^M \alpha_i e_i\|^2 = (X - \sum_{i=1}^M \alpha_i e_i)^T (X - \sum_{i=1}^M \alpha_i e_i)$$

$$= X^T X - X^T \sum_{i=1}^M \alpha_i e_i - \left(\sum_{i=1}^M \alpha_i e_i \right)^T X + \left(\sum_{i=1}^M \alpha_i e_i \right)^T \left(\sum_{i=1}^M \alpha_i e_i \right)$$

$$= X^T X - \sum_{i=1}^M \alpha_i e_i^T X - \sum_{i=1}^M \alpha_i e_i^T X + \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j e_i^T e_j$$

Let's
use to
separate
& disti-
nished

$\because e_i$ = orthonormal,

$$e_i^T e_j = \delta_{ij} \quad \text{where, } \delta_{ij} = 1 \quad \text{if } i=j$$

$$\delta_{ij} = 0 \quad \text{if } i \neq j$$

$$\therefore \sum_{i=1}^M \sum_{j=1}^M \alpha_i \alpha_j e_i^T e_j = \sum_{i=1}^M \alpha_i^2$$

minimising J ,

take derivative w.r.t α_i , $\lambda = 0$

$$\frac{\partial J}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i} \left[x^T x - 2 \sum_{i=1}^M \alpha_i x^T e_i + \sum_{i=1}^M \alpha_i^2 \right]$$

$$= -2 x^T e_i + 2 \alpha_i = 0$$

$$\Rightarrow \alpha_i = x^T e_i$$

i.e. $\alpha_i = e_i^T x$ { dot product is commutative }

\Rightarrow principal components α_i are given by projection of the data vector x onto eigenvectors e_i .

2. $P(x|w_i) \rightarrow \text{mean} = \mu_i$

Co-variance matrix Σ_i $i=1,2$

$y = w^T x$: projection.

$P(y|w_i) \rightarrow \text{mean } \mu_i$
var σ_i^2

(a) TO show:
 $J_1(w) = \frac{(\mu_1 - \mu_2)^2}{\sigma_1^2 + \sigma_2^2}$ is max by

$$w = (\Sigma_1 + \Sigma_2)^{-1} (\mu_1 - \mu_2)$$

{reference;
class notes}.

$$\mu_i = \frac{1}{\#D_i} \sum_{x \in D_i} x \quad 0 \leq \mu_i \leq 1$$

$$= \frac{1}{\#D_i} w^T \sum_{x \in D_i} x = w^T \mu_i$$

$$\therefore \text{Var}(y) = \text{Var}(w^T x)$$

$$= w^T \text{Cov}(x) w$$

$$\Rightarrow \sigma_i^2 = w^T \Sigma_i w$$

$$\text{Let } V = \mu_1 - \mu_2$$

$$\Sigma = \Sigma_1 + \Sigma_2$$

$$\therefore J_1(w) = \frac{(w^T (\mu_1 - \mu_2))^2}{w^T (\Sigma_1 + \Sigma_2) w}$$

$$= \frac{(w^T V)^2}{w^T \Sigma w}$$

to maximise, we use Lagrange multipliers:-

$$\mathcal{L}(w, \lambda) = \frac{(w^T V)^2}{w^T \Sigma w} - \lambda (w^T w - 1)$$

Considering the Rayleigh quotient:

$\frac{(w^T V)^2}{w^T \Sigma w}$ • The Rayleigh quotient is maximised when w is $\propto \Sigma^{-1} V$

$$\text{i.e. } w \propto (\Sigma_1 + \Sigma_2)^{-1} (\mu_1 - \mu_2)$$

(b) Let projection,

$$y = w^T x$$

$$\text{with } \mu_i' = w^T \mu_i$$

$$\sigma_y^{(i)} = w^T \Sigma_i w$$

→ mean & covariance matrix Σ_i

Considering contributions from each class w_1 & w_2 , weighted by prior probabilities $P(w_1)$ and $P(w_2)$

Variance of y is

$$\sigma_y^2 = P(w_1) \sigma_y^{(1)} + P(w_2) \sigma_y^{(2)}$$

$$\text{where, } \sigma_y^{(1)} = w^T \Sigma_1 w \quad \sigma_y^{(2)} = w^T \Sigma_2 w$$

$$\therefore \sigma_y^2 = P(\omega_1) \underbrace{w^T \Sigma_1 w}_{\text{scalar}} + P(\omega_2) \underbrace{w^T \Sigma_2 w}_{\text{scalar}}$$

$$\therefore \sigma_y^2 = w^T (P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2) w$$

$$\sigma_y^2 = w^T (P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2) w$$

$$\therefore \sigma_y^2 = w^T (P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2) w$$

$$J_2(w) = \frac{(w^T (\mu_1 - \mu_2))^2}{w^T (P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2) w}$$

Considering Rayleigh quotient:

$$\frac{(w^T a)^2}{w^T B w} \quad \text{where } a = \mu_1 - \mu_2$$

$$B = P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2$$

It is max when $w \propto B^{-1} a$

$$\text{i.e. } w \propto (P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2)^{-1} (\mu_1 - \mu_2)$$

$\therefore J(w)$ is maximised when

$$w = [P(\omega_1) \Sigma_1 + P(\omega_2) \Sigma_2]^{-1} (\mu_1 - \mu_2)$$

(c) The objective function for Fisher's LDA can be written as:

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \quad (\text{between-class scatter matrix})$$

$$S_W = \Sigma_1 + \Sigma_2 \quad (\text{within class scatter matrix})$$

$J_1(w)$ is closer to the criterion used by Fisher's LDA because it compares between class mean difference to the within-class variance without involving prior probabilities. Similar to original Fisher's LDA of max ratio of between-class scatter to within-class scatter.