

Notes

the means in projected space for the two classes.

$$\mu_1 = \frac{1}{N_1} \sum_{k \in C_1} z_k \quad \mu_2 = \frac{1}{N_2} \sum_{k \in C_2} z_k$$

for slide 7

goal: find  $w$  that maximizes the differences of means.

$$\operatorname{argmax}_w | \mu_2(w) - \mu_1(w) |, \|w\| = 1$$

$$= \operatorname{argmax}_w \left| \frac{1}{N_2} \sum_{k \in C_2} z_k(w) - \frac{1}{N_1} \sum_{k \in C_1} z_k(w) \right|$$

$$= \operatorname{argmax}_w \left| \frac{1}{N_2} \sum_{k \in C_2} w^T \cdot x_k - \frac{1}{N_1} \sum_{k \in C_1} w^T \cdot x_k \right|$$

$$= \operatorname{argmax}_w \left| w^T \cdot \left( \frac{1}{N_2} \sum_{k \in C_2} x_k - \frac{1}{N_1} \sum_{k \in C_1} x_k \right) \right|$$

$$= \operatorname{argmax}_w | w^T \cdot (\mu_2 - \mu_1) |$$

Conclusion: the best vector  $w$  is the one that aligns with  $(\mu_2 - \mu_1)$ , i.e.

$$w = (\mu_2 - \mu_1) / \| \mu_2 - \mu_1 \|, \quad w^T \cdot (\mu_2 - \mu_1) = \left( \frac{d}{\| d \|} \right)^T \cdot d = \frac{d^T \cdot d}{\| d \|} = \frac{\| d \|^2}{\| d \|} = \| d \|$$

proof:

$$d = \mu_2 - \mu_1$$

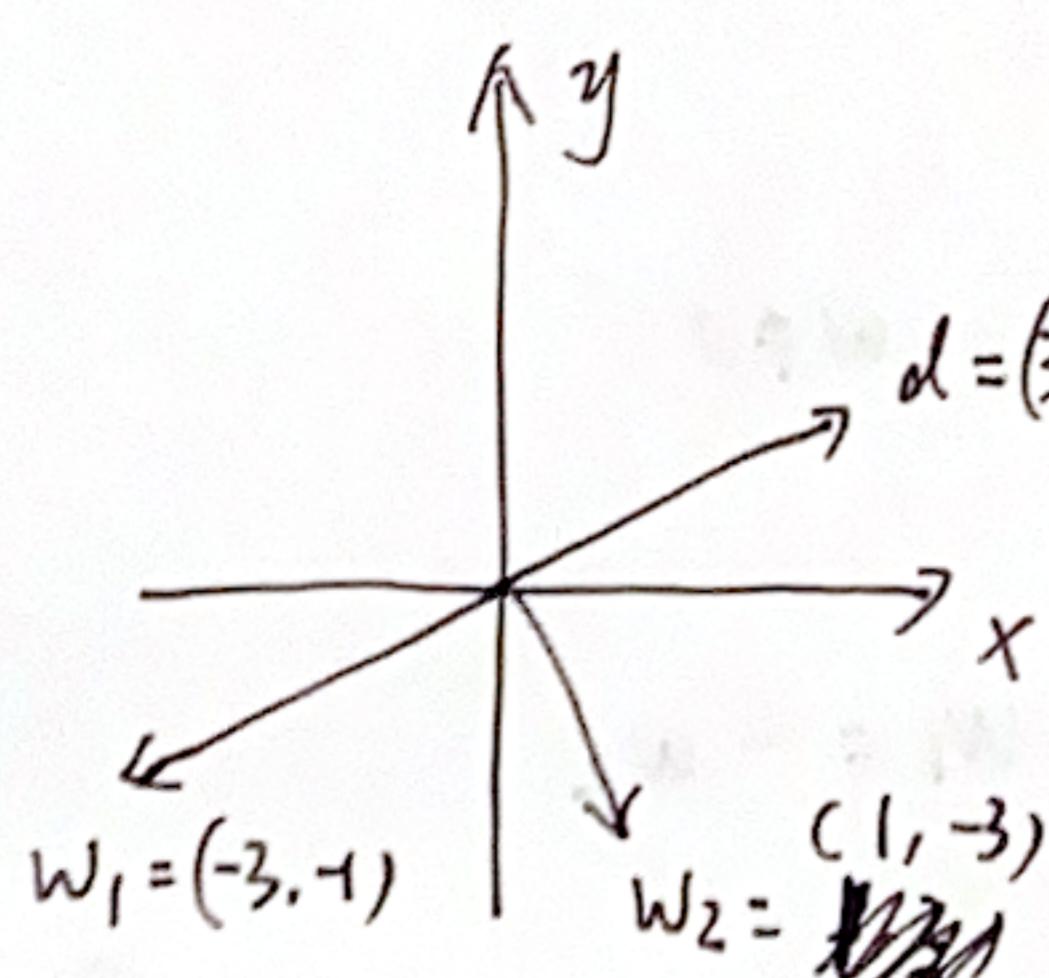
$$|w^T \cdot d| \leq \|w^T\| \cdot \|d\|, \|w^T\| = 1$$

$$\Rightarrow |w^T \cdot d| \leq 1 \cdot \|d\| = \|d\|$$

Wenn  $w$  is parallel to  $d$ , then the equation holds  $|w^T \cdot d| = \|d\|$ ,

otherwise,  $|w^T \cdot d| < \|d\|$

geometric meaning illustration



$$|w_1 \cdot d| = |(-3, -1) \cdot (3, 1)| = |-9 - 1| = 10$$

$$|w_2 \cdot d| = |(1, -3) \cdot (3, 1)| = |3 - 3| = 0$$

## Slides 8

limitation of mean separation

$$w = \frac{\mu_2 - \mu_1}{\|\mu_2 - \mu_1\|}, w \text{ points to the direction connecting the two class centroids.}$$

the class means are as far apart as possible after projection, but it not guarantee good class separability.

the reasons are :

- 1) significant class overlap in projected space
- 2) the method ignores class variances, the orientation of the clouds ; the shape of the distributions.

If we rotate the projection direction slightly (not using  $\mu_2 - \mu_1$ ) ; the projected distributions became more separated.

→ Motivation for Fischer Discriminant

## Fisher Discriminant

Slide 11

$$\mu_1 = \frac{1}{|C_1|} \sum_{k \in C_1} z_k \quad \mu_2 = \frac{1}{|C_2|} \sum_{k \in C_2} z_k$$

$$S_1 = \sum_{k \in C_1} (z_k - \mu_1)^2 \quad S_2 = \sum_{k \in C_2} (z_k - \mu_2)^2$$

$$J(w) = \underset{w}{\operatorname{argmax}} \frac{(\mu_2(w) - \mu_1(w))^2}{S_1(w) + S_2(w)}$$

$$= \frac{w^T S_B w}{w^T S_w w}$$

proof

where  $S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$

$$S_w = S_1 + S_2$$

why  $(w^T v)^2 = w^T v v^T w$

Proof:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^{d \times 1}, \quad v = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix} \in \mathbb{R}^{d \times 1}$$

$$w^T \cdot v = [w_1 \ w_2 \dots \ w_d] \cdot \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix}$$

$$= \sum_{i=1}^d w_i v_i = [w_1 v_1 \ w_2 v_2 \ \dots \ w_d v_d]$$

$$(w^T v)^2 = \left( \sum_{i=1}^d w_i v_i \right)^2$$

$$v \cdot v^T = \begin{bmatrix} v_1 \\ v_2 \\ \vdots \\ v_d \end{bmatrix} \cdot [v_1 \ v_2 \ \dots \ v_d] = \begin{bmatrix} v_1 v_1 & v_1 v_2 & \dots & v_1 v_d \\ v_2 v_1 & v_2 v_2 & \dots & v_2 v_d \\ \vdots & & & \\ v_d v_1 & v_d v_2 & \dots & v_d v_d \end{bmatrix}$$

$$w^T \cdot (v v^T) w = (w^T v) \cdot (v^T w)$$

$$v^T \cdot w = [v_1 \ v_2 \ \dots \ v_d] \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = [v_1 w_1 \ v_2 w_2 \ \dots \ v_d w_d]$$

$$\Rightarrow w^T (v v^T) w = (w^T v)^2$$

Herleitung von  $S_j$

$$S_j(w) = \sum_{k \in C_j} (z_k - \mu_j)^2$$

$$= \sum_{k \in C_j} (w^T x_k - w^T \mu_j)^2$$

$$= w^T \underbrace{\sum_{k \in C_j} (x_k - \mu_j)(x_k - \mu_j)^T}_{S_j} w = w^T S_j w$$

$S_j$  is a scatter matrix / unnormalized covariance matrix for the data of class  $j$ .

等式后转换:

$$\mu_j = \frac{1}{|C_j|} \sum_{k \in C_j} z_k = \frac{1}{|C_j|} \sum_{k \in C_j} w^T x_k = w^T \mu_j$$

$$\mu_2(w) - \mu_1(w) = w^T (\mu_2 - \mu_1)$$

$$(\mu_2(w) - \mu_1(w))^2 = (w^T (\mu_2 - \mu_1))^2$$

$$= w^T (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T w$$

\* reason see next page, detailed explanation about how to calculate.

$$S_j(w) = w^T S_j w$$

$$S_1(w) = w^T S_1 w, \quad S_2(w) = w^T S_2 w$$

$$S_1(w) + S_2(w) = w^T (S_1 + S_2) w = w^T S_w w$$

slide 13

$$\mathcal{L}(w, \lambda) = w^T S_B w + \lambda(1 - w^T S_w w)$$

$$\text{explain why } \frac{\partial}{\partial w} (w^T A w) = 2Aw,$$

where  $A$  is symmetric,  $A = A^T$

$$w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} \in \mathbb{R}^{d \times 1}, A = (a_{ij})$$

$$w^T A w = \sum_{i=1}^d \sum_{j=1}^d w_i a_{ij} w_j, \text{ because}$$

$$[w_1 \ w_2 \ \dots \ w_d] \cdot \begin{bmatrix} a_{11} & a_{12} \dots a_{1d} \\ a_{21} & a_{22} \dots a_{2d} \\ \vdots & \vdots \\ a_{d1} & \dots \ a_{dd} \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix}$$

$$= \begin{bmatrix} 1 & 0 & \dots & 0 \end{bmatrix} \cdot \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_d \end{bmatrix} = \text{a scalar}$$

$$w^T A w = \sum_{i=1}^d \sum_{j=1}^d w_i a_{ij} w_j, \text{ observe } w_k$$

case 1 :  $i=k$

$$\Rightarrow w_k a_{kj} w_j$$

$$\frac{\partial}{\partial w_k} (w_k a_{kj} w_j) = a_{kj} w_j$$

$$\text{sum over all } j \quad \sum_j a_{kj} w_j$$

case 2 :  $j=k$

$$\Rightarrow w_i a_{ik} w_k \quad \frac{\partial}{\partial w_k} (w_i a_{ik} w_k) = a_{ik} w_k$$

$$\text{sum over all } i \quad \sum_i a_{ik} w_k$$

$A$  is symmetric  $\Rightarrow a_{ij} = a_{ji}$

$$\Rightarrow \frac{\partial}{\partial w} (w^T A w) = 2Aw$$

slide 14

explain ~~why~~ How from

$$S_B w = \lambda S_w w$$

$$\cancel{S_w^{-1} S_B w} \rightarrow (S_w^{-1} S_B)_w = \lambda w,$$

where  $S_w$  is invertible.

$$S_B w = \lambda S_w w$$

$$S_w^{-1} S_B w = S_w^{-1} \lambda S_w w$$

$$= \lambda S_w^{-1} S_w w$$

$$= \lambda I w$$

$$= \lambda w$$

$\nabla_w \mathcal{E}_k(w, b)$  是如何计算出来的

$$z_k = w^T x_k + b \quad , \quad y_k = \text{sign}(z_k)$$

$$\mathcal{E}(w, b) = \frac{1}{N} \sum_{k=1}^N \underbrace{\max(0, -z_k t_k)}_{\mathcal{E}_k(w, b)}$$

$$-z_k t_k = -w^T x_k t_k - b t_k$$

$$\frac{\partial \Sigma_k(w, b)}{\partial w} = \begin{cases} -x_k t_k & , y_k \neq t_k \\ 0 & , y_k = t_k \end{cases}$$