

Σ symmetric positive-definite. (e.g. covariance matrix).

eigen-decomposition:

$$\Sigma = V D V^T.$$

$$D = \begin{pmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & \\ 0 & & & \lambda_n \end{pmatrix}$$

$\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \geq 0$ eigen-values.

$$V = \begin{pmatrix} v_1^T & v_2^T & \dots & v_n^T \end{pmatrix}$$

v_1, v_2, \dots, v_n are eigen-vectors.

Properties:

(1) $v_i \perp v_j$ ($\langle v_i, v_j \rangle = 0$).

(2) $\Sigma^{\frac{1}{2}} = V D^{\frac{1}{2}} V^T$ ($\because \Sigma^{\frac{1}{2}} \cdot \Sigma^{\frac{1}{2}} = V D^{\frac{1}{2}} \overset{I}{V^T} V D^{\frac{1}{2}} V^T = V D V^T = \Sigma$).

(3) Suppose $\text{cov}(X) = \Sigma$;

then $\text{Var}(\alpha X) = \alpha \Sigma \alpha^T$.

Consider $\max_{\alpha} \text{Var}(\alpha X) = \alpha \Sigma \alpha^T$ subject to $|\alpha| = 1$,

the solution is $\alpha^* = v_1$ and $\text{Var}(\alpha^* X) = v_1 \cdot V D V^T v_1^T = \lambda_1$

$$f(X|Y=k) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} e^{-\frac{1}{2} (X-\mu_k)^T \Sigma_k^{-1} (X-\mu_k)}.$$

Linear discriminant analysis (LDA): $\Sigma_k = \bar{\Sigma}, \forall k.$

$$C(X) = \arg \max_k \left(-\log \left((2\pi)^{p/2} |\bar{\Sigma}|^{1/2} \right) - \frac{1}{2} (X-\mu_k)^T \bar{\Sigma}^{-1} (X-\mu_k) + \log \pi_k \right)$$

$$= \arg \max_k \left(-\frac{1}{2} (X-\mu_k)^T \bar{\Sigma}^{-1} (X-\mu_k) + \log(\pi_k) \right).$$

$$\text{note: } \frac{1}{2} (X-\mu_k)^T \bar{\Sigma}^{-1} (X-\mu_k) = X^T \bar{\Sigma}^{-1} \mu_k - \frac{1}{2} \mu_k^T \bar{\Sigma}^{-1} \mu_k + \frac{1}{2} X^T \bar{\Sigma}^{-1} X$$

$$= \arg \max_k \left(X^T \bar{\Sigma}^{-1} \mu_k - \frac{1}{2} \mu_k^T \bar{\Sigma}^{-1} \mu_k + \log(\pi_k) \right).$$

$$= \arg \min_k \left(\mu_k^T \bar{\Sigma}^{-1} \mu_k - 2 X^T \bar{\Sigma}^{-1} \mu_k \right).$$

↓ uniform prior.
 $\pi_k = \pi, \forall k$

Decision boundary between class l and m :

$$(\mu_k + \mu_l)^T \bar{\Sigma}^{-1} (\mu_k - \mu_l) - 2 X^T \bar{\Sigma}^{-1} (\mu_k - \mu_l) = 0.$$

The decision boundaries are hyper-planes.

Computation: (1) estimate $\hat{\Sigma}$

(LDA) (2). $\hat{\Sigma} = U D U^T$ (eigen-decomposition).

↓
diagonal matrix
with eigen-values.

$$(3). \text{ Sphere the data } D^{-\frac{1}{2}} U^T X \rightarrow X^*$$

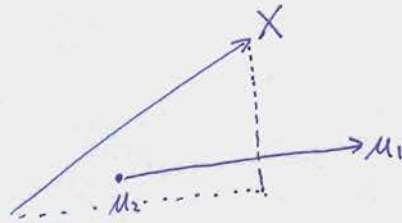
$$D^{-\frac{1}{2}} U^T \mu_k \rightarrow \mu_k^*$$

(4). Classify x^* to the closest centroid μ_k^*

LDA as a dimension reduction method.

- decision boundary: $-\frac{1}{2}(\mu_1 + \mu_2)^T \Sigma^{-1}(\mu_1 - \mu_2) + x^T \Sigma^{-1}(\mu_1 - \mu_2) = 0$.
- only the projection of X on the direction $\Sigma^{-1}(\mu_1 - \mu_2)$ matters.

Suppose $\Sigma = \sigma^2 I$ is sphered. Only the projection of X onto $\mu_1 - \mu_2$ is needed.



- $K=2 \Rightarrow 1D. H_1$
- $K=3 \Rightarrow 2D. H_2$
- $K=4 \Rightarrow 3D. H_3$
- \vdots

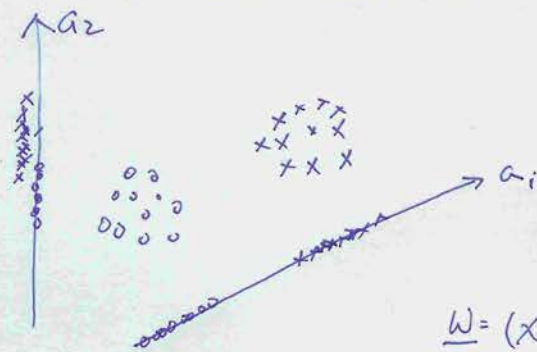
Fisher's Optimization Criteria

— When $K > 3$, we might want to find a subspace $H_L \subseteq H_{K-1}$ optimal for LDA in some sense.

Idea:

$$\max_a \frac{a^T B a}{a^T W a}$$

B : between-class variance
 W : within-class variance. } in original space



$$\begin{aligned} W &= (x - \bar{x})^T (x - \bar{x}) \\ &= (ax - a\bar{x})^T (ax - a\bar{x}) \\ &= a^T (x - \bar{x})(x - \bar{x}) a \\ &= a^T W a \end{aligned}$$

within class variance in projected space.

$$W = V D V^T \Rightarrow W^{\frac{1}{2}} = D^{\frac{1}{2}} V^T$$

Define $b = W^{\frac{1}{2}} a$.

$$\Rightarrow \max_b \frac{b^T (W^{-\frac{1}{2}})^T B W^{-\frac{1}{2}} b}{b^T b}$$

$$\text{Define } B^* = (W^{-\frac{1}{2}})^T B W^{-\frac{1}{2}}$$

$$\text{Eigen-decompose: } B^* = V^* D^* V^{*T}$$

We know maximization is achieved by

$$b_1 = v_1^*, \text{ the first eigen vector of } B^*.$$

Similarly, one can find the next direction. $b_2 = v_2^*$.

$$\text{s.t. } b_2 \perp b_1 \text{ and maximizes } \frac{b_2^T B^* b_2}{b_2^T b_2}.$$

Computation of Fisher Optimization.

— find centroids of all classes and calculate.

between-class covariance matrix \hat{B} .

within-class covariance matrix \hat{W} .

$$— \hat{W} = \hat{V} \hat{D} \hat{V}^T$$

$$— \hat{B}^* = (\hat{W}^{-\frac{1}{2}})^T \hat{B} \hat{W}^{-\frac{1}{2}}$$

$$— \hat{B}^* = \hat{V}^* \hat{D}^* \hat{V}^{*T}$$

$$— a_1 = \hat{W}^{-\frac{1}{2}} v_1^*$$

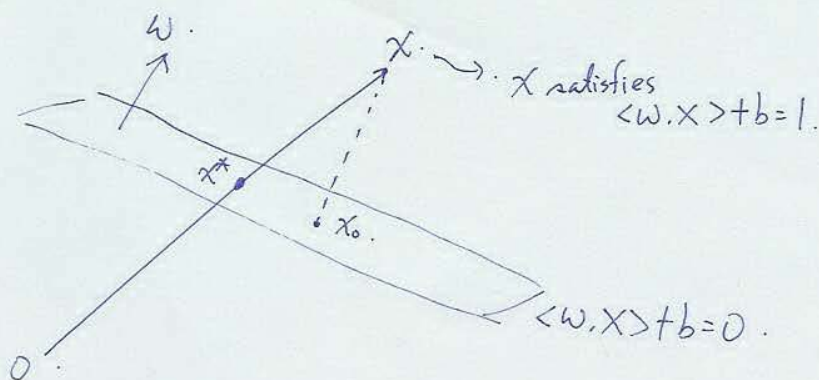
$$a_2 = \hat{W}^{-\frac{1}{2}} v_2^*$$

⋮

Support Vector Machines.

$$x \in \mathbb{R}^p, \quad f(x) = \langle w, x \rangle + b$$

P5



$$x^* = cx \quad \begin{aligned} w^T x + b &= 1 \\ w^T x^* + b &= 0 \end{aligned} \quad \begin{aligned} cw^T x + b &= 0 \\ \Rightarrow w^T x &= \frac{-b}{c} = 1-b \end{aligned}$$

$$\therefore x^* = \frac{-b}{1-b} x$$

$$\begin{aligned} |\overrightarrow{x x_0}| &= \langle \overrightarrow{x x^*}, w \rangle \frac{1}{|w|} \\ &= \langle x - \frac{-b}{1-b} x, w \rangle \frac{1}{|w|} \\ &= \langle \frac{1}{1-b} x, w \rangle \frac{1}{|w|} \\ &= \frac{1}{1-b} \cdot 1-b \cdot \frac{1}{|w|} = \frac{1}{|w|} \end{aligned}$$

maximizing margin:

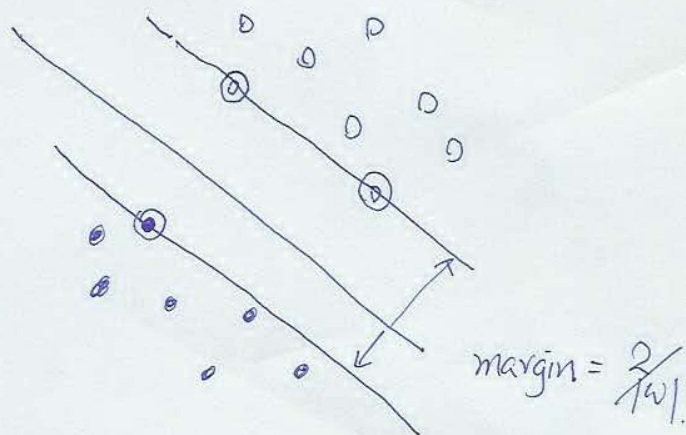
$$\min \frac{1}{2} \|w\|^2$$

under the condition that:

$$\begin{cases} w^T x_i + b \geq 1 & \text{for } x_i = +1 \\ w^T x_i + b \leq -1 & \text{for } x_i = -1 \end{cases}$$

\Downarrow

equivalent to. $x_i [(w^T x_i) + b] - 1 \geq 0, \quad i = 1, 2, \dots, n.$



Kuhn-Tucker Theorem.

Given an optimization problem:

$$L(w, \alpha) = f(w) + \sum_{\bar{i}=1}^k \alpha_{\bar{i}} g_{\bar{i}}(w)$$

$$\min f(w)$$

$$\text{subject to } g_{\bar{i}}(w) \leq 0, \quad \bar{i}=1, 2, \dots, k.$$

Then

$$\frac{\partial L(w^*, \alpha^*)}{\partial w} = 0.$$

$$\frac{\partial L(w^*, \alpha^*)}{\partial \alpha} = 0.$$

$$\alpha_{\bar{i}}^* g_{\bar{i}}(w^*) = 0, \quad \bar{i}=1, \dots, k$$

$$g_{\bar{i}}(w^*) \leq 0, \quad \bar{i}=1, \dots, k.$$

$$\alpha_{\bar{i}}^* \geq 0, \quad \bar{i}=1, \dots, k.$$

$$\min \frac{1}{2} \langle w, w \rangle.$$

$$\text{subject to } y_{\bar{i}} (\langle w, x_{\bar{i}} \rangle + b) \geq 1.$$

$$L(w, b, \alpha) = \frac{1}{2} \langle w, w \rangle - \sum_{\bar{i}=1}^n \alpha_{\bar{i}} [y_{\bar{i}} (\langle w, x_{\bar{i}} \rangle + b) - 1].$$

where $\alpha_{\bar{i}} \geq 0$ are the Lagrange multipliers.

$$\frac{\partial L(w^*, b^*, \alpha^*)}{\partial w} = w^* - \sum_{\bar{i}=1}^n y_{\bar{i}} \alpha_{\bar{i}}^* x_{\bar{i}} = 0. \quad \Rightarrow w^* = \sum_{\bar{i}=1}^n y_{\bar{i}} \alpha_{\bar{i}}^* x_{\bar{i}}$$

$$\frac{\partial L(w^*, b^*, \alpha^*)}{\partial b} = \sum_{\bar{i}=1}^n y_{\bar{i}} \alpha_{\bar{i}}^* = 0.$$

$$\alpha_{\bar{i}}^* [y_{\bar{i}} (\langle w^*, x_{\bar{i}} \rangle + b^*) - 1] = 0. \quad (1). \quad \left(\begin{array}{l} y_{\bar{i}} (\langle w^*, x_{\bar{i}} \rangle + b^*) \geq 1 \\ \Rightarrow \alpha_{\bar{i}}^* = 0 \end{array} \right)$$

classifier

$$f(x) = \langle w^*, x \rangle + b^*$$

$$= \sum_{\bar{i}=1}^n y_{\bar{i}} \alpha_{\bar{i}}^* \langle x_{\bar{i}}, x \rangle + b^*$$

$$= \sum_{\bar{i} \in SV} y_{\bar{i}} \alpha_{\bar{i}}^* \langle x_{\bar{i}}, x \rangle + b^*$$

(from (1), only the support vectors have non-zero $\alpha_{\bar{i}}^*$)