

# Pattern Classification

All materials in these slides were taken from  
Pattern Classification (2nd ed) by R. O.  
Duda, P. E. Hart and D. G. Stork, John Wiley  
& Sons, 2000  
with the permission of the authors and the  
publisher

## Chapter 2 (part 3)

## Bayesian Decision Theory (Sections 2-6, 2-9)

- Discriminant Functions for the Normal Density
- Bayes Decision Theory – Discrete Features

Quiz on coming thursday

Syllabus - lectures till coming monday

# Discriminant Functions for the Normal Density

- We saw that the minimum error-rate classification can be achieved by the discriminant function

$$g_i(x) = \ln P(x | \omega_i) + \ln P(\omega_i)$$

- Case of multivariate normal

$$g_i(x) = -\frac{1}{2}(x - \mu_i)^t \sum_i^{-1} (x - \mu_i) - \frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$



$$*) \sum_i = \sigma^2 I$$

$$g_i(x) = w^T x + b, \quad w = \frac{M_i}{\sigma^2} \quad b = -\frac{M_i^T M_i}{2\sigma^2}$$

$$*) g_1(x) = g_2(x) \rightarrow \text{for two categories}$$

$$*) \sum_i = \sum \rightarrow \text{not in the form identity}$$

$$g_i(x) = -\frac{1}{2} (x - M_i)^T \sum_i^{-1} (x - M_i) + \ln P(w_i)$$

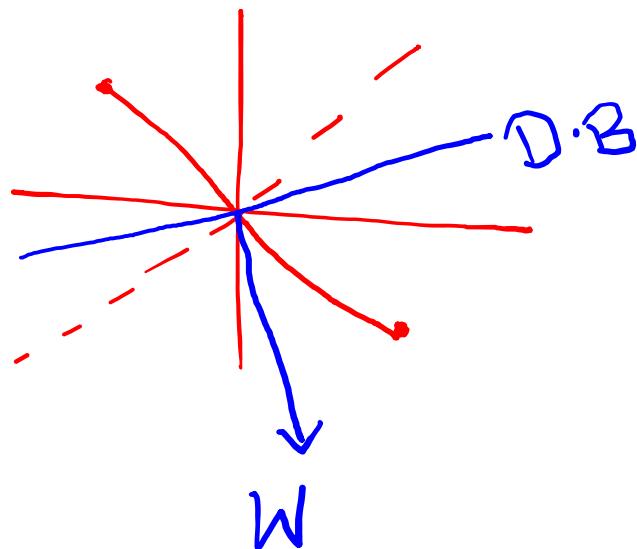
$$= -\frac{1}{2} \left\{ x^T \sum_i^{-1} x - 2 M_i^T \sum_i^{-1} x + M_i^T \sum_i^{-1} M_i + \ln P(w_i) \right\}$$

$$= M_i^T \sum_i^{-1} x - \frac{1}{2} M_i^T \sum_i^{-1} M_i + \ln P(w_i)$$

$$= w^T x + b \quad w = \sum_i^{-1} M_i \quad b = -\frac{1}{2} M_i^T \sum_i^{-1} M_i + \ln P(w_i)$$

$$D \cdot B \quad g_i(x) = g_j(x) \equiv w^T x + b$$

$$w = \sum^{-1} (m_i - m_j)$$



\*  $\Sigma_i \rightarrow$  arbitrary and different for all classes.

- Case  $\Sigma_i = \sigma^2 I$  (I stands for the identity matrix)

$$g_i(x) = w_i^t x + w_{i0} \quad (\text{linear discriminant function})$$

*where :*

$$w_i = \frac{\mu_i}{\sigma^2}; \quad w_{i0} = -\frac{1}{2\sigma^2} \mu_i^t \mu_i + \ln P(\omega_i)$$

*( $\omega_{i0}$  is called the threshold for the  $i$ th category! )*

- A classifier that uses linear discriminant functions is called “**a linear machine**”
- The decision surfaces for a linear machine are pieces of **hyperplanes** defined by:

$$g_i(x) = g_j(x)$$

- Case  $\Sigma_i = \text{arbitrary}$

- The covariance matrices are different for each category

$$g_i(x) = x^t W_i x + w_i^t x + w_{i0}$$

*where :*

$$W_i = -\frac{1}{2} \Sigma_i^{-1}$$

$$w_i = \Sigma_i^{-1} \mu_i$$

$$w_{i0} = -\frac{1}{2} \mu_i^t \Sigma_i^{-1} \mu_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)$$

(Hyperquadrics which are: hyperplanes, pairs of hyperplanes, hyperspheres, hyperellipsoids, hyperparaboloids, hyperhyperboloids)

$$g_i(\mathbf{x}) = \mathbf{x}^t \mathbf{W}_i \mathbf{x} + \mathbf{w}_i^t \mathbf{x} + w_{i0},$$

$i=1, g_1(x_0)$

$i=2, g_2(x_0)$

$$\mathbf{W}_i = -\frac{1}{2} \boldsymbol{\Sigma}_i^{-1}, \quad \text{equal}$$

$$\mathbf{w}_i = \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i$$

or  
 $\sum_i = \sigma^2 I$

$$w_{i0} = -\frac{1}{2} \boldsymbol{\mu}_i^t \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \frac{1}{2} \ln |\boldsymbol{\Sigma}_i| + \ln P(\omega_i).$$

Assume two different classes -  $w_1$  &  $w_2$

Both having multi-variate Gaussian PDF.

$$P(x|w_1) \sim N(\mu_1, \Sigma_1) \quad x \in \mathbb{R}^d$$

$$P(x|w_2) \sim N(\mu_2, \Sigma_2)$$

$$\mu_1 \neq \mu_2, \Sigma_1 \neq \Sigma_2$$

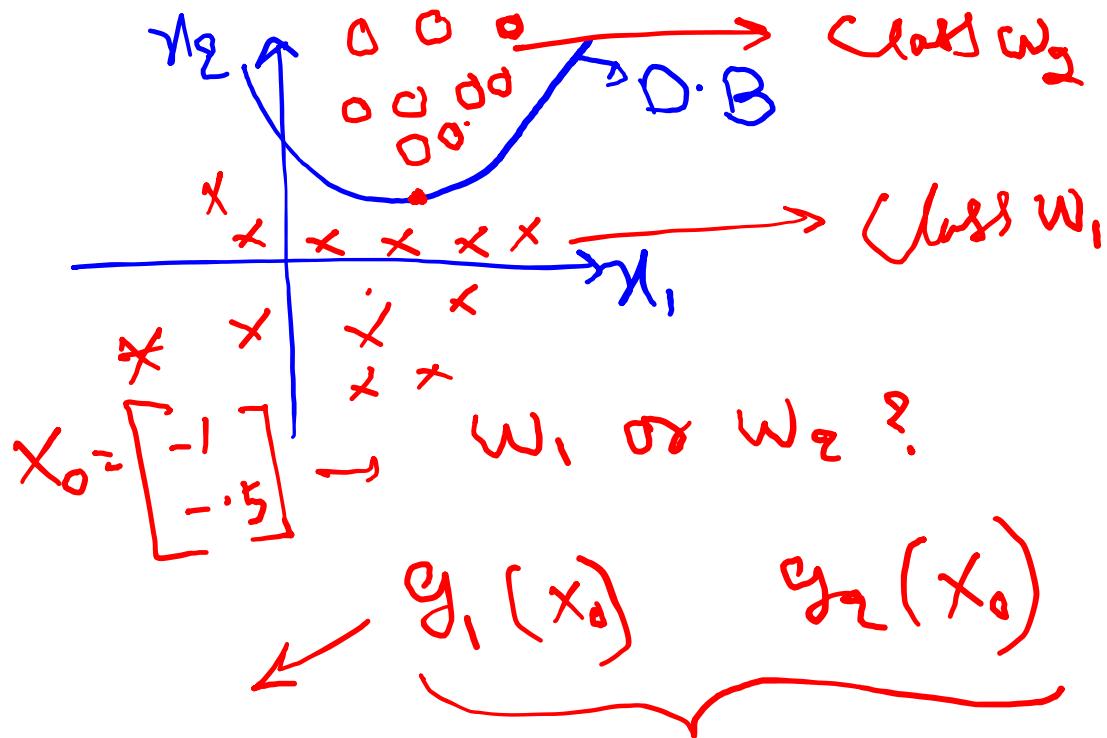
$$\boldsymbol{\mu}_1 = \begin{bmatrix} 3 \\ 6 \end{bmatrix}; \quad \boldsymbol{\Sigma}_1 = \begin{pmatrix} 1/2 & 0 \\ 0 & 2 \end{pmatrix} \text{ and } \boldsymbol{\mu}_2 = \begin{bmatrix} 3 \\ -2 \end{bmatrix}; \quad \boldsymbol{\Sigma}_2 = \begin{pmatrix} 2 & 0 \\ 0 & 2 \end{pmatrix}.$$

The inverse matrices are then,

$$\boldsymbol{\Sigma}_1^{-1} = \begin{pmatrix} 2 & 0 \\ 0 & 1/2 \end{pmatrix} \quad \text{and} \quad \boldsymbol{\Sigma}_2^{-1} = \begin{pmatrix} 1/2 & 0 \\ 0 & 1/2 \end{pmatrix}.$$

We assume equal prior probabilities,  $P(\omega_1) = P(\omega_2) = 0.5$ , and substitute these into the general form for a discriminant, Eqs. 64 – 67, setting  $g_1(\mathbf{x}) = g_2(\mathbf{x})$  to obtain the decision boundary:

$$x_2 = 3.514 - 1.125x_1 + 0.1875x_1^2. \rightarrow \textcircled{1} \cdot \textcircled{B}.$$



$$d = 10^6, 10^4 \quad g_1(x) = g_2(x)$$

$W_1 \rightarrow$  Take  $8ML \rightarrow$

$W_2 \rightarrow$  Don't Take  $8ML \rightarrow$

$\times$  Graded, I touched it, has changed.  $d=4 \rightarrow$  seek opinions b/w  
 $x_1 \quad x_2 \quad x_3 \quad x_4$  5th.

## Demo

Use LDA when cov matrices are not same?

$$\sum_{\phi} = \sum$$

$$\sum = \frac{n_1 \sum_1 + n_2 \sum_2 + \cdots + n_c \sum_c}{n_1 + n_2 + \cdots + n_c}$$

$n_1 \rightarrow$  no. of samples in class 1

:

$n_c \rightarrow$

..

:

Class c

$X \in R$   $\rightarrow$  Gaussian

## Bayes Decision Theory – Discrete Features

$$X \in \{v_1, v_2, \dots, v_m\}$$

- Components of  $x$  are binary or integer valued,  $x$  can take only one of  $m$  discrete values

$$v_1, v_2, \dots, v_m$$

$$\begin{aligned} X &\in \{v_1, v_2\} \\ X &\in \{0, 1\} \end{aligned}$$

- Case of independent binary features in 2 category problem

Let  $x = [x_1, x_2, \dots, x_d]^t$  where each  $x_i$  is either 0 or 1, with probabilities:

$$p_i = P(x_i = 1 | \omega_1)$$

$$q_i = P(x_i = 1 | \omega_2)$$

$\omega_1 \rightarrow$  Auto function.

$\omega_2 \rightarrow$  culto non function

$x_1 \rightarrow$  denotes gear

$x_2 \rightarrow$  denotes brake

:

$x_d \rightarrow$  Accelerator.

Category / classed.

$$P_i = P_{\delta}(\gamma_{i\delta}=1 | w_i)$$

$$1 - P_i = P_{\delta}(\gamma_{i\delta}=0 | w_i)$$

$$\cancel{P(x_i)} = P_i^{\gamma_{i\delta}} * (1-P_i)^{1-\gamma_{i\delta}} \checkmark$$

$$P(\gamma_{i\delta} | w_i) = P_i \delta(\gamma_{i\delta}) + (1-P_i) \delta(1-\gamma_{i\delta}) \times$$

$$= P_i \delta(1-\gamma_{i\delta}) + (1-P_i) \delta(\gamma_{i\delta}) \checkmark$$

$\delta \rightarrow$  didicates:  $\delta(0) = 1$

$$\delta(\neq 0) = 0$$

$$P(x_i | w_i) = \beta_i^{x_i} (1 - \beta_i)^{1 - x_i}$$

$$P(x | w_i) = P\left(x = [x_1, x_2 \dots x_i, \dots x_d]^\top | w_i\right)$$
$$x = [x_1, x_2 \dots x_d]^\top$$

\* Assumption  $\rightarrow$  Conditional independence.

$$P(x_1, x_2 | w) = P(x_1 | w) P(x_2 | w)$$

$$P(x | w_i) = \prod_{j=1}^d P(x_j | w_i)$$
$$= \prod_{j=1}^d \beta_j^{x_j} (1 - \beta_j)^{1 - x_j}$$

If the classifier assumes this type of  
Conditional independence.  $\rightarrow$  Naive Bayes class.

$$P(X|w_2) = \prod_{i=1}^d q_i^{x_i} (1-q_i)^{1-x_i}$$

We have ~~total~~ likelihood for both classes.

$$g_1(x) = ? = P(w_2) P(X|w_2) + P(w_1) P(X|w_1)$$

↓ discriminant:

$$\ln P(X|w_1) \\ + \ln P(w_1)$$

$$g_2(x) = \underline{-\ln P(w_2)} \\ \ln P(w_2) \rightarrow$$

$$g_1(x) = \ln \prod_{j=1}^N p_j^{x_j} (1-p_j)^{1-x_j} + \ln p(w)$$

$$= \sum_{j=1}^N x_j \ln p_j + (1-x_j) \ln (1-p_j) + \ln p(w)$$

$g_2(x)$

$$g_1(x) = g_2(x) \rightarrow Q \cdot \beta$$

Usually for two category class.

We define:  $\underbrace{g(x)}_{\text{discriminizer}} = g_1(x) - g_2(x)$

$w$

- The discriminant function in this case is:

$$g(x) = \sum_{i=1}^d w_i x_i + w_0$$

*where :*

$$w_i = \ln \frac{p_i(1-q_i)}{q_i(1-p_i)} \quad i = 1, \dots, d$$

*and :*

$$w_0 = \sum_{i=1}^d \ln \frac{1-p_i}{1-q_i} + \ln \frac{P(\omega_1)}{P(\omega_2)}$$

*decide  $\omega_1$  if  $g(x) > 0$  and  $\omega_2$  if  $g(x) \leq 0$*

Ex.

- Suppose two categories consist of independent binary features in three dimensions with known feature probabilities. Let us construct the Bayesian decision boundary if  $P(\omega_1) = P(\omega_2) = 0.5$  and the individual components obey:
- $p_i = 0.8$  and  $q_i = 0.5$  for  $i = 1, 2, 3$ .
- $w_i = \ln .8(1 - .5)/.5(1 - .8) = 1.3863$

$$\begin{aligned}g(x) &= g_1(x) - g_2(x) \\&= 0 \rightarrow \text{D.B.}\end{aligned}$$

$$P_i = P_{\omega_1}(x_i=1|\omega_1)$$

$$q_i = P_{\omega_2}(x_i=1|\omega_2)$$

$$w_i = \ln \frac{0.8 \cdot 0.5}{0.5 \cdot 0.2} = \ln 4 = w_2 = w_3$$

$$w_0 = \sum_{i=1}^3 \ln \frac{2}{5} .$$

$$g(x) = 0$$

$$\left[ \sum_{i=1}^3 w_i x_i + w_0 = 0 \right]$$

$$w_1 (x_1 + x_2 + x_3) + 3 \ln \frac{2}{5} = 0$$

$$ax + by + cz + d = 0$$

Find the multiclass discriminant corresponding to  
min probability of error for independent binary  
valued features.

$$\omega_1, \omega_2, \dots, \omega_c$$

$$g_i(x) \rightarrow \forall i=1, \dots, c$$

$$P(x|\omega_i) = \prod_{j=1}^d \rho_i^{x_j} (1-\rho_i)^{1-x_j}$$

$$\rho_{ij} = P(x_i=1 | \omega_j)$$

$$\rho_{i1} = P(x_i=1 | \omega_1) \quad 1 - \rho_{i1} = 1 - P(x_i=1 | \omega_1)$$

$$P(x|w_j) = \prod_{i=1}^d P_{ij}^{x_i} (1-P_{ij})^{1-x_i}$$

$$P(x_i|w_j) = P_{ij}^{x_i} (1-P_{ij})^{1-x_i}$$

\*  $g_j(x) = \ln P(x|w_j) + \ln P(w_j)$   
 $\forall j = 1, 2, \dots <$

$\max_j g_j(x) \rightarrow$  map  $x$  to  $j$