

# Discriminant Analysis

## (1) Background:

For  $m$  population  $G_i \sim F_i(x)$ . or we have samples from it. If we're given another sample  $x$  to be tested. How can we discriminate which one population  $x$  comes from?

## (2) Distance:

The ideal of discrimination by distance is: choose  $G_i$  where  $i = \arg \min d(x, G_i)$  and regard  $x$  is from  $G_i$  which is the nearest.

### ① Mahalanobis Distance:

Def: If  $\vec{x}, \vec{y}$  are samples from population  $G$  whose mean is  $\vec{\mu}$  and covariance  $\Sigma = (\sigma_{ij})_{exp}$

Define:  $d_m(\vec{x}, \vec{y}) = d(\vec{x}, \vec{y}) = (\vec{x} - \vec{y})^T \Sigma^{-1} (\vec{x} - \vec{y})$

$d_m(\vec{x}, G_i) = d(\vec{x}, G_i) = (\vec{x} - \vec{\mu}^i)^T \Sigma^{-1} (\vec{x} - \vec{\mu}^i)$

where  $\mu^i$  is mean of  $G_i$ .

Remark:  $d_m$  normalizes the variables so that it's indept with unit of measure.



## ② Discrimination:

i)  $m^i$ 's are known.  $\Sigma^i = \Sigma$  known:

Rule: 
$$\begin{cases} \eta \in h_1 & \text{if } \lambda^2(\eta, h_1) < \lambda^2(\eta, h_2) \\ \eta \in h_2 & \text{if } \lambda^2(\eta, h_2) < \lambda^2(\eta, h_1) \\ \text{undetermined} & \text{otherwise.} \end{cases}$$

Note that:  $\lambda^2(\eta, h_1) - \lambda^2(\eta, h_2)$

$$= 2 \left( \eta - \frac{m^1 + m^2}{2} \right)^T \Sigma^{-1} (m^1 - m^2)$$

Denote:  $\alpha = \Sigma^{-1} (m^1 - m^2)$ ,  $\bar{m} = (m^1 + m^2)/2$  ( $\vec{\alpha}$  is projection

$$\Rightarrow \lambda^2(\eta, h_1) - \lambda^2(\eta, h_2) = 2 \alpha^T (\eta - \bar{m}) \quad \text{direction!})$$

ii)  $m^i, \Sigma^i$  known:

$$\lambda^2(\eta, h_1) - \lambda^2(\eta, h_2) = (\eta - m^1)^T \Sigma_1^{-1} (\eta - m^1) - (\eta - m^2)^T \Sigma_2^{-1} (\eta - m^2)$$

iii)  $m^i, \Sigma^i$  unknown:

Consider we have samples  $X_1 = (x_{ij}^{(1)})_{n \times p}$

and  $X_2 = (x_{ij}^{(2)})_{n \times p}$  from  $h_1, h_2$ .

Denote:  $\bar{X}_i^{(t)} = \frac{\sum_{k=1}^{n_t} x_{ki}^{(t)}}{n_t}$ ,  $t=1,2$

$$S_{ij}^{(t)} = \frac{1}{n_t - 1} \sum_k (x_{ki}^{(t)} - \bar{x}_i^{(t)}) (x_{kj}^{(t)} - \bar{x}_j^{(t)})$$

$$S^{(t)} = (S_{ij}^{(t)})_{p \times p}, \quad \bar{X}^{(t)} = (\bar{x}_1^{(t)} \dots \bar{x}_p^{(t)})$$

Def:  $\lambda^2(x, h^{(t)}) = (x - \bar{X}^{(t)})^T S^{(t)-1} (x - \bar{X}^{(t)})$

$$\Rightarrow \text{Find } i = \arg \min_k \lambda^2(x, h_k)$$



### (3) Bayesian:

- Distance Discriminant doesn't take prior prob. of occurrence in  $G_i$  into account and consider the loss of incorrect discrimination.

$\Rightarrow$  Suppose we have  $m$  populations  $G_i \sim f_i(x)$  whose occurrence prob is  $\pi_i$  based on the past data.

#### ① Maximum Posterior:

Note that  $P(G_i | x) = \pi_i f_i(x) / \sum_{k=1}^m \pi_k f_k(x)$

Find  $l = \operatorname{argmax}_k P(G_k | x) \Rightarrow x \in G_l$ .

Rmk: When  $f_i(x) \sim N(\mu^i, \Sigma^i)$ .

Since  $\max_k P(G_k | x) \Leftrightarrow \max_k \pi_k f_k(x)$

$$\pi_k f_k(x) = \frac{\pi_k}{(2\pi)^{\frac{p}{2}} |\Sigma_k|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} (x - \mu^k)^T \Sigma_k^{-1} (x - \mu^k)\right)$$

#### ② Minimum ECM:

Suppose  $D = \cup D_i$  a partition, if  $x$  falls into

$D_i$ . Then regard  $x \in G_i$ .

Next, find a partition  $\cup D_i$  st. minimize the loss of incorrect discrimination.

Def:  $p_{cjl}(i) = p(x \in D_j | G_i) = \int_{D_j} f_i(x) dx$ .

$L_{cjl}(i)$  is loss coefficient, st.  $L_{cjl}(i) = 0$ .



Expected Cost of Misclassification of  $D$ :

$$E(CM(D)) = \sum_{i=1}^m \alpha_i \left( \sum_{j=1}^m L(c_j|i) p(c_j|i) \right)$$

We call  $D^* = \arg \min_D E(CM(D))$ , the partition is a Bayesian Rule.

Thm. Denote  $h_j(x) = \sum_{i=1}^m \alpha_i L(c_j|i) f_i(x)$ . Then:  
the Bayesian Rule is  $D_i^* = \{x | h_i(x) \leq h_j(x), \forall j \neq i\}$

Pf: If  $D = \bigcup_i D_i$  is another partition.

$$\Rightarrow E(CM(D^*)) - E(CM(D))$$

$$= \sum_{i=1}^m \sum_{j=1}^m \int_{D_i^* \cap D_j} [h_i(x) - h_j(x)] dx \leq 0$$

Remark: i) It means: find  $i = \arg \min_i h_i(x)$ .

$\Rightarrow$  Discriminate  $x \in C_i$ .

ii) Note  $h_i(x) = \sum_{k=1}^m \alpha_k f_k(x) - \alpha_i f_i(x)$

$\Leftrightarrow$  maximize  $\alpha_i f_i(x)$ . It's eqn.

with criteria ①, if  $L(c_j|i) = 1 - \delta_{ij}$

#### (4) Fisher Discrimination:

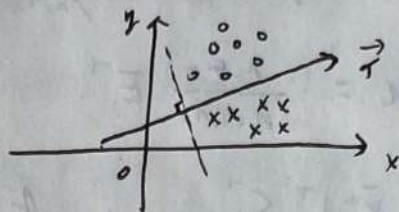
The ideal is projecting the data into several direction to separate each class significantly.

Then choose a criteria to classify.



Def:

We want to project the data to some direction which can separate different kinds of data as most as possible.



Link: As we did before. When

$\Sigma^1 = \Sigma^2 = \Sigma$ . Then To minimize  $ECD$ :

We classify  $D$ , if  $(m_1^T - m_2^T) \Sigma^{-1} x - k \geq \ln \frac{L(1)}{L(2)} \frac{p_2}{p_1}$

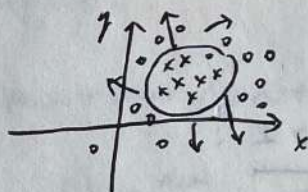
It projects to  $(m_1 - m_2)^T \Sigma^{-1}$ . (The direction)

But when  $\Sigma^1 \neq \Sigma^2$ . The allocation is:

$D$ , if  $-\frac{1}{2} x^T (\Sigma_1^{-1} - \Sigma_2^{-1}) x + m_1^T \Sigma_1^{-1} x - m_2^T \Sigma_2^{-1} x \geq k + \ln \square$

It's quadratic form:

(i.e. separate by ellipse)



proceed: Set  $Y = C^T X$ . st. minimize the Variance in the same group. and maximize Variance between different groups. i.e.

$$\text{Set } E_0 = \sum_{i=1}^m V_{y_i}^i = C^T \sum_{i=1}^m V_x^i C = C^T E C$$

$$\text{where } E = \sum V_x^i = \sum_{i=1}^m (X^{(i)} - \bar{X}^{(i)}) (X^{(i)} - \bar{X}^{(i)})^T$$

$$B_0 = \sum_{i=1}^m n_i (\bar{Y}^{(i)} - \bar{Y})^2 = C^T \sum (\bar{X}^{(i)} - \bar{X}) (\bar{X}^{(i)} - \bar{X})^T n_i C$$

$$= C^T B C$$

$$\text{maximize: } A^2(C) = \frac{C^T B C}{C^T E C} \quad \text{i.e. } \max = \max \{ \lambda \mid \lambda \in \sigma_{E^{-1}B} \}$$

Link: If number of group is large. Then we will find  $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_k$ . eigenvalues of  $E^{-1}B$ .



i) When  $m=2$ :

$$\text{Then } B = \frac{n_1 n_2}{n_1 + n_2} (\bar{X}^1 - \bar{X}^2) (\bar{X}^1 - \bar{X}^2)^T, \text{rank}(B) = 1$$

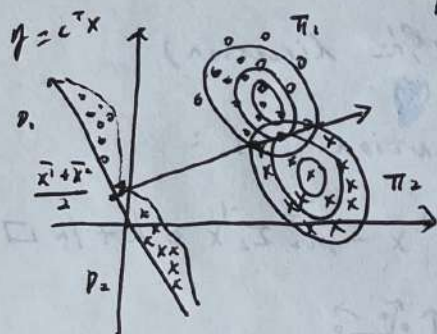
$$E^{-1}B \text{ has max eigenvalue: } \frac{1}{\frac{1}{n_1} + \frac{1}{n_2}} (\bar{X}^1 - \bar{X}^2)^T E^{-1} (\bar{X}^1 - \bar{X}^2)$$

$$\text{Denote } S_{pooled}^{-1} = \frac{n_1 n_2}{n_1 + n_2} E^{-1}.$$

$$\Rightarrow C^T = (\bar{X}^1 - \bar{X}^2)^T S_{pooled}^{-1}, \quad \eta = C^T X. \text{ Suppose } \bar{X}^1 > \bar{X}^2.$$

Allocate  $X$  to  $D_1$  if  $\eta > \frac{1}{2} C^T (\bar{X}^1 + \bar{X}^2)$

$D_2$  if  $\eta < \frac{1}{2} C^T (\bar{X}^1 + \bar{X}^2)$



Rmk: When  $\eta = \frac{1}{2} C^T (\bar{X}^1 + \bar{X}^2)$

It doesn't work.

ii) When  $m > 2$ :

$E^{-1}B$  may have different eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots$

suppose each corresponds  $C_i$ , the project direction

$$\text{Criteria Functions: } U_i(X) = C_i^T X, \quad \bar{U}_i^t = C_i^T \bar{X}^t$$

If  $\exists$  unique  $i_0$  st.

$$i_0 = \arg \min_t |U_i(X) - \bar{U}_i^t| / \hat{\sigma}_t^2, \quad \hat{\sigma}_t^2 = C_i^T S_t C_i$$

Then allocate  $X$  to  $\pi_{i_0}$ .

Otherwise, consider using  $U_2(X), U_3(X), \dots$

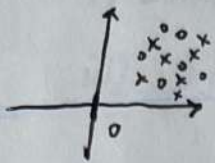
until exists a unique  $\pi_{i_0}$ .

Rmk: We may choose  $\{U_i(X)\}_i$  st.  $\frac{\sum_i \lambda_i}{\sum_i \lambda_i} > p_0$ .

$p_0$  is the expected efficient.



① Classification sometimes may not be a good idea for the data:



So, before separating the data.

We can apply  $T^2$ 's test:

$$H_0: \mu_1 = \mu_2 \quad \text{V.S.} \quad H_1: \mu_1 \neq \mu_2$$

Remark: Significant separation  $\Rightarrow$  good Classification

② Evaluating classification function:

i) Total prob of misclassification  $TPM = \sum P_i \int_{R/D_i} f_i(x) dx$

ii) Actual error rate  $AER = \sum P_i \int_{R/\hat{D}_i} \hat{f}_i(x) dx$ .  $\hat{f}_i = f_i(x|\hat{\theta})$   
 $\hat{D}_i$  is separation based on samples.

iii) Apparent error rate  $APER = n_{1m} + n_{2m} / (n_1 + n_2)$

	$\pi_1$	$\pi_2$	
Actual $\pi_1$	$n_{1c}$	$n_{1m}$	$= n_1$
Actual $\pi_2$	$n_{2c}$	$n_{2m}$	$= n_2$
	prediction		

$$= \frac{n_1}{n_1+n_2} \cdot \frac{n_{1m}}{n_1} + \frac{n_2}{n_1+n_2} \cdot \frac{n_{2m}}{n_2}$$

$\hookrightarrow n_i / (n_1 + n_2)$  can estimate the rate of class  $\pi_i$  occurs

Remark: i) need the information of populations ( $P_i, f_i, \dots$ )

ii) need the information from samples and pdf's.

iii) Apparent depends on density. But it may under-estimate the error rate. Since it totally bases on the data — which separates them as most.

$\Rightarrow$  Cross-validation method: (Leave one out)

Start with  $\pi_1$  group. omit one observation.  $\Rightarrow$  Develop classification on remain  $n_1 - 1 + n_2$  ob's.  $\Rightarrow$  Classify the

omit ob.  $\xRightarrow{\text{rep. on } \pi_1}$  Calculate misclassification  $n_{1m}^{(1)} / n_1 \Rightarrow$  On  $\pi_2$   $\xRightarrow{\text{rep. on } \pi_2}$  Calculate misclassification  $n_{2m}^{(1)} / n_2 \Rightarrow \hat{E}(ARE) = \frac{n_{1m}^{(1)} + n_{2m}^{(1)}}{n_1 + n_2}$