

Discrimination amongst several populations

We want to determine if an observation vector

$$\mathbf{X} = x_1, \dots, x_p$$

comes from one of the g populations:

$$\pi_1 : f_1(x_1, \dots, x_p)$$

...

$$\pi_g : f_g(x_1, \dots, x_p)$$

Usually, the densities f_1, \dots, f_g will be assumed to be multivariate normal.

For this purpose we need to partition p -dimensional space into g regions R_1, R_2, \dots, R_g

We will make the decision $D_i = \{X \text{ came from } \pi_i\}$
if X belongs to R_i

Misclassification probabilities

$$\begin{aligned} P(k|i) &= P[\text{classify the case in } \pi_k \text{ when case is from } \pi_i] \\ &= P(X \in R_k | \pi_i) = \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} \end{aligned}$$

Cost of Misclassification

$c_{k|i}$ = Cost classifying the case in π_k when case is from π_i

Prior probabilities of inclusion

$P(i) = P[\text{classify the case is from } \pi_i \text{ initially}]$

Expected Cost of Misclassification of a case from population i

We assume that we know the case came from π_i

$$\begin{aligned} ECM(i) = & c_{1|i}P[1|i] + \dots + c_{i-1|i}P[i-1|i] + c_{i+1|i}P[i+1|i] + \\ & + \dots + c_{k|i}P[g|i] \end{aligned}$$

$$= \sum_{j \neq i} c_{j|i}P[j|i]$$

Total Expected Cost of Misclassification

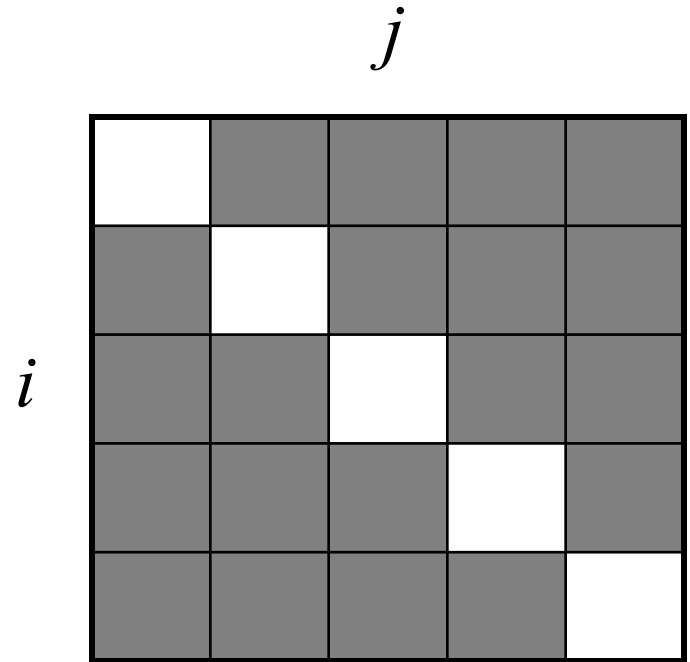
$$ECM = P[1]ECM(1) + P[2]ECM(2) + \dots + P[g]ECM(g)$$

$$= \sum_i P[i] ECM(i)$$

$$= \sum_i P[i] \sum_{j \neq i} c_{j|i} P[j|i]$$

$$= \sum_i P[i] \sum_{j \neq i} c_{j|i} \int_{C_j} f_i(\vec{x}) d\vec{x}$$

$$= \sum_j \int_{C_j} \sum_{i \neq j} P[i] f_i(\vec{x}) c_{j|i} d\vec{x}$$



Optimal Classification Rule

The optimal classification rule will find the regions C_j that will minimize:

$$\begin{aligned} ECM &= \sum_j \int_{C_j} \sum_{i \neq j} P[i] f_i(\vec{x}) c_{j|i} d\vec{x} \\ &= c \sum_j \int_{C_j} \sum_{i \neq j} P[i] f_i(\vec{x}) d\vec{x} \quad \text{if } c_{j|i} = c \\ &= c \sum_j \int_{C_j} \left[\sum_{i=1}^k P[i] f_i(\vec{x}) - P[j] f_j(\vec{x}) \right] d\vec{x} \end{aligned}$$

ECM will be minimized if C_j is chosen where the term that is omitted:

$$P[j] f_j(\vec{x})$$

is the largest

Optimal Regions when misclassification costs are equal

Allocate X to π_k if

$$P[k]f_k(X) > P[i]f_i(X) \text{ for all } i \neq k$$

Or, equivalently if

$$\ln[P[k]f_k(X)] > \ln[P[i]f_i(X)] \text{ for all } i \neq k$$

Optimal Regions when misclassification costs are equal and distributions are p -variate Normal with common covariance matrix Σ

$$R_k = \{\mathbf{X}: P[k]f_k(\mathbf{X}) > P[i]f_i(\mathbf{X}) \text{ for all } i \neq k\}$$

$$= \{\mathbf{X}: \ln[P[k]f_k(\mathbf{X})] > \ln[P[i]f_i(\mathbf{X})] \text{ for all } i \neq k\}$$

In the case of normality

$$f_i(\vec{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} e^{-\frac{1}{2}(\vec{x} - \vec{\mu}_i)' \Sigma^{-1} (\vec{x} - \vec{\mu}_i)}$$

$$\ln P[i] f_i(\vec{x})$$

$$= \ln P[i] - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\vec{x} - \vec{\mu}_i)' \Sigma^{-1} (\vec{x} - \vec{\mu}_i)$$

and $\ln P[j] f_j(\vec{x}) > \ln P[i] f_i(\vec{x})$ if:

$$\begin{aligned} \ln P[j] - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\vec{x} - \vec{\mu}_j)' \Sigma^{-1} (\vec{x} - \vec{\mu}_j) \\ > \ln P[i] - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\vec{x} - \vec{\mu}_i)' \Sigma^{-1} (\vec{x} - \vec{\mu}_i) \end{aligned}$$

that is

$$\vec{\mu}_j' \Sigma^{-1} \vec{x} - \frac{1}{2} \vec{\mu}_j' \Sigma^{-1} \vec{\mu}_j + \ln P[j] > \vec{\mu}_i' \Sigma^{-1} \vec{x} - \frac{1}{2} \vec{\mu}_i' \Sigma^{-1} \vec{\mu}_i + \ln P[i]$$

or $\vec{a}_j' \vec{x} + b_j > \vec{a}_i' \vec{x} + b_i$

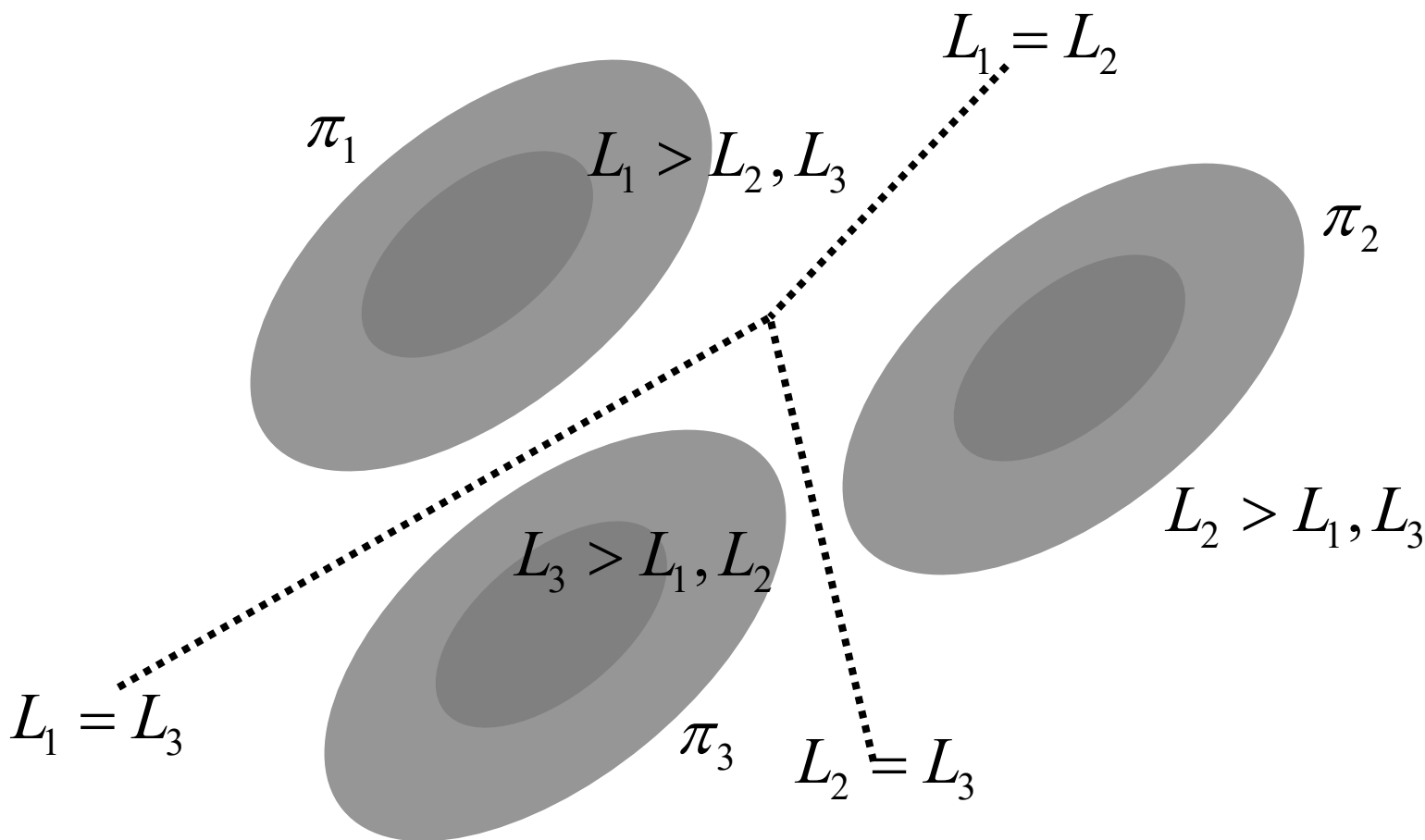
where $\vec{a}_i = \Sigma^{-1} \vec{\mu}_i$ and $b_i = \ln P[i] - \frac{1}{2} \vec{\mu}_i' \Sigma^{-1} \vec{\mu}_i$

Summarizing

We will classify the observation vector in population π_j

if: $L_j = \vec{a}_j' \vec{x} + b_j = \max_i L_i = \max_i (\vec{a}_i' \vec{x} + b_i)$

where $\vec{a}_i = \Sigma^{-1} \vec{\mu}_i$ and $b_i = \ln P[i] - \frac{1}{2} \vec{\mu}_i' \Sigma^{-1} \vec{\mu}_i$



Classification with Normal Populations with Different Covariance Matrices

Let

$$f_i(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\boldsymbol{\Sigma}_i|^{1/2}} e^{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)' \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)}, i = 1, \dots, g$$

Assume further that $c(i|i) = 0$ and $c(k|i) = 0$ for $k \neq i$

Rule: Allocate \mathbf{x} to π_k if:

$$\begin{aligned} & \ln P[k] f_k(\mathbf{x}) \\ &= \ln P[k] - \left(\frac{p}{2}\right) \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \\ &= \max_i \ln P[i] f_i(\mathbf{x}) \end{aligned}$$