



Discrimination and classification (I)

- Objectives
- Two populations
- Minimum expected cost of misclassification rule
- Minimum total probability of misclassification rule
- Two multivariate normal populations
 - Equal population variances
 - Unequal population variances



Objectives

- Assume the existence of clusters which are known a priori.
- Construct discrimination rules to allocate new observations into these known groups.



Two populations

Two populations: π_1 and π_2

Classification measurements $\mathbf{x} = (x_1, x_2, \dots, x_p)'$

An object with measurements \mathbf{x} , must be assigned into either π_1 or π_2

Let Ω be the sample space, and

R_1 = set of \mathbf{x} for which the object is being classified into π_1

R_2 = set of \mathbf{x} for which the object is being classified into π_2

In addition, $R_1 \cup R_2 = \Omega$, $R_1 \cap R_2 = \emptyset$

Classification rule (Region R_1 and R_2): minimizes the chances of making misclassification error



Two populations

Important concepts

1. Prior

p_1 = prior probability of π_1

p_2 = prior probability of π_2

2. Misclassification

$P(2|1)$ = Conditional probability of misclassification of an object as π_2 when, in fact, it is from π_1

$P(1|2)$ = Conditional probability of misclassification of an object as π_1 when, in fact, it is from π_2

3. Cost

Cost matrix

True Population	Classify as	
	π_1	π_2
π_1	0	$c(2 1)$
π_2	$c(1 2)$	0



Two populations

Expected cost of misclassification (ECM)

$$\text{ECM} = c(2|1)P(2|1)p_1 + c(1|2)P(1|2)p_2$$

Minimum Expected Cost of Misclassification

The Regions R_1 and R_2 that minimize the ECM are defined by the values \mathbf{x} for which the following inequalities hold.

$$R_1: \quad \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \geq \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$

$$R_2: \quad \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} < \left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right)$$



Two populations

Special Cases of Minimum Expected Cost of Misclassification Rule

1. Equal priors: $p_1 = p_2$

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \left(\frac{c(1|2)}{c(2|1)} \right); \quad R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{c(1|2)}{c(2|1)} \right)$$

2. Equal misclassification costs: $c(1|2) = c(2|1)$

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \left(\frac{p_2}{p_1} \right); \quad R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{p_2}{p_1} \right)$$

3. Equal priors and equal misclassification costs : $p_1 = p_2, c(1|2) = c(2|1)$

$$R_1: \frac{f_1(x)}{f_2(x)} \geq 1; \quad R_2: \frac{f_1(x)}{f_2(x)} < 1$$



Two populations

TMP rule: Minimize the total probability of misclassification

$$\begin{aligned}\text{TMP} &= P(\text{misclassifying a } \pi_1 \text{ observation or misclassifying a } \pi_2 \text{ observation}) \\ &= P(\text{observation from } \pi_1 \text{ and is misclassified}) + P(\text{observation from } \pi_2 \text{ and is misclassified}) \\ &= p_1 \int_{R_2} f_1(\mathbf{x}) d\mathbf{x} + p_2 \int_{R_1} f_2(\mathbf{x}) d\mathbf{x}\end{aligned}$$

Allocate observations in order to minimize TPM. Mathematically, equivalent to minimizing the expected cost of misclassification when the costs of misclassification are equal.



Two populations (Normal populations: $(\mu_1, \Sigma_1), (\mu_2, \Sigma_2)$)

Equal population variances $\Sigma_1 = \Sigma_2 = \Sigma$

Minimum ECM rule:

Allocate a new observation \mathbf{x}_0 to π_1 if

$$(\mu_1 - \mu_2)' \Sigma^{-1} \mathbf{x}_0 - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln \left[\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \right]$$

Allocate \mathbf{x}_0 to π_2 otherwise.



Two populations (Normal populations: $(\mu_1, \Sigma_1), (\mu_2, \Sigma_2)$)

Equal population variances $\Sigma_1 = \Sigma_2 = \Sigma$

Estimated Minimum ECM rule:

Allocate a new observation \mathbf{x}_0 to π_1 if

$$(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{S}_{pooled}^{-1} \mathbf{x}_0 - \frac{1}{2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)' \mathbf{S}_{pooled}^{-1} (\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2) \geq \ln \left[\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \right]$$

Allocate \mathbf{x}_0 to π_2 otherwise.



Two populations (Normal populations: $(\mu_1, \Sigma_1), (\mu_2, \Sigma_2)$)

Unequal population variances $\Sigma_1 \neq \Sigma_2$

Minimum ECM rule:

$$\text{Let } k = \frac{1}{2} \ln \left(\frac{|\Sigma_1|}{|\Sigma_2|} \right) + \frac{1}{2} (\mu_1' \Sigma_1^{-1} \mu_1 - \mu_2' \Sigma_2^{-1} \mu_2)$$

Allocate a new observation \mathbf{x}_0 to π_1 if

$$-\frac{1}{2} \mathbf{x}_0' (\Sigma_1^{-1} - \Sigma_2^{-1}) \mathbf{x}_0 + (\mu_1' \Sigma_1^{-1} - \mu_2' \Sigma_2^{-1}) \mathbf{x}_0 - k \geq \ln \left[\left(\frac{c(1|2)}{c(2|1)} \right) \left(\frac{p_2}{p_1} \right) \right]$$

Allocate \mathbf{x}_0 to π_2 otherwise.

Estimated Minimum ECM rule: same as above except the unknown parameters are replaced by sample estimates. Hence, $\mu_1, \Sigma_1, \mu_2, \Sigma_2$ are replaced by $\bar{\mathbf{x}}_1, \mathbf{S}_1, \bar{\mathbf{x}}_2, \mathbf{S}_2$.