

4. Repeat Steps 1 through 3 for the π_2 observations. Let $n_{2M}^{(H)}$ be the number of holdout observations misclassified in this group.

$$\hat{P}(2|1) = \frac{n_{1M}^{(H)}}{n_1} \quad \hat{E}(\text{AER}) = \frac{n_{1M}^{(H)} + n_{2M}^{(H)}}{n_1 + n_2}$$

$$\hat{P}(1|2) = \frac{n_{2M}^{(H)}}{n_2}$$

❖ **Reading:** textbook, 11.1, 11.2, 11.3, 11.4

• Classification with Several Populations

- generalization of classification procedure from 2 to $g \geq 2$ groups
- minimum expected cost of misclassification method
 - Let $f_i(\mathbf{x})$ be the density associated with population $\pi_i, i = 1, 2, \dots, g$.
 - Let p_i = the prior probability of population $\pi_i, i = 1, 2, \dots, g$
 - Let $c(k|i)$ = the cost of allocating an item to π_k when, in fact, it belongs to $\pi_i, \text{ for } k, i = 1, 2, \dots, g$
 - Let R_k be the set of \mathbf{x} 's classified as π_k
 - Then,

$$P(k|i) = P(\text{classifying item as } \pi_k | \pi_i) = \int_{R_k} f_i(\mathbf{x}) d\mathbf{x} \text{ for } k, i = 1, 2, \dots, g$$

NTHU STAT 5191, 2010, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- The conditional expected cost of misclassifying an \mathbf{x} from π_1 into π_2 , or π_3, \dots , or π_g is

$$\text{ECM}(1) = P(2|1)c(2|1) + P(3|1)c(3|1) + \dots + P(g|1)c(g|1) = \sum_{k=2}^g P(k|1)c(k|1)$$

- In a similar manner, we can obtain the conditional expected costs of misclassification $\text{ECM}(2), \dots, \text{ECM}(g)$.
- The overall ECM is:

$$\begin{aligned} \text{ECM} &= p_1 \text{ECM}(1) + p_2 \text{ECM}(2) + \dots + p_g \text{ECM}(g) \\ &= p_1 \left(\sum_{k=2}^g P(k|1)c(k|1) \right) + p_2 \left(\sum_{k=1, k \neq 2}^g P(k|2)c(k|2) \right) + \dots + p_g \left(\sum_{k=1}^{g-1} P(k|g)c(k|g) \right) \\ &= \sum_{i=1}^g p_i \left(\sum_{k=1, k \neq i}^g P(k|i)c(k|i) \right) \end{aligned}$$

- **Result 11.5.** The classification regions that minimize the ECM are defined by allocating \mathbf{x} to that population $\pi_k, k = 1, 2, \dots, g$, for which

$$\sum_{\substack{i=1 \\ i \neq k}}^g p_i f_i(\mathbf{x}) c(k|i)$$

is smallest. If a tie occurs, \mathbf{x} can be assigned to any of the tied populations.

- ◆ Suppose all the misclassification costs are equal. The minimum ECM is the minimum total probability of misclassification. In which case, we would allocate \mathbf{x} to that population $\pi_k, k = 1, 2, \dots, g$, for which

$$\sum_{\substack{i=1 \\ i \neq k}}^g p_i f_i(\mathbf{x})$$

is smallest It will be smallest when the omitted term, $p_k f_k(\mathbf{x})$, is largest.

- ◆ Minimum ECM Classification Rule with equal misclassification costs

Allocate \mathbf{x}_0 to π_k if $p_k f_k(\mathbf{x}) > p_i f_i(\mathbf{x})$ for all $i \neq k$

- ◆ Notice that the classification rule is identical to the one that maximize the posterior probability

$$\begin{aligned} P(\pi_k | \mathbf{x}) &= P(\mathbf{x} \text{ comes from } \pi_k \text{ given that } \mathbf{x} \text{ was observed}) \\ &= \frac{p_k f_k(\mathbf{x})}{\sum_{i=1}^g p_i f_i(\mathbf{x})} = \frac{(\text{prior}) \times (\text{likelihood})}{\sum [(\text{prior}) \times (\text{likelihood})]} \quad \text{for } k = 1, 2, \dots, g \end{aligned}$$

➤ Classification with Normal Populations

- Under normality assumption,

$$f_i(\mathbf{x}) = \frac{1}{(2\pi)^{p/2} |\Sigma_i|^{1/2}} \exp \left[-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)' \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) \right], \quad i = 1, 2, \dots, g$$

NTHU STAT 5191, 2010, Lecture Notes
made by S.-W. Cheng (NTHU, Taiwan)

- Allocate \mathbf{x} to π_k if

$$\begin{aligned} \ln p_k f_k(\mathbf{x}) &= \ln p_k - \left(\frac{p}{2} \right) \ln(2\pi) - \frac{1}{2} \ln |\Sigma_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)' \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \\ &= \max_i \ln p_i f_i(\mathbf{x}) \end{aligned}$$

- define quadratic discrimination score for i th population

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)' \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i) + \ln p_i \quad i = 1, 2, \dots, g$$

- Minimum Total Probability of Misclassification (TPM) rule for normal populations with unequal Σ_i

Allocate \mathbf{x} to π_k if the quadratic score $d_k^Q(\mathbf{x}) = \text{largest of } d_1^Q(\mathbf{x}), d_2^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x})$

- In practice, the $\boldsymbol{\mu}_i$ and Σ_i are unknown \Rightarrow replaced by their sample quantities

$$\hat{d}_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\mathbf{S}_i| - \frac{1}{2} (\mathbf{x} - \bar{\mathbf{x}}_i)' \mathbf{S}_i^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i) + \ln p_i, \quad i = 1, 2, \dots, g$$

- Estimated Minimum TPM rule for normal population with unequal Σ_i

Allocate \mathbf{x} to π_k if the quadratic score $\hat{d}_k^Q(\mathbf{x}) = \text{largest of } \hat{d}_1^Q(\mathbf{x}), \hat{d}_2^Q(\mathbf{x}), \dots, \hat{d}_g^Q(\mathbf{x})$

- When $\Sigma_i = \Sigma$, for $i = 1, 2, \dots, g$,

$$d_i^Q(\mathbf{x}) = -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} \mathbf{x}' \Sigma^{-1} \mathbf{x} + \boldsymbol{\mu}_i' \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i' \Sigma^{-1} \boldsymbol{\mu}_i + \ln p_i$$

The first two terms are the same for $d_1^Q(\mathbf{x}), d_2^Q(\mathbf{x}), \dots, d_g^Q(\mathbf{x})$,

- define the linear discriminant score

$$d_i(\mathbf{x}) = \boldsymbol{\mu}_i' \boldsymbol{\Sigma}^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i' \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_i + \ln p_i \quad \text{for } i = 1, 2, \dots, g$$

An estimate $\hat{d}_i(\mathbf{x})$ of the linear discriminant score $d_i(\mathbf{x})$ is based on the pooled

$$\mathbf{S}_{\text{pooled}} = \frac{1}{n_1 + n_2 + \dots + n_g - g} ((n_1 - 1) \mathbf{S}_1 + (n_2 - 1) \mathbf{S}_2 + \dots + (n_g - 1) \mathbf{S}_g)$$

and is given by

$$\hat{d}_i(\mathbf{x}) = \bar{\mathbf{x}}_i' \mathbf{S}_{\text{pooled}}^{-1} \mathbf{x} - \frac{1}{2} \bar{\mathbf{x}}_i' \mathbf{S}_{\text{pooled}}^{-1} \bar{\mathbf{x}}_i + \ln p_i \quad \text{for } i = 1, 2, \dots, g$$

- Estimated Minimum TPM rule for normal populations with equal covariance

Allocate \mathbf{x} to π_k if

the linear discriminant score $\hat{d}_k(\mathbf{x}) = \text{the largest of } \hat{d}_1(\mathbf{x}), \hat{d}_2(\mathbf{x}), \dots, \hat{d}_g(\mathbf{x})$

- ◆ An equivalent classifier for the equal-covariance case is to use

$$D_i^2(\mathbf{x}) = (\mathbf{x} - \bar{\mathbf{x}}_i)' \mathbf{S}_{\text{pooled}}^{-1} (\mathbf{x} - \bar{\mathbf{x}}_i)$$

It measure the squared distances

from \mathbf{x} to the sample mean vector $\bar{\mathbf{x}}_i$.

The allocatory rule is then

Assign \mathbf{x} to the population π_i for which $-\frac{1}{2} D_i^2(\mathbf{x}) + \ln p_i$ is largest

- ◆ If the prior probabilities are unknown, the usual procedure is to set $p_1 = p_2 = \dots = p_g = 1/g$. An observation is then assigned to the closest population.

