

Chapter 6

Discriminant Analysis: Classification of Observations

Classification and discriminant analysis in multivariate analysis are useful for separating distinct sets of observations or objects and for allocating new observations or objects to previously defined groups.

Suppose an investigator makes several measurements on an individual and wants to classify the individual into one of several categories on the basis of these observations. It is not possible to directly identify and classify an individual into a category.

The difference between the techniques of cluster analysis and discriminant analysis is that the categories in discriminant analysis, in which the observation is to be classified, are known a priori, whereas the categories in which the observation is to be classified in cluster analysis are unknown a priori.

It is assumed that there is a finite number of categories or populations, and each category or population is characterized by a probability distribution of measurements. It is also assumed that the individual belongs to one of such known categories or populations. Thus, an individual is considered as a random observation from this population. The question to be answered is, given an individual with certain measurements, from which population or category did it arise?

For example, suppose that there are two areas or fields in which a mineral is found. The minerals in the two fields differ in their chemical properties, i.e., the minerals from both fields have the same chemical compounds, but their proportions differ. The chemical composition of the minerals from the two fields is supposed to be known. Now, a sample of mineral arrives, and it is to be found in which of the fields it belongs to. A possible solution is to determine the chemical properties of the new sample and compare them with those of the two known fields. Based on such a comparison, the new sample can be classified into one of the two known fields. The classification procedure is said to be good if the sample is correctly classified to the originating field, i.e., if the sample originating from, say, field 1 is classified to field 1 or the sample originating from, say,

field 2 is classified to field 2. On the other hand, the misclassification occurs when the sample is classified to the field from which it does not originate, i.e., if the sample arose from field 1 and is classified to field 2 and vice versa. Each field can be characterized by a probability distribution of measurements. So the observations on the chemical analysis of the sample can be considered as a random observation from one of the two mineral populations in the two fields.

Finally, the problem is that the measurements about the individual are given, and it is to be determined from which population they come. A good solution is to make a decision such that the chances of misclassification, or equivalently, the probability of misclassification, are minimum.

Classification as a Problem of Decision Functions

The classification problem can be considered a problem of “statistical decision functions”. Suppose there are a number of hypotheses, and each hypothesis is that the distribution of the observation is a given one. One of the hypotheses must be accepted and the other hypothesis rejected. If there are only two populations, there are only two hypotheses. The test of a hypothesis is conducted to test one hypothesis about a specified distribution against another.

Since a population or a category is specified by a probability distribution, the categories are specified beforehand in the sense that the probability distributions of the measurements are completely known. Another possibility is that the form of each probability distribution is known, but the corresponding distribution parameters are unknown; hence, they must be estimated from a sample from that population.

Considerations for Classification

It is desired to have a classification procedure that is constructed such that the probability of misclassification is minimum or the bad effects of misclassification are, on average,

minimized.

We consider a simple case of two populations. The two populations are denoted as π_1 and π_2 , and an observation is from either π_1 or π_2 . The observation is a measurement denoted as $\tilde{x} = (x_1, x_2, \dots, x_p)'$, which is a point in p -dimensional space, and the classification depends upon \tilde{x} . The aim is to set up a rule that if an individual is characterized by a set of measurement values x_1, x_2, \dots, x_p , then classify it into π_1 , i.e., the given observation is coming from the population π_1 , and for all other values, it is classified into π_2 .

As the observation \tilde{x} is considered as a point in the p -dimensional space, the entire p -dimensional space is divided into two regions, say R_1 and R_2 . Then the classification rule is as follows: If the observation falls in R_1 , it is classified as coming from π_1 ; if it falls in R_2 , it is classified as coming from π_2 . The error of misclassifying the observation is inevitable in this procedure. There are two types of misclassifications that can occur: If the observation is actually coming from π_1 , it is classified to π_2 or If the observation is coming from π_2 , it is classified to π_1 .

Both types of misclassification yield undesirable results, and it is necessary to know the relative undesirability of each misclassification. Such undesirable outcomes can be described by some associated costs.

Let the cost of misclassifying an individual from π_2 classified into π_1 be $C(1|2)$ and the cost of misclassifying an individual from π_1 classified into π_2 be $C(2|1)$. Both $C(1|2)$ and $C(2|1)$ are assumed to be positive and measured in some units. The statistician may not always know these costs exactly, but may have a rough idea about them. The misclassification may lead to a loss in a statistical sense.

The following notations are used for an observation \tilde{x} and populations π_1 and π_2 .

- (i) $\pi_1|\pi_1$: $\tilde{x} \in \pi_1$ and actually $\tilde{x} \in \pi_1$.

So the observation is from π_1 and is correctly classified into π_1 . This is the correct decision. So the associated cost of misclassification is $C(1|1) = 0$.

(ii) $\pi_2|\pi_2 : \mathcal{X} \in \pi_2$ and actually $\mathcal{X} \in \pi_2$.

The observation is from π_2 and is correctly classified into π_2 . This is the correct decision. So the associated cost of misclassification is $C(2|2) = 0$.

(iii) $\pi_1|\pi_2 : \mathcal{X} \in \pi_2$ and actually $\mathcal{X} \in \pi_1$.

The observation is from π_2 but is incorrectly classified into π_1 . This is an incorrect decision. The associated cost of misclassification is $C(1|2) > 0$.

(iv) $\pi_2|\pi_1 : \mathcal{X} \in \pi_1$ and actually $\mathcal{X} \in \pi_2$.

The observation comes from π_1 but is incorrectly classified as π_2 . This is an incorrect decision. The associated cost of misclassification is $C(2|1) > 0$.

A good classification procedure is one that minimizes, in some sense or other, the cost of misclassification. The statistical loss due to classification is considered, and the “minimum cost” is defined for two cases - when the a priori probabilities of the two populations are

(i) known and

(ii) unknown.

We first consider the case in which the a priori probabilities of the two populations are known.

Case 1: When A Priori Probabilities are Known

Suppose the a priori probabilities that an observation comes from π_1 and π_2 are q_1 and q_2 respectively, i.e.,

$$P(\mathcal{X} \in \pi_1) = q_1 ,$$

$$P(\mathcal{X} \in \pi_2) = q_2 ,$$

$$q_1 + q_2 = 1 .$$

The probability properties of populations are specified by a distribution function which, for convenience, is assumed to have a density. Let $p_1(\mathcal{X})$ and $p_2(\mathcal{X})$ are the densities of π_1 and π_2 , respectively.

Suppose R_1 and R_2 are the regions of classification of observation from π_1 and π_2 , respectively.

The probability that an observation from π_1 is correctly classified into π_1 is denoted as

$$P(1|1, R) = \int_{R_1} p_1(\underline{x}) d\underline{x}$$

where $d\underline{x} = dx_1 dx_2 \dots dx_p$

The probability that an observation from π_1 is misclassified into π_2 is denoted as

$$P(2|1, R) = \int_{R_2} p_1(\underline{x}) d\underline{x}.$$

Similarly, the probability that an observation from π_2 is correctly classified into π_2 is denoted as

$$P(2|2, R) = \int_{R_2} p_2(\underline{x}) d\underline{x}.$$

The probability that an observation from π_2 is misclassified into π_1 is denoted as

$$P(1|2, R) = \int_{R_1} p_2(\underline{x}) d\underline{x}.$$

Since $P(\underline{x} \in \pi_1) = q_1$, i.e., the probability of drawing an observation from π_1 is q_1 , so the probability of drawing an observation from π_1 and correctly classifying it is

$$q_1 \int_{R_1} p_1(\underline{x}) d\underline{x} = q_1 P(1|1, R).$$

Similarly, the probability of drawing an observation from π_1 and misclassifying it is $q_1 P(2|1, R)$.

The probability of drawing an observation from π_2 and correctly classifying it is $q_2 P(2|2, R)$.

The probability of drawing an observation from π_2 and misclassifying it is $q_2 P(1|2, R)$.

The average or expected loss from costs of misclassification is given by

$$EL = C(2|1) \cdot P(2|1, R) \cdot q_1 + C(1|2) \cdot P(1|2, R) \cdot q_2.$$

It is an aim to minimize this average loss. In other words, divide the space into regions R_1 and R_2 such that the expected loss EL is as small as possible.

A procedure that minimizes the expected loss EL for a given q_1 and q_2 is called a Bayes procedure, and the average expected loss is called the Bayes risk.

We have assumed so far that the errors of misclassification are equally important to avoid hanging an innocent person as compared to letting a murderer free.

Here $C(1|2)$ and $C(2|1)$ are the penalties of misclassification. For example, if a potentially good candidate for admission to a medical programme is rejected, then the nation will suffer a shortage of good doctors. If a bad candidate is admitted to a medical programme, the candidate may not be able to complete it successfully, the college resources used on the candidate will be wasted, and the nation will suffer a shortage of medical doctors.

Case 2: When A Priori Probabilities are Unknown

When the priori probabilities are not known, the expected loss if the observation is from π_1 is

$$C(2|1)P(2|1, R) = r(1, R)$$

and the expected loss if the observation is from π_2 is

$$C(1|2)P(1|2, R) = r(2, R).$$

We do not know whether the observation comes from π_1 or π_2 , nor the probabilities of these two instances.

A procedure is at least as good as R^* if

$$r(1, R) \leq r(1, R^*) \quad \text{and} \quad r(2, R) \leq r(2, R^*);$$

the procedure R is better than R^* if at least one of these inequalities is a strict inequality.

Usually, there is no procedure that is better than all others, or at least as good as all of them.

A procedure R is called admissible if there is no procedure better than R . Under certain conditions, the class of admissible procedures is the same as the class of Bayes procedures (obtained when the a priori probabilities are known).

The minimax principle leads to a unique procedure. A procedure is minimax if the maximum expected loss $r(i, R)$ is a minimum. From a conservative point of view, this may be considered an optimum procedure.

Procedure of Classification into One of Two Populations with Known Probability Distributions and A Priori Probabilities are Known:

The aim is now to choose the regions R_1 and R_2 such that the expected loss

$$C(2|1) \cdot P(2|1, R) \cdot q_1 + C(1|2) \cdot P(1|2, R) \cdot q_2$$

is minimized.

Since the a priori probabilities are assumed to be known, the joint probabilities of the population and the observed set of variables can be defined. The probability that an observation comes from π_1 and that each variate is less than the corresponding component in \underline{y} is

$$\int_{-\infty}^{y_p} \int_{-\infty}^{y_{p-1}} \dots \int_{-\infty}^{y_1} q_1 p_1(\underline{x}) dx_1 dx_2 \dots dx_p.$$

The conditional probability that an observation came from a given population, given the observed values of the variates, is defined as follows.

The conditional probability that an observation is coming from π_1 , given an observation \underline{x} , is

$$\frac{q_1 p_1(\underline{x})}{q_1 p_1(\underline{x}) + q_2 p_2(\underline{x})}.$$

Recall the following definition of conditional probability for two sets A and B

$$P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|A^c)P(A^c)}$$

which is used to find the conditional probability, and A^c is the complement of set A .

Suppose for a moment that the costs of misclassifications $C(1|2) = C(2|1) = 1$. Then the expected loss is

$$q_1 \int_{R_2} p_1(\underline{x}) d\underline{x} + q_2 \int_{R_1} p_2(\underline{x}) d\underline{x}.$$

This is also the probability of a misclassification. Hence, we wish to minimize the probability of misclassification.

Similarly, the conditional probability that an observation is coming from π_2 , given an observation \underline{x} , is

$$\frac{q_2 p_2(\underline{x})}{q_1 p_1(\underline{x}) + q_2 p_2(\underline{x})}.$$

For a given observed point \underline{x} , we minimize the probability of misclassification by assigning the population that has the higher conditional probability. If

$$\frac{q_1 p_1(\underline{x})}{q_1 p_1(\underline{x}) + q_2 p_2(\underline{x})} \geq \frac{q_2 p_2(\underline{x})}{q_1 p_1(\underline{x}) + q_2 p_2(\underline{x})},$$

we choose population π_1 . Otherwise, choose population π_2 . Since the probability of misclassification is minimized at each point, we minimize it over the whole space. Thus, the rule is as follows:

When q_1 and q_2 are known, then an individual is classified into π_1 , if

$$R_1 : q_1 p_1(\underline{x}) \geq q_2 p_2(\underline{x})$$

and the individual is classified into π_2 , if

$$R_2 : q_1 p_1(\underline{x}) < q_2 p_2(\underline{x}).$$

If $q_1 p_1(\underline{x}) = q_2 p_2(\underline{x})$, then the individual could be classified as either π_1 or π_2 . We have arbitrarily classified into R_1 . If $q_1 p_1(\underline{x}) + q_2 p_2(\underline{x}) = 0$ for a given \underline{x} , the individual may be classified into π_1 or π_2 . Note that R_2 is the complement of R_1 .

This is the Bayes procedure against a priori q_1 and q_2 .

Next we prove that the Bayes procedure given by $R_1 : q_1 p_1(\underline{x}) \geq q_2 p_2(\underline{x})$ and $R_2 : q_1 p_1(\underline{x}) < q_2 p_2(\underline{x})$ for assigning an individual into π_1 and π_2 , respectively is the best procedure.

Consider any other procedure $R^* = (R_1^*, R_2^*)$. Assume that the individual must belong to either of the two populations. The average expected loss or the probability of misclassification is

$$\begin{aligned} & q_1 \cdot \int_{R_2^*} p_1(\underline{x}) d\underline{x} + q_2 \cdot \int_{R_1^*} p_2(\underline{x}) d\underline{x} \\ &= q_1 \cdot \int_{R_2^*} p_1(\underline{x}) d\underline{x} + q_2 \cdot \left[\int_{R_1^*} p_2(\underline{x}) d\underline{x} + \int_{R_2^*} p_2(\underline{x}) d\underline{x} - \int_{R_2^*} p_2(\underline{x}) d\underline{x} \right] \\ &= \int_{R_2^*} (q_1 p_1(\underline{x}) - q_2 p_2(\underline{x})) d\underline{x} + q_2 \cdot \int_R p_2(\underline{x}) d\underline{x} \end{aligned}$$

where $R = R_1^* \cup R_2^*$.

Note that the second term $q_2 \int_R p_2(\underline{x}) d\underline{x}$ is a given number and a constant, as we are integrating over the whole space. So minimize only the first term

$$\int_{R_2^*} (q_1 p_1(\underline{x}) - q_2 p_2(\underline{x})) d\underline{x}.$$

If we define

$$R_2^* : q_1 p_1(\underline{x}) - q_2 p_2(\underline{x}) < 0$$

then the term $\int_{R_2^*} (q_1 p_1(\underline{x}) - q_2 p_2(\underline{x})) d\underline{x}$ can be positive, negative or equal to zero. So $\int_{R_2^*} (q_1 p_1(\underline{x}) - q_2 p_2(\underline{x})) d\underline{x}$ will be minimized if R_2^* includes the point \underline{x} such that

$$q_1 p_1(\underline{x}) - q_2 p_2(\underline{x}) < 0$$

and excludes the points for which

$$q_1 p_1(\underline{x}) - q_2 p_2(\underline{x}) > 0.$$

Hence, the average expected loss is minimum when

$$R_2^* : q_1 p_1(\underline{x}) < q_2 p_2(\underline{x})$$

or

$$R_2^* : \frac{p_1(\underline{x})}{p_2(\underline{x})} < \frac{q_2}{q_1}.$$

Obviously, then

$$R_1^* : q_1 p_1(\underline{x}) \geq q_2 p_2(\underline{x})$$

or

$$R_1^* : \frac{p_1(\underline{x})}{p_2(\underline{x})} \geq \frac{q_2}{q_1},$$

which is the same as the regions R_1 and R_2 obtained in the Bayes procedure.

If

$$P\left(\frac{p_1(\underline{x})}{p_2(\underline{x})} = \frac{q_2}{q_1} \middle| \pi_i\right) = 0, i = 1, 2,$$

then the Bayes procedure is unique except for sets of probability zero.

Now the problem of classification can be stated mathematically as follows:

Given nonnegative constants q_1 and q_2 and nonnegative functions $p_1(\underline{x})$ and $p_2(\underline{x})$, choose regions R_1 and R_2 so as to minimize

$$q_1 \int_{R_2} p_1(\underline{x}) d\underline{x} + q_2 \int_{R_1} p_2(\underline{x}) d\underline{x}.$$

The solution is

$$R_1 : q_1 p_1(\underline{x}) \geq q_2 p_2(\underline{x})$$

$$R_2 : q_1 p_1(\underline{x}) < q_2 p_2(\underline{x}).$$

So far, we have assumed that the costs of misclassification are equal. If the costs $C(1|2)$ and $C(2|1)$ are unequal, then the expected loss function is given by

$$C(2|1) \cdot q_1 \cdot \int_{R_2} p_1(\underline{x}) d\underline{x} + C(1|2) \cdot q_2 \cdot \int_{R_1} p_2(\underline{x}) d\underline{x}$$

which is minimized and the R_1 and R_2 are chosen as follows

$$R_1 : [C(2|1) \cdot q_1] p_1(\underline{x}) \geq [C(1|2) \cdot q_2] p_2(\underline{x})$$

$$R_2 : [C(2|1) \cdot q_1] p_1(\underline{x}) < [C(1|2) \cdot q_2] p_2(\underline{x}),$$

since $[C(2|1)q_1]$ and $[C(1|2)q_2]$ are nonnegative constants.

In another way,

$$\begin{aligned} R_1 : \frac{p_1(\underline{x})}{p_2(\underline{x})} &\geq \frac{C(1|2)q_2}{C(2|1)q_1} \\ R_2 : \frac{p_1(\underline{x})}{p_2(\underline{x})} &< \frac{C(1|2)q_2}{C(2|1)q_1} \end{aligned}$$

and the Bayes procedure will be unique, except for sets of probability zero, if

$$P \left[\frac{p_1(\underline{x})}{p_2(\underline{x})} = \frac{C(1|2)q_2}{C(2|1)q_1} \right] = 0.$$

Here we note that $\frac{p_1(\underline{x})}{p_2(\underline{x})}$ is the likelihood ratio for a sample of size one.

When q_1 and q_2 are Unknown

When q_1 and q_2 are unknown, we look for any other procedure, since the Bayes procedure is based on an a priori distribution. So we use the minimax rule.

When q_1 and q_2 are unknown, then the minimax rule can be obtained by equating the two losses due to misclassification. Thus

$$C(1|2) \int_{R_1} p_2(\underline{x}) d\underline{x} = C(2|1) \int_{R_2} p_1(\underline{x}) d\underline{x}.$$

If costs of classification are equal, i.e., $C(1|2) = C(2|1)$, then

$$\int_{R_1} p_2(\underline{x}) d\underline{x} = \int_{R_2} p_1(\underline{x}) d\underline{x}.$$

Another approach to minimise the average expected loss is due to the Neyman-Pearson theory. A linear combination of α and β , say $a\alpha + b\beta$ is minimized such that $a + b = 1$. This will lead to Bayes' decision rule against a priori distribution (a, b) . So there is no loss of generality in setting $C(1|2) = C(2|1)$. Also, α and β are the errors of the first and second type in the context of testing a hypothesis.

Classification into One of Two Known Multivariate Normal Populations

Now we use the Bayes procedure in the case of two multivariate normal populations with equal covariance matrices, say $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$. Then for $i = 1, 2$, the i^{th} density is

$$p_i(\underline{x}) = \frac{1}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp \left[-\frac{1}{2}(\underline{x} - \mu_i)' \Sigma^{-1} (\underline{x} - \mu_i) \right].$$

The ratio of densities is

$$\begin{aligned} \frac{p_1(\underline{x})}{p_2(\underline{x})} &= \frac{\frac{1}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp \left[-\frac{1}{2}(\underline{x} - \mu_1)' \Sigma^{-1} (\underline{x} - \mu_1) \right]}{\frac{1}{(2\pi)^{p/2}|\Sigma|^{1/2}} \exp \left[-\frac{1}{2}(\underline{x} - \mu_2)' \Sigma^{-1} (\underline{x} - \mu_2) \right]} \\ &= \exp \left[-\frac{1}{2} \left\{ (\underline{x} - \mu_1)' \Sigma^{-1} (\underline{x} - \mu_1) - (\underline{x} - \mu_2)' \Sigma^{-1} (\underline{x} - \mu_2) \right\} \right] \\ &= \exp \left[-\frac{1}{2} \left\{ -2\underline{x}' \Sigma^{-1} \mu_1 + 2\underline{x}' \Sigma^{-1} \mu_2 + \mu_1' \Sigma^{-1} \mu_1 + \mu_2' \Sigma^{-1} \mu_2 - \mu_2' \Sigma^{-1} \mu_1 + \mu_1' \Sigma^{-1} \mu_2 \right\} \right] \\ &= \exp \left[\underline{x}' \Sigma^{-1} (\mu_1 - \mu_2) - \frac{1}{2} (\mu_1 + \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2) \right]. \end{aligned}$$

The region of classification into π_1 is R_1 and R_1 is the set of \underline{x} 's for which

$$\frac{p_1(\underline{x})}{p_2(\underline{x})} \geq k$$

for suitably chosen k . Using the monotonicity property of the logarithmic function, we consider

$$\log \left(\frac{p_1(\underline{x})}{p_2(\underline{x})} \right) \geq \log k$$

or

$$\left[\underline{x}' \Sigma^{-1} (\mu_1 - \mu_2) - \frac{1}{2} (\mu_1 + \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2) \right] \geq \log k.$$

Here $\Sigma^{-1}(\mu_1 - \mu_2)$ is a scalar as μ_1, μ_2 and Σ are known quantities. Thus $\underline{x}' \Sigma^{-1} (\mu_1 - \mu_2)$ is a linear function of x_1, x_2, \dots, x_p . It is called as discriminant function. It is a linear function of the component of the observation vector.

Thus, the best regions of classification are given by

$$R_1 : \tilde{x}'\Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2) \geq \log k$$

$$R_2 : \tilde{x}'\Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2) < \log k.$$

If a priori probabilities q_1 and q_2 are known, then

$$k = \frac{q_2 C(1|2)}{q_1 C(2|1)}.$$

Note that the same k was obtained using the Bayes rule as well.

If the two populations are equally likely and the costs are equal, then $k = 1$ and $\log k = 0$. Then the regions of classification are

$$R_1 : \tilde{x}'\Sigma^{-1}(\mu_1 - \mu_2) \geq \frac{1}{2}(\mu_1 + \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2)$$

$$R_2 : \tilde{x}'\Sigma^{-1}(\mu_1 - \mu_2) < \frac{1}{2}(\mu_1 + \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2).$$

If we do not have a priori probabilities, we may choose $\log k = c$, say, on the basis of making the expected losses due to misclassification equal.

Let \tilde{X} be a random observation. There are now two cases.

Case (i): If $\tilde{X} \in \pi_1$, then $\tilde{X} \sim N_p(\mu_1, \Sigma)$.

Case (ii): If $\tilde{X} \in \pi_2$, then $\tilde{X} \sim N_p(\mu_2, \Sigma)$.

Let

$$U = \tilde{X}'\Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)'\Sigma^{-1}(\mu_1 - \mu_2).$$

Now we find the distribution of U under two cases - when $\tilde{X} \sim N_p(\mu_1, \Sigma)$ and when $\tilde{X} \sim N_p(\mu_2, \Sigma)$.

Case (i): When $\tilde{X} \in N_p(\mu_1, \Sigma)$

Since U is a linear function of normally distributed random variables, so U is also normally distributed. The mean of U under $N_p(\mu_1, \Sigma)$ is

$$\begin{aligned} E_1(U) &= \mu_1' \Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= \left[\mu_1 - \frac{1}{2}(\mu_1 + \mu_2) \right]' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= \frac{1}{2}(\mu_1 - \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= \frac{\Delta^2}{2} \end{aligned}$$

where Δ^2 is the Mahalanobis distance between $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$.

The variance of U under $N_p(\mu_1, \Sigma)$ is

$$\begin{aligned} V_1(U) &= E[(\mu_1 - \mu_2)' \Sigma^{-1}(\tilde{X} - \mu_1)(\tilde{X} - \mu_1)' \Sigma^{-1}(\mu_1 - \mu_2)] \\ &= (\mu_1 - \mu_2)' \Sigma^{-1} \Sigma \Sigma^{-1}(\mu_1 - \mu_2) \\ &= (\mu_1 - \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= \Delta^2. \end{aligned}$$

Thus if $\tilde{X} \in \pi_1$, then $U \sim N(\frac{\Delta^2}{2}, \Delta^2)$.

Case (ii): When $\tilde{X} \in N_p(\mu_2, \Sigma)$

Since U is a linear function of normally distributed random variables, so U is also normally distributed.

The mean of U under $N_p(\mu_2, \Sigma)$ is

$$\begin{aligned} E_2(U) &= \mu_2' \Sigma^{-1}(\mu_1 - \mu_2) - \frac{1}{2}(\mu_1 + \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= -\frac{1}{2}(\mu_1 - \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2) \\ &= -\frac{\Delta^2}{2}. \end{aligned}$$

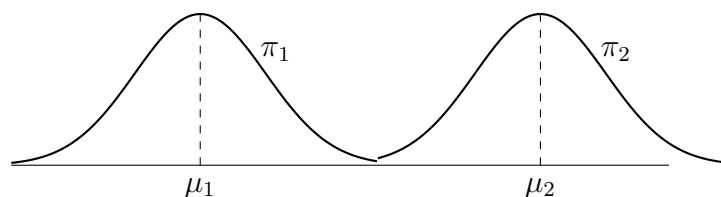
The variance of U under $N_p(\mu_2, \Sigma)$ is

$$\begin{aligned} V_2(U) &= (\mu_1 - \mu_2)' \Sigma^{-1} E(\tilde{X} - \mu_2)(\tilde{X} - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2) \\ &= (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 - \mu_2) \\ &= \Delta^2. \end{aligned}$$

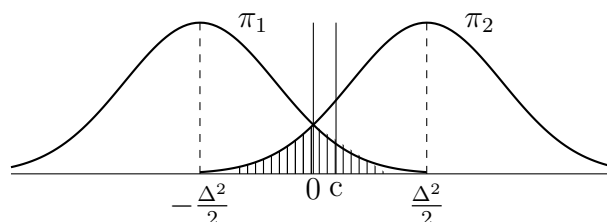
Thus if $\tilde{X} \in \pi_2$, then $U \sim N(-\frac{\Delta^2}{2}, \Delta^2)$.

Probability of Classification

Two well classified normal populations π_1 and π_2 with no misclassification with means μ_1 and μ_2 with same variances looks as follows:



The same populations with misclassification look like follows, and we want to find the probabilities of misclassification.



If any observation lies in the shaded region, then the observation can belong to π_1 or π_2 with certain probabilities of misclassification.

If $\tilde{X} \in \pi_1$, then the probability of misclassification is

$$\begin{aligned} P(2|1) &= \int_{R_2} N\left(\frac{\Delta^2}{2}, \Delta^2\right) \\ &= \int_{-\infty}^c \frac{1}{\sqrt{2\pi\Delta^2}} \exp\left[-\frac{1}{2}\left(\frac{z - \frac{\Delta^2}{2}}{\Delta^2}\right)^2\right] dz \\ &= \int_{-\infty}^{\frac{c - \frac{\Delta^2}{2}}{\Delta}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy. \end{aligned}$$

If $\tilde{X} \in \pi_2$, then the probability of misclassification is

$$\begin{aligned} P(1|2) &= \int_{R_1} N\left(-\frac{\Delta^2}{2}, \Delta^2\right) \\ &= \int_c^{\infty} \frac{1}{\sqrt{2\pi\Delta^2}} \exp\left[-\frac{1}{2}\left(\frac{z + \frac{\Delta^2}{2}}{\Delta^2}\right)^2\right] dz \\ &= \int_{\frac{c + \frac{\Delta^2}{2}}{\Delta}}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy. \end{aligned}$$

The probabilities $P(1|2)$ and $P(2|1)$ can be calculated with the help of a normal probability table. Here, c is known because we are considering the Bayes' solution in which a priori probabilities and costs are known. Also, given the probability, we can find the value of c .

When prior probabilities are unknown, we use the minimax rule. If the costs are known, then the expected loss is

$$C(1|2) \cdot P(1|2) + C(2|1) \cdot P(2|1).$$

For the minimax solution, we choose c such that

$$C(1|2) \int_{\frac{c + \frac{\Delta^2}{2}}{\Delta}}^{\infty} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy = C(2|1) \int_{-\infty}^{\frac{c - \frac{\Delta^2}{2}}{\Delta}} \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{y^2}{2}\right) dy.$$

After finding c , we can easily find R_1 and R_2 .

Fisher's Approach to find the Linear Discriminant Function

The linear discriminant function $\tilde{x}'\Sigma^{-1}(\mu_1 - \mu_2)$ was also derived by Fisher with a different approach.

Suppose there are two populations $N(\mu_1, \sigma^2)$ and $N(\mu_2, \sigma^2)$. These populations, if they are well classified, will look like those shown in Figure 1.

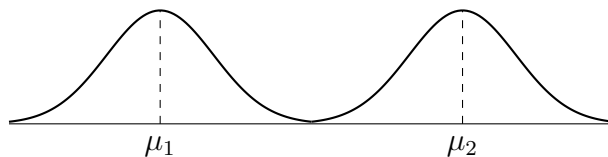


Figure 1: No misclassification

When misclassification is possible, the same populations will appear as in Figure 2.

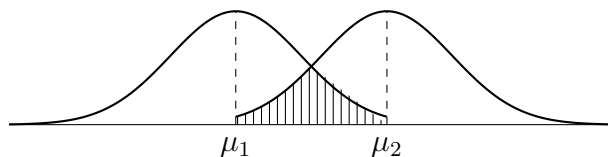


Figure 2: The misclassification area is indicated by the shaded region

In such a case, the misclassification is possible if an individual lies in the shaded region. So the question arises: what criterion should we use to classify the observations into either population? It is to be observed that in such a situation, the populations are differing with respect to the centre of location only, i.e., the means μ_1 and μ_2 . So, a criterion where the distance between μ_1 and μ_2 is maximum, i.e., $(\mu_1 - \mu_2)$ is maximum, can be used to decide the misclassification. Since σ^2 also has an important role as a scale parameter, considering $(\frac{\mu_1 - \mu_2}{\sigma})$ will make the distance independent of σ . Thus the choice for the criterion becomes now “ $(\frac{\mu_1 - \mu_2}{\sigma})$ is maximum”. Fisher used this idea in a multivariate normal population.

It is intuitively clear that an observation \tilde{x} will be classified to $\pi_1 \equiv N(\mu_1, \sigma^2)$ if \tilde{x} is nearer to μ_1 rather than μ_2 . So the risk of misclassification is smaller if $(\frac{\mu_1 - \mu_2}{\sigma})^2$, which is unit-free, is large because the two normal density curves are farther apart and the overlap is smaller. The reverse is true otherwise. Also note that maximizing $(\frac{\mu_1 - \mu_2}{\sigma})$ is the same as to maximize $(\frac{\mu_1 - \mu_2}{\sigma})^2$.

So, suppose two multivariate normal populations are given as $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$ which differ in mean vectors μ_1 and μ_2 and having the same covariance matrix Σ . We

want a linear function of \underline{x} , viz., $\underline{a}'\underline{x}$ where \underline{a} is unknown and the choice of \underline{a} has to be made such that

$$\frac{[E_1(\underline{a}'\underline{x}) - E_2(\underline{a}'\underline{x})]^2}{\underline{a}'\Sigma\underline{a}} \quad (1)$$

is maximum. Fisher used this idea in finding such an \underline{a} . Now, similar to (1), we can write

$$\frac{(\underline{a}'\underline{\mu}_1 - \underline{a}'\underline{\mu}_2)^2}{\underline{a}'\Sigma\underline{a}} = \frac{\underline{a}'(\underline{\mu}_1 - \underline{\mu}_2)(\underline{\mu}_1 - \underline{\mu}_2)'\underline{a}}{\underline{a}'\Sigma\underline{a}} \quad (2)$$

So we have to maximize (2) with an aim to obtain \underline{a} , using the following lemma:

Lemma 1: Let B be a $p \times p$ positive definite matrix and \underline{d} is a $p \times 1$ given vector. Then for any arbitrary nonzero vector \underline{x} of order $p \times 1$,

$$\max_{\underline{x} \neq 0} \frac{(\underline{x}'\underline{d})^2}{\underline{x}'B\underline{x}} = \underline{d}'B^{-1}\underline{d}$$

with the maximum attained when $\underline{x} = cB^{-1}\underline{d}$ for any constant $c \neq 0$.

Using Lemma 1, we get

$$\max_{\underline{a} \neq 0} \frac{[\underline{a}'(\underline{\mu}_1 - \underline{\mu}_2)]^2}{\underline{a}'\Sigma\underline{a}} = (\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}(\underline{\mu}_1 - \underline{\mu}_2) = \Delta^2$$

with the maximum attained with $\underline{a} = c\Sigma^{-1}(\underline{\mu}_1 - \underline{\mu}_2)$ for any constant $c \neq 0$. Thus the linear discriminant function $\underline{a}'\underline{x}$ is now the same as earlier $\underline{a}'\underline{x} = (\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}\underline{x}$. Any linear function which is proportional to $\underline{a}'\underline{x}$ can be used as a discriminant function.

Note that the Mahalanobis distance Δ^2 is given by

$$\Delta^2 = (\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}(\underline{\mu}_1 - \underline{\mu}_2).$$

The same result can also be alternatively proved by maximizing the function using calculus

$$\underline{a}'(\underline{\mu}_1 - \underline{\mu}_2)(\underline{\mu}_1 - \underline{\mu}_2)'\underline{a} - \lambda(\underline{a}'\Sigma\underline{a} - 1)$$

where without loss of generality, we can take $\underline{a}'\Sigma\underline{a} = 1$. This gives $\underline{a}'\underline{x} = (\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}\underline{x}$ which is again the same result as obtained earlier.

Performance of Classification Rules

The performance of a discriminant function depends, in principle, on how well it correctly classifies observations into the correct populations.

The performance of a sample classification rule is evaluated by calculating the Actual Error Rate (*AER*) defined as

$$AER = q_1 \int_{R_2} f_1(\underline{x}) d\underline{x} + q_2 \int_{R_1} f_2(\underline{x}) d\underline{x}$$

where R_1 and R_2 represent the classification regions determined by samples of sizes N_1 and N_2 . For example,

$$R_1 : (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' S^{-1} \underline{x} - \frac{1}{2} (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' S^{-1} (\bar{\underline{x}}_1 + \bar{\underline{x}}_2) \geq \log_e \left(\frac{C(1|2)}{C(2|1)} \cdot \frac{q_2}{q_1} \right)$$

$$R_2 : (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' S^{-1} \underline{x} - \frac{1}{2} (\bar{\underline{x}}_1 - \bar{\underline{x}}_2)' S^{-1} (\bar{\underline{x}}_1 + \bar{\underline{x}}_2) < \log_e \left(\frac{C(1|2)}{C(2|1)} \cdot \frac{q_2}{q_1} \right).$$

AER indicates how the sample classification function will perform in future samples, but it depends on the unknown density functions $f_1(\underline{x})$ and $f_2(\underline{x})$.

However, an estimate of a quantity related to *AER*, called Apparent Error Rate (*APER*), is defined as the fraction of observations in the training sample that are misclassified by the sample classification function. The *APER* does not depend on the form of the parent populations and can be calculated for any classification procedure. The *APER* can be calculated from the **confusion matrix**, which shows the actual versus predicted group membership.

The confusion matrix for N_1 observations from population π_1 and N_2 observations from population π_2 has the following form:

| | | Predicted membership | | |
|-------------------|---------|-------------------------|-------------------------|-------|
| | | π_1 | π_2 | |
| Actual membership | π_1 | n_{1c} | $n_{1m} = n_1 - n_{1c}$ | n_1 |
| | π_2 | $n_{2m} = n_2 - n_{2c}$ | n_{2c} | n_2 |

Here

n_{1c} : Number of π_1 items correctly classified as π_1 items.

n_{1m} : Number of π_1 items misclassified as π_2 items.

n_{2c} : Number of π_2 items correctly classified as π_2 items.

n_{2m} : Number of π_2 items misclassified as π_1 items.

The *APER* is defined as

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

which is the proportion of items in the training set that are misclassified.

Violation of Assumptions

Now, we consider situations in which the assumptions underlying the earlier discriminant function are violated.

When the Populations do not have the Same Covariance Matrices

Instead of having two multivariate normal populations with equal covariance matrices, suppose there are multivariate normal populations $N_p(\mu_1, \Sigma_1)$ and $N_p(\mu_2, \Sigma_2)$ having different covariance matrices Σ_1 and Σ_2 . In such a case, we will not have a linear discriminant function but a quadratic function in \tilde{x} is obtained. This can be seen from the following

$$\frac{f_1(\tilde{x})}{f_2(\tilde{x})} = \exp \left[(\tilde{x} - \mu_2)' \Sigma_2^{-1} (\tilde{x} - \mu_2) - (\tilde{x} - \mu_1)' \Sigma_1^{-1} (\tilde{x} - \mu_1) \right]$$

which is a quadratic function in \tilde{x} .

Hence, the earlier defined R_1 and R_2 will change accordingly. Anderson and Bahadur (1962) considered a linear discriminant function $R_1 : \tilde{b}'\tilde{x} < c$ and $R_2 : \tilde{b}'\tilde{x} > c$, but c is difficult to determine.

2. When the Parameters μ_1, μ_2 and Σ are Unknown

Consider two multivariate normal populations $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$. If μ_1, μ_2 and Σ are unknown, then they are replaced by their respective maximum likelihood estimates,

i.e., μ_1, μ_2 and Σ are replaced by \bar{x}_1, \bar{x}_2 and S , respectively and the discriminant function is obtained as follows:

Let the two samples of sizes N_1 and N_2 are obtained from $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$ respectively and these samples are called as training samples.

Let $x_1^{(1)}, x_2^{(1)}, \dots, x_{N_1}^{(1)} \sim N_p(\mu_1, \Sigma)$ and $x_1^{(2)}, x_2^{(2)}, \dots, x_{N_2}^{(2)} \sim N_p(\mu_2, \Sigma)$. Then the sample means based on N_1 and N_2 observations are obtained as $\bar{x}^{(1)}$ and $\bar{x}^{(2)}$ which are the maximum likelihood estimators of μ_1 and μ_2 , respectively. The maximum likelihood estimator of Σ is based on pooled covariance matrix as

$$S = \frac{1}{N_1 + N_2 - 2} \left[\sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}_1)(x_{\alpha}^{(1)} - \bar{x}_1)' + \sum_{\alpha=1}^{N_2} (x_{\alpha}^{(2)} - \bar{x}_2)(x_{\alpha}^{(2)} - \bar{x}_2)' \right].$$

Our aim is to classify the observation x into $\pi_1 \equiv N_p(\mu_1, \Sigma)$ or $\pi_2 \equiv N_p(\mu_2, \Sigma)$.

Substitute \bar{x}_1, \bar{x}_2 and S in place of μ_1, μ_2 and Σ , respectively in the earlier obtained linear discriminant function U and we obtain

$$V(x) = x'S^{-1}(\bar{x}_1 - \bar{x}_2) - \frac{1}{2}(\bar{x}_1 + \bar{x}_2)'S^{-1}(\bar{x}_1 - \bar{x}_2).$$

This V statistic is called as **Anderson statistic**.

Now the regions R_1 and R_2 are defined accordingly as in U :

$$R_1 : V > c$$

$$R_2 : V < c$$

where c is a constant. The V statistic can be used as a classification criterion, as we did with U earlier.

Note that V is no longer a linear discriminant function. Both the terms in V are functions of random variables. Since \bar{x}_1 and \bar{x}_2 are normally distributed whereas S has a Wishart distribution, so the exact distribution of V is difficult to find but approximate results for the distribution of V are available in literature.

When the population is known, we can argue that the classification criterion is best, in the sense that it minimises expected loss when a priori probabilities are known, and

generates the class of admissible procedures when a priori probabilities are not known. We cannot justify V in the same way. However, intuitively, V should perform well.

Classification of Whole Sample

Suppose we have a sample x_1, x_2, \dots, x_N from either π_1 or π_2 and we wish to classify the sample as a whole, then we proceed as follows:

Define S as

$$(N_1 + N_2 + N - 3)S = \sum_{\alpha=1}^{N_1} (x_{\alpha}^{(1)} - \bar{x}_1)(x_{\alpha}^{(1)} - \bar{x}_1)' + \sum_{\alpha=1}^{N_2} (x_{\alpha}^{(2)} - \bar{x}_2)(x_{\alpha}^{(2)} - \bar{x}_2)' + \sum_{\alpha=1}^N (x_{\alpha} - \bar{x})(x_{\alpha} - \bar{x})'$$

where $\bar{x} = \frac{1}{N} \sum_{\alpha=1}^N x_{\alpha}$. Then the criterion is

$$\left[\bar{x} - \frac{1}{2}(\bar{x}_1 + \bar{x}_2) \right]' S^{-1}(\bar{x}_1 - \bar{x}_2).$$

The larger N is, the smaller the probabilities of misclassification.

On the Distribution of V Statistic

Now we discuss the distribution of V statistic given as

$$V = \bar{X}' S^{-1}(\bar{X}_1 - \bar{X}_2) - \frac{1}{2}(\bar{X}_1 + \bar{X}_2)' S^{-1}(\bar{X}_1 - \bar{X}_2)$$

for random \bar{X} , \bar{X}_1 and \bar{X}_2 . The distribution of V is very complicated and depends on the sample sizes and unknown Δ^2 . We attempt, as follows, to gain some information about the distribution of V . Let

$$Y_1 = C_1 \left[\bar{X} - \frac{(N_1 \bar{X}_1 + N_2 \bar{X}_2)}{N_1 + N_2} \right]$$

$$Y_2 = C_2 [\bar{X}_1 - \bar{X}_2]$$

where $C_1 = \sqrt{\frac{N_1+N_2}{N_1+N_2+1}}$ and $C_2 = \sqrt{\frac{N_1N_2}{N_1+N_2}}$, Y_1 and Y_2 are independently and normally distributed with the following mean vectors and covariance matrices:

$$\begin{aligned} E(Y_1) &= C_1 \left[E(\tilde{X}) - \frac{N_1\mu_1 + N_2\mu_2}{N_1 + N_2} \right] \\ &= \begin{cases} C_1 \left(\frac{N_2}{N_1+N_2} \right) (\mu_1 - \mu_2) & \text{if } \tilde{X} \in \pi_1 \\ -C_1 \left(\frac{N_1}{N_1+N_2} \right) (\mu_1 - \mu_2) & \text{if } \tilde{x} \in \pi_2 \end{cases} \\ E(Y_2) &= C_2(\mu_1 - \mu_2) \text{ irrespective of either } \tilde{X} \in \pi_1 \text{ or } \tilde{X} \in \pi_2. \end{aligned}$$

where $E(\bar{X}_1) = \mu_1$, $E(\bar{X}_2) = \mu_2$.

$$\begin{aligned} Var(Y_1) &= C_1^2 \left[Var(\tilde{X}) + \frac{1}{(N_1 + N_2)^2} \{N_1^2 Var(\bar{X}_1) + N_2^2 Var(\bar{X}_2)\} \right] \\ &= C_1^2 \left[\Sigma + \frac{1}{(N_1 + N_2)^2} \left(\frac{N_1^2}{N_1} + \frac{N_2^2}{N_2} \right) \Sigma \right] \\ &= C_1^2 \left(\frac{N_1 + N_2 + 1}{N_1 + N_2} \right) \Sigma \\ &= \Sigma \text{ irrespective of either } \tilde{X} \in \pi_1 \text{ or } \tilde{X} \in \pi_2. \\ Var(Y_2) &= C_2^2 [Var(\bar{X}_1) + Var(\bar{X}_2)] \\ &= C_2^2 \left(\frac{\Sigma}{N_1} + \frac{\Sigma}{N_2} \right) \\ &= \Sigma \text{ irrespective of either } \tilde{X} \in \pi_1 \text{ or } \tilde{X} \in \pi_2. \end{aligned}$$

Further, also

$$Cov(Y_1, Y_2) = 0.$$

Thus Y_1 and Y_2 are independently distributed. Both Y_1 and Y_2 follow multivariate normal distribution with mean vectors given by $E(Y_1)$ and $E(Y_2)$ respectively and the same covariance matrix Σ .

Let

$$\begin{aligned} Y &= (Y_1, Y_2)' \\ M &= Y' S^{-1} Y = \begin{pmatrix} m_{11} & m_{12} \\ m_{21} & m_{22} \end{pmatrix}. \end{aligned}$$

Then V can be expressed as

$$V = \sqrt{\frac{N_1 + N_2 + 1}{N_1 N_2}} \cdot m_{12} + \frac{N_1 - N_2}{2N_1 N_2} \cdot m_{22}.$$

The density of M has been studied by Sitgreaves (1952) and Anderson (1951). The density of V has been studied by Wald (1944).

If $N_1 = N_2$, then the distribution of V for \tilde{X} from π_1 is the same as that of $-V$ for \tilde{X} from π_2 . Then, if $V \geq 0$ is the region of classification as π_1 , then the probability of misclassifying \tilde{X} when it is from π_1 is equal to the probability of misclassifying it when it is from π_2 .

The Asymptotic Distribution of V

In the case of large samples from $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$, we can apply limiting distribution theory. We denote $\hat{\theta}_n \xrightarrow{Pr} \theta$, i.e., $\hat{\theta}_n$ converges to θ in probability, denoted as $\text{plim}_{n \rightarrow \infty} \hat{\theta}_n = \theta$.

Since $\bar{\tilde{X}}_1$ and $\bar{\tilde{X}}_2$ are the means of N_1 and N_2 observations, respectively from $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$, respectively, so $\bar{\tilde{X}}_1 \xrightarrow{Pr} \mu_1$ and $\bar{\tilde{X}}_2 \xrightarrow{Pr} \mu_2$ or $\text{plim}_{N_1 \rightarrow \infty} \bar{\tilde{X}}_1 = \mu_1$, $\text{plim}_{N_2 \rightarrow \infty} \bar{\tilde{X}}_2 = \mu_2$. So $\text{plim } S = \Sigma$ as $N_1 \rightarrow \infty$ or $N_2 \rightarrow \infty$ or both $N_1, N_2 \rightarrow \infty$.

Also $\text{plim } S^{-1} = \Sigma^{-1}$. Thus

$$\text{plim}_{N_1, N_2 \rightarrow \infty} S^{-1}(\bar{\tilde{X}}_1 - \bar{\tilde{X}}_2) = \Sigma^{-1}(\mu_1 - \mu_2),$$

$$\text{plim}_{N_1, N_2 \rightarrow \infty} (\bar{\tilde{X}}_1 + \bar{\tilde{X}}_2)' S^{-1}(\bar{\tilde{X}}_1 - \bar{\tilde{X}}_2) = (\mu_1 + \mu_2)' \Sigma^{-1}(\mu_1 - \mu_2).$$

Thus $V \xrightarrow{Pr} U$.

For sufficiently large samples from π_1 and π_2 , we can use this criterion as if we knew the population exactly and make only a small error.

So the limiting distribution V as $N_1 \rightarrow \infty$ and $N_2 \rightarrow \infty$ is

- $N(\frac{\Delta^2}{2}, \Delta^2)$ if \tilde{X} is distributed according to $N_p(\mu_1, \Sigma)$
- $N(-\frac{\Delta^2}{2}, \Delta^2)$ if \tilde{X} is distributed according to $N_p(\mu_2, \Sigma)$.

Testing of Hypothesis and Likelihood Ratio Test

The classification problem can also be viewed as a hypothesis-testing problem. The main aim of classification is to assign an observation to one of the two populations, which can be expressed in terms of the null and alternative hypotheses.

Suppose we have two samples of sizes N_1 and N_2 as follows:

$$\mathbf{x}_1^{(1)}, \mathbf{x}_2^{(1)}, \dots, \mathbf{x}_{N_1}^{(1)} \sim N_p(\mu_1, \Sigma)$$

$$\mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}, \dots, \mathbf{x}_{N_2}^{(2)} \sim N_p(\mu_2, \Sigma).$$

An individual observation \mathbf{x} is coming from either $N_p(\mu_1, \Sigma)$ or $N_p(\mu_2, \Sigma)$. This can be framed as a problem of testing of hypothesis as follows:

$$H_0 : \mathbf{x}, \mathbf{x}_1^{(1)}, \mathbf{x}_2^{(1)}, \dots, \mathbf{x}_{N_1}^{(1)} \sim N_p(\mu_1, \Sigma) \quad \text{and} \quad \mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}, \dots, \mathbf{x}_{N_2}^{(2)} \sim N_p(\mu_2, \Sigma)$$

$$H_1 : \mathbf{x}_1^{(1)}, \mathbf{x}_2^{(1)}, \dots, \mathbf{x}_{N_1}^{(1)} \sim N_p(\mu_1, \Sigma) \quad \text{and} \quad \mathbf{x}, \mathbf{x}_1^{(2)}, \mathbf{x}_2^{(2)}, \dots, \mathbf{x}_{N_2}^{(2)} \sim N_p(\mu_2, \Sigma).$$

We know that $\bar{\mathbf{x}}_1$ and $\bar{\mathbf{x}}_2$ are the maximum likelihood estimators of μ_1 and μ_2 , respectively.

The maximum likelihood estimators under H_0 are as follows:

$$\begin{aligned} \hat{\mu}_1 &= \frac{\mathbf{x} + \sum_{\alpha=1}^{N_1} \mathbf{x}_\alpha^{(1)}}{N_1 + 1} = \frac{\mathbf{x} + N_1 \bar{\mathbf{x}}_1}{N_1 + 1} \\ \hat{\mu}_2 &= \frac{\sum_{\alpha=1}^{N_2} \mathbf{x}_\alpha^{(2)}}{N_2} = \bar{\mathbf{x}}_2 \\ \hat{\Sigma}_1 &= \frac{1}{N_1 + N_2 + 1} \left[\sum_{\alpha=1}^{N_1} (\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)(\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)' \right. \\ &\quad \left. + \sum_{\alpha=1}^{N_2} (\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)(\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)' + (\mathbf{x} - \hat{\mu}_1)(\mathbf{x} - \hat{\mu}_1)' \right]. \end{aligned}$$

The maximum likelihood estimators under H_1 are as follows:

$$\begin{aligned} \hat{\mu}_1 &= \bar{\mathbf{x}}_1 \\ \hat{\mu}_2 &= \frac{\mathbf{x} + \sum_{\alpha=1}^{N_2} \mathbf{x}_\alpha^{(2)}}{N_2 + 1} = \frac{\mathbf{x} + N_2 \bar{\mathbf{x}}_2}{N_2 + 1} \\ \hat{\Sigma}_2 &= \frac{1}{N_1 + N_2 + 1} \left[\sum_{\alpha=1}^{N_1} (\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)(\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)' \right. \\ &\quad \left. + \sum_{\alpha=1}^{N_2} (\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)(\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)' + (\mathbf{x} - \hat{\mu}_2)(\mathbf{x} - \hat{\mu}_2)' \right]. \end{aligned}$$

First we simplify the terms of $\hat{\Sigma}_1$ and $\hat{\Sigma}_2$ as follows:

Let

$$\begin{aligned} A_{11} &= \sum_{\alpha=1}^{N_1} (\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1)(\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1)' \\ A_{22} &= \sum_{\alpha=1}^{N_2} (\mathcal{X}_{\alpha}^{(2)} - \bar{x}_2)(\mathcal{X}_{\alpha}^{(2)} - \bar{x}_2)' \\ A &= A_{11} + A_{22}. \end{aligned}$$

Consider

$$\begin{aligned} &\sum_{\alpha=1}^{N_1} (\mathcal{X}_{\alpha}^{(1)} - \hat{\mu}_1)(\mathcal{X}_{\alpha}^{(1)} - \hat{\mu}_1)' \\ &= \sum_{\alpha=1}^{N_1} \left[(\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1) + (\bar{x}_1 - \hat{\mu}_1) \right] \left[(\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1) + (\bar{x}_1 - \hat{\mu}_1) \right]' \\ &= \sum_{\alpha=1}^{N_1} (\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1)(\mathcal{X}_{\alpha}^{(1)} - \bar{x}_1)' + N_1(\bar{x}_1 - \hat{\mu}_1)(\bar{x}_1 - \hat{\mu}_1)' \\ &= A_{11} + N_1 \left[\left(\bar{x}_1 - \frac{N_1 \bar{x}_1 + \bar{x}}{N_1 + 1} \right) \left(\bar{x}_1 - \frac{N_1 \bar{x}_1 + \bar{x}}{N_1 + 1} \right)' \right] \\ &= A_{11} + \frac{N_1}{(N_1 + 1)^2} (N_1 \bar{x}_1 + \bar{x}_1 - N_1 \bar{x}_1 - \bar{x})(N_1 \bar{x}_1 + \bar{x}_1 - N_1 \bar{x}_1 - \bar{x})' \\ &= A_{11} + \frac{N_1}{(N_1 + 1)^2} (\bar{x} - \bar{x}_1)(\bar{x} - \bar{x}_1)'. \end{aligned}$$

So

$$\begin{aligned} &\sum_{\alpha=1}^{N_1} (\mathcal{X}_{\alpha}^{(1)} - \hat{\mu}_1)(\mathcal{X}_{\alpha}^{(1)} - \hat{\mu}_1)' + (\bar{x} - \hat{\mu}_1)(\bar{x} - \hat{\mu}_1)' \\ &= A_{11} + \frac{N_1}{(N_1 + 1)^2} (\bar{x} - \bar{x}_1)(\bar{x} - \bar{x}_1)' + \left(\bar{x} - \frac{N_1 \bar{x}_1 + \bar{x}}{N_1 + 1} \right) \left(\bar{x} - \frac{N_1 \bar{x}_1 + \bar{x}}{N_1 + 1} \right)' \\ &= A_{11} + \frac{N_1}{(N_1 + 1)^2} (\bar{x} - \bar{x}_1)(\bar{x} - \bar{x}_1)' + N_1^2 (\bar{x} - \bar{x}_1)(\bar{x} - \bar{x}_1)' \\ &= A_{11} + \frac{N_1 (\bar{x} - \bar{x}_1)(\bar{x} - \bar{x}_1)'}{N_1 + 1}. \end{aligned}$$

Thus, under H_0

$$\begin{aligned}\hat{\Sigma}_1 &= \frac{1}{N_1 + N_2 + 1} \left[\sum_{\alpha=1}^{N_1} (\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)(\mathbf{x}_\alpha^{(1)} - \hat{\mu}_1)' + (\mathbf{x} - \hat{\mu}_1)(\mathbf{x} - \hat{\mu}_1)' \right. \\ &\quad \left. + \sum_{\alpha=1}^{N_2} (\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)(\mathbf{x}_\alpha^{(2)} - \hat{\mu}_2)' \right] \\ &= \frac{1}{N_1 + N_2 + 1} \left[A_{11} + \left(\frac{N_1}{N_1 + 1} \right) (\mathbf{x} - \bar{\mathbf{x}}_1)(\mathbf{x} - \bar{\mathbf{x}}_1)' + A_{22} \right] \\ &= \frac{1}{N_1 + N_2 + 1} \left[A + \left(\frac{N_1}{N_1 + 1} \right) (\mathbf{x} - \bar{\mathbf{x}}_1)(\mathbf{x} - \bar{\mathbf{x}}_1)' \right]\end{aligned}$$

where $A = A_{11} + A_{22}$.

Similarly, under H_1 ,

$$\hat{\Sigma}_2 = \frac{1}{N_1 + N_2 + 1} \left[A + \left(\frac{N_2}{N_2 + 1} \right) (\mathbf{x} - \bar{\mathbf{x}}_2)(\mathbf{x} - \bar{\mathbf{x}}_2)' \right]$$

The likelihood ratio test criterion is

$$\begin{aligned}\lambda &= \frac{\sup_{H_0} L}{\sup_{H_1} L} = \frac{|\hat{\Sigma}_2|^{\frac{N_1 + N_2 + 1}{2}}}{|\hat{\Sigma}_1|^{\frac{N_1 + N_2 + 1}{2}}} \\ \text{or} \quad \lambda^{\frac{2}{N_1 + N_2 + 1}} &= \frac{|\hat{\Sigma}_2|}{|\hat{\Sigma}_1|}.\end{aligned}$$

Now we first simplify

$$\begin{aligned}\frac{|\hat{\Sigma}_2|}{|\hat{\Sigma}_1|} &= \frac{|A + \frac{N_2}{N_2 + 1}(\mathbf{x} - \bar{\mathbf{x}}_2)(\mathbf{x} - \bar{\mathbf{x}}_2)'|}{|A + \frac{N_1}{N_1 + 1}(\mathbf{x} - \bar{\mathbf{x}}_1)(\mathbf{x} - \bar{\mathbf{x}}_1)'|} \\ &= \frac{1 + \frac{N_2}{N_2 + 1}(\mathbf{x} - \bar{\mathbf{x}}_2)'A^{-1}(\mathbf{x} - \bar{\mathbf{x}}_2)}{1 + \frac{N_1}{N_1 + 1}(\mathbf{x} - \bar{\mathbf{x}}_1)'A^{-1}(\mathbf{x} - \bar{\mathbf{x}}_1)}\end{aligned}$$

which is obtained by using the following result:

Result 2: For a vector \mathbf{y} ,

$$\begin{vmatrix} C & \mathbf{y} \\ \mathbf{y}' & 1 \end{vmatrix} = |C - \mathbf{y}\mathbf{y}'| = \begin{vmatrix} 1 & \mathbf{y}' \\ \mathbf{y} & C \end{vmatrix} = |C|(1 - \mathbf{y}'C^{-1}\mathbf{y})$$

and then

$$\frac{|\hat{\Sigma}_2|}{|\hat{\Sigma}_1|} = \frac{n + \frac{N_2}{N_2 + 1}(\mathbf{x} - \bar{\mathbf{x}}_2)'S^{-1}(\mathbf{x} - \bar{\mathbf{x}}_2)}{n + \frac{N_1}{N_1 + 1}(\mathbf{x} - \bar{\mathbf{x}}_1)'S^{-1}(\mathbf{x} - \bar{\mathbf{x}}_1)}$$

where $n = N_1 + N_2 - 2$ and $S = \frac{A}{n}$. The region of classification into π_1 consists of those points for which

$$\lambda > K_n \text{ (where } K_n \text{ is given)}$$

or

$$R_1 : n + \frac{N_2}{N_2 + 1}(\bar{x} - \bar{x}_2)'S^{-1}(\bar{x} - \bar{x}_2) > K_n \left[n + \frac{N_1}{N_1 + 1}(\bar{x} - \bar{x}_1)'S^{-1}(\bar{x} - \bar{x}_1) \right].$$

If $K_n = 1$, then the rule is the maximum likelihood rule. Let

$$Z = \frac{1}{2} \left[\frac{N_2}{N_2 + 1}(\bar{x} - \bar{x}_2)'S^{-1}(\bar{x} - \bar{x}_2) - \frac{N_1}{N_1 + 1}(\bar{x} - \bar{x}_1)'S^{-1}(\bar{x} - \bar{x}_1) \right].$$

Then the maximum likelihood rule becomes:

Classify to π_1 if $Z > 0$

Classify to π_2 if $Z < 0$

as $\log K_n = 0$.

Roughly speaking, assign \bar{x} to π_1 or π_2 according to whether the distance to \bar{x}_1 is less or greater than the distance to \bar{x}_2 .

The difference between V and Z is

$$V - Z = \frac{1}{2} \left[\frac{1}{N_2 + 1}(\bar{x} - \bar{x}_2)'S^{-1}(\bar{x} - \bar{x}_2) - \frac{1}{N_1 + 1}(\bar{x} - \bar{x}_1)'S^{-1}(\bar{x} - \bar{x}_1) \right]$$

for which

$$\text{plim}(V - Z) = 0$$

as $N_1, N_2 \rightarrow \infty$. The probability of misclassification with V is equivalent asymptotically to that with Z for large samples.

Classification into One of Several Populations when A Priori Probabilities are Known

Now we consider the case when there are more than two populations and an observation is to be classified into one of them when a priori probabilities are known.

Let $\pi_1, \pi_2, \dots, \pi_m$ be m populations with probability density functions $p_1(\mathcal{X}), p_2(\mathcal{X}), \dots, p_m(\mathcal{X})$, respectively. Then the observation space is divided into m mutually exclusive and exhaustive regions R_1, R_2, \dots, R_m .

We say that the observation comes from π_i if the observation falls in R_i .

Let

- $C(j|i)$: Cost of misclassification when an observation coming from π_j is misclassified into π_i .
- $P(j|i, R) = \int_{R_j} p_i(\mathcal{X}) d\mathcal{X}$,
- q_1, q_2, \dots, q_m are the a priori probabilities of the populations $\pi_1, \pi_2, \dots, \pi_m$,

respectively.

The expected loss is given as

$$\sum_{i=1}^m q_i \left[\sum_{\substack{j=1 \\ i \neq j}}^m C(j|i) P(j|i, R) \right]$$

and R_1, R_2, \dots, R_m are chosen such that this expected loss is minimum. Since the a priori probabilities are known, the conditional probabilities of an observation coming from a population given the values of \mathcal{X} can be defined.

The probability of an observation coming from π_i given \mathcal{X} is

$$P[\text{Observation coming from } \pi_i | \mathcal{X}] = \frac{q_i p_i(\mathcal{X})}{\sum_{k=1}^m q_k p_k(\mathcal{X})}.$$

If the observation is classified as from π_j , the expected loss is

$$EL = \sum_{\substack{i=1 \\ i \neq j}}^m \left[\frac{q_i p_i(\mathcal{X})}{\sum_{k=1}^m q_k p_k(\mathcal{X})} C(j|i) \right].$$

We minimize the expected loss at this point if we choose j so as to minimize the expected loss EL, i.e., we consider

$$\sum_{\substack{i=1 \\ i \neq j}}^m q_i p_i(\mathcal{X}) C(j|i)$$

for all j and select that j which gives the minimum.

If two different indices give the minimum, it is irrelevant which index is selected. Follow this procedure for each \mathbf{x} and define R_1, R_2, \dots, R_m . The classification rule is to classify an observation as coming from π_j if it falls into the region R_j .

Classification into One of Several Populations when A Priori Probabilities are Unknown

Now we consider the case when there are more than two populations and an observation is to be classified into one of them when a priori probabilities are unknown.

The unconditional expected loss for a classification procedure cannot be defined in such a case.

The conditional expected loss if the observation is coming from π_i is

$$\sum_{\substack{j=1 \\ j \neq i}}^m P(j|i, R) = r(i, R).$$

A procedure R is at least as good as another procedure R^* if

$$r(i, R) \leq r(i, R^*), \quad i = 1, 2, \dots, m.$$

The procedure R is **better** if at least one inequality is strict. R is **admissible** if there is no procedure R^* that is better than R . A class of procedures is **complete** if for every procedure R outside the class, there is a procedure in the class that is better.

Bayes procedure is admissible.

Classification Into one of Several Multivariate Normal Populations

Suppose there are m multivariate normal populations $\pi_i : N_p(\mu_i, \Sigma)$, $i = 1, 2, \dots, m$ and an observation \mathbf{x} is to be classified into one of these populations.

Recall the details of classifying an observation into one of the two multivariate normal populations $N_p(\mu_1, \Sigma)$ and $N_p(\mu_2, \Sigma)$. Based on that classification rule, we use

$$\begin{aligned} u_{jk} &= \log \frac{p_j(\underline{x})}{p_k(\underline{x})} \\ &= \left[\underline{x} - \frac{1}{2}(\mu_j + \mu_k) \right]' \Sigma^{-1}(\mu_j - \mu_k) \end{aligned}$$

for classifying an observation \underline{x} between $N_p(\mu_i, \Sigma)$ and $N_p(\mu_j, \Sigma)$. If a priori probabilities are known, R_j is defined by those \underline{x} satisfying

$$R_j : u_{jk}(\underline{x}) > \log \frac{q_k}{q_j}, k = 1, 2, \dots, m, k \neq j.$$

Such regions R_1, R_2, \dots, R_m minimizes the expected cost and

$$u_{jk}(\underline{x}) = -u_{kj}(\underline{x}).$$

Based on earlier developments about the U statistic, we can find the following:

$$\begin{aligned} U_{ji} &= \left[\underline{X} - \frac{1}{2}(\mu_i + \mu_j) \right]' \Sigma^{-1}(\mu_j - \mu_i) \\ U_{ji} &= -U_{ij} \\ U_{ji} &\sim N\left(\frac{\Delta_{ji}^2}{2}, \Delta_{ji}^2\right) \\ \Delta_{ji}^2 &= (\mu_j - \mu_i)' \Sigma^{-1}(\mu_j - \mu_i) \\ Cov(U_{ji}, U_{jk}) &= \Delta_{jk,ji} = (\mu_j - \mu_k)' \Sigma^{-1}(\mu_j - \mu_i). \end{aligned}$$

To determine the constant c_j , consider the integral

$$P(j|j, R) = \int_{c_j - c_m}^{\infty} \dots \int_{c_j - c_1}^{\infty} f_j du_{j1} \dots du_{j,j-1} du_{j,j+1} \dots du_m$$

where f_j is the density of $U_{ji}, i = 1, 2, \dots, m, i \neq j$.

When the parameters are unknown, then they are estimated from the training samples $\underline{x}_1^{(i)}, \underline{x}_2^{(i)}, \dots, \underline{x}_{N_i}^{(i)}$ drawn from $N_p(\mu_i, \Sigma), i = 1, 2, \dots, m$ as

$$\begin{aligned} \hat{\mu}_i &= \bar{\underline{x}}_i = \frac{1}{N_i} \sum_{\alpha=1}^{N_i} \underline{x}_{\alpha}^{(i)} \\ \hat{\Sigma} &= S \end{aligned}$$

where $(\sum_{i=1}^m N_i - m)S = \sum_{i=1}^m \sum_{\alpha=1}^{N_i} (\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}_i)(\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}_i)'$.

An analogue of $U_{ij}(\mathbf{X})$ when the parameters are unknown and estimated from the training sample is

$$W_{ij}(\mathbf{X}) = \left[\mathbf{X} - \frac{1}{2}(\bar{\mathbf{x}}_i + \bar{\mathbf{x}}_j) \right]' S^{-1}(\bar{\mathbf{x}}_i - \bar{\mathbf{x}}_j).$$

As the variables in W_{ij} are random, so the distribution of W_{ij} differs from the distribution of U_{ij} . However, as $N_i \rightarrow \infty$, the joint distributions approach those of U_{ij} . So for sufficiently large samples, W_{ij} can be used.

Classification Into One of Two Known Multivariate Populations with Unequal Covariance Matrices

Let there be two multivariate normal populations with unequal mean vectors and unequal covariance matrices as

$$\pi_1 : N_p(\boldsymbol{\mu}_1, \Sigma_1), \pi_2 : N_p(\boldsymbol{\mu}_2, \Sigma_2),$$

where $\boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$ and $\Sigma_1 \neq \Sigma_2$. When the parameters are known, the likelihood ratio for classifying an observation \mathbf{x} into one of the populations is

$$\begin{aligned} \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} &= \frac{|\Sigma_2|^{1/2}}{|\Sigma_1|^{1/2}} \cdot \frac{\exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)' \Sigma_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right]}{\exp \left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)' \Sigma_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right]} \\ &= |\Sigma_2|^{1/2} |\Sigma_1|^{-1/2} \exp \left[\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)' \Sigma_2^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)' \Sigma_1^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) \right]. \end{aligned}$$

Here $\log \frac{p_1(\mathbf{x})}{p_2(\mathbf{x})}$ is a quadratic function in \mathbf{x} . So use a transformation of \mathbf{x} such that the covariance matrix is an identity matrix and matrix of quadratic form is diagonal; then $\log \left(\frac{p_1(\mathbf{x})}{p_2(\mathbf{x})} \right)$ is distributed as linear combination of noncentral χ^2 distribution and a constant in addition.

When the parameters are unknown, then under the same null and alternative hy-

potheses about the observation \mathbf{x} , as earlier, use the likelihood ratio test criterion as

$$\begin{aligned} & \frac{|\hat{\Sigma}_1(2)|^{\frac{N_1}{2}} |\hat{\Sigma}_2(2)|^{\frac{N_2+1}{2}}}{|\hat{\Sigma}_1(1)|^{\frac{N_1+1}{2}} |\hat{\Sigma}_2(1)|^{\frac{N_2}{2}}} \\ &= \frac{[1 + (\mathbf{x} - \bar{\mathbf{x}}_2)' A_2^{-1} (\mathbf{x} - \bar{\mathbf{x}}_2)]^{\frac{N_2+1}{2}}}{[1 + (\mathbf{x} - \bar{\mathbf{x}}_1)' A_1^{-1} (\mathbf{x} - \bar{\mathbf{x}}_1)]^{\frac{N_1+1}{2}}} \cdot \frac{|A_2|^{1/2}}{|A_1|^{1/2}} \cdot \frac{(N_1 + 1)^{\frac{(N_1+1)p}{2}} N_2^{\frac{N_2 p}{2}}}{(N_2 + 1)^{\frac{(N_2+1)p}{2}} N_1^{\frac{N_1 p}{2}}} \end{aligned}$$

where

$$\begin{aligned} \hat{\Sigma}_1(1) &= \frac{1}{N_1 + 1} \left[A_1 + \frac{N_1}{N_1 + 1} (\mathbf{x} - \bar{\mathbf{x}}_1)(\mathbf{x} - \bar{\mathbf{x}}_1)' \right] \\ \hat{\Sigma}_2(1) &= \frac{A_2}{N_2} \\ \hat{\Sigma}_1(2) &= \frac{A_1}{N_1} \\ \hat{\Sigma}_2(2) &= \frac{1}{N_2 + 1} \left[A_2 + \frac{N_2}{N_2 + 1} (\mathbf{x} - \bar{\mathbf{x}}_2)(\mathbf{x} - \bar{\mathbf{x}}_2)' \right] \\ A_i &= \sum_{\alpha=1}^{N_i} (\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}_i)(\mathbf{x}_\alpha^{(i)} - \bar{\mathbf{x}}_i)', i = 1, 2. \end{aligned}$$

Tests Associated with Discriminant Function

1. Goodness of Fit of a Hypothetical Discriminant Function

Suppose we want to test whether the discriminant function is good enough. So the null hypothesis can be framed as

$$H_0 : \text{A given function } \alpha' \mathbf{x} \text{ is good enough for discriminating between } \pi_1 \text{ and } \pi_2.$$

Transform \mathbf{X} to Z_1 as

$$Z_1 = \alpha' \mathbf{X}$$

and Z_2, Z_3, \dots, Z_p are $(p - 1)$ other suitable linear functions of \mathbf{X} .

First, we prove that a discriminant function is invariant under a linear transformation as follows:

Let $\mathbf{Z} = (Z_1, Z_2, \dots, Z_p)' = A\mathbf{X}$, then $\text{Var}(\mathbf{Z}) = A\Sigma A'$ and difference in the means

of \underline{Z} in π_1 and π_2 is $A(\underline{\mu}_1 - \underline{\mu}_2)$. The discriminant function based on \underline{Z} is therefore

$$[A(\underline{\mu}_1 - \underline{\mu}_2)]'(A\Sigma A')^{-1}\underline{Z} = (\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}\underline{X}$$

as A is non-singular and chosen such that the first column is $\underline{\alpha}$ and rest are arbitrary.

If H_0 is true, then $(\underline{\mu}_1 - \underline{\mu}_2)'\Sigma^{-1}\underline{X}$ must be proportional to $\underline{\alpha}'\underline{X}$, i.e., Z_1 . So

$$[A(\underline{\mu}_1 - \underline{\mu}_2)]'(A\Sigma A')^{-1}\underline{Z} = \gamma Z_1$$

where γ is some constant. The coefficients Z_2, Z_3, \dots, Z_p in $[A(\underline{\mu}_1 - \underline{\mu}_2)]'(A\Sigma A')^{-1}\underline{Z}$ are thus all zero.

So H_0 is equivalent to the hypothesis that the variables Z_2, Z_3, \dots, Z_p do not occur in the discriminant function based on \underline{Z} .

Using the hypothesis related to $\Delta_p^2 = \Delta_k^2$, putting $k = 1$, the hypothesis reduces to in this case as $H_0 : \Delta_p^2 = \Delta_1^2$ which can be tested by Rao's U -statistic and under H_0 ,

$$\frac{N_1 + N_2 - p - 1}{p - 1} \cdot \frac{N_1 N_2 (D_p^2 - D_1^2)}{(N_1 + N_2)(N_1 + N_2 - 2) + N_1 N_2 D_1^2} \sim F_{p-1, N_1 + N_2 - p - 1}$$

where D_p^2 is the distance (D^2) based on Z_1, Z_2, \dots, Z_p and D_1^2 is distance (D^2) based on Z_1 alone. But D^2 is invariant for a non-singular linear transformation like $\underline{Z} = A\underline{X}$ as $(A\underline{d})'(ASA')^{-1}(A\underline{d}) = \underline{d}'S^{-1}\underline{d}$. So D_p^2 is distance (D^2) based on \underline{X} , viz., $(N_1 + N_2)\underline{d}'S^{-1}\underline{d}$ where $\underline{d} = \bar{\underline{X}} - \underline{\bar{y}}$. However, as $\underline{Z}_1 = \underline{\alpha}'\underline{X}$,

$$\frac{D_1^2}{N_1 + N_2} = (\underline{\alpha}'\underline{d})(\underline{\alpha}'S\underline{\alpha})^{-1}(\underline{\alpha}'\underline{d}) = \frac{(\underline{\alpha}'\underline{d})^2}{\underline{\alpha}'S\underline{\alpha}}.$$

If $H_0 : \Delta_p^2 = \Delta_1^2$ is accepted, then we conclude that the first observation, i.e., Z_1 is sufficient to discriminate between the two populations. Thus, we test the adequacy of the discriminant function.

2. Testing Discriminating Ability of Some Variables

To test the discriminating ability of some variables in a discriminant function, we consider the following three types of null hypotheses, which can be tested by Rao's U -statistic.

$$(I) H_{01} : a_{k+1} = a_{k+2} = \dots = a_p = 0$$

where $\underline{a} = \Sigma^{-1}(\underline{\mu}_1 - \underline{\mu}_2)$ is coefficient vector of discriminating function $\underline{a}'\underline{X}$ and $\underline{a} = (\underline{a}_1 \ \underline{a}_2)'$, $\underline{X} = (\underline{X}_1 \ \underline{X}_2)'$ and $S = \begin{pmatrix} S_{11} & S_{12} \\ S_{21} & S_{22} \end{pmatrix}$ where \underline{X}_1 and \underline{a}_1 are of order $(k \times 1)$; and \underline{X}_2 and \underline{a}_2 are of order $(p-k) \times 1$ and orders of S_{11} , S_{12} and S_{22} are suitably determined. So H_{01} can be expressed as

$$H_{01} : \underline{a}_2 = \underline{0}$$

.

$$(II) H_{02} : \Delta_p^2 = \Delta_q^2$$

$$(III) H_{03} : \underline{\delta}_2 = B\underline{\delta}_1$$

where $\underline{\delta} = (\underline{\delta}_1 \ \underline{\delta}_2)'$, $\underline{\delta} = \underline{\mu}_1 - \underline{\mu}_2$ and $\underline{\delta}_1$ and $\underline{\delta}_2$ are of orders $(k \times 1)$ and $(p-k) \times 1$. So H_{03} can be expressed as

$$H_{03} : E(\underline{X}_2 | \underline{X}_1)$$

is the same in both π_1 and π_2 .

All three null hypotheses (I), (II) and (III) can be tested by Rao's U statistic as follows:

$$\frac{N_1 + N_2 - p - 1}{p - 1} \cdot \frac{N_1 N_2 (D_p^2 - D_1^2)}{(N_1 + N_2)(N_1 + N_2 - 2) + N_1 N_2 D_1^2} \sim F_{p-1, N_1 + N_2 - p - 1}$$

which is obtained using

$$\frac{N_1 + N_2 - p - 1}{p - 1} \cdot \left(\frac{1}{U} - 1 \right).$$

When H_0 is accepted, it means that the variables $x_{k+1}, x_{k+2}, \dots, x_p$ do not have additional discriminating ability, once x_1, x_2, \dots, x_k have already been considered.