# UNIT 3

## The Problem of classification

The Problem of classification arises when an investigator makes a number of measurements on an individuals and wishes to classify the individual into one of several categories on the basis of these measurements. An individual is considered as a random observation from the population, the question is : ~~Given an individual is considered as a random observation from the population~~ Given an individual with certain measurements, from which population did the person arise?

This Statistical technique developed for estimation, hypothesis testing and confidence statement based on exact specification of the response variate. In the applied source one kind of multivariate problem frequently occur in which on observations must be assign in some optimum fashion to one of several population. This kind of multivariate technique is called as problem of classification.

For example, a banking officer may ~~which work~~ wish to classify the loan application has low or high credit risk on the basis of the elements of certain accounting statements.

## Standards of good classification

In constructing ~~the~~ a procedure of classification, it is desired to minimize the probability of misclassification. For convenience we shall now consider the case of only two categories.

Suppose an individual is an observation from either population $\Pi_1$ or population $\Pi_2$. The classification of an observation depends on the vector of measurements
$$X = (X_1, X_2, \ldots X_p)'$$ on that individual.

we set up a rule that if an individual is characterised by certain sets of values of $x_1, x_2, \ldots x_p$ that person will be classified as from $\Pi_1$, if other values as from $\Pi_2$.

We can think of an observation as a point in a $p$-dimensional space. We divide this space into two regions. If the observation falls in $R_1$, we classify it as comming from population $\Pi_1$, and if it falls in $R_2$, we classify it as comming from population $\Pi_2$.

In the classification procedure, the statistician can make two kinds of errors in classification. If the individual is actually from $\Pi_1$, the statistician can classify him or as comming from population $\Pi_2$. If the individual is actually from $\Pi_2$, the statistician can classify him as comming from $\Pi_1$.

Let the cost of the first type of misclassification be $c\left(\frac{2}{1}\right)$ $(>0)$ and let the cost of misclassifying an individual from $\Pi_2$ as from $\Pi_1$ be $c\left(\frac{1}{2}\right)(>0)$. These costs may be measured in any kind of units. There is no rewards for correct classification.

The following table indicates the costs of correct and incorrect classifications. Clearly a good classification procedure is one that minimizes the cost of misclassification.

Statistician's Decision

|  |  | $\Pi_1$ | $\Pi_2$ |
|---|---|---|---|
| Population | $\Pi_1$ | $0$ | $c\left(\frac{2}{1}\right)$ |
|  | $\Pi_2$ | $c\left(\frac{1}{2}\right)$ | $0$ |

Baye's Procedure for minimizing expected loss of misclassification –
Two cases of two populations

Let the Probability that the observation, comes from population $\pi_1$ be $q_1$ and from $\pi_2$ be $q_2$. Such that $q_1 + q_2 = 1$.

Let the density of the population $\pi_1$ be $P_1(x)$ and for the population $\pi_2$ be $P_2(x)$.

If we have a region $R_1$ of classification as from $\pi_1$, the probability of correctly classifying an observation that actually is drawn from population $\pi_1$ is

$$P(1/1, R) = \int_{R_1} P_1(x) \, dx \longrightarrow ①$$

where $dx = dx_1 \ldots dx_p$

and the probability of misclassification of an observation from $\pi_1$ is

$$P(2/1, R) = \int_{R_2} P_1(x) \, dx \longrightarrow ②$$

III$^y$ the probability of correctly classifying an observation from $\pi_2$ is

$$P(2/2, R) = \int_{R_2} P_2(x) \, dx \longrightarrow ③$$

and the probability of misclassifying such an observation is

$$P(1/2, R) = \int_{R_1} P_2(x) \, dx \longrightarrow ④$$

Since the probability of drawing an observation from $\pi_1$ is $q_1$, the probability of drawing an observation from $\pi_1$ and correctly classifying it is

$$q_1 \cdot P(1/1, R)$$

III$^y$ the probability of drawing an observation from $\pi_2$ and correctly classifying it is $q_1 P(2/1, R)$.

On the other hand, the probability of drawing an observation from $\Pi_2$ and correctly classifying it is $q_2 \, P(2/2, R)$. Similarly, the probability of drawing an observation from $\Pi_2$ and misclassifying it is $q_2 \cdot P(1/2, R)$.

Therefore, the expected loss of misclassification costs is the sum of the products of costs of misclassifications with their respective probabilities of occurrence.

ie
$$C(2/1) \; P(2/1, R) \, q_1 + C(1/2) \; P(1/2, R) \, q_2$$

$$\longrightarrow \text{⑤}$$

We wish to minimize this average loss, to divide our space into region $R_1$ and $R_2$ such that the expected loss is as small as possible. A procedure that minimizes equation ⑤ for given $q_1$ and $q_2$ is called a Bayes procedure.

## Procedure for classification of into one of two populations with known probability distributions

When a prior probabilities are known, we can define joint probabilities of the population and the observation set of variables.

The probability that an observation comes from the population $\pi_1$ and that each variate is less than the corresponding component in $Y$ is

$$\int_{-\infty}^{y_p} \cdots \int_{\infty}^{y_1} q_1 \, p_1(x) \, dx_1 \cdots dx_p$$

The conditional probability of for the population comming from $\pi_1$ given an observation $x$ is

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

The expected loss of misclassification is

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

Suppose $C(2/1) = C(1/2) = 1$ then the expected loss is

$$q_1 \, P(2/1, R) + q_2 \, P(1/2, R)$$

ie $q_1 \int_{R_2} P_1(x) \, dx + q_2 \int_{R_1} P_2(x) \, dx \longrightarrow ②$

This is also the probability of misclassification. For a given observation observed point $x$, we minimize the probability of misclassification by assigning the population that has the higher conditional probability.

If $\dfrac{q_1 \, p_1(x)}{q_1 \, p_1(x) + q_2 \, p_2(x)} \geqslant \dfrac{q_2 \, p_2(x)}{q_1 \, p_1(x) + q_2 \, p_2(x)}$ we choose population $\pi_1$

Otherwise we choose population $\pi_2$.

Since we minimize the probability of misclassification at each point, we minimize it over the whole space. Thus the rule is,

$$R_1 : q_1 P_1(x) \geqslant q_2 P_2(x)$$
$$R_2 : q_1 P_1(x) < q_2 P_2(x) \Bigg\} \longrightarrow ③$$

If $P_r \left\{ \dfrac{P_1(x) = q_2}{P_2(x) = q_1} \Big/ \pi_i \right\} = 0, \qquad i = 1, 2$

Then the Bayes Procedure is unique except for sets of probability zero.

∴ choose $R_1$ and $R_2$ so as to minimize equation ② the solution is equation ③.

we wish to minimize equation ①, which can be written as

$$c(2/1) \, q_1 \int_{R_2} P_1(x) \, dx + c(1/2) \, q_2 \int_{R_1} P_2(x) \, dx$$

we choose $R_1$ and $R_2$ according to

$$R_1 : c(2/1) \, q_1 \, P_1(x) \geqslant c(1/2) \, q_2 \, P_2(x)$$
$$R_2 : c(2/1) \, q_1 \, P_1(x) < c(1/2) \, q_2 \, P_2(x) \Bigg\} \longrightarrow ④$$

~~Since~~ Since, $c(2/1) \, q_1$ and $c(1/2) \, q_2$ are non negative constants.

Another way of writing equation ④ is

$$R_1 : \frac{P_1(x)}{P_2(x)} \geqslant \frac{c(1/2) \, q_2}{c(2/1) \, q_1}$$

$$R_2 : \frac{P_1(x)}{P_2(x)} < \frac{c(1/2) \, q_2}{c(2/1) \, q_1}$$

## Classification into one of two known multivariate normal Populations.

Let us Consider in the Case of two multivariate normal populations with equal Covariance matrices ie $(\Sigma_1 = \Sigma_2 = \Sigma)$ namely $N(\mu^{(1)}, \Sigma)$ and $N(\mu^{(2)}, \Sigma)$. where $\mu^{(i)} = (\mu_1^{(i)}, \mu_2^{(i)} \dots \mu_p^{(i)})'$ is the vector of means of the $i^{th}$ population. $i = 1, 2$ and $\Sigma$ is the matrix of variances and Covariances of the each population. Then the $i^{th}$ density is

$$P_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{\frac{1}{2}}} \exp\left\{-\frac{1}{2}(x - \mu^{(i)})' \Sigma^{-1} (x - \mu^{(i)})\right\} \longrightarrow ①$$

The ratio of the densities is

$$\frac{P_1(x)}{P_2(x)} = \frac{\exp\left\{-\frac{1}{2}(x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)})\right\}}{\exp\left\{-\frac{1}{2}(x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)})\right\}}$$

$$= \exp\left\{-\frac{1}{2}\left[(x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)}) - (x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)})\right]\right\} \longrightarrow ②$$

Let $R_1$, is the region of classification into $\pi_1$ is the set of values $x$ for which equation ② is greater than or equal to $k$. (Choose $k$ arbitrarily) then the inequality can be written as (using log.)

$$-\frac{1}{2}\left[(x - \mu^{(1)})' \Sigma^{-1} (x - \mu^{(1)}) - (x - \mu^{(2)})' \Sigma^{-1} (x - \mu^{(2)})\right] \geqslant \log k \longrightarrow ③$$

The LHS of ③ Can be expressed as

$$-\frac{1}{2}\left[x'\Sigma^{-1}x - x'\Sigma^{-1}\mu^{(1)} - \mu^{(1)'}\Sigma^{-1}x + \mu^{(1)'}\Sigma^{-1}\mu^{(1)} - x'\Sigma^{-1}x \right.$$
$$\left. + x'\Sigma^{-1}\mu^{(2)} + \mu^{(2)'}\Sigma^{-1}x - \mu^{(2)'}\Sigma^{-1}\mu^{(2)}\right]$$

$$= -\frac{1}{2}\left[-2x'\Sigma^{-1}\mu^{(1)} + 2x'\Sigma^{-1}\mu^{(2)}\right] - \frac{1}{2}\left[\mu^{(1)'}\Sigma^{-1}\mu^{(1)} - \mu^{(2)'}\Sigma^{-1}\mu^{(2)}\right]$$

$$= x'\Sigma^{-1}\mu^{(1)} - x'\Sigma^{-1}\mu^{(2)} - \frac{1}{2}\left[\mu^{(1)'}\Sigma^{-1}\mu^{(1)} - \mu^{(1)'}\Sigma^{-1}\mu^{(2)} + \mu^{(2)'}\Sigma^{-1}\mu^{(1)} - \mu^{(2)'}\Sigma^{-1}\mu^{(2)}\right]$$

$$= x'\Sigma^{-1}\mu^{(1)} - x'\Sigma^{-1}\mu^{(2)} - \frac{1}{2}\left[\mu^{(1)} + \mu^{(2)'}\Sigma^{-1}\mu^{(1)} - \mu^{(2)}\right]$$

$$= x'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right) - \frac{1}{2}\left(\mu^{(1)} + \mu^{(2)}\right)'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right)$$
$$\hookrightarrow ④$$

The first term of equation ④ is well known discriminant function.

It is a linear function of the components of the observation vector.

If $\pi_i$ has the density ① $i = 1, 2$ the best region of classification are

$$R_1 : x'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right) - \frac{1}{2}\left(\mu^{(1)} + \mu^{(2)}\right)'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right)$$
$$\geqslant \log k$$

$$R_2 : x'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right) - \frac{1}{2}\left(\mu^{(1)} + \mu^{(2)}\right)'\Sigma^{-1}\left(\mu^{(1)} - \mu^{(2)}\right) < \log k$$

If a prior probabilities $q_1$ and $q_2$ are known then $k$ is given by

$$k = \frac{q_2 \, C(1/2)}{q_1 \, C(2/1)}$$

In a Particular case of the two population being equally likely and the Cost being equal, $k=1$ and $k=0$

$\log k = 0$. Then the region of classification are

$R_1$: $x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) \geqslant \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})$

$R_2$: $x' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)}) < \frac{1}{2} (\mu^{(1)} + \mu^{(2)})' \Sigma^{-1} (\mu^{(1)} - \mu^{(2)})$

If we donot have a prior probabilities we may select $\log k = c$ (say) on the basis of the expected loss. due to misclassification equal.

# Classification ~~into~~ with Several populations

Let $f_i(x)$ be the density associated with population $\pi_i$

$i = 1, 2 \cdots g$   $(g > 2)$.

Let $P_i$ is the prior probability of population $\pi_i$

$c(k/i)$ is the cost of allocating an item to $\pi_k$

In fact it belongs to $\pi_i$   for $k, i = 1, 2 \cdots g$

for $k = i$   $c(i/i) = 0$

Finally let $R_k$ be the set of $x$'s classified as $\pi_k$

and $P(k/i) = P(\text{classify item as } \pi_k/\pi_i) = \int_{R_k} f_i(x) \, dx$

for $k = i = 1, 2, \cdots g$ with $P(i/i) = 1 - \sum_{\substack{k=1 \\ k \neq i}}^{g} P(k/i)$

The Conditional expected cost of misclassifying an $x$ from $\pi_1$ into $\pi_2$ or $\pi_3 \cdots \pi_g$ is

$$ECM(1) = P(2/1) \, c(2/1) + P(3/1) \, c(3/1) + \cdots + P(g/1) \, c(g/1)$$

$$= \sum_{k=2}^{g} P(k/1) \cdot c(k/1)$$

This Conditional expected cost occurs with prior probability $P_1$, the population of $\pi_1$.

In a similar manner, we can obtain the conditional expected costs of misclassification, $ECM(2), \cdots ECM(g)$.

Multiplying each Conditional ECM by its prior probability and ~~sum~~ summing gives the overall ECM.

$\therefore ECM = P_1 \, ECM(1) + P_2 \, ECM(2) + \cdots + P_g \, ECM(g)$

$$= P_1 \left[ \sum_{k=2}^{g} P(k/1) \, c(k/1) \right] + P_2 \left[ \sum_{\substack{k=1 \\ k \neq 2}}^{g} P(k/2) \, c(k/2) \right] + \cdots$$

$$+ P_g \left[ \sum_{k=1}^{g-1} P(k/g) \, c(k/g) \right]$$

$$\therefore \; ECM = \sum_{i=1}^{g} P_i \left[ \sum_{\substack{k=1 \\ k \neq i}}^{g} P(k/i) \, c(k/i) \right] \longrightarrow \text{①}$$

Determining an optimal classification procedure amounts to choosing the mutually exclusive and exhaustive classification region $R_1, R_2, \cdots R_g$ such that equation ① is a minimum.

The classification regions that minimize the ECM of equation ① are defined ~~as follows~~ by allocating x to that population $\pi_k$ : $k = 1, 2 \cdots g$ for which

$$\sum_{\substack{i=1 \\ k \neq i}}^{g} P_i \, f_i(x) \, c(k/i) \longrightarrow \text{②}$$

is smallest.

Suppose all the misclassification cost are equal, in which case the minimum expected cost of misclassification rule is the minimum total probability of misclassification.

using the argument leading to equation ② we would allocate x to that population $\pi_k$ : where $k = 1, 2 \cdots g$

for which $\sum\limits_{i=1}^{g} p_i \, f_i(x) \longrightarrow ③$ is minimum,

when the misclassification cost are same,
the minimum ECM rule has the following form

Allocate x to $\Pi_K$ if

$$p_K \, f_K(x) > p_i \, f_i(x) \qquad \text{for all} \quad i \neq k$$

or ~~the~~ equivalently,

Allocate x to $\Pi_K$ if

$$\ln \left[ p_K \, f_K(x) \right] > \ln \left[ p_i \, f_i(x) \right] \text{ for all } i \neq k$$

<u>Classification into one of Several Multivariate</u>
<u>Normal Population.</u>

we know that

$$f_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma_i|^{1/2}} \exp\left\{-\frac{1}{2}(x-\mu_i)' \Sigma_i^{-1}(x-\mu_i)\right\} \qquad i=1,2\cdots g$$

①

are multivariate normal densities with mean vector $\mu_i$
and Covariance matrices $\Sigma_i$.

If $c(i/i) = 0$ and $c(k/i) = 1$ for $k \neq i$
the minimum expected cost misclassification (ECM) rule becomes,
allocate $x$ to $\pi_k$ if

~~$\frac{}{}$~~

$$P_k f_x(x) = P_k (2\pi)^{-p/2} |\Sigma_k|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2}(x-\mu_k)' \Sigma_k^{-1}(x-\mu_k)\right\}$$

Taking log on both sides

$$\log\left[P_k f_x(n)\right] = \log P_k - \frac{p}{2}\log(2\pi) - \frac{1}{2}\log|\Sigma_k| - \frac{1}{2}(x-\mu_k)' \Sigma_k^{-1}(x-\mu_k)$$

$$= \text{Max}_i \log P_i f_i(x) \qquad\longrightarrow ②$$

Here the Constant term $\frac{p}{2}\log(2\pi)$ can be ignored in ②
Since it is fixed for all the ~~terms~~ populations

∴ we define the quadratic discrimination Score
for the $i^{th}$ population

$$d \ Q_i(x) = -\frac{1}{2}\log|\Sigma_i| - \frac{1}{2}(x-\mu_i)' \Sigma_i^{-1}(x-\mu_i) + \log P_i$$

$$i=1,2\cdots g$$

$$\longrightarrow ③$$

The quadratic Score $d Q_i(x)$ is
Composed of Contributions from the generalised variance
$|\Sigma_i|$, the Prior Probability $P_i$ and the Squared distance

from the $X$ to the population $\Pi_i$

using discriminant Score to the classification

rule in the equation ② , ~~it becomes~~

Allocate $X_i$ to $\Pi_k$ if

$$d\, Q_k(x) = \text{maximum of } d\,Q_1(x),\ d\,Q_2(x) \ldots \ d\,Q_g(x) \longrightarrow ④$$

where $d\,Q_k(x)$ is given in equation ③, i.e.,

In Particular, $\mu_i$ and $\Sigma_i$ are unknown we use the relevant estimates $\bar{x}_i$ and $S_i$ then ③ becomes

$$d\,Q_i(x) = -\tfrac{1}{2} \log |S_i| - \tfrac{1}{2} (x-\bar{x}_i)' S_i^{-1} (x-\bar{x}_i) + \log P_i \longrightarrow ⑤ \quad \text{if } \Sigma_i = \Sigma$$

$$d\,Q_i(x) = -\tfrac{1}{2} \log |\Sigma| - \tfrac{1}{2} (x-\mu_i)' \Sigma^{-1} (x-\mu_i) + \log P_i$$

$$= -\tfrac{1}{2} \log |\Sigma| - \tfrac{1}{2} x' \Sigma^{-1} x + \tfrac{1}{2} x' \Sigma^{-1} \mu_i + \tfrac{1}{2} \mu_i' \Sigma^{-1} x$$

$$\qquad\qquad -\tfrac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \log P_i$$

$$= -\tfrac{1}{2} \log |\Sigma| - \tfrac{1}{2} x' \Sigma^{-1} x + \mu_i' \Sigma^{-1} x - \tfrac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \log P_i$$

The first two terms are ignored, since ~~they~~ they are Common for $d\,Q_1(x),\ d\,Q_2(x) \ldots d\,Q_g(x)$

∴ The linear discriminant Score.

$$d\,Q_i(x) = \mu_i' \Sigma^{-1} x - \tfrac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \log P_i \longrightarrow ⑥$$

∴ The classification rule becomes; allocate $x$ to $\Pi_k$ if

$$d\,Q_k(x) = \text{Max} \left\{ d\,Q_1(x),\ d\,Q_2(x) \ldots d\,Q_g(x) \right\}$$

where $d\,Q_i(x)$ is given in ⑥ for $i = 1, 2, \ldots g$

# Discriminant Analysis

Discriminant analysis and classification are the multivariate techniques concerned with separating distinct set of objects. or observations and with allocating new objects to previously defined groups.

    Discriminant analysis is rather exploratory in nature. As a separatory procedure, it is often employed on a one-time basis inorder to investigate observed differences when caused relationships are not well understood.

    The basic idea of discriminant analysis consist of assigning an individual or a group of individual to one or several known or ~~unknown~~ unknown distinct population on the basis of the observations on several characteristics of the individual. or the group.

    In scientific literature, discriminant analysis has many synonyms such as classification, Identification, prediction. and selection depending on the types of scientific area inwhich is used.

    Thus the immediate goal of discriminant analysis is to describe either graphically or algebraically, the differential features of objectives from several known Collections. we try to find 'discriminants' whose numerical values are such that the Collections are separated as much as possible.

# Fisher Discriminant function — Separation of population

Fisher idea was to transform the multivariate observations $x$ to univariate observation $Y$ such that the $Y$'s derived from population $\pi_1$ and $\pi_2$ were separated as much as possible

If we let $\mu_{1Y}$ be the mean of $Y$ obtained from $x$ belonging to $\pi_1$ and $\mu_{2Y}$ be the mean of $Y$ obtained from $x$ belonging to $\pi_2$.

Fisher selected the linear combination $Y = \ell' x$, to maximize the distance between $\mu_{1Y}$ and $\mu_{2Y}$. by defining

$$\mu_1 = E(x/\pi_1) = \text{Expected value of multivariate observation from } \pi_1$$

$$\mu_2 = E(x/\pi_2) = \text{Expected value of multivariate observation from } \pi_2$$

and the Covariance matrix

$$\Sigma = E\left[(x - \mu_i)(x - \mu_i)'\right] \quad : \quad i = 1, 2$$

is same for the populations $\pi_1$ and $\pi_2$

Consider the linear Combination

$$Y = \ell' x$$
$$= (\mu_1 - \mu_2)' \Sigma^{-1} x$$
$$= (\bar{x}_1 - \bar{x}_2)' s^{-1} x$$

using this transformation, Y has a mean

$$\mu_{1Y} = E(Y/\pi_1) = E(\ell'x/\pi_1) = \ell'\mu_1$$

$$\mu_{2Y} = E(Y/\pi_2) = E(\ell'x/\pi_2) = \ell'\mu_2$$

and its variances are same,

ie, 
$$\sigma^2_Y = var(\ell'x)$$
$$= \ell' \, covar(x) \, \ell \qquad \{\because x \text{ is a vector}$$
$$= \ell' \, \Sigma \, \ell$$

The best linear combination $Y = \ell'x$, maximizes the ratio

$$\frac{\text{square distance between mean of } Y}{\text{variance of } Y} = \frac{(\mu_{1Y} - \mu_{2Y})^2}{\sigma^2_Y}$$

$$= \frac{(\ell'\mu_1 - \ell'\mu_2)^2}{\ell'\Sigma\ell}$$

$$= \frac{\ell'.(\mu_1-\mu_2)(\mu_1-\mu_2)'\ell}{\ell'\Sigma\ell}$$

$$= \frac{(\ell'(\mu_1-\mu_2))^2}{\ell'\Sigma\ell}$$

$$= \frac{\left[\ell'(\mu_1-\mu_2)\right]\left[\ell'(\mu_1-\mu_2)\right]'}{\ell'\Sigma\ell}$$

$$= \frac{(\ell's)^2}{\ell'\Sigma\ell}$$

where $S = \mu_1 - \mu_2$ and $\ell' = (\ell_1, \ell_2 \cdots \ell_p)$ is the fisher linear combinations of co-efficients.

Scanned by CamScanner

Maximizing the ratio, fisher introduce the linear Combination $Y = (M_1 - M_2)' \Sigma^{-1} x$ which is known as fisher linear discriminant function.

let $Y_0 = (M_1 - M_2)' \Sigma^{-1} x$ be the value of the discriminant function for a new observation $X_0$ and let

$$m = \tfrac{1}{2} (M_1 - M_2)' \Sigma^{-1} (M_1 + M_2)$$ be the mid point

between the two population mean. Therefore the Classification rule is,

allocate $X_0$ to $\pi_1$ if $Y_0 = (M_1 - M_2)' \Sigma^{-1} X_0 \geq m$

allocate $X_0$ to $\pi_2$, if $Y_0 = (M_1 - M_2)' \Sigma^{-1} X_0 < m$

alternatively, Subtract $m$ from $Y_0$ and Compare the results with zero, in Such Case, the rule becomes,

Allocate $X_0$ to $\pi_1$ if $Y_0 - m \geq 0$

allocate $X_0$ to $\pi_2$ if $Y_0 - m < 0$

---