

# STATISTICAL INFERENCE APPLIED TO CLASSIFICATORY PROBLEMS

By C. RADHAKRISHNA RAO  
*Statistical Laboratory, Calcutta*

## PART II: THE PROBLEM OF SELECTING INDIVIDUALS FOR VARIOUS DUTIES IN A SPECIFIED RATIO

### 1. INTRODUCTION

In a discussion of the author's paper on the 'Utilization of multiple measurements in problems of biological classification' read before the Royal Statistical Society in 1948, Mr. Patrick Slater posed the following problem.

"Could Mr. Rao give a rule of procedure to be followed when out of an undistributed population with parameters  $X_1, X_2, \dots, X_m$  and  $[M]$ ,  $n_a$  have to be selected for duty  $A$ ,  $n_b$  for duty  $B$  etc., given the means for each duty? This supposes that the whole population must be distributed among the duties and the quota and the desired standards for each duty are laid down."

In reply the author has given a complete solution without proof. The object of this paper is to study this problem in its generality, supply the necessary proofs and suggest some practical methods of arriving at the desired classification. Although the main results are deducible from the general theory of decision functions as developed by Prof. Abraham Wald in his book on 'Statistical decision functions' and in his lectures at the Indian Statistical Institute, the computational procedure presents some difficulty. The two lemmas given in the next section are useful for this purpose. It is also shown that the solution derived by the method of maximum likelihood possesses some important properties.

### 2. TWO LEMMAS

Consider an array of elements

$$\begin{array}{cccc} a_{11} & a_{21} & \dots & a_{p1} \\ a_{12} & a_{22} & \dots & a_{p2} \\ \cdot & \cdot & \dots & \cdot \\ a_{1n} & a_{2n} & \dots & a_{pn} \end{array}$$

consisting of  $p$  columns and  $n$  rows. Let  $P$  denote the product and  $S$  the sum of  $n$  elements chosen one from each row such that the total number of elements coming from the first column is equal to a specified value  $n_1$ , from the second  $n_2$ , and so on from the  $p$ -th column  $n_p$ . Obviously  $n \geq p$  if no  $n_i$  is zero and  $n = n_1 + \dots + n_p$ .

Lemma 1: If the elements  $a_{ij}$  are not negative and if there exist quantities  $\lambda_1, \lambda_2, \dots, \lambda_p$  such that each element  $a_{ik}$  of  $n_i$  elements chosen from the  $i$ -th column satisfies the relationships

$$\lambda_i a_{ik} \geq \lambda_j a_{jk} \quad \text{for } j = 1, 2, \dots, p \quad \dots (2.1.1)$$

and if similar relationships are satisfied for all  $i = 1, 2, \dots, p$  with the same set of  $\lambda$ 's, then for this choice of  $n_1, n_2, \dots, n_p$  elements the product  $P$  defined above is a maximum.

To prove this consider any other choice of  $n_i$  elements

$$a_{1r}, a_{1s}, \dots$$

from the  $i$ -th column. Remembering that the second subscript refers to the row number, the elements from the  $r$ th,  $s$ th, ... rows occurring in the former selection may be represented by

$$a_{br}, a_{bs}, \dots$$

from which by definition it follows that the product

$$a_{br} a_{cs} \dots$$

is not less than

$$\frac{\lambda_1}{\lambda_b} a_{1r} \frac{\lambda_1}{\lambda_c} a_{1s} \dots$$

The  $\lambda$ 's are positive by definition so that division by a  $\lambda$  does not change the inequality sign. Considering all the groups of elements in the second selection we find

$$\begin{aligned} \Pi a_{br} a_{cs} \dots &\geq \Pi \frac{\lambda_1}{\lambda_b} \frac{\lambda_1}{\lambda_c} \dots a_{1r} a_{1s} \dots \\ &\geq \Pi a_{1r} a_{1s} \dots \end{aligned}$$

since

$$\Pi \frac{\lambda_1}{\lambda_b} \frac{\lambda_1}{\lambda_c} \dots = 1$$

all the  $\lambda$ 's in the numerator cancelling with those in the denominator in the final product.

Corollary 1.1: If the object is to maximise the product without the restriction on the number of elements coming from each column then, obviously, the best procedure is to choose the biggest element from each row.

Corollary 1.2.: The product  $P$  will be a minimum if the inequality relationship (2.1.1) is reversed.

## STATISTICAL INFERENCE APPLIED TO CLASSIFICATORY PROBLEMS

Corollary 1.3: If  $p=2$  the method described in lemma 1 reduces to evaluating the ratios

$$\frac{a_{11}}{a_{21}}, \frac{a_{12}}{a_{22}}, \dots, \frac{a_{1n}}{a_{2n}}$$

and arranging them in descending order of magnitude and choosing the elements in the numerator from the first  $n_1$  ratios and the elements in the denominator from the second  $n_2$  ratios.

Lemma 2: If there exist quantities  $\mu_1, \mu_2, \dots, \mu_p$  such that each element  $a_{ik}$  of the  $n_1$  elements chosen from the  $i$ th column satisfies the relationships

$$a_{ik} + \mu_i \geq a_{jk} + \mu_j \text{ for all } j = 1, 2, \dots, p$$

and if similar relationships hold for all  $i = 1, 2, \dots, p$  with the same set of  $\mu$ 's, then for this choice of  $n_1, \dots, n_p$  elements the sum  $S$  defined above is a maximum.

The result follows from lemma 1 by considering  $\exp(\sum a_{ik})$  and maximising the product. The existence of  $\mu$ 's leads to the existence of positive quantities  $\lambda$ 's used in lemma 1. The sum  $S$  is minimised when the reverse relationships hold good.

### 3. PROBLEMS OF TWO GROUPS

Let us consider two groups characterised by probability densities  $f_1(x|\theta_1)$  and  $f_2(x|\theta_2)$  where  $x$  stands for all the available set of measurements and  $\theta$  for all the parameters. These may be simply denoted by  $f_1(x)$  and  $f_2(x)$ .

Two types of problems arise. Firstly a sample of size  $n_1 + n_2$  may be drawn from a group with two populations mixed in the ratio  $n_1 : n_2$  and a division of the sample in two groups of  $n_1$  and  $n_2$  individuals is required. Secondly two samples of sizes  $n_1$  and  $n_2$  drawn independently from the first and second groups may get mixed. In this case what is the best method of separating the individuals belonging to two different groups? The difference between the first and second problems is that in the latter every sample is known to consist of  $n_1$  individuals from the first group and  $n_2$  from the second while in the former no such information is available, the sample being drawn at random from a mixed population.

#### 3.1. Solution to the first problem

Let  $x_1, \dots, x_n$  represent the measurements on  $n$  individuals. An observed earlier  $x_i$  will stand for all the available set of measurement on the  $i$ th individual. The probability of the set is

$$\Pi [n_1 f_1(x_i) + n_2 f_2(x_i)]$$

Consider the following set of functions

$$\delta_i(x_1, \dots, x_n), \delta'_i(x_1, \dots, x_n); i = 1, \dots, n$$

which can be represented simply as  $\delta_1, \delta_1'$  satisfying the conditions

$$\delta_1 = 0 \text{ or } 1, \quad \delta_1 + \delta_1' = 1$$

and 
$$\sum_1^n \delta_1 = n_1, \quad \sum_1^n \delta_1' = n_2$$

where  $n_1$  and  $n_2$  are the specified numbers. If the individual with measurements  $x_1$  is assigned to the first group when  $\delta_1 = 1$  and to the second when  $\delta_1' = 1$  then the above set of functions constitute a decision rule for the problem (Wald, 1950). The problem is then to construct the above functions such that the expected risk associated with this decision rule is a minimum. To calculate the expected risk we need know as a datum of the problem the loss of assigning an individual of one group to another. If  $r_{12}$  represents the loss in assigning an individual of the first to the second and  $r_{21}$  in the other case then the quantity to be minimised is the expected value of

$$\begin{aligned} & \sum_1^n \frac{r_{21} \delta_1 n_1 f_1(x_1) + r_{12} \delta_1' n_2 f_2(x_1)}{n_1 f_1(x_1) + n_2 f_2(x_1)} \\ &= \sum_1^n (\delta_1 a_{11} + \delta_1' a_{12}) \end{aligned} \quad \dots (3.1.1)$$

where  $a_{11} = r_{21} n_2 f_2 / (n_1 f_1 + n_2 f_2), a_{12} = r_{12} n_1 f_1 / (n_1 f_1 + n_2 f_2)$

The expected loss will be a minimum if  $\delta_1$  and  $\delta_1'$  are chosen such that the expression (3.1.1) has the least value. The problem is the same as that treated in lemma 2 leading to the solution

$$\begin{aligned} \delta_1 &= 1 & \text{if } a_{11} + \mu_1 < a_{12} + \mu_2 \\ &= 0 & \text{if } a_{11} + \mu_1 \geq a_{12} + \mu_2 \end{aligned}$$

where  $\mu_1$  and  $\mu_2$  are suitably chosen to satisfy the condition  $\sum \delta_1 = n_1$ . Now  $a_{11} - a_{12} < \mu_2 - \mu_1$  implies that

$$\frac{r_{21} n_2 f_2(x_1) - r_{12} n_1 f_1(x_1)}{n_1 f_1(x_1) + n_2 f_2(x_1)} < \lambda$$

or 
$$\frac{f_1(x_1)}{f_2(x_1)} > \lambda$$

so that the decision rule reduces to the evaluation of the likelihood ratios  $\lambda_1 = f_1(x_1)/f_2(x_1)$  and assigning, all the individuals with highest  $n_1$  values of the ratios to the first group and the rest to the second. Or if the differences  $a_1 = \log f_1(x_1) - \log f_2(x_1)$  are first calculated then the highest  $n_1$  values will form the basis for selection to the first group. Fortunately the decision rule is independent of the a-priori probabilities and also the loss function.

## STATISTICAL INFERENCE APPLIED TO CLASSIFICATORY PROBLEMS

Corresponding to every decision rule we can set up a density function of the observations by considering all individuals assigned to the first group as having been drawn at random from the first group and similarly for the second. By using lemma 1 we find that the best decision rule found above maximises the corresponding probability density. As shown below this forms the basis on which the solution of the second problem depends.

### 3.2. Solution to the second problem

In this problem the mixture is known to consist of  $n_1$  individuals drawn from the first group and  $n_2$  from the second. The observations  $x_1, x_2, \dots, x_n$  could have arisen in  ${}^nC_{n_1}$  ways, any subset of  $n_1$  observations belonging to the first group. The probability density of the observations is equal to the sum of the densities associated with  ${}^nC_{n_1}$  ways of splitting the sample. If  $x_p, x_q, \dots$  represents a division into two groups of sizes  $n_1$  and  $n_2$  then the probability density of the observations can be written as

$$P(x) = \sum f_1(x_p) f_1(x_q) \dots f_2(x_r) f_2(x_s) \dots$$

where the summation is over  ${}^nC_{n_1}$  such terms. Corresponding to any one of  ${}^nC_{n_1}$  possible decision rules the loss relative to the given set of observations is  $1/P(x)$  times

$$\sum l(a, b, \dots, p, q, \dots) f_1(x_p) f_1(x_q) \dots f_2(x_r) f_2(x_s) \dots \quad \dots \quad (3.2.1)$$

where  $l(a, b, \dots, p, q, \dots)$  is the loss incurred in adopting an assigned decision rule when in fact  $x_p, x_q, \dots$  come from the first group and  $x_r, x_s, \dots$  from the second. The loss will generally be a function of the number of wrong classifications only. That decision rule for which the expression (3.2.1) is the least is the best in the sense that it minimises the expected loss. The solution depends on the evaluation of the expression (3.2.1) for each of the  ${}^nC_{n_1}$  decision rules and this makes its application a little difficult in practice even when  $n$  is small.

As an alternative we may try to minimise the maximum loss incurred by following a decision rule. If we suppose that the loss function is proportional to the number of wrong classifications then the maximum loss occurs when all the individuals in the smaller group are wrongly classified. With this assumption it can be shown that the division of the sample corresponding to the maximum probability density supplies the best possible solution to the problem. This solution may be referred to as the maximum likelihood solution; we consider the  ${}^nC_{n_1}$  ways of splitting the sample as associated with  ${}^nC_{n_1}$  different hypotheses concerning the individuals in the sample and choose that hypothesis which has the maximum likelihood.

To prove the property referred above, consider any other decision rule leading to a division of the sample

$$x_{a_1}, x_{a_2}, \dots, x_{b_1}, x_{b_2}, \dots \quad \text{and} \quad x_{c_1}, x_{c_2}, \dots, x_{d_1}, x_{d_2}, \dots$$

into two groups of sizes  $n_1$  and  $n_2$  and compare with the division

$$x_{11}, x_{12}, \dots, x_{1r_1}, x_{1s_1}, \dots \text{ and } x_{21}, x_{22}, \dots, x_{2r_2}, x_{2s_2}, \dots$$

associated with the maximum density. The measurements classified in the same way by the two decision rules are represented by  $x_{11}, x_{12}, \dots$  for the first group and by  $x_{21}, x_{22}, \dots$  for the second group. By definition

$$\frac{f_1(x_{1j})}{f_2(x_{1j})} \geq \frac{f_1(x_{2j})}{f_2(x_{2j})} \quad \dots (3.2.2)$$

and the same is true for the product of a number of ratios involving  $x_r$  and the product of the same number of ratios involving  $x_b$ .

Let  $n_2 < n_1$  without loss of generality. In this case maximum loss occurs, by following the first decision rule when

$$x_{c1}, x_{c2}, \dots, x_{c1}, x_{c2}, \dots \quad \dots (3.2.3)$$

arise from the first group in which case a subset  $n_2$  out of

$$x_{c1}, x_{c2}, \dots, x_{b1}, x_{b2}, \dots \quad \dots (3.2.4)$$

arise from the second group. Let this subset be

$$x_{c1}, x_{c2}, \dots, x_{bp}, x_{bq}, \dots \quad \dots (3.2.5)$$

By replacing the subscript  $c$  by  $b$  we obtain the corresponding situation for the proposed maximum likelihood decision rule. The difference between the two probability densities associated with maximum errors for the two rules is, apart from a common multiplier, equal to

$$f_1(x_{cp})f_1(x_{c1}) \dots f_2(x_{cp})f_2(x_{bq}) \dots - f_1(x_{1p})f_1(x_{bq}) \cdot f_2(x_{cp})f_2(x_{c1}) \dots$$

which is not less than zero according to result (3.2.2). By considering all subsets of  $n_2$  observations out of (3.2.4) we exhaust all possible ways in which the observations leading to a maximum error according to the first rule can arise. To each such case there is a corresponding division leading to the maximum error for the proposed rule. But this division leads to a smaller probability density. The total chance of maximum error relative to the given set of observations is thus a minimum for the maximum likelihood decision rule.

#### 4. THE PROBLEM OF THREE GROUPS

As in the case of two groups we consider two situations, firstly when the sample consists of  $n$  individuals observed at random from a mixed population and secondly when the sample is a mixture of  $n_1$  individuals drawn from the first group,  $n_2$  from the second and  $n_3$  from the third. The problem in either case is to select  $n_1$  individuals for the first group,  $n_2$  for the second and  $n_3$  for the third, where  $n_1 + n_2 + n_3 = n$ .

STATISTICAL INFERENCE APPLIED TO CLASSIFICATORY PROBLEMS

Let  $f_1(x)$ ,  $f_2(x)$  and  $f_3(x)$  represent the probability densitoc of  $x$  for the three groups and  $\pi_1$ ,  $\pi_2$ ,  $\pi_3$  the proportions of mixture in the general population. The loss in assigning a person to the  $i$ -th group when, in fact, he belongs to the  $j$ -th group is denoted by  $r_{ij}$ . The aposteriori risks in assigning an individual with measurements  $x_1$  to the first, second and third groups are respectively equal to

$$a_{11} = \frac{\pi_1 f_1(x_1) r_{11} + \pi_2 f_2(x_1) r_{21}}{\pi_1 f_1(x_1) + \pi_2 f_2(x_1) + \pi_3 f_3(x_1)}$$

$$a_{21} = \frac{\pi_1 f_1(x_1) r_{12} + \pi_2 f_2(x_1) r_{22}}{\pi_1 f_1(x_1) + \pi_2 f_2(x_1) + \pi_3 f_3(x_1)}$$

and

$$a_{31} = \frac{\pi_1 f_1(x_1) r_{13} + \pi_2 f_2(x_1) r_{23}}{\pi_1 f_1(x_1) + \pi_2 f_2(x_1) + \pi_3 f_3(x_1)}$$

Consider a set of functions

$$\delta_i = 0 \text{ or } 1, \delta'_i = 0 \text{ or } 1, \delta''_i = 0 \text{ or } 1, \delta_i + \delta'_i + \delta''_i = 1$$

such that

$$\sum \delta_i = n_1, \sum \delta'_i = n_2, \sum \delta''_i = n_3$$

They define a decision rule if the individual with measurements  $x_1$  is assigned to the first group when  $\delta_1 = 1$  to the second when  $\delta'_1 = 1$  and to the third when  $\delta''_1 = 1$ . The aposteriori risk for such a selection procedure is

$$\sum_1 (\delta_1 a_{11} + \delta'_1 a_{21} + \delta''_1 a_{31}) \quad \dots (4.1)$$

The best decision rule is one which minimises the above expression. This is exactly the problem solved in lemma 2. The best solution is

$$\delta_1 = 1 \text{ if } a_{11} + \mu_1 \leq a_{21} + \mu_2, \quad a_{11} + \mu_1 \leq a_{31} + \mu_3$$

$$\delta'_1 = 1 \text{ if } a_{21} + \mu_2 \leq a_{11} + \mu_1, \quad a_{21} + \mu_2 \leq a_{31} + \mu_3$$

$$\delta''_1 = 1 \text{ if } a_{31} + \mu_3 \leq a_{11} + \mu_1, \quad a_{31} + \mu_3 \leq a_{21} + \mu_2$$

where  $\mu_1$ ,  $\mu_2$  and  $\mu_3$  are determined such that  $\sum \delta_i = n_1$ ,  $\sum \delta'_i = n_2$  and  $\sum \delta''_i = n_3$ . As it stands the problem of determination of  $\mu_1$ ,  $\mu_2$ ,  $\mu_3$  appears to be complicated. There is a geometrical device which is helpful in the solution of the problem. In higher-dimensional cases involving four or more groups the geometrical method cannot be applied. For three groups we replace  $a_{11}$ ,  $a_{21}$ ,  $a_{31}$  by two coordinates

$$X_1 = a_{11} - a_{21}, \quad Y_1 = a_{11} - a_{31}$$

and represent the  $n$  points  $(X_1, Y_1)$  on a two dimensional chart with rectangular axes. The problem is to determine a point  $(X_0, Y_0)$  on this chart such that the regions formed by the lines  $X = X_0$ ,  $Y = Y_0$  and  $Y - Y_0 = X - X_0$  contain the requisite number of points. This can be done by moving three thin rods  $OX'$ ,  $OY'$ ,  $OZ'$  fixed at the point

*O* as shown in the figure below, with *OX'* and *OY'* parallel to the *X* and *Y* axes and arrive at the required division by trial and error. It will help in this process if the numbers falling in the three regions are recorded for a few positions of the frame with the frame marked on the chart.

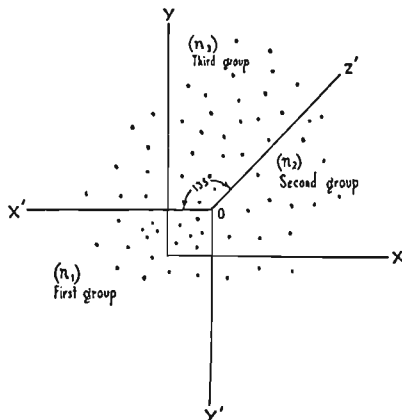


Figure 1

To solve the second problem when the sample consists of a mixture of  $n_1$  individuals drawn from the first group,  $n_2$  from the second and  $n_3$  from the third, one can set up the total risk relative to the given set of observations for each of  $n!/(n_1!n_2!n_3!)$  possible decision rules and choose that rule for which the risk is a minimum. This is very difficult in practice so that a simplified procedure is needed. As before we may choose that decision rule which leads to a division of the sample with the maximum probability density. This possesses an important property that the probability of the maximum number of wrong classifications for any one group is as small as possible.

The method of arriving at the required division is first to obtain the quantities

$$a_{1i} = \log f_1(x_i), \quad a_{2i} = \log f_2(x_i), \quad a_{3i} = \log f_3(x_i) \\ i = 1, 2, \dots, n$$

and plot the points

$$X_i = a_{2i} - a_{1i}, \quad Y_i = a_{3i} - a_{1i}$$



## STATISTICAL INFERENCE APPLIED TO CLASSIFICATORY PROBLEMS

and proceed geometrically as in figure 1. For this division the probability density will be a maximum.

### 5. ON A PROBLEM OF OPTIMUM SELECTION

Birnbaum and Chapman (1950) considered a problem of selecting candidates on the basis of  $p$  admission scores  $y_1, y_2, \dots, y_p$ . The object is to select those whose performance is expected to be better in the final test. The offered solution does not refer to a case where the scores of a number of individuals  $N$  have been observed but to a hypothetical set of individuals offering for the admission test. The former problem is often met with because the question asked is whom out of a number of individuals whose scores are available should be admitted. Let the scores of  $N$  individuals be represented by

$$\begin{array}{cccc} y_{11}, & y_{12}, & \dots, & y_{1p} \\ y_{21}, & y_{22}, & \dots, & y_{2p} \\ \dots & \dots & \dots & \dots \\ y_{N1}, & y_{N2}, & \dots, & y_{Np} \end{array}$$

To answer this problem we need the expected performance in the final test of an individual with the initial scores  $y_{11}, \dots, y_{1p}$ . Let this expected performance be

$$x_1 = \phi(y_{11}, y_{12}, \dots, y_{1p})$$

which actually stands for the regression equation of the final performance on the initial scores. The regression function which may be of any complicated type supplies us with the expected performances  $x_1, x_2, \dots, x_N$  of the candidates and these latter scores form the basis for selection. The regression function can be estimated on the basis of the previous information. No assumption of multivariate normality is needed.

For instance if  $k$ , a given number of seats, are available then the best plan is to admit  $k$  candidates corresponding to the  $k$  largest  $x$ 's because this maximises the expected performance under the condition that  $k$  have to be chosen.

A second alternative may be to admit as many as possible with the restriction that the expected average performance is not less than an assigned number  $x_0$ . The best plan is then to arrange the  $x$  scores in a decreasing order and find the cumulative averages from the top and admit all those for whom the cumulative average is greater than or equal to  $x_0$ . Obviously under such a selection procedure the maximum number is admitted subject to the condition that the expected average performance of the selected candidates is not less than  $x_0$ .

If the restriction is that the average performance of the chosen candidates should exceed a given value  $x_0$  with a probability greater than  $\beta$ , then again we start with the highest score of  $x$  and go on adding the others in the decreasing order of  $x$  till the required probability remains greater than  $\beta$ . The calculations are not however, simple.

If we consider a hypothetical set of candidates, a situation which may arise when the statistician is asked to give a uniform rule for independent recruitment at various places without specifying the numbers to be selected from each place, then what is needed is the determination of a score  $x_0$  leading to the selection of all individuals with expected  $x$  score (calculated on the basis of the admission  $y$  score) greater than or equal to  $x_0$ . For this the distribution of expected score  $x$  as a function of  $y$  has to be studied. Let this be  $f(x)$ . If the criterion is that the maximum number of candidates have to be admitted subject to the condition that the expected average performance is not less than  $x_0$  then  $x_0$  is determined from the formula

$$\int_{x_0}^{\infty} x f(x) dx = x_0 \int_{x_0}^{\infty} f(x) dx$$

in which case the expected proportion admitted is

$$\int_{x_0}^{\infty} f(x) dx$$

## REFERENCES

- DIRNBACH, Z. W. and CHAPMAN, D. G. (1950): On optimum selections from multinormal populations  
*Ann. Math. Stat.* 21, 443.
- RAO, C. R. (1948): The utilization of multiple measurements in problems of biological classification  
*J.R.S.S., Series B*, 10, 150.
- WALD, A. (1950): *Statistical decision functions*. John Wiley & Sons, New York.

*Paper received: January, 1951.*

This treatment of the problem then presupposes that the individual agent has stable probability estimates of finding someone to trade with, of that person being costly to bargain with, of that party being inclined to cheat on the terms of the agreement. In the case of externalities this implies that some unwanted side effects remain because of the uncertainty associated with undertaking the transaction that would eliminate them. The original problem was how to perform the market basket analysis which attempted to find all the interesting relationships between products bought in a given context. Association rule mining has been applied to LMS in order to reveal which contents students tend to access together, or which combination of tools they use. Subgroup discovery is usually seen as being different from classification, as it addresses different goals.