

ON A DISTRIBUTION-FREE TWO-WAY CLASSIFICATION

HIROSI HUDIMOTO

1. Introduction

Suppose that $X^{(1)}$ and $X^{(2)}$ are random variables with continuous distribution functions $F_1(x^{(1)})$ and $F_2(x^{(2)})$, respectively. When $F_2(x^{(2)})$ is completely known and $F_1(x^{(1)})$ can be estimated from a random sample, Z. W. Birnbaum [3] proposed the statistic

$$(1.1) \quad \hat{p} = \int_{-\infty}^{\infty} \hat{F}_1(t) dF_2(t)$$

as an unbiased estimator of

$$(1.2) \quad p = P\{X^{(1)} < X^{(2)}\} = \int_{-\infty}^{\infty} F_1(t) dF_2(t),$$

where $\hat{F}_1(t)$ is the empirical distribution function of $F_1(t)$. The statistic proposed by H. B. Mann and D. R. Whitney [6] is applicable as an unbiased estimator of p in the case when both $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are unknown and can be estimated from samples. In this paper, we shall consider the use of such statistics for classifying a sample (X_1, \dots, X_n) into either of two populations Π_1 and Π_2 characterized by $F_1(x^{(1)})$ and $F_2(x^{(2)})$, respectively, and give an evaluation of their success rates. In section 2, we treat the case where $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are known and in section 3, the case where both are unknown.

A fairly general consideration for classifying an individual into one of two or more populations has been presented by T. W. Anderson [1], [2], and C. R. Rao [7]. The particular case of two p -variate normal populations with equal covariance matrices has been fully discussed by A. Wald [10], T. W. Anderson [2], C. R. Rao [8], R. Sitgreaves [9] and the others. However, approaches in these papers cannot be applied to the non-parametric case where distributions are not specified but must be estimated from samples. Our method will be useful for such a case. Although it is assumed that $F_1(t) > F_2(t)$ for all t in the following sections, the method is applicable so far as it can be expected that $A = \int_{-\infty}^{\infty} [F_1(t) -$

$F_2(t)]d\frac{F_1(t)+F_2(t)}{2} \neq 0$ and may be favorable in cases where Δ is fairly different from zero.

2. The case of known distribution functions

When (X_1, \dots, X_n) is a random sample of size n which has come from either of two populations Π_1 and Π_2 , we wish to classify it correctly. Suppose Π_1 and Π_2 are characterized by univariate distribution functions $F_1(x^{(1)})$ and $F_2(x^{(2)})$, and $X^{(1)}$ and $X^{(2)}$ stand for their random variables, respectively. Then, it will be desirable for our procedure to have higher probability of correct classification.

As a preliminary for the next section, we shall first deal with the case where the distribution functions $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are completely known, and assume that $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are continuous and $F_1(t) > F_2(t)$ for all real t , i.e. $X^{(2)}$ is stochastically larger than $X^{(1)}$, for convenience. Consider statistics

$$(2.1) \quad \hat{p}_1 = \frac{1}{n} \sum_{k=1}^n F_1(X_k),$$

and

$$(2.2) \quad \hat{p}_2 = 1 - \frac{1}{n} \sum_{k=1}^n F_2(X_k).$$

Their expectations are

$$(2.3) \quad \mathcal{E}(\hat{p}_1 | \Pi_1) = \int_{-\infty}^{\infty} F_1(t) dF_1(t) = \frac{1}{2},$$

and

$$(2.4) \quad \begin{aligned} \mathcal{E}(\hat{p}_2 | \Pi_1) &= 1 - \int_{-\infty}^{\infty} F_2(t) dF_1(t) \\ &= \int_{-\infty}^{\infty} F_1(t) dF_2(t) = p, \end{aligned}$$

where $\mathcal{E}(\hat{p}_u | \Pi_v)$ denotes the conditional expectation of \hat{p}_u under the condition that (X_1, \dots, X_n) has come from Π_v ($u=1, 2$ and $v=1, 2$). Similarly, we have

$$(2.5) \quad \mathcal{E}(\hat{p}_1 | \Pi_2) = \int_{-\infty}^{\infty} F_1(t) dF_2(t) = p,$$

and

$$(2.6) \quad \mathcal{E}(\hat{p}_2 | \Pi_2) = 1 - \int_{-\infty}^{\infty} F_2(t) dF_1(t) = \frac{1}{2}.$$

The statistic of this kind was proposed by Z. W. Birnbaum [3]. p must be larger than $1/2$ because of the assumption that $X^{(2)}$ is stochastically larger than $X^{(1)}$. Thus, we shall take the following as our procedure.

Decision rule 1. If the continuous distribution functions $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are completely known and $F_1(t) > F_2(t)$ for all real t , classify (X_1, \dots, X_n) into Π_1 or Π_2 according to $\hat{p}_1 < \hat{p}_2$ or $\hat{p}_1 > \hat{p}_2$, respectively, where \hat{p}_1 and \hat{p}_2 are given by (2.1) and (2.2).

Let $H(x) = \frac{1}{2} \{F_1(x) + F_2(x)\}$. Then, $\hat{p}_1 \cong \hat{p}_2$ is equivalent to $\frac{1}{n} \sum_{k=1}^n H(X_k) \cong \frac{1}{2}$. The conditional expectations of $\frac{1}{n} \sum_{k=1}^n H(X_k)$ are

$$(2.7) \quad \mathcal{E} \left(\frac{1}{n} \sum_{k=1}^n H(X_k) \mid \Pi_1 \right) = \frac{1}{2} \left\{ 1 - \left(p - \frac{1}{2} \right) \right\} = \frac{1}{2} - \frac{D}{2}$$

and

$$(2.8) \quad \mathcal{E} \left(\frac{1}{n} \sum_{k=1}^n H(X_k) \mid \Pi_2 \right) = \frac{1}{2} \left\{ 1 + \left(p - \frac{1}{2} \right) \right\} = \frac{1}{2} + \frac{D}{2},$$

where $D = p - \frac{1}{2}$. Thus, if $\frac{1}{n} \sum_{k=1}^n [H(X_k) - \mathcal{E}\{H(X_k) \mid \Pi_1\}]$ is smaller than $\frac{D}{2}$, $\frac{1}{n} \sum_{k=1}^n H(X_k)$ will be smaller than $\frac{1}{2}$ when (X_1, \dots, X_n) comes from Π_1 , and if $\frac{1}{n} \sum_{k=1}^n [\mathcal{E}\{H(X_k) \mid \Pi_2\} - H(X_k)]$ is smaller than $\frac{D}{2}$, $\frac{1}{n} \sum_{k=1}^n H(X_k)$ will be larger than $\frac{1}{2}$ when (X_1, \dots, X_n) comes from Π_2 .

We shall use above relations to evaluate the success rate of the rule 1. W. Hoeffding [5] has given the certain probability inequalities for sums of bounded random variables. Since $H(X_1), \dots, H(X_n)$ are mutually independent and $0 \leq H(X_k) \leq 1$ ($k=1, \dots, n$), theorem 1 of his paper [5] gives us

$$(2.9) \quad P \left\{ \frac{1}{n} \sum_{k=1}^n H(X_k) < \frac{1}{2} \mid \Pi_1 \right\} \geq 1 - e^{-2n(D/2)^2}$$

and

$$(2.10) \quad P \left\{ \frac{1}{n} \sum_{k=1}^n H(X_k) > \frac{1}{2} \mid \Pi_2 \right\} \geq 1 - e^{-2n(D/2)^2},$$

where $P\{\dots | \prod_v\}$ denotes the conditional probability under the condition that each X_1, \dots, X_n have come from $\prod_v (v=1, 2)$. Thus, we obtain the following:

The success rate of the rule 1 is greater than or equal to $1 - \exp\left(-\frac{n\Delta^2}{2}\right)$,

where $\Delta = \int_{-\infty}^{\infty} F_1(t)dF_2(t) - \frac{1}{2}$ and n is the size of sample to be classified.

It will be found that this result is useful in the next section.

3. The case of unknown distribution functions

In this section, we shall consider the case where the distribution functions $F_1(x^{(1)})$ and $F_2(x^{(2)})$ are not known, but random samples of sizes n_1 and n_2 obtained from these populations are available. When $(X_1^{(1)}, \dots, X_{n_1}^{(1)})$ and $(X_1^{(2)}, \dots, X_{n_2}^{(2)})$ are the random samples from \prod_1 and \prod_2 , respectively, and (X_1, \dots, X_n) is a random sample to be classified into either of them, we shall use the statistics \hat{p}_1 and \hat{p}_2 instead of p_1 and p_2 :

$$(3.1) \quad \hat{p}_1 = \frac{U_{n_1 n}}{n_1 n} \quad \text{and} \quad \hat{p}_2 = \frac{U_{n n_2}}{n n_2},$$

where $U_{n_1 n}$ is the number of pairs $(X_i^{(1)}, X_k)$ such that $X_i^{(1)} < X_k$, and $U_{n n_2}$ is the number of pairs $(X_k, X_j^{(2)})$ such that $X_k < X_j^{(2)}$, $i=1, \dots, n_1$, $j=1, \dots, n_2$ and $k=1, \dots, n$, respectively. Such a statistic U has been proposed by Mann and Whitney [6] for the two-sample problem. In this case, we have also

$$(3.2) \quad \mathcal{E}(\hat{p}_1 | \prod_1) = \frac{1}{2} \quad \text{and} \quad \mathcal{E}(\hat{p}_2 | \prod_1) = p$$

and

$$(3.3) \quad \mathcal{E}(\hat{p}_1 | \prod_2) = p \quad \text{and} \quad \mathcal{E}(\hat{p}_2 | \prod_2) = \frac{1}{2}.$$

Thus, a procedure analogous to the rule 1 in the previous section can be used.

Decision rule 2. If it can be assumed that $F_1(t) > F_2(t)$ for all real t , classify (X_1, \dots, X_n) into \prod_1 or \prod_2 according to $\hat{p}_1 < \hat{p}_2$ or $\hat{p}_1 > \hat{p}_2$, where \hat{p}_1 and \hat{p}_2 are given by (3.1).

We shall assume here that $n \leq n_1$ and $n \leq n_2$. Let $g(X_{i(k)}^{(1)}, X_k)$ and $g(X_{j(k)}^{(2)}, X_k)$ be two sets of random variables defined as

$$(3.4) \quad g(X_{i(k)}^{(1)}, X_k) = \begin{cases} \frac{1}{2} & \text{if } X_{i(k)}^{(1)} < X_k \\ 0 & \text{if } X_{i(k)}^{(1)} \geq X_k \end{cases}$$

$$g(X_{j(k)}^{(2)}, X_k) = \begin{cases} \frac{1}{2} & \text{if } X_{j(k)}^{(2)} < X_k \\ 0 & \text{if } X_{j(k)}^{(2)} \geq X_k \end{cases}$$

$$(k=1, \dots, n; i(k)=1, \dots, n_1 \text{ and } j(k)=1, \dots, n_2)$$

and let

$$(3.5) \quad V(X_{i(1)}^{(1)}, \dots, X_{i(n)}^{(1)}; X_{j(1)}^{(2)}, \dots, X_{j(n)}^{(2)}) \\ = \frac{1}{n} \sum_{k=1}^n \{g(X_{i(k)}^{(1)}, X_k) + g(X_{j(k)}^{(2)}, X_k)\} \\ (k=1, \dots, n, i(1) \neq \dots \neq i(n) \text{ and } j(1) \neq \dots \neq j(n)).$$

Then, we have also

$$(3.6) \quad \mathcal{E} \left\{ V(X_{i(1)}^{(1)}, \dots, X_{i(n)}^{(1)}, X_{j(1)}^{(2)}, \dots, X_{j(n)}^{(2)}) \mid \Pi_1 \right\} = \frac{1}{2} - \frac{d}{2}$$

and

$$(3.7) \quad \mathcal{E} \left\{ V(X_{i(1)}, \dots, X_{i(n)}, X_{j(1)}, \dots, X_{j(n)}) \mid \Pi_2 \right\} = \frac{1}{2} + \frac{d}{2}.$$

Consider

$$(3.8) \quad V = \frac{(n!)^2}{n_1! n_2!} \sum_{i(1) \neq \dots \neq i(n)}^{n_1 P_n} \sum_{j(1) \neq \dots \neq j(n)}^{n_2 P_n} V(X_{i(1)}^{(1)}, \dots, X_{i(n)}^{(1)}, X_{j(1)}^{(2)}, \dots, X_{j(n)}^{(2)}),$$

where sums $\sum_{i(1) \neq \dots \neq i(n)}^{n_1 P_n}$ and $\sum_{j(1) \neq \dots \neq j(n)}^{n_2 P_n}$ are taken over all n -tuples $\{i(1), \dots, i(n)\}$ of distinct positive integers $\leq n_1$ and all n -tuples $\{j(1), \dots, j(n)\}$ of distinct positive integers $\leq n_2$.

Then V will become a statistic corresponding to $\frac{1}{n} \sum_{k=1}^n H(X_k)$ in the previous section and $\hat{p}_1 \cong \hat{p}_2$ is equivalent to $V \cong \frac{1}{2}$ except possibly for $X_k = X_{i(k)}^{(1)}$ and $X_k = X_{j(k)}^{(2)}$ ($k=1, \dots, n, i(k)=1, \dots, n_1$ and $j(k)=1, \dots, n_2$). Thus, the situation is similar as before and the result of section 5 in Hoeffding's paper [5] is applicable in this case, too.

If we rewrite (3.8) as

$$V = \alpha_1 V_1 + \dots + \alpha_N V_N,$$

where $\alpha_1 = \dots = \alpha_N = \frac{(n!)^2}{n_1! n_2!}$ and $N = \frac{n_1! n_2!}{(n!)^2}$,

we have

$$\begin{aligned} (3.9) \quad P \left\{ V \geq \frac{1}{2} \mid \Pi_1 \right\} &= P \left\{ V - \mathcal{E}_{\Pi_1} V \geq \frac{\Delta}{2} \mid \Pi_1 \right\} \\ &\leq \mathcal{E} \left[\exp \left\{ h \left(V - \mathcal{E}_{\Pi_1} V - \frac{\Delta}{2} \right) \right\} \mid \Pi_1 \right] \\ &= \mathcal{E} \left[\exp \left\{ h \sum_{i=1}^N \alpha_i \left(V_i - \mathcal{E}_{\Pi_1} V_i - \frac{\Delta}{2} \right) \right\} \mid \Pi_1 \right] \end{aligned}$$

for an arbitrary positive constant h . Since the exponential function $\exp \{ \ } is convex,$

$$(3.10) \quad P \left\{ V \geq \frac{1}{2} \mid \Pi_1 \right\} \leq \sum_{i=1}^N \alpha_i \mathcal{E} \left[\exp \left\{ h \left(V_i - \mathcal{E}_{\Pi_1} V_i - \frac{\Delta}{2} \right) \right\} \mid \Pi_1 \right].$$

Since $g(X_{i(1)}^{(1)}, X_1) + g(X_{j(1)}^{(2)}, X_1), \dots, g(X_{i(n)}^{(1)}, X_n) + g(X_{j(n)}^{(2)}, X_n)$ are mutually independent and $0 \leq g(X_{i(k)}^{(1)}, X_k) + g(X_{j(k)}^{(2)}, X_k) \leq 1$, it can be seen from [5] that the minimum of $\mathcal{E} \left[\exp \left\{ h \left(V_i - \mathcal{E}_{\Pi_1} V_i - \frac{\Delta}{2} \right) \right\} \mid \Pi_1 \right]$ with respect

to $h (> 0)$ is not larger than $\exp \left(-\frac{n\Delta^2}{2} \right)$. Thus,

$$P \left\{ V < \frac{1}{2} \mid \Pi_1 \right\} \geq 1 - e^{-n\Delta^2/2},$$

and similarly

$$P \left\{ V > \frac{1}{2} \mid \Pi_2 \right\} \geq 1 - e^{-n\Delta^2/2}.$$

Summarizing the above, we have the following:

The success rate of the rule 2 is greater than or equal to $1 - \exp \left(-\frac{n\Delta^2}{2} \right)$, if it can be assumed that $F_1(t) > F_2(t)$ for all real t , and $\Delta = p - \frac{1}{2}$ is given.

Although the above assumptions are not always fulfilled in actual

applications, Δ can be estimated from samples $(X_1^{(1)}, \dots, X_{n_1}^{(1)})$ and $(X_1^{(2)}, \dots, X_{n_2}^{(2)})$, and we may remove the assumption of stochastic orderings, if Δ is fairly larger than zero.

THE INSTITUTE OF STATISTICAL MATHEMATICS

REFERENCES

- [1] T. W. Anderson, *An Introduction to Multivariate Statistical Analysis*, Wiley and Sons, New York, 1958.
- [2] T. W. Anderson, "Classification by multivariate analysis," *Psychometrika*, 16 (1951), 31-50.
- [3] Z. W. Birnbaum, "On a use of Mann-Whitney statistic," *Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability*, Univ. of California Press, I (1956) 13-17.
- [4] R. A. Fisher, "The use of multiple measurements in taxonomic problems," *Ann. Eugen.*, 7 (1936), 179-188.
- [5] W. Hoeffding, "Probability inequalities for sums of bounded random variables," *Amer. Statist. Assoc.*, 58 (1963), 13-30.
- [6] H. B. Mann and D. R. Whitney, "On a test of whether one of two random variables is stochastically larger than the other," *Ann. Math. Statist.*, 18 (1947), 50-60.
- [7] C. R. Rao, *Advanced Statistical Methods in Biometric Research*, John Wiley and Sons, New York, 1952.
- [8] C. R. Rao, "On a general theory of discrimination when the information on alternative hypotheses is based on samples," *Ann. Math. Statist.*, 25 (1954) 651-669.
- [9] R. Sitgreaves, "On the distribution of two random matrices used in classification procedures," *Ann. Math. Statist.*, 23 (1952), 263-270.
- [10] A. Wald, "On a statistical problem arising in the classification of an individual into one of two groups," *Ann. Math. Statist.*, 15 (1944) 145-163.