



---

Discriminatory Analysis. Nonparametric Discrimination: Consistency Properties

Author(s): Evelyn Fix and J. L. Hodges, Jr.

Source: *International Statistical Review / Revue Internationale de Statistique*, Dec., 1989, Vol. 57, No. 3 (Dec., 1989), pp. 238-247

Published by: International Statistical Institute (ISI)

Stable URL: <https://www.jstor.org/stable/1403797>

---

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact [support@jstor.org](mailto:support@jstor.org).

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <https://about.jstor.org/terms>



JSTOR

*International Statistical Institute (ISI)* is collaborating with JSTOR to digitize, preserve and extend access to *International Statistical Review / Revue Internationale de Statistique*

## References

- Agrawala, A.K. (Ed). (1977). *Machine Recognition of Patterns*. New York: IEEE Press.
- Cacoullos, T. (1966). Estimation of a multivariate density. *Ann. Inst. Statist. Math.* **18**, 179–189.
- Coomans, D. & Broeckaert, I. (1986). *Potential Pattern Recognition in Chemical and Medical Decision Making*. Letchworth: Research Studies Press.
- Cover, T.M. & Hart, P.E. (1967). Nearest neighbour pattern classification. *IEEE Trans. Info. Theory* **IT-13**, 21–27.
- das Gupta, S. (1973). Theories and methods in classification: a review. In *Discriminant Analysis and Applications*, Ed. T. Cacoullos, pp. 77–137. New York: Academic Press.
- Devijver, P.A. & Kittler, J. (1982). *Pattern Recognition: a Statistical Approach*. London: Prentice-Hall.
- Devroye, L. & Györfi, L. (1985). *Nonparametric Density Estimation: the  $L_1$  View*. New York: Wiley.
- Fisher, R.A. (1936). The use of multiple measurements in taxonomic problems. *Ann. Eugen.* **7**, 179–188.
- Fix, E. & Hodges, J.L. (1951). Discriminatory analysis. Nonparametric discrimination; consistency properties. Report Number 4, Project Number 21-49-004, USAF School of Aviation Medicine, Randolph Field, Texas. (Reprinted as pp 261–279 of Agrawala, 1977).
- Fix, E. & Hodges, J.L. (1952). Discriminatory analysis. Nonparametric discrimination: small sample performance. Report Number 11, Project Number 21-49-004, USAF School of Aviation Medicine, Randolph Field, Texas. (Reprinted as pp 280–322 of Agrawala, 1977).
- Hand, D.J. (1981). *Discrimination and Classification*. Chichester: Wiley.
- Hand, D.J. (1982). *Kernel Discriminant Analysis*. Chichester: Research Studies Press.
- Hoel, P.G. & Peterson, R.P. (1949). A solution to the problem of optimum classification. *Ann. Math. Statist.* **20**, 433–438.
- Loftsgaarden, D.O. & Quesenberry, C.P. (1965). A nonparametric estimate of a multivariate density function. *Ann. Math. Statist.* **36**, 1049–1051.
- Marron, J.S. (1989). Automatic smoothing parameter selection: a survey. *Empirical Econ.* **13**, 187–208.
- Parzen, E. (1962). On estimation of a probability density function and mode. *Ann. Math. Statist.* **33**, 1065–1076.
- Prakasa Rao, B.L.S. (1983). *Nonparametric Functional Estimation*. Orlando, Florida: Academic Press.
- Rosenblatt, M. (1956). Remarks on some nonparametric estimates of a density function. *Ann. Math. Statist.* **27**, 832–837.
- Silverman, B.W. (1986). *Density Estimation for Statistics and Data Analysis*. London: Chapman and Hall.

## Résumé

En 1951, E. Fix et J.L. Hodges, Jr. ont écrit un rapport technique prophétique sur l'analyse non-paramétrique de discrimination et l'estimation de la densité de probabilité, mais celui-ci ne fut jamais publié par ses auteurs. Ce rapport introduit plusieurs idées nouvelles et importantes. Il nous intéresse non seulement pour des raisons historiques, mais aussi parce qu'il contient des concepts qui sont encore importants de nos jours. Nous le publions ici en entier, accompagné d'un commentaire qui l'interprète d'un point de vue plus moderne.

[Received June 1988, accepted May 1989]

## Discriminatory Analysis—Nonparametric Discrimination: Consistency Properties

**Evelyn Fix† and J.L. Hodges, Jr.**

*University of California, Berkeley*

This paper originally appeared as Report Number 4, Project Number 21-49-004, USAF School of Aviation Medicine, Randolph Field, Texas, in February 1951. Enquiries should be addressed to the authors of the accompanying commentary, B.W. Silverman and M.C.

† Evelyn Fix died on 30 December 1965.

Jones, School of Mathematical Sciences, University of Bath, Bath BA2 7AY, UK, who prepared the paper for publication.

## 1 Introduction

The discrimination problem (two population case) may be defined as follows: a random variable  $Z$ , of observed value  $z$ , is distributed over some space (say,  $p$ -dimensional) either according to distribution  $F$ , or according to distribution  $G$ . The problem is to decide, on the basis of  $z$ , which of the two distributions  $Z$  has.

The problem may be classified in various ways into subproblems. One pertinent method of classification is according to the amount of information assumed to be available about  $F$  and  $G$ . We may distinguish three stages:

- (i)  $F$  and  $G$  are completely known;
- (ii)  $F$  and  $G$  are known except for the values of one or more parameters;
- (iii)  $F$  and  $G$  are completely unknown, except possibly for assumptions about existence of densities, etc.

Subproblem (i) has been, in a sense, completely solved. The solution is implicit in the Neyman–Pearson lemma (Neyman & Pearson, 1936), and was made explicit by Welch (1939). We may without loss of generality assume the existence of density functions, say  $f$  and  $g$ , corresponding to  $F$  and  $G$ , since  $F$  and  $G$  are absolutely continuous with respect to  $F + G$ . If  $f$  and  $g$  are known, the discrimination should depend only on  $f(z)/g(z)$ . An appropriate (positive) constant  $c$  is chosen, and the following rule is observed:

- if  $f(z)/g(z) > c$ , we decide in favor of  $F$ ;
- if  $f(z)/g(z) < c$ , we decide in favor of  $G$ ;
- if  $f(z)/g(z) = c$ , the decision may be made in an arbitrary manner.

These procedures are known to have optimum properties with regard to control of probability of misclassification (probability of wrong decision). We shall refer to this as the ‘likelihood ratio procedure’, and denote it by  $L(c)$ .

For simplicity, we shall assume throughout the paper that the borderline case  $f(z) = cg(z)$  can be neglected. Formally, we postulate that

$$P\{f(Z) = cg(Z)\} = 0$$

regardless of whether  $Z$  comes from  $F$  or  $G$ . Since the classification is arbitrary when  $f(z) = cg(z)$ , it hardly seems worth while to introduce complications into the methods to allow for it. However, it is not difficult to extend our methods to take care of the situation which arises when

$$P\{f(Z) = cg(Z)\} > 0.$$

The choice of  $c$  depends on considerations relating to the relative importance of the two possible errors: saying  $Z$  is distributed according to  $G$  when in fact it is distributed according to  $F$ , and conversely. Two choices of  $c$  have been widely advocated:

- (a) take  $c = 1$ ;
- (b) choose  $c$  so that the two probabilities of error are equal.

Choice (a) has been called ‘logical’; choice (b) yields the minimax procedure. In this paper we shall not concern ourselves with the choice of  $c$ , but shall assume that a given positive  $c$  is a datum of the problem.

The usual approach to subproblem (ii) is as follows. We assume there are available

samples from the two distributions, say

$$\begin{aligned} X_1, X_2, \dots, X_m &: \text{sample from } F, \\ Y_1, Y_2, \dots, Y_n &: \text{sample from } G. \end{aligned}$$

We assume further that  $F$  and  $G$  are known in form: that is, that we know them except for the values of some real parameters, which may be denoted collectively by  $\theta$ . We may denote the distributions corresponding to a given  $\theta$  by  $F_\theta, G_\theta$ . The procedure currently employed is to use the  $X$ 's and  $Y$ 's to estimate  $\theta$ , by, say,  $\hat{\theta}$ , and then to proceed as under (i), using the distributions  $F_{\hat{\theta}}, G_{\hat{\theta}}$  as though they were known to be correct.

The most familiar example of this process is the linear discriminant function (Fisher, 1936). There, it is (tacitly) assumed that  $F$  and  $G$  are  $p$ -variate normal distributions having the same (unknown) covariance matrix, and unknown expectation vectors. The two expectation vectors and the covariance matrix are estimated from the samples, and the likelihood ratio procedure is then employed, using the estimated values as though they were known to be correct.

Not much is known about the desirability of the usual method of attack on (ii). We give in §3 a theorem concerning asymptotic properties of the method. Undoubtedly, this procedure is reasonable *provided* the assumed parametric form is correct. But the validity of the use of the linear discriminant function with data obviously not normal or, if normal, with obviously unequal covariance matrices has been of general concern. Presumably, very bad results may ensue if a procedure is used, based on certain assumptions about parametric form, when those assumptions are not even approximately correct.

There seems to be a need for discrimination procedures whose validity does not require the amount of knowledge implied by the normality assumption, the homoscedastic assumption, or any assumption of parametric form. The present paper is, as far as the authors are aware, the first one to attack subproblem (iii): can reasonable discrimination procedures be found which will work even if no parametric form can be assumed?

It is not to be expected that any procedure can be guaranteed to give good results without any restriction whatsoever on the distributions  $F$  and  $G$ . To clarify this point, we need to state a precise meaning for 'good results'. This is done in §2, with the introduction of the concept of 'consistency'. We then proceed in §4 to prove, under weak restrictions on the densities  $f$  and  $g$ , the consistency of a class of nonparametric procedures there proposed. A modification of these procedures is then considered in §5.

It may be noted that all of the methods and results of this paper can be extended without difficulty to the situation in which there are more than two populations to be discriminated.

The authors are engaged in further work along the lines here laid down. Specifically, some sampling experiments are being conducted, intended to throw some light on the performance of the procedures for moderate sample sizes; and asymptotic properties of a class of sequential nonparametric discriminatory procedures is being investigated. It is intended to prepare further reports setting forth the results.

## 2 The Notion of Consistency

In setting out to define an optimum property in statistical inference, it is useful to have in mind the limit of excellence beyond which it is not possible to go. The procedures  $L(c)$  described in §1 provide such a limit in the case of nonparametric discrimination: we cannot, with any nonparametric classification procedure, expect to do better than the best which is possible when the densities themselves are assumed to be known. This fact is

intuitively obvious, but if desired an exact proof is easily given. When  $f$  and  $g$  are known,  $Z$  is sufficient for the classification, with respect to  $(Z; X_1, X_2, \dots, X_m; Y_1, Y_2, \dots, Y_n)$ , and we may (by using randomization) exactly duplicate (with a procedure based on  $Z$ ) the performance characteristic of any procedure based on  $(Z; X_1, X_2, \dots, X_m; Y_1, Y_2, \dots, Y_n)$ .

Thus, no nonparametric procedure can have probabilities of error less than those of a likelihood ratio procedure. On the other hand, we shall propose in §§ 4 and 5 classes of (sequences of) nonparametric procedures which, in the limit as  $m$  and  $n$  tend to infinity, have the same probabilities of error as the procedures  $L(c)$ . We may therefore reasonably say that our procedures are *consistent with* the likelihood ratio procedures.

There are two different notions of consistency for sequences of statistical decision functions, and it may be worthwhile to distinguish them. Suppose that the decision space is finite (as is the case in discriminatory analysis when there are finitely many populations). Let the possible decisions be denoted by  $\delta_1, \delta_2, \dots, \delta_r$ . Now suppose we are considering two sequences of decision functions, say  $\{\Delta'_n\}$  and  $\{\Delta''_n\}$ . How should we define the notion that these two sequences tend to agree with each other, or be consistent with each other, as  $n \rightarrow \infty$ ? On the one hand, we might require that in the limit there should be close agreement between the probabilities of decision; on the other hand we might require that in the limit there be high probability of agreement of decision. The former requirement relates to the performance characteristics of the decision functions; the latter requirement relates to the decision functions themselves. We have then two definitions.

**Definition 1.** We shall say that the sequences  $\{\Delta'_n\}$  and  $\{\Delta''_n\}$  are *consistent in the sense of performance characteristics* if, whatever be the true distributions, and whatever be  $\varepsilon > 0$ , there exists a number  $N$  such that whenever  $m > N$  and  $n > N$ ,

$$|P(\Delta'_m = \delta_i) - P(\Delta''_n = \delta_i)| < \varepsilon$$

for every decision  $\delta_i$ .

**Definition 2.** We shall say that the sequences  $\{\Delta'_n\}$  and  $\{\Delta''_n\}$  are *consistent in the sense of decision functions* if, whatever be the true distributions, and whatever be  $\varepsilon > 0$ , there exists a number  $N$  such that whenever  $m > N$  and  $n > N$ ,

$$P(\Delta'_m = \Delta''_n) > 1 - \varepsilon.$$

We observe that consistency in the second sense implies that in the first, since  $P(\Delta'_m \neq \Delta''_n)$  is not less than each of the quantities

$$P(\Delta'_m = \delta_i \text{ and } \Delta''_n \neq \delta_i) \geq P(\Delta'_m = \delta_i) - P(\Delta''_n = \delta_i).$$

The definitions are not equivalent however, as the following trivial example shows. If  $\Delta'$  and  $\Delta''$  each denotes (for any  $m, n$ ) the process of choosing between two alternatives  $\delta_1$  and  $\delta_2$  by tossing a coin, then  $P(\Delta' = \Delta'') = \frac{1}{2}$ , while

$$P(\Delta' = \delta_i) = P(\Delta'' = \delta_i) = \frac{1}{2} \quad \text{for } i = 1, 2.$$

Inasmuch as it is customary to evaluate decision functions solely in terms of their performance characteristics, Definition 1 is the more natural. However, all proofs of consistency given in this paper provide consistency in the stronger sense of the second definition, and consequently we shall adopt it.

Since our procedures are based on two samples, we must consider a double limit process as both  $m$  and  $n$  tend to infinity. To avoid difficulties which would otherwise arise in § 5, we shall assume throughout that  $m$  and  $n$  approach infinity at the same speed.

Precisely, we assume  $m/n$  and  $n/m$  are both bounded away from 0 as  $n, m \rightarrow \infty$ . Whenever we write ' $m, n \rightarrow \infty$ ' this restriction should be understood. Our restriction has the effect of reducing the limiting process from a double to a single one.

In the sequel we shall be comparing certain discriminatory procedures with procedures of the type  $L(c)$ . It is convenient to introduce the following definition.

**Definition 3.** A sequence  $\{\Delta_{m,n}\}$  of discriminatory procedures, based on  $Z$  and on samples  $X_1, X_2, \dots, X_m$  from  $F$  and  $Y_1, Y_2, \dots, Y_n$  from  $G$ , is said to be *consistent with*  $L(c)$  if, whatever be the distributions  $F$  and  $G$ , regardless of whether  $Z$  is distributed according to  $F$  or according to  $G$ , and whatever be  $\varepsilon > 0$ , we can assure

$$P\{\Delta_{m,n} \text{ and } L(c) \text{ yield the same classification of } Z\} > 1 - \varepsilon$$

provided only that  $m$  and  $n$  are sufficiently large.

We may also define a corresponding notion of uniform consistency. If, in Definition 3, the bound on probability of agreement can be assured for all  $F$  and  $G$  with a single size specification on  $m$  and  $n$ , we say that  $\{\Delta_{m,n}\}$  is uniformly consistent with  $L(c)$ .

### 3 Consistency for the Parametric Case

We shall now demonstrate that the analogy of the notion of consistency just introduced with the like-named notion in point estimation is more than formal. Consider the problem of parametric discrimination (subproblem (ii)) of § 1.

We shall from time to time have occasion to consider probabilities computed under the assumption that  $Z$  is distributed according to  $F$ , or according to  $G$ . It is convenient to let  $P_1$  and  $P_2$  denote probabilities computed under these respective assumptions.

Let  $\mathcal{F}$  and  $\mathcal{G}$  be classes of densities parametrized by parameters denoted collectively by  $\theta$ . Let there be a notion of convergence introduced in the space  $\Theta$  of parameter values. Suppose there is given a sequence  $\{\hat{\theta}_{m,n}\}$  of estimates for  $\theta$ ,  $\hat{\theta}_{m,n}$  being a function of  $X_1, X_2, \dots, X_m$  and  $Y_1, Y_2, \dots, Y_n$ .

**THEOREM 1.** *If*

(a) *the estimates  $\{\hat{\theta}_{m,n}\}$  are consistent,*

(b) *for every  $\theta$ ,  $f_\theta(z)$  and  $g_\theta(z)$  are continuous functions of  $\theta$  for every  $z$  except perhaps for  $z \in Z_\theta$ , where  $P_i(Z_\theta) = 0$ ,  $i = 1, 2$ ,*

*then the sequence of discrimination procedures  $\{\hat{L}_{m,n}(c)\}$  obtained by applying the likelihood ratio principle with critical value  $c > 0$  to  $f_{\hat{\theta}_{m,n}}(z)$  and  $g_{\hat{\theta}_{m,n}}(z)$  is consistent with  $L(c)$ .*

**Proof.** The idea of the proof is very simple: since  $\hat{\theta}_{m,n}$  is consistent,  $\hat{\theta}_{m,n}$  will probably be near  $\theta$  if  $m$  and  $n$  are large. But since  $f_\theta$  and  $g_\theta$  are continuous, this means that  $f_{\hat{\theta}_{m,n}}$  will probably be near  $f_\theta$ , and  $cg_{\hat{\theta}_{m,n}}$  will probably be near  $cg_\theta$ . Therefore, it is not likely to make much difference whether we compare  $f_\theta$  and  $cg_\theta$  or  $f_{\hat{\theta}_{m,n}}$  and  $cg_{\hat{\theta}_{m,n}}$ .

Fix  $c > 0$ ,  $\varepsilon > 0$ , and  $\theta \in \Theta$ . Find  $\delta > 0$  so small that

$$P_i\{|f_\theta(Z) - cg_\theta(Z)| \leq \delta\} < \frac{1}{2}\varepsilon \quad (i = 1, 2).$$

(This is possible since  $P_i\{|f_\theta(Z) - cg_\theta(Z)| \leq u\}$  is the cumulative function of the random variable  $|f_\theta(Z) - cg_\theta(Z)|$  and hence is continuous on the right, and by assumption takes on the value 0 when  $u = 0$ ). We now assume that  $z$  does not lie in  $Z_\theta$ , thus excluding an event of zero probability. Since  $f_\theta(z)$  is a continuous function of  $\theta$  for all  $z$ , we can associate with every  $z$  a quantity  $\eta_1(z) > 0$  such that

$$|f_{\hat{\theta}}(z) - f_\theta(z)| < \frac{1}{2}\delta \quad \text{whenever} \quad |\hat{\theta} - \theta| < \eta_1(z).$$



A like function  $\eta_2(z)$  arises if  $f$  is replaced by  $cg$ . Let  $\eta(z) = \min \{\eta_1(z), \eta_2(z)\}$  and find  $\eta > 0$  such that  $P_i(\eta(Z) < \eta) < \frac{1}{4}\epsilon$ ,  $i = 1, 2$ .

Using finally the consistency of the estimates, choose  $M$  and  $N$  large enough so that, whenever  $m > M$  and  $n > N$ ,  $P\{|\hat{\theta}_{m,n} - \theta| > \eta\} < \frac{1}{4}\epsilon$ . Combining the above, a disagreement between  $L(c)$  and  $\hat{L}_{m,n}(c)$  will arise with probability less than  $\epsilon$ .  $\square$

*Remark 1.* The dependence of the discontinuity sets  $Z_\theta$  on  $\theta$  is important. Were we to demand the stronger property that  $f_\theta(z)$  and  $g_\theta(z)$  be continuous in  $\theta$  for all  $z \notin Z$ ,  $Z$  a fixed set,  $P_i(Z) = 0$ ,  $i = 1, 2$ , we should exclude many cases which are included under the theorem as given.

*Remark 2.* Two notions of convergence in  $\Theta$  are involved: that with respect to which the estimates are consistent, and that with respect to which the densities are continuous. These need not be the same, provided the former implies the latter.

*Remark 3.* If uniformity is added to the hypotheses of Theorem 1, it may also be added to the conclusions. Specifically, if the estimates  $\hat{\theta}_{m,n}$  are uniformly consistent, if the densities  $f$  and  $g$  are uniformly continuous functions of  $\theta$ , uniformly in  $z$ , and if the  $\delta$  of the proof of Theorem 1 may be fixed independently of  $\theta$ , then that proof goes through for all  $\theta$  using the same value of  $\epsilon$ . We can then conclude the uniform consistency of  $\{\hat{L}_{m,n}(c)\}$ .

#### 4 Nonparametric Discrimination and its Consistency

Let us next consider the discrimination problem of the third kind delineated in § 1. We admit the possibility that the densities  $f$  for  $X$  and  $g$  for  $Y$  may be any in certain classes  $\mathcal{F}$  and  $\mathcal{G}$  of densities which are too large to be characterized by a finite number of parameters. Thus,  $\mathcal{F}$  and  $\mathcal{G}$  may consist of all uniformly continuous densities, or of all continuous densities, or of all densities continuous save at most at countably many points. Can we have any discrimination procedures which are reasonable to use when so little is assumed about the populations being discriminated?

Recall that, once  $c$  has been selected and  $Z$  has been observed to have the value  $z$ , the only information needed to carry out the procedure  $L(c)$  are the two real numbers  $f(z)$  and  $g(z)$ . In the procedure  $\hat{L}_{m,n}(c)$ , we employed the estimate for  $\theta$  as a means of obtaining estimates for  $f_\theta(z)$  and  $g_\theta(z)$ . In the nonparametric case there is no  $\theta$  to be estimated, but we may instead proceed to estimate the numbers  $f(z)$  and  $g(z)$  directly. Once estimates have been obtained, we may apply the procedure  $L(c)$ , using these estimates instead of  $f(z)$  and  $g(z)$ . We shall designate such procedures by  $L^*(c, \hat{f}, \hat{g})$ , where  $\hat{f}$  and  $\hat{g}$  are the estimates for  $f$  and  $g$ .

Before considering the problem of estimating the densities, let us note the properties which such estimates should have if we are to be able to prove the consistency of  $L^*(c, \hat{f}, \hat{g})$  with  $L(c)$ .

**THEOREM 2.** *If  $\hat{f}_{m,n}(z)$  and  $\hat{g}_{m,n}(z)$  are consistent estimates for  $f(z)$  and  $g(z)$  for all  $z$  except possibly  $z \in Z_{f,g}$ , where  $P_i(Z_{f,g}) = 0$ ,  $i = 1, 2$ , then  $\{L^*_{m,n}(c, \hat{f}, \hat{g})\}$  is consistent with  $L(c)$ .*

The proof follows lines similar to that of Theorem 1, and will be omitted.

Our problem is now to find consistent estimates for  $f(z)$  and  $g(z)$ . We shall for brevity consider  $f(z)$  only, as analogous remarks apply to  $g(z)$ . We fix  $z$ , since the argument is the same for each value. Our basic idea is this: the proportion of the  $m$   $X$ 's which fall in a stated (small) neighborhood of  $z$  may be used to estimate the  $X$ -probability in that

neighborhood. The ratio of this estimated probability to the measure of the neighborhood is then an estimate of the average value of  $f(x)$  near  $z$ . This is in turn an estimate of  $f(z)$  itself if we make some assumption about the smoothness of  $f$ . To obtain consistency, we may let the neighborhood shrink down to  $z$  as  $m \rightarrow \infty$ , so that the average of  $f(x)$  over the neighborhood will approach  $f(z)$ ; but we will take care to have the neighborhood shrink slowly enough so that the proportion of the  $X$ 's therein will have a positive expectation. This will assure that the proportion of  $X$ 's in the neighborhood is a consistent estimate of the probability.

It is obvious that we cannot hope to estimate  $f(z)$  from  $X_1, X_2, \dots, X_m$  unless some continuity assumption is made. For, otherwise we could alter  $f(z)$  arbitrarily without in any way changing the distribution of  $X_1, X_2, \dots$ , and thus without changing the distribution of any sequence of estimates based on  $X_1, X_2, \dots$ .

Now let  $\mu$  denote Lebesgue measure in our ( $p$ -dimensional) sample space, and let  $|x - y|$  denote the (Euclidean) distance between points  $x$  and  $y$  of this space.

LEMMA 3. *If  $f(x)$  is continuous at  $x = z$ , and if  $\{\Delta_m\}$  is a sequence of sets such that*

$$\lim_{m \rightarrow \infty} \sup_{d \in \Delta_m} |z - d| = 0,$$

and

$$\lim_{m \rightarrow \infty} m\mu(\Delta_m) = \infty,$$

and if  $M$  is the number of  $X_1, X_2, \dots, X_m$  which lie in  $\Delta_m$ , then  $M/\{m\mu(\Delta_m)\}$  is a consistent estimate for  $f(z)$ .

*Proof.* Observe that  $P(\Delta_m)/\mu(\Delta_m) \rightarrow f(z)$  as  $m \rightarrow \infty$ . If  $f(z) > 0$ ,  $mP(\Delta_m) \rightarrow \infty$ . Since  $\mu(\Delta_m) \rightarrow 0$ ,  $P(\Delta_m) \rightarrow 0$  and we conclude  $M/\{mP(\Delta_m)\} \rightarrow 1$  in probability. Combining,  $M/\{m\mu(\Delta_m)\} \rightarrow f(z)$  in probability, as was to be shown. If  $f(z) = 0$ ,

$$E[M/\{m\mu(\Delta_m)\}] = P(\Delta_m)/\mu(\Delta_m) \rightarrow 0$$

and the Markoff lemma completes the proof.  $\square$

We have in Lemma 3 a class of estimates, any of which, by virtue of Theorem 2, will provide consistent discrimination of any (nonparametric) classes  $\mathcal{F}$  and  $\mathcal{G}$  whose members are continuous (except possibly for a set of values of zero measure).

## 5 Alternative Procedures

While the procedures

$$L^*(c, M/\{m\mu(\Delta_m)\}, N/\{n\mu(\Delta_n)\})$$

of the last section provide consistent discrimination, the question of their applicability when  $m$  and  $n$  are not large remains open. (Like criticism may of course be applied to any asymptotic theorem.) We shall in the present section suggest some alternative estimates for  $f(z)$  and  $g(z)$ , which seem on intuitive grounds more likely to give good results than the estimates proposed before. The former estimates are the natural ones when thinking of the simplicity of consistency proofs, but need not be desirable in practice.

The main practical difficulty in using the former estimates lies in the choice of the regions  $\{\Delta_m\}$ , (and the corresponding regions for  $g$ , say  $\{\Delta_n\}$ ). If these regions are made too small, the numbers  $M$  and  $N$  of sample points falling into them will be too small, so that the proportions  $M/m$  and  $N/n$  will not be accurate estimates for the corresponding



probabilities  $P_1(\Delta_m)$ ,  $P_2(\Lambda_n)$ . On the other hand, if the regions are made too large, these probabilities will not be good approximations for  $f(z)\mu(\Delta_m)$  and  $g(z)\mu(\Lambda_n)$ . We are between twin perils and must steer a middle course. We might, for example, decide the smallest values of  $M$  and  $N$  we could tolerate, and choose  $\Delta_m$  and  $\Lambda_n$  just big enough to include the chosen number of points. But to do so alters the probabilistic properties; now  $M$  and  $N$  are fixed and  $\Delta$  and  $\Lambda$  are random. Are the results of Lemma 3 still valid?

Even if they are we may still be in difficulties. It may happen that near  $z$  there are numerous  $X$ 's, but few  $Y$ 's; but by going a little further we find the situation reversed. The indication is clearly for  $\pi_1$ , but if we take separate  $\Delta$  and  $\Lambda$  the estimated  $f$  and  $g$  may be close. To avoid this difficulty the following idea is suggested: choose a number  $k$ , and take in the neighborhood of  $z$  a single region,  $\Delta_{m,n}$ , containing a total of  $k$  points of either sample. Intuitively this procedure seems sound, but since  $M + N = k$  we have introduced dependence of our estimates and further altered the probabilistic properties. The question which now arises is whether or not estimates for  $f(z)$  and  $g(z)$  based on  $M$  and  $N$ , when so determined, are still consistent.

As a first step in answering these questions, observe that we may by means of a preliminary transformation reduce our space from  $p$  dimensions to one. Let  $\rho(x, y)$  denote a non-negative real valued function of pairs  $(x, y)$  of points in the sample space. Suppose  $\rho$  is so constructed that when  $x_n \rightarrow x$ ,  $\rho(x_n, x) \rightarrow 0$ , and suppose further that for each  $z$ , except perhaps for  $z \in Z_{f,g}$ , where  $P_i(Z_{f,g}) = 0$ ,  $i = 1, 2$ ,  $\rho(X, z)$  and  $\rho(Y, z)$  are random variables possessing densities, say  $f_z(x)$  and  $g_z(x)$ , continuous and not both 0 at 0. (These properties are satisfied, for example, by  $\rho(x, y) = |x - y|^{1/p}$ .) We now replace the problem of deciding whether  $f(z)$  or  $g(z)$  is the larger, by the problem of deciding whether  $f_z(0)$  or  $g_z(0)$  is larger; and further replace the samples  $X_1, X_2, \dots, X_m$  and  $Y_1, Y_2, \dots, Y_n$  by

$$\rho(X_1, z), \rho(X_2, z), \dots, \rho(X_m, z), \quad \text{and} \quad \rho(Y_1, z), \rho(Y_2, z), \dots, \rho(Y_n, z),$$

respectively. We may now, without real loss of generality, assume that  $f$  and  $g$  are densities of non-negative univariate random variables, and that  $z = 0$ .

**THEOREM 4.** *Let  $X$  and  $Y$  be non-negative. Let  $f$  and  $g$  be positive and continuous at 0. Let  $k(m, n)$  be a positive-integer-valued function such that  $k(m, n) \rightarrow \infty$ ,  $m^{-1}k(m, n) \rightarrow 0$ , and  $n^{-1}k(m, n) \rightarrow 0$ , as  $m$  and  $n \rightarrow \infty$ . (This tendency being restricted so that  $m/n$  is bounded away from 0 and  $\infty$ .) Define:*

*$U$ ,  $k$ th smallest value of combined samples of  $X$ 's and  $Y$ 's;*

*$M$ , number of  $X$ 's  $\leq U$ ;*

*$N$ , number of  $Y$ 's  $\leq U$ .*

*Then  $M/(mU)$  is a consistent estimate for  $f(0)$  and  $N/(nU)$  is a consistent estimate for  $g(0)$ .*

*Proof.* Fix  $\varepsilon > 0$  and  $\delta > 0$ . Define  $k_1(m, n)$  and  $k_2(m, n)$  by

$$k_1(m, n) + k_2(m, n) = k(m, n), \quad k_1(m, n)/k_2(m, n) = mf(0)/\{ng(0)\}.$$

Define

$$v(m, n) = k_1(m, n)/\{mf(0)(1 + \delta)^2\}, \quad w(m, n) = k_1(m, n)/\{mf(0)(1 - \delta)^2\}.$$

Observe

$$v(m, n) = k_2(m, n)/\{ng(0)(1 + \delta)^2\}, \quad w(m, n) = k_2(m, n)/\{ng(0)(1 - \delta)^2\}.$$

Define:

- $M_{m,n}^v$ , number of  $X$ 's  $< v(m, n)$ ;
- $M_{m,n}^w$ , number of  $X$ 's  $< w(m, n)$ ;
- $N_{m,n}^v$ , number of  $Y$ 's  $< v(m, n)$ ;
- $N_{m,n}^w$ , number of  $Y$ 's  $< w(m, n)$ .

Using the continuity and positiveness of  $f$  and  $g$  at 0, find  $q > 0$  so small that when  $0 \leq x \leq q$ ,

$$|f(x)/f(0) - 1| < \delta, \quad |g(x)/g(0) - 1| < \delta.$$

Find  $m_1, n_1$  such that when  $m > m_1$  and  $n > n_1$ ,  $w(m, n) < q$ , and make these restrictions. Observe

$$\mathbf{E}(M_{m,n}^v) = m \int_0^{v(m,n)} f(x) dx$$

and hence

$$mf(0)v(m, n)(1 - \delta) < \mathbf{E}(M_{m,n}^v) < mf(0)v(m, n)(1 + \delta).$$

Similarly observe

$$\begin{aligned} mf(0)w(m, n)(1 - \delta) &< \mathbf{E}(M_{m,n}^w) < mf(0)w(m, n)(1 + \delta), \\ ng(0)v(m, n)(1 - \delta) &< \mathbf{E}(N_{m,n}^v) < ng(0)v(m, n)(1 + \delta), \\ ng(0)w(m, n)(1 - \delta) &< \mathbf{E}(N_{m,n}^w) < ng(0)w(m, n)(1 + \delta). \end{aligned}$$

Thus

$$\begin{aligned} \mathbf{E}(M_{m,n}^v) &< k_1(m, n)/(1 + \delta), \quad \mathbf{E}(M_{m,n}^w) > k_1(m, n)/(1 - \delta), \\ \mathbf{E}(N_{m,n}^v) &< k_2(m, n)/(1 + \delta), \quad \mathbf{E}(N_{m,n}^w) > k_2(m, n)/(1 - \delta). \end{aligned}$$

The random variables involved are binomials, whose expectations tend to  $\infty$ , but more slowly than the numbers of trials, as  $m, n \rightarrow \infty$ . Therefore, if we take  $m_2, n_2$  large enough, we can assure

$$\begin{aligned} P\{M_{m,n}^v < k_1(m, n)\} &> 1 - \varepsilon, \quad P\{N_{m,n}^v < k_2(m, n)\} > 1 - \varepsilon, \\ P\{M_{m,n}^w > k_1(m, n)\} &> 1 - \varepsilon, \quad P\{N_{m,n}^w > k_2(m, n)\} > 1 - \varepsilon, \end{aligned}$$

as soon as  $m > m_2$  and  $n > n_2$ , which restrictions we now make. Combining, using the fact that  $U$  will exceed  $v(m, n)$  if  $M_{m,n}^v + N_{m,n}^v < k(m, n)$ , we have

$$P\{U > v(m, n)\} > 1 - 2\varepsilon, \quad P\{U < w(m, n)\} > 1 - 2\varepsilon.$$

The event  $U > v(m, n)$  implies the event that all  $X$ 's  $< v(m, n)$  are among the first  $k$   $X$ 's and  $Y$ 's and hence the event  $M_{m,n}^v \leq M$ . Therefore,

$$P(M_{m,n}^v \leq M) \geq P\{U > v(m, n)\} > 1 - 2\varepsilon.$$

Similarly,  $P(M_{m,n}^w \geq M) > 1 - 2\varepsilon$ . Restricting  $m > m_3, n > n_3$ , we can further assure

$$\begin{aligned} P\{M_{m,n}^v > mf(0)v(m, n)(1 - \delta)^2\} &> 1 - \varepsilon, \\ P\{M_{m,n}^w < mf(0)w(m, n)(1 + \delta)^2\} &> 1 - \varepsilon. \end{aligned}$$

Combining,

$$P\{M/m < f(0)w(m, n)(1 + \delta)^2\} > 1 - 3\varepsilon,$$

$$P\{M/m > f(0)v(m, n)(1 - \delta)^2\} > 1 - 3\varepsilon.$$

Hence

$$P\left\{f(0)\frac{v(m, n)}{w(m, n)}(1 - \delta)^2 < \frac{M}{mU} < f(0)\frac{w(m, n)}{v(m, n)}(1 + \delta)^2\right\} > 1 - 10\varepsilon.$$

Since

$$v(m, n)/w(m, n) = (1 - \delta)^2/(1 + \delta)^2,$$

the conclusion  $M/(mU) \rightarrow f(0)$  in probability is at hand. A similar argument shows  $N/(nU) \rightarrow g(0)$  in probability.  $\square$

A situation in which one of the densities is 0 at 0 can be dealt with by a corresponding but simpler argument which we omit. The effect of Theorem 4 is to assure us of satisfactory large sample results if we employ procedures of the following kind:

Choose  $k$ , a positive integer which is large but small compared to the sample sizes. Specify a metric in the sample space, for example ordinary Euclidean distance. Pool the two samples and find, of the  $k$  values in the pooled samples which are nearest to  $z$ , the number  $M$  which are  $X$ 's. Let  $N = k - M$  be the number which are  $Y$ 's. Proceed with the likelihood ratio discrimination, using however  $M/m$  in place of  $f(z)$  and  $N/n$  in place of  $g(z)$ . That is, assign  $Z$  to  $F$  if and only if

$$M/m > cN/n.$$

## References

- Fisher, R.A. (1936). The use of multiple measurements in taxonomic problems. *Ann. Eugen.* **7**, 179–188.  
 Neyman, J. & Pearson, E.S. (1936). Contributions to the theory of testing statistical hypotheses. *Statist. Res. Mem.* **1**, 1–37.  
 Welch, B.L. (1939). Note on discriminant functions. *Biometrika* **31**, 218–220.