

PROCEEDINGS *of the* THIRD  
BERKELEY SYMPOSIUM ON  
MATHEMATICAL STATISTICS  
AND PROBABILITY

*Held at the Statistical Laboratory  
University of California  
December, 1954  
July and August, 1955*

VOLUME I

CONTRIBUTIONS TO THE THEORY OF STATISTICS

EDITED BY JERZY NEYMAN

For further effective support of the Symposium thanks must be given the National Science Foundation, the United States Air Force Research and Development Command, the United States Army Office of Ordnance Research, and the United States Navy Office of Naval Research.

UNIVERSITY OF CALIFORNIA PRESS  
BERKELEY AND LOS ANGELES  
1956

# THE ROLE OF ASSUMPTIONS IN STATISTICAL DECISIONS

WASSILY Hoeffding  
UNIVERSITY OF NORTH CAROLINA

## 1. Introduction

In order to obtain a good decision rule for some statistical problem we start by making assumptions concerning the class of distributions, the loss function, and other data of the problem. Usually these assumptions only approximate the actual conditions, either because the latter are unknown, or in order to simplify the mathematical treatment of the problem. Hence the assumptions under which a decision rule is derived are ordinarily not satisfied in a practical situation to which the rule is applied. It is therefore of interest to investigate how the performance of a decision rule is affected when the assumptions under which it was derived are replaced by another set of assumptions.

We shall confine ourselves to the consideration of assumptions concerning the class of distributions. Investigations of particular problems of this type are numerous in the literature. There are many studies of the performance of "standard" tests under "non-standard" conditions, for example [3], where further references are given. Most of them considered only the effect of deviations from the assumptions on the significance level of the test. The relatively few studies of the effect on the power function include several papers by David and Johnson, the latest of which is [6]. For some problems tests have been proposed whose significance level is little affected by certain deviations from standard assumptions, for instance R. A. Fisher's randomization tests (see section 3; see also Box and Andersen [4]). Some other relevant work will be mentioned later.

In sections 2, 3, and 4 we shall be concerned with problems of the following type. Let  $P$  denote the joint distribution of the random variables under observation. Suppose that we contemplate making the assumption that  $P$  belongs to a class  $\rho_1$ , but we admit the possibility that actually  $P$  is contained in another class,  $\rho_2$ . The performance of a decision rule (decision function)  $d$  is assumed to be expressed by the given risk function  $r(P, d)$ , defined for all  $P \in \rho_1 + \rho_2$  and all  $d$  in  $D$ , the class of decision rules available to the statistician. Let  $d_i$  be a decision rule which is optimal in some specified sense (for example, minimax) under the assumption  $P \in \rho_i$ ,  $i = 1, 2$ . Suppose first that the optimal rule  $d_i$  is unique except for equivalence in  $\rho_1 + \rho_2$ , for  $i = 1, 2$ , that is, if  $d'_i$  is also optimal for  $P \in \rho_i$  then  $r(P, d'_i) = r(P, d_i)$  for all  $P \in \rho_1 + \rho_2$ . Then we may assess the consequences of the assumption  $P \in \rho_1$  when actually  $P \in \rho_2$  by comparing the values  $r(P, d_1)$  and  $r(P, d_2)$  for  $P \in \rho_2$ . If the optimal rules are not unique, we may pick out from the class of rules which are optimal for  $P \in \rho_1$  a subclass of rules which come closest to optimality under the assumption  $P \in \rho_2$ , and compare their performance with that of the rules which are optimal under the latter assumption. In

This research was supported by the United States Air Force through the Office of Scientific Research of the Air Research and Development Command, and by the Office of Ordnance Research, U.S. Army, under Contract DA-04-200-ORD-355.

some situations other ways of approaching the problem may be more adequate (see, for example, section 2).

In section 2 the consequences of assuming that a distribution is continuous are discussed. Problems involved in comparing assumptions of varying generality are considered in section 3. Section 4 is concerned with cases where decision rules derived under assumptions of normality retain their optimal properties when these assumptions are relaxed.

The last three sections deal with distinguishable sets of distributions, a concept related to the problem of the existence of unbiased or consistent tests under given assumptions. Criteria for the distinguishability of two sets by means of a test based on finitely many observations and by a sequential test are considered and their uses illustrated in sections 5 and 7. An example where two sets are indistinguishable by a nonrandomized test, but distinguishable by a randomized test, is discussed in section 6.

## 2. The assumption of a continuous distribution

The assumption that we are dealing with a class of continuous distributions is usually made when actually the observations are integer multiples of the unit of measurement  $h$ , a (small) positive constant. Suppose that a sample  $x = (x_1, \dots, x_n)$  is a point in  $R^n$ , and let  $\rho_1$  be a class of distributions (probability measures) which are absolutely continuous with respect to  $n$ -dimensional Lebesgue measure. Let  $S$  be the set of all points in  $R^n$  whose coordinates are integer multiples of  $h$ . Let us suppose that when we say that the distribution is  $P_1 \in \rho_1$ , we "have in mind" that the distribution is  $P_2 = f(P_1)$ , where the probability measure  $P_2$  is defined by

$$(1) \quad P_2(\{y\}) = P_1\left(\left\{x: y_j - \frac{h}{2} < x_j \leq y_j + \frac{h}{2}, j = 1, \dots, n\right\}\right)$$

for all  $y = (y_1, \dots, y_n)$  in  $S$ . Let  $\rho_2 = \{f(P): P \in \rho_1\}$ . Thus we are interested in the consequences of assuming  $P \in \rho_1$  when actually  $P \in \rho_2$ .

Let  $d$  be a decision function which is optimal in some sense under the assumption  $P \in \rho_1$ . Then any decision rule which differs from  $d$  only on the set  $S$  is equivalent to  $d$  for  $P \in \rho_1$ . Since  $P(S) = 1$  for all  $P \in \rho_2$ , the mere fact that a rule is optimal for  $P \in \rho_1$  does not tell us anything about its performance when  $P \in \rho_2$ ; indeed, it can be as bad as we please under the latter assumption.<sup>1</sup> Of course, in general there are rules which are optimal under either assumption. But the main reason for making the simplifying assumption of continuity is that we do not want to bother with rules which are optimal for  $P \in \rho_2$ . Now it is clear that if there is a determination  $d'$  of  $d$  which is sufficiently regular, its risk at  $P_2 = f(P_1)$  will differ arbitrarily little from the risk at  $P_1$  if  $h$  is small enough; also,  $d'$  may not be much worse than an optimal rule for  $P \in \rho_2$ . We shall not investigate here under what conditions such a regular decision rule exists or how small  $h$  has to be in order that the assumption of continuity cause little harm. These questions may deserve attention. Fortunately, when a statistician applies a decision rule, he is likely to choose the most regular determination available anyway. However, the theoretical statistician might do well to be careful when he neglects sets of measure zero.

<sup>1</sup> The author's attention was drawn to situations of this kind by H. Robbins some years ago.

### 3. Assumptions of varying generality

Suppose we consider making one of two assumptions,  $P \in \rho_1$  and  $P \in \rho_2$ , where  $\rho_1 \subset \rho_2$ . The second assumption is safer, but with the first assumption we may achieve a smaller risk.

The consequences of making the broader assumption when actually the narrower assumption is justified may be called serious if any decision rule which is "good" under the broader assumption is "bad" under the narrower assumption. Thus the consequences will depend on what we mean by a good decision rule. But even with a given definition of "good" or "best" the consequences may depend on the class of decision rules at our disposal. For example, suppose we require a minimax estimator of the mean  $\mu$  of a normal distribution when the loss function is the squared deviation from  $\mu$ , and we assume that the variance  $\sigma^2$  does not exceed a given number  $A$ . If we are restricted to estimators based on a sample of fixed size, the minimax estimator is the sample mean  $\bar{x}$  and does not depend on  $A$ . On the other hand, if we are permitted to choose the sample size in advance, and the cost of sampling is taken into account, the minimax estimator will depend on  $A$ . If  $A_2$  is substantially larger than  $A_1$ , the assumption  $\sigma^2 \leq A_2$  will give us a unique minimax estimator whose performance is poor under the assumption  $\sigma^2 \leq A_1$ .

Sometimes a considerable broadening of the assumption does not lead to serious consequences when the narrower assumption is justified. Thus in the standard problems concerning the variance of a normal distribution we need, when the mean is completely unknown, just one more observation to obtain the same expected loss as when the mean is known. Somewhat similar results have been obtained in certain cases where a parametric class of distributions is enlarged to a nonparametric class. Several examples can be found in [9]. For instance, consider the problem of testing whether two distributions are equal (and not otherwise specified) against the alternative that the distributions are normal with common variance and means  $\mu_1 < \mu_2$ . The uniformly most powerful similar test, based on two random samples of fixed size, is asymptotically as powerful in large samples (in a sense explained in [9]) as the corresponding standard test for testing the equality of the means of two normal distributions. (The former test is of the randomization type introduced by R. A. Fisher; its optimal properties were proved by Lehmann and Stein [12].) Here we assumed that the class of alternatives is the same under both assumptions. Actually the test retains its property of being uniformly most powerful similar even when the class of alternatives is enlarged to a nonparametric class of distributions of an exponential type (see Lehmann and Stein [12]). If the class is further extended, a uniformly most powerful similar test will in general not exist, and it will be necessary to specify against what types of alternatives the power of a test should be large. This can be done in many ways, and an optimal test and its performance in the class of normal distributions will depend on this specification.

### 4. Nonparametric justifications of assumptions of normality

Given a decision rule  $d$  which is optimal in a specified sense under the assumption that  $P$  is in a class  $\rho_1$ , it is of interest to determine other classes  $\rho$  such that  $d$  is optimal (in the same or a suitably extended sense) under the assumption  $P \in \rho$ . If optimal means minimax, an obvious sufficient condition for  $d$  to remain a minimax rule in  $\rho \supset \rho_1$  is that the risk of  $d$  in  $\rho$  attain its maximum in  $\rho_1$ . Situations of this type were considered by Hodges and Lehmann [8].

In certain cases we find that a decision rule derived under the assumption of a normal distribution retains its optimal character in a large, nonparametric class of distributions. One result of this type, concerning the minimax character of Markov estimators, can be found in [8]. Similar though weaker results can be obtained in certain testing problems.

As an example consider the following extension of Student's problem. Let  $\mathcal{Q}$  be the class of distributions  $F$  with finite mean  $\mu(F)$ , positive variance  $\sigma^2(F)$  and such that

$$(2) \quad \int_{-\infty}^{\infty} |x - \mu(F)|^2 dF(x) \leq M\sigma^2(F),$$

where  $M$  is fixed. Let  $\mathcal{Q}_\delta$  be the subclass of  $\mathcal{Q}$  with  $\sqrt{n} \mu(F)/\sigma(F) = \delta$ . We want to test the hypothesis  $F \in \mathcal{Q}_\delta, \delta \leq 0$ , against the alternative  $F \in \mathcal{Q}_\delta, \delta > 0$ . We restrict ourselves to the class  $D$  of tests  $d$  based on  $n$  independent observations from  $F$ , with critical region  $W = W(d)$ . We choose the risk function

$$(3) \quad r(F, d) = \begin{cases} aP(W|F) & \text{if } F \in \mathcal{Q}_\delta, \delta \leq 0, \\ b[1 - P(W|F)] & \text{if } F \in \mathcal{Q}_\delta, \delta \geq \delta_1, \\ 0 & \text{elsewhere,} \end{cases}$$

where  $P(W|F)$  denotes the probability of  $(X_1, \dots, X_n) \in W$  when the  $X_j$  are independent with the common distribution  $F$ , and  $a, b$ , and  $\delta_1$  are positive constants.

Let  $d_0$  be the test with critical region  $W_0 = \{t > c\}$ , where

$$(4) \quad t = \frac{n^{1/2}\bar{x}}{s}, \quad \bar{x} = n^{-1} \sum_{j=1}^n x_j, \quad s^2 = (n-1)^{-1} \sum_{j=1}^n (x_j - \bar{x})^2,$$

and the constant  $c$  is determined by

$$(5) \quad a[1 - S_{n-1}(c, 0)] = bS_{n-1}(c, \delta_1);$$

here  $S_{n-1}(x, \delta)$  denotes the noncentral Student distribution function with  $n-1$  degrees of freedom and noncentrality parameter  $\delta$ . It can be shown by standard methods that  $d_0$  is the minimax test in the subclass  $\mathcal{Q}^0$  of  $\mathcal{Q}$  which consists of the normal distributions. By an inequality of Berry and Esseen (see, for example, [7]) the distribution function  $F_n(y)$  of  $n^{1/2}[\bar{X} - \mu(F)]/\sigma(F)$  converges to the standard normal distribution function  $\Phi(y)$  uniformly for  $F \in \mathcal{Q}$  (and uniformly in  $y$ ) as  $n \rightarrow \infty$ . Also, for any  $\epsilon > 0$ ,  $P[|s/\sigma(F) - 1| < \epsilon|F] \rightarrow 1$  uniformly for  $F \in \mathcal{Q}$ . Hence it can be shown that for any real  $\delta$  and for all  $F \in \mathcal{Q}_\delta$  we have

$$(6) \quad |P(t \leq y|F) - \Phi(y - \delta)| \leq C_n(\delta), \quad -\infty < y < \infty,$$

where  $C_n(\delta)$  depends on  $n, \delta$ , and  $M$  only and tends to 0 as  $n \rightarrow \infty$ , for  $\delta$  fixed. It follows that

$$(7) \quad |P(t \leq y|F) - S_{n-1}(y|\delta)| \leq 2C_n(\delta), \quad -\infty < y < \infty$$

for all  $F \in \mathcal{Q}_\delta$ .

Now if  $\mathcal{Q}^*$  denotes the subclass of  $\mathcal{Q}$  with  $\mu(F) = 0$  and  $\sigma^2(F) = 1$ , we have

$$(8) \quad \sup_{F \in \mathcal{F}_\delta} P(t > c|F) = \sup_{F \in \mathcal{F}^*} P\left[\frac{n^{1/2}\bar{x} + \delta}{s} > c|F\right],$$

which is a nondecreasing function of  $\delta$ . The same is true of the infimum in  $\mathcal{Q}_\delta$ . Hence we obtain

$$(9) \quad \begin{aligned} \sup_{F \in \mathcal{F}} r(F, d_0) &\leq \sup_{F \in \mathcal{F}} r(F, d_0) + \epsilon, \\ &\leq \inf_{d \in \mathcal{D}} \sup_{F \in \mathcal{F}} r(F, d) + \epsilon \end{aligned}$$

where  $\epsilon = 2 \max \{aC_n(0), bC_n(\delta_1)\}$ . Thus the maximum risk in  $\mathcal{Q}$  of Student's test  $d_0$  exceeds the minimax risk in  $\mathcal{Q}$  by at most  $\epsilon$ , where  $\epsilon$  is arbitrarily small for  $n$  sufficiently large. (Note that the minimax risk is bounded away from zero as  $n \rightarrow \infty$ .)

In the corresponding problem with  $\sigma(F) = \sigma$  fixed we find in a similar way the stronger result that the maximum risk of the  $\bar{x}$ -test in  $\mathcal{Q}_\delta$  [the class with  $\mu(F) = \delta\sigma/\sqrt{n}$ ] lies within a small  $\epsilon$  of its "normal" risk, uniformly in  $\delta$ . The argument which was used above does not permit us to decide whether an analogous result is true when  $\sigma(F)$  is unrestricted.

The explanation for the near-optimal behavior of the "normal" decision rules in these cases is, of course, the distribution-free character of the central limit theorem, combined with the fact that the class  $\mathcal{Q}$  was so chosen as to make the approach to the normal distribution uniform.

### 5. Distinguishable sets of distributions

If we relax the assumptions more and more, the minimax risk will in general increase, and eventually we may reach a point where the maximum risk of any decision rule is not smaller than the risk of a rule which does not depend on the observations. We shall consider criteria for recognizing when this or a similar situation occurs in testing problems.

Consider a testing (or two decision) problem such that one or the other decision is definitely preferred according as the distribution  $P$  belongs to  $\rho^1$  or  $\rho^2$ , two disjoint subsets of the given class  $\rho$ . Unless otherwise stated we assume that each  $P$  in  $\rho$  is a probability measure on  $(\mathcal{X}, \mathcal{A})$ , where  $\mathcal{X}$  is the space of infinite sequences  $x = (x_1, x_2, \dots)$  of real numbers and  $\mathcal{A}$  is Kolmogorov's extension to  $\mathcal{X}$  of the ordinary Borel field.

A test will be called finite if it depends only on a finite number of coordinates (observations)  $x_j$ . By the critical function of a finite test we mean a measurable function  $\psi$  from  $\mathcal{X}$  to the interval  $[0, 1]$  such that  $1 - \psi(x)$  [ $\psi(x)$ ] is the probability of taking the decision corresponding to  $P \in \rho^1$  [ $P \in \rho^2$ ] when  $x$  is the sequence of observations.

Let  $D$  be any class of finite tests, and let  $\Psi$  be the class of the critical functions of tests in  $D$ . We shall say that the sets  $\rho^1$  and  $\rho^2$  are distinguishable in  $D$  if there exists a  $\psi$  in  $\Psi$  such that

$$(10) \quad \sup_{P \in \rho^1} E(\psi | P) + \sup_{P \in \rho^2} E(1 - \psi | P) < 1,$$

where  $E(f|P)$  is the expected value of  $f(X)$  when  $X$  has the distribution  $P$ . Otherwise  $\rho^1$  and  $\rho^2$  are said to be indistinguishable in  $D$ . [The property of  $\psi$  expressed in (10) has an obvious relation with unbiasedness.]

Let  $D_f$  denote the class of all finite tests, and let  $D_n, n = 1, 2, \dots$ , be the class of all fixed sample size tests based on the observations  $(x_1, \dots, x_n)$ . Two sets which are distinguishable in  $D_f$  will be called finitely distinguishable. We observe that two sets are finitely distinguishable if and only if they are distinguishable in  $D_n$  for some  $n$ .

Berger and Wald [2] gave conditions under which two sets of distributions are dis-

tinguishable in the class of all nonrandomized tests in  $D_n$  if and only if they are disjoint. (Their theorem 3.1 is stated in a slightly more special form.)

A sufficient condition for two sets to be indistinguishable in  $D_n$  can be stated as follows. Let  $\mathcal{X}_n$  be the space of points  $(x_1, \dots, x_n)$ , and let  $\mathcal{A}_n$ ,  $\rho_n$  and  $\rho_n^i$  be the  $\sigma$ -field of subsets of  $\mathcal{X}_n$  and the classes of distributions on  $\mathcal{A}_n$  which are determined by  $\mathcal{A}$ ,  $\rho$  and  $\rho^i$  in an obvious way. For any two distributions  $P_1$  and  $P_2$  on  $\mathcal{A}_n$  we denote by  $\nu$  any measure on  $\mathcal{A}_n$  relative to which  $P_1$  and  $P_2$  are absolutely continuous, and by  $p_1$  and  $p_2$  the respective densities (Radon-Nikodym derivatives). With this notation, the sets  $\rho^1$  and  $\rho^2$  (or, equivalently, the sets  $\rho_n^1$  and  $\rho_n^2$ ) are indistinguishable in  $D_n$  if for any  $\epsilon > 0$  there exist two distributions,  $P_1 \in \rho_n^1$  and  $P_2 \in \rho_n^2$ , such that

$$(11) \quad \int |p_1 - p_2| d\nu < \epsilon,$$

where the integral extends over  $\mathcal{X}_n$ . This follows from the inequality

$$(12) \quad \inf_{P \in \rho^1} \int \psi dP - \sup_{P \in \rho^2} \int \psi dP \leq \int \psi (p_2 - p_1) d\nu.$$

The statement of the condition remains true in the more general form where  $P_i$  is any mixture of distributions in  $\rho_n^i$  with respect to some probability measure  $\xi_i$  on a  $\sigma$ -field of subsets of  $\rho_n^i$ , subject to an obvious measurability condition. The proof is similar and uses theorem 3 of Robbins [13]. A theorem of Le Cam (see Kraft [15]) implies that if the distributions in  $\rho_n^1$  and  $\rho_n^2$  are absolutely continuous with respect to a fixed measure, the condition expressed in (11), with  $P_1$  and  $P_2$  mixtures, is also necessary for the indistinguishability of  $\rho_n^1$  and  $\rho_n^2$ .

With  $S = \{x: p_1(x) > p_2(x)\}$  we have

$$(13) \quad \frac{1}{2} \int |p_1 - p_2| d\nu = \sup_{A \in \mathcal{A}_n} |P_1(A) - P_2(A)| = P_1(S) - P_2(S).$$

The first equation (13) shows that condition (11) is independent of the choice of  $\nu$ . The last expression in (13) is often convenient when applying this condition.

It follows from an earlier remark that two sets  $\rho^1$  and  $\rho^2$  are finitely indistinguishable if the condition expressed in (11) is satisfied for every  $n$ .

We shall say that  $\rho^1$  and  $\rho^2$  are finitely absolutely distinguishable if for any  $\epsilon > 0$  there exists a finite test with critical function  $\psi$  such that

$$(14) \quad \sup_{P \in \rho_n^1} E(\psi | P) + \sup_{P \in \rho_n^2} E(1 - \psi | P) < \epsilon.$$

This property has also been expressed by saying that there exists a uniformly consistent sequence of tests [1].

Now suppose that each  $P$  in  $\rho$  is the distribution of a sequence of independent, identically distributed random variables. Then if two sets are finitely distinguishable, they are finitely absolutely distinguishable. This is a simple partial extension of a theorem of Berger [1]; the theorem gives a necessary and sufficient condition for the existence of a uniformly consistent sequence of nonrandomized tests. Further interesting results on the existence of a uniformly consistent sequence of tests were recently obtained by Kraft [15].

We now give three examples of finitely indistinguishable sets.

*Example 5.1.* If  $P$  is the distribution of independent, normal random variables with mean  $\mu$  and variance  $\sigma^2$ , and  $\rho^i$  is the set with  $\mu = \mu_i$ ,  $0 < \sigma^2 < \infty$ , then  $\rho^1$  and  $\rho^2$  are finitely indistinguishable. Condition (11) is satisfied for every  $n$  if  $P_i$  is the distribution

with  $\mu = \mu_i$  and  $\sigma$  sufficiently large. The corresponding result for tests with constant power in  $\rho^1$  and  $\rho^2$  was proved by Dantzig [5] in 1940.

*Example 5.2.* If  $P$  is the distribution of independent, normal random variables with means  $\mu_1, \mu_2, \dots$  and common variance  $\sigma^2$ , and  $\rho^i$  the set with  $\sigma = \sigma_i, -\infty < \mu_j < \infty, j = 1, 2, \dots$ , then  $\rho^1$  and  $\rho^2$  are finitely indistinguishable. Here we can apply the general form of condition (11). For if  $P_i$  is the mixture of the  $P$  in  $\rho_n^i$  according to  $\xi_i$ , where under  $\xi_i$  the means  $\mu_1, \dots, \mu_n$  are independent normal with zero mean and variance  $\tau_i^2$ , such that  $\sigma_1^2 + \tau_1^2 = \sigma_2^2 + \tau_2^2$ , then  $P_1 = P_2$ .

*Example 5.3.* This is a further extension of Student's problem (see section 4). Let  $\mathcal{Q}^i$  be the class of all distributions  $F$  on the real line with finite mean  $\mu(F)$  and positive variance  $\sigma^2(F)$  such that  $\mu(F)/\sigma(F) = \gamma_i, \gamma_1 < \gamma_2$ . Let  $\rho^i$  be the class of distributions of independent random variables with common distribution  $F \in \mathcal{Q}^i$ . Then  $\rho^1$  and  $\rho^2$  are finitely absolutely distinguishable if  $\gamma_1 < 0 < \gamma_2$ , and finitely indistinguishable if  $\gamma_2 \leq 0$  or  $\gamma_1 \geq 0$ .

If  $\gamma_1 < 0 < \gamma_2$ , it is easy to show with the aid of Chebyshev's inequality that the tests with critical functions  $\psi_n(x) = 0$  or 1 according as  $\sum_{j=1}^n x_j \leq 0$  or  $> 0$  form a uniformly consistent sequence.

If  $\gamma_1 \geq 0$ , condition (11) is satisfied for every  $n$  if  $P_i$  is the distribution with  $F = F_i$ , where  $F_i$  ascribes probabilities  $1 - \pi_i$  and  $\pi_i = (1 + t_i^2)^{-1}$  to the respective points  $\gamma_2 - t_2^{-1}$  and  $\gamma_2 + t_2$ ; here  $t_2 > 0, t_1 = f(t_2)$  is the positive root (unique for  $t_2$  small) of

$$(15) \quad (1 - \gamma_2 t_2) t_1^2 + \gamma_1 (1 + t_2^2) t_1 - \gamma_2 t_2 - t_2^2 = 0,$$

and  $t_2 \rightarrow 0$ . The case  $\gamma_2 \leq 0$  can be reduced to this case.

### 6. Sets distinguishable only by randomized tests: An example

Some results of Lehmann [11] suggest that two sets may be distinguishable in  $D_n$  but indistinguishable in the class  $D'_n$  of nonrandomized tests in  $D_n$ . We shall consider a problem where this situation occurs. We denote by  $\Psi_n(\Psi'_n)$  the class of critical functions of the tests in  $D_n(D'_n)$ . Thus if  $\psi \in \Psi'_n, \psi(x) = 0$  or 1 for all  $x$ .

Let  $\mathcal{Q}_\mu$  be a class of distributions  $F$  on the real line with mean  $\mu$  and variance 1, which contains all distributions with this property which assign probability 1 to at most three points. Let  $\rho_{\mu, n}$  be the class of all distributions of  $n$  independent random variables with a common distribution in  $\mathcal{Q}_\mu$ . We shall show that  $\rho_{\lambda, n}$  and  $\rho_{\mu, n}$  are distinguishable in  $D_n$  for all  $\lambda \neq \mu$  and all  $n = 1, 2, \dots$ , but indistinguishable in  $D'_n$  for any  $n$  unless  $|\lambda - \mu|$  exceeds a positive constant (which depends on  $n$ ). It is clearly sufficient to take  $\lambda = -h, \mu = h > 0$ . We denote by  $E(f|F)$  the expected value of  $f(X)$  when the components of  $X$  are independent with the common distribution  $F$ .

We first prove the second part of the statement in the stronger form: For any  $n$  and for any  $\alpha \in (0, 1)$  the inequalities

$$(16) \quad \sup_{F \in \mathcal{F}_{-h}} E(\psi|F) \leq \alpha \leq \inf_{F \in \mathcal{F}_h} E(\psi|F)$$

cannot both be satisfied with  $\psi \in \Psi'_n$  unless  $h$  exceeds a positive number which depends only on  $n$  (and is of order  $n^{-1/2}$ ). If  $\psi$  is in  $\Psi'_n$  and satisfies the first inequality (16), we must have

$$(17) \quad \psi(y, \dots, y) = 0 \quad \text{if} \quad \alpha [1 + (y + h)^2]^n < 1,$$



for all real  $y$ . For if  $t = y + h \neq 0$ , let  $F'$  be the distribution (in  $\mathcal{Q}_{-h}$ ) which assigns the probabilities  $(1 + t^2)^{-1}$  and  $1 - (1 + t^2)^{-1}$  to the respective points  $t - h$  and  $-t^{-1} - h$ . Then  $a \geq E(\psi | F') \geq \psi(t - h, \dots, t - h)(1 + t^2)^{-n}$ . This implies (17) for  $y + h \neq 0$ . If  $y + h = 0$ , we use a similar argument with  $F'$  any distribution in  $\mathcal{Q}_{-h}$  which assigns to the point  $-h$  a probability arbitrarily close to 1.

Similarly, for any  $\psi \in \Psi'_n$  which satisfies the second inequality (16) we must have

$$(18) \quad \psi(y, \dots, y) = 1 \quad \text{if} \quad (1 - a)[1 + (y - h)^2]^n < 1,$$

for all real  $y$ . Taking  $y = -h$  and  $y = h$ , we find that a  $\psi \in \Psi'_n$  cannot satisfy both inequalities (16) if  $[1 + (2h)^2]^n < \max [a^{-1}, (1 - a)^{-1}]$ ; and hence cannot satisfy them for any  $a$  if  $[1 + (2h)^2]^n < 2$ . [This is not the best bound which can be obtained from (17) and (18).]

We now show that for any  $h > 0$ , any  $n \geq 1$ , and any  $a \in (0, 1)$  condition (16), with at least one strict inequality, can be satisfied by a randomized test in  $D_n$ . Let  $a = hn^{1/2}$ ,  $-a < c < a$ ,

$$(19) \quad k(c) = a + c + (a - c)^{-1}, \quad b = \frac{-k(-c)}{2}, \quad d = \frac{k(c)}{2},$$

$$(20) \quad \phi(y) = \begin{cases} 0 & \text{if } y \leq b, \\ \frac{y - b}{d - b} & \text{if } b < y < d, \\ 1 & \text{if } d \leq y. \end{cases}$$

If we let  $\psi(x) = \phi\left(n^{-1/2} \sum_1^n x_j\right)$ , we have

$$(21) \quad \sup_{F \in \mathcal{F}_{-h}} E(\psi | F) \leq \frac{1}{2ak(c)} < 1 - \frac{1}{2ak(-c)} \leq \inf_{F \in \mathcal{F}_h} E(\psi | F)$$

for  $|c| < a$ . As  $c$  increases from  $-a$  to  $a$ , either side of (21) decreases continuously from 1 to 0.

We sketch the proof of (21). Let  $f(y)$  be any polynomial of the second degree such that  $\phi(y) \leq f(y)$  for all real  $y$ . If  $g(x) = f(n^{-1/2} \sum x_j)$ , then  $E(\psi | F) \leq E(g | F)$ , and  $E(g | F)$  is constant in  $\mathcal{Q}_\mu$  for each  $\mu$ . Now choose  $f$  so as to minimize  $E(g | F)$ ,  $F \in \mathcal{Q}_\mu$ .

### 7. Sequentially distinguishable sets of distributions

We shall restrict ourselves to sequences of independent random variables with a common distribution  $F$ . Suppose that  $F \in \mathcal{Q}$ , and let  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  be two disjoint subsets of  $\mathcal{Q}$ . Let  $D_s = D_s(\mathcal{Q})$  be the class of all sequential tests for taking one of two decisions,  $a_1$  and  $a_2$ , which terminate with probability one for all  $F \in \mathcal{Q}$ . We denote by  $Pr\{a_i | F, d\}$  and  $E(n | F, d)$ , respectively, the probability of the decision  $a_i$  and the expected number of observations required to reach a decision when the distribution is  $F$  and test  $d$  is used.

The sets  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  will be called sequentially distinguishable (indistinguishable) at  $F$  if there exists (does not exist) a  $d$  in  $D_s$  such that  $E(n | F, d) < \infty$  and

$$(22) \quad \sup_{F \in \mathcal{F}^1} Pr\{a_2 | F, d\} + \sup_{F \in \mathcal{F}^2} Pr\{a_1 | F, d\} < 1.$$

If the left side of (22) is arbitrarily small for some  $d$  in  $D$ , with  $E(n|F, d) < \infty$ , then  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  are said to be sequentially absolutely distinguishable at  $F$ . If  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  are sequentially [absolutely] distinguishable (indistinguishable) at every  $F$  in a class  $\mathcal{Q}^*$ , then  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  will be said to be sequentially [absolutely] distinguishable (indistinguishable) in  $\mathcal{Q}^*$ .

Note that these definitions are stated in terms of the sets  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  rather than in terms of the corresponding sets of distributions of sequences. Statements such as  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  are finitely indistinguishable will have an obvious meaning in this context.

A sufficient condition for two sets to be sequentially indistinguishable is implied by an inequality proved in [10]. Let  $F_1 \in \mathcal{Q}^1, F_2 \in \mathcal{Q}^2, F \in \mathcal{Q}$ , and let  $\nu$  be a measure relative to which these three distributions are absolutely continuous, with respective densities  $f_1, f_2$ , and  $f$ . By a trivial extension of equation (4) in [10], if  $d$  is any test in  $D$ , such that

$$(23) \quad \sup_{F \in F^1} Pr \{ a_2 | F, d \} \leq a_1, \quad \sup_{F \in F^2} Pr \{ a_1 | F, d \} \leq a_2$$

where  $a_1 > 0, a_2 > 0, a_1 + a_2 < 1$ , then

$$(24) \quad E(n|F, d) \geq \frac{-\log [ a_1^c (1 - a_2)^{1-c} + (1 - a_1)^c a_2^{1-c} ]}{c \int f \log \frac{f}{f_1} d\nu + (1 - c) \int f \log \frac{f}{f_2} d\nu}$$

for  $0 < c < 1$ , where the integrals are taken over the entire space. If, in particular,  $F \in \mathcal{Q}^1$ , the right side of (24) is maximized with  $F_1 = F$  and  $c \rightarrow 0$ , and we obtain

$$(25) \quad E(n|F, d) \geq \frac{a_1 \log \frac{a_1}{1 - a_2} + (1 - a_1) \log \frac{1 - a_1}{a_2}}{\int f \log \frac{f}{f_2} d\nu} \quad \text{if } F \in \mathcal{Q}^1.$$

We note that the numerators and denominators in (24) and (25) are positive; the denominators may be infinite.

Hence if for any positive number  $M$  and any two positive numbers  $a_1$  and  $a_2$  with  $a_1 + a_2 < 1$  the distributions  $F_1 \in \mathcal{Q}^1$  and  $F_2 \in \mathcal{Q}^2$  and the number  $c$  can be so chosen that the right side of (24) exceeds  $M$ , the sets  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  are sequentially indistinguishable at  $F$ . If  $F \in \mathcal{Q}^1$ , the two sets are sequentially indistinguishable at  $F$  if for any  $\epsilon > 0$  we can find an  $F_2 \in \mathcal{Q}^2$  such that

$$(26) \quad \int f \log \frac{f}{f_2} d\nu < \epsilon.$$

By example 5.1 two sets of normal distributions with fixed means and unrestricted variances are finitely indistinguishable. On the other hand, by a well-known result of Stein [14], these sets are sequentially absolutely distinguishable in the class of all normal distributions. However, if the requirement  $E(n|F, d) < \infty$  is replaced by the stronger condition that  $E(n|F, d) = E(n|\mu, \sigma; d)$  be bounded in  $\sigma$  for  $\mu$  fixed, inequality (24) easily implies that condition (22) cannot be satisfied.

As an application of condition (26) we shall show that the sets  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  of example 5.3, with  $\gamma_1 = 0 < \gamma_2$ , are sequentially indistinguishable in  $\mathcal{Q}^2$ . Let  $F$  be any distribution in  $\mathcal{Q}^2$ , so that  $\mu(F)/\sigma(F) = \gamma_2$ . Let  $F_1 = (1 - t)F + tG$ , where  $0 < t < 1$  and  $G$  is the distribution which assigns probability one to the point  $a = -\mu(F)(1 - t)/t$ . Then

$F_1 \in \mathcal{Q}^1$ . Both  $F_1$  and  $F$  are absolutely continuous relative to  $\nu = F_1$ , with respective densities  $f_1(x) = 1$  and

$$(27) \quad f(x) = \begin{cases} \frac{1}{1-t} & \text{if } x \neq a, \\ \frac{b}{b+t-bt} & \text{if } x = a, \end{cases}$$

where  $b = b(t)$  is the  $F$ -probability of the point  $a$ . Hence

$$(28) \quad \int f \log \frac{f}{f_1} d\nu = -(1-b) \log(1-t) + b \log \frac{b}{b+t-bt},$$

where the last term is to be omitted if  $b = 0$ . The right side of (28) tends to 0 as  $t \rightarrow 0$ . Thus condition (26) (with  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  interchanged) can be satisfied for any  $\epsilon > 0$ .

The proof shows that this result still holds if  $\mathcal{Q}^1$  and  $\mathcal{Q}^2$  consist only of the mixtures  $H = (1-t)F + tG$  of a normal distribution  $F$  and an arbitrary distribution  $G$ , where  $0 \leq t < \epsilon$  and  $\epsilon$  is positive and as small as we please. The distributions  $H$  are, in a sense, very close to normal distributions.

#### REFERENCES

- [1] A. BERGER, "On uniformly consistent tests," *Annals of Math. Stat.*, Vol. 22 (1951), pp. 289-293.
- [2] A. BERGER and A. WALD, "On distinct hypotheses," *Annals of Math. Stat.*, Vol. 20 (1949), pp. 104-109.
- [3] G. E. P. BOX, "Non-normality and tests on variances," *Biometrika*, Vol. 40 (1953), pp. 318-335.
- [4] G. E. P. BOX and S. L. ANDERSEN, "The robust tests for variances and effect of non-normality and variance heterogeneity on standard tests," *Institute of Statistics Mimeo Series, No. 101*, Consolidated University of North Carolina, 1954.
- [5] G. B. DANTZIG, "On the non-existence of tests of 'Student's' hypothesis having power functions independent of  $\sigma$ ," *Annals of Math. Stat.*, Vol. 11 (1940), pp. 186-192.
- [6] F. N. DAVID and N. L. JOHNSON, "Extension of a method of investigating the properties of analysis of variance tests to the case of random and mixed models," *Annals of Math. Stat.*, Vol. 23 (1952), pp. 594-601.
- [7] B. V. GNEDENKO and A. N. KOLMOGOROV, *Limit Distributions for Sums of Independent Random Variables*, Cambridge, Mass., Addison-Wesley, 1954.
- [8] J. L. HODGES, JR. and E. L. LEHMANN, "Some problems in minimax point estimation," *Annals of Math. Stat.*, Vol. 21 (1950), pp. 182-197.
- [9] W. Hoeffding, "The large-sample power of tests based on permutations of observations," *Annals of Math. Stat.*, Vol. 23 (1952), pp. 169-192.
- [10] ———, "A lower bound for the average sample number of a sequential test," *Annals of Math. Stat.*, Vol. 24 (1953), pp. 127-130.
- [11] E. L. LEHMANN, "Consistency and unbiasedness of certain nonparametric tests," *Annals of Math. Stat.*, Vol. 22 (1951), pp. 165-179.
- [12] E. L. LEHMANN and C. STEIN, "On the theory of some non-parametric hypotheses," *Annals of Math. Stat.*, Vol. 20 (1949), pp. 28-45.
- [13] H. ROBBINS, "Mixture of distributions," *Annals of Math. Stat.*, Vol. 19 (1948), pp. 360-369.
- [14] C. STEIN, "A two-sample test for a linear hypothesis whose power is independent of the variance," *Annals of Math. Stat.*, Vol. 16 (1945), pp. 243-258.
- [15] C. KRAFT, "Some conditions for consistency and uniform consistency of statistical procedures," *Univ. Calif. Publ. Stat.*, Vol. 2 (1956), pp. 125-142.