

# DECISION THEORY FOR SOME NONPARAMETRIC MODELS

KJELL A. DOKSUM  
UNIVERSITY OF CALIFORNIA, BERKELEY

## 1. Introduction and summary

The problem considered in this paper is that of obtaining optimal decision rules when a parametric form of the distribution of the observations is not known exactly. Thus we assume that the underlying distribution function  $F$  of the  $X_i$  in the random sample  $X = (X_1, \dots, X_n)$  is in a class  $\Omega$  of distribution functions, and  $\Omega$  is not indexed in a natural way by a parameter  $\theta$  in  $m$  dimensional Euclidean space  $R^m$ . Let  $R(F, d)$  denote the risk of the decision rule  $d = d(X)$  when  $F$  is the true distribution. Minimax procedures that minimize the maximum risk  $\sup \{R(F, d); F \in \Omega\}$  have been obtained in special cases by Hoeffding [8], Ruist [14], Huber [9], [10], [11], and Doksum [3]. In particular, Huber was able to show that if  $\Omega$  is the class of all distributions in a neighborhood of a normal distribution, then the minimax procedures are based on statistics that are, approximately, trimmed means. Most stringent procedures that minimize the maximum shortcoming  $\sup_F \{R(F, d) - \inf_d R(F, d)\}$  have been considered by Schaafsma [15].

Another approach would be to define a probability (weight function)  $P$  on  $\Omega$  and then minimize the average (Bayes) risk  $\int_{\Omega} R(F, d) P(dF)$ , thereby obtaining what is called the Bayes solution. This approach has been taken by Kraft and van Eeden [13], Ferguson [6], and Antoniak [1], who were able to obtain explicit Bayes solutions for some probabilities  $P$ . Their work is closely related to the work of Fabius [5], who considered properties of posterior distributions for a class of probability measures  $P$  that essentially contains those of Kraft and van Eeden and of Ferguson. Fabius' work in turn is related to that of Freedman [7], who considered properties of Bayes procedures in the case where the  $X_i$  are discrete random variables. The relationship between these papers will be discussed further in Section 5.

In this paper, we introduce a criterion which involves minimizing a quantity between the maximum risk and the average risk: This criterion is appropriate when the probability  $P$  on  $\Omega$  is not fully specified, but only the distribution of  $F(t_1), \dots, F(t_k)$  is known for some  $t_1 < \dots < t_k$ . Thus past records may

Prepared with the partial support of the National Science Foundation under Grant GP-15283, and with the partial support of the Office of Naval Research under Contract N00014-69-A-0200-1036. Completed while the author was at the University of Oslo.

be available for  $F(1)$ ,  $F(2)$ , and so on, but not for  $F(e)$ ,  $F(\pi)$  and so on. The criterion is to minimize the average maximum risk, where the average is computed with respect to the distribution  $\lambda$  of  $F(t_1), \dots, F(t_k)$ . More precisely, let  $t_1 < \dots < t_k$  be fixed and let  $\Omega(q, k)$  be the class of distribution functions in  $\Omega$  that pass through  $(t_1, q_1), (t_2, q_2), \dots, (t_k, q_k)$ ,  $0 \leq q_1 \leq \dots \leq q_k \leq 1$ . We define the average maximum risk (or mixed risk) as  $\int_{R^k} [\sup_{F \in \Omega(q, k)} R(F, d)] \lambda(dq)$ . The decision rule that minimizes this risk is called the mixed Bayes-minimax rule, or mixed rule. It will be shown in Sections 2 and 3 that, under certain conditions on the risk function, the mixed rule can be obtained by computing the posterior distribution of a multinomial parameter  $p = (p_1, \dots, p_{k-1})$  having as prior distribution the distribution  $\mu$  of  $F(t_2) - F(t_1), F(t_3) - F(t_2), \dots, F(t_k) - F(t_{k-1})$ .

The mixed Bayes-minimax rule  $d_k$  can be thought of as an approximation to the Bayes rule. In Section 4, Prohorov's theorem is used to show that if  $\Omega$  is contained in  $C[0, 1]$  or  $D[0, 1]$ , and if a Bayes solution  $d$  exists, then  $d_k$  converges to  $d$  in the sense of convergence of Bayes risks. Thus in situations where  $P$  is known, but the Bayes rule is hard to compute, one can use the mixed rule as an approximation. If the limit  $\lim_{k \rightarrow \infty} d_k$  can be computed, then it gives a method of obtaining the Bayes solution. Note that  $k = 0$  corresponds to the minimax problem.

## 2. The mixed Bayes-minimax problem

Let  $X_1, \dots, X_n$  be independent, identically distributed random variables with distribution function  $F$ , where  $F$  belongs to some specified class  $\Omega$  of distribution functions. It will be convenient to assume that  $\Omega$  is a measurable subset of some larger class  $\Gamma$  of functions with a  $\sigma$ -field  $\mathcal{S}$ . We further assume that there is a probability  $P$  on  $(\Gamma, \mathcal{S})$ , with  $P(\Omega) = 1$ .

Let  $L(F, d)$  be a real valued function that denotes the loss of the real valued decision rule  $d = d(X_1, \dots, X_n)$  when  $F \in \Omega$  is the true distribution. Let  $X = (X_1, \dots, X_n)$ . Then the risk of  $d$  is

$$(2.1) \quad R(F, d) = E_F L(F, d(X)).$$

If there is a rule (procedure) that minimizes the maximum risk

$$(2.2) \quad R(d) = \sup_{F \in \Omega} R(F, d),$$

then it is called a minimax procedure. Similarly, if there is a  $d$  that minimizes the average risk

$$(2.3) \quad r(P, d) = \int_{\Omega} R(F, d) P(dF),$$

then it is a Bayes procedure. Here  $F$  is thought of as a random distribution function with distribution  $P$ ; that is,  $F$  is a stochastic process with sample paths  $F(t)$ ,  $t \in R$ . Computing Bayes procedures will involve computing the posterior distribution of  $F$  given  $X$ .

The mixed Bayes-minimax procedure minimizes a function that is between  $R(d)$  and  $r(P, d)$ . This function is the average maximum risk when a finite dimensional distribution corresponding to  $P$  is known, and the average is taken with respect to this distribution. We now proceed with the definition. The carrier of a given distribution is in general the smallest compact set whose probability under the given distribution is one. Let  $C(F)$  denote the carrier of  $F \in \Omega$ ; we define the support of  $\Omega$  to be  $S(\Omega) = \cup_{F \in \Omega} C(F)$ . Let  $t_1 < \dots < t_k$  be  $k$  points in  $S(\Omega)$ . The distribution of  $F(t_1), \dots, F(t_k)$  under  $P$  will be denoted by  $\lambda$ , or  $\lambda(\cdot; P, k)$ . Thus, if we write  $q = (q_1, \dots, q_k)$ ,  $0 \leq q_1 \leq \dots \leq q_k \leq 1$ , then

$$(2.4) \quad \lambda(q; P, k) = P(F: F(t_1) \leq q_1, \dots, F(t_k) \leq q_k).$$

The class of distributions in  $\Omega$  whose value at  $t_i$  is  $q_i$ ,  $i = 1, \dots, k$ , will be denoted by  $\Omega(q, k)$ , that is,

$$(2.5) \quad \Omega(q, k) = \{F \in \Omega: F(t_i) = q_i, i = 1, \dots, k\}.$$

The average maximum risk of a decision rule  $d$  is now defined as

$$(2.6) \quad r_k(P, d) = \int_{R^k} [\sup_{F \in \Omega(q, k)} R(F, d)] \lambda(dq; P, k).$$

If there is a  $d_0$  that minimizes  $r_k(P, d)$ ,  $d_0$  is called a mixed Bayes-minimax (or mixed) procedure.

As we are dealing with functions of  $X \in R^n$  and  $F \in \Omega$ , we need a joint distribution for  $X$  and a random element  $F$  of  $\Omega$  with distribution  $P$ . If  $\mathcal{B}^n$  denotes the  $\sigma$ -field of Borel sets in  $R^n$ , we define the probability  $\tilde{P}$  on  $(R^n \times \Gamma, \mathcal{B}^n \times \mathcal{S})$  by

$$(2.7) \quad \tilde{P}(B \times S) = \int_S \mu_F(B) P(dF), \quad B \in \mathcal{B}^n, \quad S \in \mathcal{S},$$

where  $\mu_F$  is the probability in  $R^n$  corresponding to the distribution of  $X$ , and it is assumed that  $\mu_F$  is  $\mathcal{S}$  measurable. We also assume that a conditional distribution of  $X$  given  $F$  exists and satisfies

$$(2.8) \quad \tilde{P}(X_1 \leq x_1, \dots, X_n \leq x_n | F) = \prod_1^n F(x_i).$$

Furthermore, we assume that  $F$  has a conditional distribution given  $F(t_i) = q_i$ ,  $i = 1, \dots, k$ ; we denote this conditional distribution by  $P^q$ . For a further discussion of these definitions, see Fabius [5]. The assumptions involved in the definitions are satisfied for the complete, separable metric spaces considered in Section 4.

The following inequalities follow at once from the definitions. Note that when more than one set of  $t_1, \dots, t_k$  is considered, we will use double subscripts and write  $t_{m,1}, \dots, t_{m,k_m}$ .

LEMMA 2.1. *The risks defined above satisfy the following relations:*

(i)  $R(d) = r_0(P, d) \geq r_k(P, d) \geq r(P, d), \quad k \geq 1;$

(ii) *if  $\{\prod_m: t_{m,1} < \dots < t_{m,k_m}\}$ ,  $m = 1, 2, \dots$ , is a sequence of partitions such that each partition is a refinement of the previous one, then*

$$(2.9) \quad r_{k_m}(P, d) \geq r_{k_\ell}(P, d) \quad \text{for } m < \ell.$$

We next give a parametric example in which the mixed risk  $r_k$  equals the Bayes risk  $r$ .

**EXAMPLE 2.1.** Let  $\Gamma = \Omega$  be the class of normal distribution functions  $F_\theta$  with mean  $\theta$  and variance unity. Suppose that  $P$  is the measure for which  $\theta$  has a normal distribution with mean  $\xi$  and variance unity. Let  $\mathcal{S}$  be the class of sets of the form  $\{F_\theta: \theta \in B\}$ , where  $B$  is a Borel subset of the reals. All the quantities of this section are now defined. Moreover, since  $F(t_1)$  determines  $\theta$  and  $\theta$  determines  $F(t_1)$ , then  $r_k(P, d) = r(P, d)$ ,  $k \geq 1$ .

Next we consider a "discrete" example in which the mixed risk eventually equals the Bayes risk.

**EXAMPLE 2.2.** Let  $\Omega = \Gamma$  be a countable class  $\{F_1, F_2, \dots\}$  and let  $\mathcal{S}$  be the collection of all subsets of  $\Omega$ . If  $\Omega$  is discrete (that is,  $\Omega$  has no limit points for the sup norm), and if  $\{t_{m,j}\}$  of Lemma 2.1 (ii) becomes dense in  $S(\Omega)$  as  $m \rightarrow \infty$ , then there exists  $m_1 \geq 1$  such that  $r_{k_m}(P, d) = r(P, d)$  for all  $m \geq m_1$ . To see this, note that, by our assumptions, there exists  $m_1$  such that

$$(2.10) \quad (F_i(t_{m,1}), \dots, F_i(t_{m,k_m})) \neq (F_j(t_{m,1}), \dots, F_j(t_{m,k_m}))$$

for  $m \geq m_1$  and all  $i \neq j$ .

We now define a decision rule  $d_k$  that, in some cases, will be the mixed Bayes-minimax procedure. Let  $t_1 = \inf \{t: t \in S(\Omega)\}$  and  $t_k = \sup \{t: t \in S(\Omega)\}$ . It follows that  $F(t_k) = 1$ . We will assume that  $F(t_1) = 0$  a.s. ( $P$ ),  $-\infty < t_1 < t_k < \infty$ , and that  $k \geq 3$ . Let  $q = \{q_1, \dots, q_k\}$  with  $0 = q_1 \leq \dots \leq q_k = 1$ . Let  $F_{q,k}$  be the polygonal distribution function that equals  $q_i$  at  $t_i$ ,  $i = 1, \dots, k$ , and is linear over each interval  $[t_i, t_{i+1}]$ ,  $i = 1, \dots, k - 1$ . Let  $F_k$  denote the random distribution function obtained by letting  $q$  in  $F_{q,k}$  have distribution  $\lambda = \lambda(\cdot: P, k)$ . We assume that  $F_k$  is a measurable function on some measure space to  $(\Gamma, \mathcal{S})$ . Let  $P_k$  denote the distribution of  $F_k$ . Finally,  $d_k$  will denote the Bayes solution for  $F_{q,k}$  when  $q$  has prior  $\lambda$ , that is,  $d_k$  minimizes

$$(2.11) \quad r(P_k, d) = \int_{R_k} R(F_{q,k}, d)\lambda(dq).$$

**THEOREM 2.1.** *If  $F_{q,k} \in \Omega$  for almost all  $q$  in the carrier  $C_\lambda$  of  $\lambda$ , if  $d_k$  minimizing (2.11) exists, and if*

$$(2.12) \quad L(F, d_k) = L(G, d_k)$$

*for all  $F, G \in \Omega(q, k)$  and almost all  $q$  in  $C_\lambda$  then  $d_k$  is the mixed Bayes-minimax procedure.*

**PROOF.** Let  $N_i$  denote the number of  $X$  in  $(t_i, t_{i+1}]$ ,  $i = 1, \dots, k - 1$ . Then  $\mathbf{N} = (N_1, \dots, N_{k-1})$  is sufficient for  $F_{q,k}$ , and  $d_k$  depends on the  $X$  only through  $\mathbf{N}$ . Thus the distribution of  $d_k$  as a function of the  $X$  is the same for all  $F$  in  $\Omega(q, k)$ . This and (2.12) imply that  $R(F, d_k) = R(G, d_k)$  for all  $F, G \in \Omega(q, k)$ . Thus  $r_k(P, d_k) = r(P_k, d_k)$ , and since  $d_k$  is optimal for  $P_k$ , then  $r(P_k, d_k) \leq r(P_k, d)$  for all other rules  $d$ . Finally, since  $F_{q,k} \in \Omega$  for almost all  $q$ ,  $r(P_k, d) \leq r_k(P, d)$  and the results follow.

REMARK 2.1. From the above proof, it is clear that (2.12) can be replaced by the condition

$$(2.13) \quad r_k(P, d_k) = r(P_k, d_k)$$

and that (2.13) is weaker than (2.12).

EXAMPLE 2.3. If  $0 \in \{t_1, \dots, t_k\}$  and  $L(F, d) = [d - F(0)]^2$ , then (2.12) is satisfied. Generalizing this, we have that (2.12) is satisfied for any loss function depending on  $F$  only through  $F(t_1), \dots, F(t_k)$ ; that is, the loss is defined through those points where we have information about  $F$ . Such a loss function corresponding to squared error when estimating the mean of  $F$  would be  $[d - \mu(F, k)]^2$ , where

$$(2.14) \quad \mu(F, k) = \frac{1}{2} \sum_{i=1}^{k-1} (t_{i+1} + t_i) [F(t_{i+1}) - F(t_i)]$$

In the next section, we consider a testing problem for which (2.12) is satisfied.

### 3. A testing problem

Suppose  $\Omega = \Omega_0 \cup \Omega_1$  with  $\Omega_0 \cap \Omega_1$  empty,  $\Omega_0, \Omega_1 \in \mathcal{S}$ , and we want to test  $H_0: F \in \Omega_0$  against  $H_1: F \in \Omega_1$ .  $\varphi = \varphi(X)$  will denote a test function, and  $L(F, \varphi)$  will be the usual loss for the testing problem; that is,  $L(F, \varphi) = L_i$ , a positive constant, if  $H_i$  is falsely rejected,  $i = 0, 1$ , and  $L(F, \varphi) = 0$  otherwise. The Bayes risk is

$$(3.1) \quad r(P, \varphi) = L_0 \int_{\Omega_0} E_F(\varphi) P(dF) + L_1 \int_{\Omega_1} [1 - E_F(\varphi)] P(dF)$$

and the average maximum risk is

$$(3.2) \quad r_k(P, \varphi) = L_0 \int_{Q_0} \left[ \sup_{F \in \Omega(q, k)} E_F(\varphi) \right] \lambda(dq; P, k) \\ + L_1 \int_{Q_1} \left[ \sup_{F \in \Omega(q, k)} [1 - E_F(\varphi)] \right] \lambda(dq; P, k),$$

where  $Q_i = \{(F(t_1), \dots, F(t_k)): F \in \Omega_i\}$ ,  $i = 1, 2$ ; that is,  $Q_0$  and  $Q_1$  are the sets in  $R^k$  corresponding to  $\Omega_0$  and  $\Omega_1$ . We assume that  $Q_0 \cap Q_1$  has  $\lambda$  probability zero. This assumption is needed to obtain (3.2) above. For this loss function it is clear that (2.12) of Theorem 2.1 is satisfied; if in addition  $F_{q, k} \in \Omega$ , then the result can be applied. If  $F_{q, k} \notin \Omega$ , then  $L(F_{q, k}, \varphi)$  is not defined. However, there is a natural way of defining  $L(F_{q, k}, \varphi)$  and making  $F_{q, k}$  a member of  $\Omega(q, k)$ : for a given  $q$  not in  $Q_0 \cap Q_1$ ,  $L(F, \varphi)$  has the same value for each  $F \in \Omega(q, k)$ ; define  $L(F_{q, k}, \varphi)$  to be this value.

In what follows, it will be assumed that either  $F_{q, k} \in \Omega$  or  $L(F_{q, k}, \varphi)$  has been defined as above. Let  $f(x|q)$  denote the density of  $F_{q, k}$ . Then from (3.2) and Theorem 2.1 we can conclude that the mixed Bayes-minimax procedure  $\varphi_k$  rejects  $H_0$  when

$$(3.3) \quad L_1 \int_{Q_1} \prod_{i=1}^{k-1} (q_{i+1} - q_i)^{N_i} \lambda(dq) \geq L_0 \int_{Q_0} \prod_{i=1}^{k-1} (q_{i+1} - q_i)^{N_i} \lambda(dq), \quad i = 1, 2.$$

Let  $p_i = q_{i+1} - q_i, i = 1, \dots, k - 1$  and

$$(3.4) \quad A_i = \{(F(t_2) - F(t_1)), \dots, (F(t_k) - F(t_{k-1})) : F \in \Omega_i\}.$$

Then (3.3) becomes

$$(3.5) \quad L_1 \int_{A_1} \prod_{i=1}^{k-1} p_i^{N_i} \pi(dp) \geq L_0 \int_{A_0} \prod_{i=1}^{k-1} p_i^{N_i} \pi(dp),$$

where  $\pi$  is the distribution of  $(F(t_2) - F(t_1), \dots, F(t_k) - F(t_{k-1}))$  when  $F$  has distribution  $P$ . Note that (3.5) is the solution to a Bayesian multinomial testing problem, in which  $(N_1, \dots, N_{k-1} | P)$  is a multinomial variable with parameters  $n$  and  $p$ , and we are testing  $H_0 : p \in A_0$  versus  $H_1 : p \in A_1$ . The solution is the test  $\varphi_k$  that rejects  $H_0$  when the ratio of the posterior probability of  $H_1$  to the posterior probability of  $H_0$  exceeds the ratio of the losses  $L_0/L_1$ ; that is, if  $p$  has the posterior density  $g(p|N)$ ,  $N = (N_1, \dots, N_{k-1})$ , then the test  $\varphi_k$  rejects  $H_0$  when

$$(3.6) \quad \frac{\int_{A_1} g(p|N) dp}{\int_{A_0} g(p|N) dp} \geq \frac{L_0}{L_1}.$$

**EXAMPLE 3.1.** Consider the goodness of fit problem where  $\Omega_0$  contains only the uniform distribution  $F_0$  on  $(0, 1)$  and

$$(3.7) \quad \Omega_1 = \{F : C_F \subset [0, 1], F \leq F_0\} - \Omega_0.$$

If  $P$  assigns probability  $\frac{1}{2}$  to  $\Omega_0$ , and if  $t_1 = 0, t_2 = \frac{1}{2}, t_3 = 1$ , then from (3.5)

$$(3.8) \quad \varphi_k = \begin{cases} 1 & \text{if } 2^{-n} \int_0^{1/2-0} p_1^{N_1} (1 - p_1)^{n-N_1} \pi_1(dp_1) \geq \frac{L_1}{L_0} \\ 0 & \text{otherwise} \end{cases}$$

where  $\pi_1(p_1) = P(F(\frac{1}{2}) \leq p_1)$ , so that  $\pi_1(\frac{1}{2} - 0) = \frac{1}{2}$  and  $\pi_1(\frac{1}{2}) = 1$ . Thus  $\varphi_k$  is based on a decreasing function of  $N_1 =$  number of  $X$  less than or equal to  $\frac{1}{2}$ . When  $2\pi_1$  is the beta distribution, then  $\varphi_k$  can be obtained from the tables of the incomplete beta function.

**REMARK 3.1.** The loss function of this section is not the only one for which optimal mixed tests can be obtained. For instance, if  $\Omega_0$  and  $\Omega_1$  are as in Example 3.1, then the loss for deciding  $H_0$  when  $F \in \Omega_1$  could be defined to be  $L_1(\frac{1}{2} - F(\frac{1}{2}))$ . The conditions of Theorem 2.1 would then still be satisfied and the optimal test can be obtained from the corresponding multinomial problem.

#### 4. Convergence of the mixed solutions to the Bayes solution

In this section, we will consider classes of distribution functions on  $[0, 1]$ , that is,  $\Omega$  is such that its support  $S(\Omega)$  is contained in  $[0, 1]$ . Then  $\Omega$  is a subset of the class  $D[0, 1]$  of right continuous functions on  $[0, 1]$  with limits from the left. On  $D[0, 1]$  we use the Skorohod topology with the modified Skorohod

metric (see for example [2], p. 113) so that  $D[0, 1]$  is a complete separable metric space. In the notation of Section 2,  $\Gamma = D[0, 1]$  and  $\mathcal{S}$  is the  $\sigma$ -field generated by the open sets,  $F$  is a random function (a measurable function on some measure space to  $(\Gamma, \mathcal{S})$ ) with distribution  $P$ . We assume that the probability  $P$  on  $(\Gamma, \mathcal{S})$  satisfies  $P(F \in \Omega) = 1$ . If  $\{F_k\}$  is a sequence of random functions with distributions  $\{P_k\}$ , then  $F_k$  is said to *converge in distribution* to  $F$  if

$$(4.1) \quad \int_{\Gamma} h(F)P_k(dF) \rightarrow \int_{\Gamma} h(F)P(dF)$$

for each continuous, bounded, real valued function  $h$  on  $\Gamma$ . If (4.1) holds, then  $P_k$  is also said to *converge weakly* to  $P$ .

If  $F_k$  is the polygonal random distribution function of Section 2, then  $F_k$  does not necessarily converge in distribution to  $F$ . However, we will show that the Bayes risk of the Bayes solution for the prior  $P_k$  converges to the Bayes risk of the Bayes solution for the prior  $P$ . To do this, we will make use of the random distribution function  $G_k$  that is constant over each interval  $[t_i, t_{i+1})$ , and for which the joint distribution of  $G_k(t_1), \dots, G_k(t_k)$  equals that of  $F(t_1), \dots, F(t_k)$ . The symbol  $Q_k$  will denote the distribution of  $G_k$ . We will need double subscripts on the  $t$  and again write  $t_{m,1}, \dots, t_{m,k_m}$  instead of  $t_1, \dots, t_k$ . We assume that  $\{\Pi_m: t_{m,1} < \dots < t_{m,k_m}\}$ ,  $m = 1, \dots$ ,  $k_m$  is a sequence of partitions of  $[0, 1]$  such that  $\Pi_{m+1}$  is a refinement of  $\Pi_m$  and  $\max_i |t_{m,(i+1)} - t_{m,i}| \rightarrow 0$  as  $m \rightarrow \infty$ .

It is now easy to show using Prohorov's theorem that:

LEMMA 4.1. *If  $\Omega \subset D[0, 1]$ , then  $G_k$  converges in distribution to  $F$ .*

PROOF. For each  $F_0 \in \Omega$  and  $\delta \in (0, 1)$  define

$$(4.2) \quad v(F_0, \delta) = \sup \min [F_0(t) - F_0(s), F_0(u) - F_0(t)],$$

where the sup is over  $s \leq t \leq u$ ,  $u - s = \delta$ . Similarly, define

$$(4.3) \quad w_0(F_0, \delta) = \sup [F_0(u) - F_0(s)],$$

where the sup is over  $1 - \delta \leq s < u < 1$ . By Prohorov's theorem applied to  $D[0, 1]$  it is enough to show that (i) the finite dimensional distributions of  $Q_k$  converge weakly to the finite dimensional distributions of  $P$  at points  $s_1, \dots, s_r$  in the set  $\{t: P[F(t) \neq F(t^-)] = 0\} \cup \{0, 1\}$ , and that (ii) for each  $\varepsilon, \eta > 0$ , there exists  $\delta \in (0, 1)$  and an integer  $k_0$  such that

$$(4.4) \quad Q_k(v(G_k, \delta) < \varepsilon) > 1 - \eta \quad \text{for } k \geq k_0,$$

and

$$(4.5) \quad Q_k(w_0(G_k, \delta) < \varepsilon) > 1 - \eta \quad \text{for } k \geq k_0.$$

(See for example [2], pp. 125-126.)

The convergence of the indicated finite dimensional distributions follows from the definition of  $G_k$ . The inequalities (4.4) and (4.5) are easy to establish using the

tightness of  $P$ . We omit the first subscript on the  $t$  to simplify the notation. Let  $k_0$  be such that  $t_{(i+1)} - t_i < \delta$ ,  $i = 1, \dots, k$ , for all  $k \geq k_0$ . Then since  $G_k$  is constant between points  $t_i$ ,

$$(4.6) \quad v(G_k, \delta) \leq \max \min [F(t_i) - F(t_j), F(t_j) - F(t_\ell)],$$

where the max is over  $t_\ell \leq t_j \leq t_i$ ,  $t_i - t_\ell \leq 2\delta$ . The right side is bounded above by  $v(F, 2\delta)$ . Thus

$$(4.7) \quad v(G_k, \delta) \leq v(F, 2\delta) \quad \text{for } k \geq k_0.$$

Similarly,  $w_0(G_k, \delta) \leq w_0(F, 2\delta)$ . Since  $P$  is tight, we can choose  $\delta$  so that  $P(v(F, 2\delta) < \varepsilon) > 1 - \eta$  and  $P(w_0(F, 2\delta) < \varepsilon) > 1 - \eta$ . This implies the result.

Let  $G_{q,k}$  denote the distribution function that is constant on  $[t_i, t_{i+1})$  and whose value at  $t_i$  is  $q_i$ , where  $0 = q_1 \leq \dots \leq q_k = 1$ . Thus  $Q_k$  is the distribution of  $G_{q,k}$  when  $q$  has distribution  $\lambda$ , and  $G_k$  is the random distribution function obtained by letting  $q$  in  $G_{q,k}$  have distribution  $\lambda$ .

Recall that  $d_k$  is the decision rule that minimizes the Bayes risk  $r(P_k, d)$ .

**THEOREM 4.1.** *If  $G_{q,k}$  is an element of  $\Omega$  for almost all  $q$  in  $C_\lambda$ , if the conditions of Theorem 2.1 are satisfied, if a Bayes solution  $d$  exists, and if  $d$  has a continuous (in  $F$ ) bounded risk  $R(F, d)$ , then for  $\Omega \subset D[0, 1]$ ,  $d_k$  converges to  $d$  in the sense that the Bayes risk  $r(P, d_k)$  of  $d_k$  converges to the Bayes risk  $r(P, d)$  of  $d$ .*

**PROOF.** Consider the Bayes solution  $\hat{d}_k$  for the prior  $Q_k$ . Since  $G_{q,k}$  is constant between points  $t_i$ ,  $\hat{d}_k$  will depend on the  $X$  only through  $\mathbf{S} = (S_1, \dots, S_{k-1})$ , where  $S_i =$  number of  $X$  equal to  $t_{i+1}$ ,  $i = 1, \dots, k-1$ . (Recall that  $F(t_1) = F(0) = 0$  a.s. ( $P$ ) by assumption.) Note that  $\mathbf{S}$  and  $\mathbf{N}$  have the same distribution under  $G_{q,k}$  and that this is the distribution  $\mathbf{N}$  has under  $F_{q,k}$ . This implies that  $d_k$  is also a Bayes solution for the prior  $Q_k$ . We now have

$$(4.8) \quad r_k(P, d_k) = r(P_k, d_k) = r(Q_k, d_k) = r(Q_k, \hat{d}_k) \leq r(Q_k, d).$$

By Lemma 4.1,

$$(4.9) \quad \lim_{k \rightarrow \infty} r(Q_k, d) = r(P, d).$$

Equations (4.8) and (4.9) yield

$$(4.10) \quad \limsup_{k \rightarrow \infty} r_k(P, d_k) \leq r(P, d).$$

On the other hand, by Lemma 2.1,

$$(4.11) \quad r_k(P, d_k) \geq r(P, d_k).$$

Since  $d$  is the Bayes solution for  $P$ ,

$$(4.12) \quad r(P, d_k) \geq r(P, d).$$

Putting these inequalities together, we get



$$(4.13) \quad \lim_{k \rightarrow \infty} r_k(P, d_k) = r(P, d), \quad \lim_{k \rightarrow \infty} r(P, d_k) = r(P, d).$$

If  $\Omega$  is a class of continuous distribution functions on  $[0, 1]$ , then it is a subset of the class  $C[0, 1]$  of continuous functions on  $[0, 1]$ . On  $C[0, 1]$  we use the sup norm and the  $\sigma$ -algebra generated by the open sets.

LEMMA 4.2. *If  $\Omega \subset C[0, 1]$ , then  $F_k$  converges in distribution to  $F$ .*

PROOF. The convergence of finite dimensional distributions follows from the definition of  $F_k$ . For each of  $F_0 \in \Omega$  and  $\delta \in (0, 1)$ , let

$$(4.14) \quad w(F_0, \delta) = \sup_{t-s=\delta} [F_0(t) - F_0(s)].$$

Let  $k_0$  be such that  $t_{i+1} - t_i < \delta$ ,  $i = 1, \dots, k$ , for all  $k \geq k_0$ . Now the tightness of  $\{P_k\}$  follows from the inequality  $w(F_k, \delta) \leq w(F, 3\delta)$  and the tightness of  $P$ . Thus the result follows from Prohorov's theorem applied to  $C[0, 1]$ .

We can now prove that  $d_k$  converges to  $d$  under fewer conditions than when  $\Omega \subset D[0, 1]$ .

THEOREM 4.2. *If  $F_{q,k} \in \Omega$  for almost all  $q$  in  $C_\lambda$ , if  $d_k$  is a mixed Bayes-minimax solution, if a Bayes solution  $d$  exists, and if  $d$  has a continuous bounded risk  $R(F, d)$ , then for  $\Omega \subset C[0, 1]$ ,  $\lim_{k \rightarrow \infty} r(P, d_k) = r(P, d)$ .*

PROOF. Since  $d_k$  is Bayes for  $P_k$ , then

$$(4.15) \quad r_k(P, d_k) = r(P_k, d_k) \leq r(P_k, d).$$

By Lemma 4.2,

$$(4.16) \quad \lim_{k \rightarrow \infty} r(P_k, d) = r(P, d).$$

The rest of the proof now follows on the lines of the proof of Theorem 4.1.

### 5. Examples of random distribution functions

In order to obtain the mixed Bayes-minimax solutions, we have to specify a distribution for the random distribution function  $F_k$  which is linear between the points  $(t_1, q_1), \dots, (t_k, q_k)$ ,  $0 = q_1 \leq \dots \leq q_k = 1$ . Equivalently, we have to define a probability  $\lambda$  on

$$(5.1) \quad A_k = \{q \in R^k : 0 = q_1 \leq \dots \leq q_k = 1\},$$

or a probability  $\pi$  on

$$(5.2) \quad B_k = \left\{ p \in R^{k-1} : 0 \leq p_i \leq 1, \sum_{i=1}^{k-1} p_i = 1 \right\}.$$

Here  $q_i$  is thought of as  $F(t_i)$  and  $p_i$  as  $F(t_{i+1}) - F(t_i)$ . We say that  $p = (p_1, \dots, p_{k-1})$  has distribution  $\pi$ .

One way to obtain a class of distributions  $\pi$  of  $p$  (Freedman [7], Fabius [5], Connor and Mosimann [16]), is to let  $p$  have the same distribution as the

vector whose  $i$ th coordinate is

$$(5.3) \quad Z_i \prod_{j=1}^{i-1} (1 - Z_j), \quad i = 1, \dots, k-1,$$

where  $Z_1, \dots, Z_{k-1}$  are independent random variables satisfying

$$(5.4) \quad 0 < Z_i \leq 1, \quad Z_{k-1} = 1.$$

For each choice of distributions  $H_1, \dots, H_{k-2}$  of the  $Z$ , we obtain a probability  $\pi$  on  $B_k$ . For this class of probabilities, it is easy ([7], p. 1401 and [5], p. 848) to compute posterior probabilities  $\pi(p|\mathbf{N})$  of  $p$  given  $\mathbf{N} = (N_1, \dots, N_{k-1})$ . Such probabilities  $\pi$  are called *tailfree* by Freedman [7] and Fabius [5] and *neutral* by Connor and Mosimann [16]. If we let each  $Z_i$  have a beta distribution  $B(r_i, s_i)$  with parameters  $r_i$  and  $s_i$ , then  $\pi$  is called the *generalized Dirichlet* distribution [16]. If in addition,

$$(5.5) \quad s_i = \sum_{j=i+1}^{k-1} r_j, \quad i = 1, \dots, k-2, \quad s_{k-1} = 0,$$

then  $\pi$  is called the *Dirichlet* distribution with parameters  $r_1, \dots, r_{k-1}$ .

Extensions of the definition of  $\pi$  on  $B_k$  to

$$(5.6) \quad B_\infty = \left\{ p \in R^\infty : 0 \leq p_i \leq 1, \sum_{i=1}^{\infty} p_i = 1 \right\}$$

are obtained by replacing (5.4) by

$$(5.7) \quad 0 \leq Z_i \leq 1, \quad \lim_{r \rightarrow \infty} \prod_{i=1}^r (1 - Z_i) = 0 \quad \text{a.s.}$$

If the  $Z$  have beta distributions, then the resulting  $\pi$  on  $B_\infty$  is called the *infinite dimensional generalized Dirichlet* distribution (Freedman [7]).

Note that, if  $p$  has a Dirichlet distribution, then  $\text{Cov}(p_i, p_j) < 0$ . However, for the generalized Dirichlet distribution, it is possible to have  $\text{Cov}(p_i, p_j) > 0$  (see [16], p. 198).

**EXAMPLE 5.1.** Consider the problem of estimating the mean  $\mu(F) = E_F(X_1)$  when  $p = [F(t_2) - F(t_1), \dots, F(t_k) - F(t_{k-1})]$  has a generalized Dirichlet distribution  $\pi$  with parameters  $(r_1, s_1), \dots, (r_{k-1}, s_{k-1})$ . If we consider the squared error loss function, then the Bayes estimate of

$$(5.8) \quad \mu(F_k) = \frac{1}{2} \sum_{i=1}^{k-1} [F(t_{i+1}) - F(t_i)](t_{i+1} + t_i)$$

is

$$(5.9) \quad \hat{\mu}_k = E_\pi(\mu(F_k)|X) = \frac{1}{2} \left( \sum_{i=1}^{k-1} [t_{i+1} + t_i] E_\pi(p_i|\mathbf{N}) \right),$$

where (see [5])

$$(5.10) \quad E_{\pi}(p_i|\mathbf{N}) = \frac{r_i + N_i}{r_1 + s_1 + n} \prod_{\ell=1}^{i-1} \frac{s_{\ell} + n - \sum_{j=1}^{\ell} N_j}{r_{\ell+1} + s_{\ell+1} + n - \sum_{j=1}^{\ell} N_j}.$$

Following Ferguson [6], let  $\alpha$  be a finite, finitely additive measure on  $R$  and let

$$(5.11) \quad r_i = \alpha(t_i, t_{i+1}], \quad s_i = \sum_{j=i+1}^{k-1} r_j, \quad s_{k-1} = 0.$$

Assume that  $\alpha$  assigns measure zero to the region outside  $(t_1, t_k]$ . Then

$$(5.12) \quad E_{\pi}(\mu(F_k)|X) = (\alpha(R) + n)^{-1} \sum_{i=1}^{k-1} \frac{1}{2}(t_{i+1} + t_i)(\alpha(t_i, t_{i+1}] + N_i) \\ = \alpha_n E(\mu(F_k)) + (1 - \alpha_n)\bar{X}',$$

where  $\alpha_n = \alpha(R)[\alpha(R) + n]^{-1}$  and  $\bar{X}'$  is the average of the random variables  $X'$  obtained by replacing each  $X$  in the interval  $(t_i, t_{i+1}]$  by the midpoint  $\frac{1}{2}(t_{i+1} + t_i)$ . Note that if the  $t$  become dense in  $(t_1, t_k]$  as in Section 4, then from (5.12)

$$(5.13) \quad \lim_{k \rightarrow \infty} E_{\pi}(\mu(F_k)|X) = \alpha_n \mu_0 + (1 - \alpha_n)\bar{X},$$

where  $\mu_0 = \alpha(R)^{-1} \int x\alpha(dx)$ . This is the estimate obtained by Ferguson [6]. Note that in addition to being the Bayes estimate of  $\mu(F_k)$ , the estimate (5.9) is the mixed Bayes-minimax estimate of  $\mu(F)$  for the loss function  $[d - \mu(F, k)]^2$  of Example 2.3.

EXAMPLE 5.2. Suppose again that  $p$  has a generalized Dirichlet distribution  $\pi$ . If  $s \in \{t_1, \dots, t_k\}$ , the mixed Bayes-minimax estimate of  $F(s)$  using squared error loss  $[d - F(s)]^2$  is

$$(5.14) \quad \hat{F}(s) = \sum_{j=1}^{i-1} E_{\pi}(p_j|\mathbf{N}),$$

where  $s = t_i$  and  $E_{\pi}(p_j|\mathbf{N})$  is given by (5.10). If in addition (5.11) is satisfied, then (5.14) becomes

$$(5.15) \quad \alpha_n \alpha_0(s) + (1 - \alpha_n)F_n(s),$$

where  $\alpha_0(s) = \alpha(-\infty, s]/\alpha(R)$  and  $F_n(s)$  is the empirical distribution function of the sample. This is exactly the estimate obtained by Ferguson [6].

Next we consider the problem of defining a probability  $P$  on a set  $\Omega$  of distribution functions  $F$  in such a way that it is possible to compute the posterior of  $F$  given  $X$  under the prior  $P$ . Ferguson [6] shows that for each finite, finitely additive measure  $\alpha$  on  $R$ , it is possible to define a *Dirichlet process*  $F$  in such a way that  $P_F(A_1), \dots, P_F(A_m)$  has a Dirichlet distribution with parameters  $\alpha_1, \dots, \alpha_m$ , where  $d_i = \alpha(A_i)$ , and  $A_1, \dots, A_m$  is a measurable partition of  $R$ .

He shows that the posterior of  $F$  given  $X$  is again a Dirichlet process with  $\alpha$  replaced by  $\alpha + \sum_{i=1}^n \delta(x_i)$ , where  $\delta(x)$  is the measure giving mass one to the point  $x$ . This makes it possible to compute the Bayes procedure for this prior. The estimate (5.15) in Example 5.2 above is both the Bayes and the mixed Bayes-minimax estimate for this prior.

Fabius ([5], p. 853) gives a general construction of probabilities on the set  $\Omega$  of all distribution functions on  $[0, 1]$  that include the Dirichlet process on  $[0, 1]$ , the processes of Kraft [12], Kraft and van Eeden [13], and those special cases of the processes of Dubins and Freedman that are contained in [13]. Kraft and van Eeden [13] compute the Bayes estimate for one of these processes for a problem in bioassay. Ferguson shows that if  $F$  is the Dirichlet process and  $\Omega_1$  is the class of discontinuous distribution functions, then  $P(F \in \Omega_1) = 1$ . Kraft [12] shows that it is possible to use the construction of Fabius to obtain a process  $F$  such that  $P(F \in \Omega^*) = 1$ , where  $\Omega^*$  is the class of absolutely continuous distribution functions.

Using definitions (2.2) and (2.3) of Fabius [5], it is possible to check that the Dirichlet process is tailfree for all trees of partitions. Thus one can use expression (2.4) of [5] for the posterior distribution of a tailfree process to conclude that the posterior of a Dirichlet process is again Dirichlet.



I am indebted to Lucien LeCam, Thomas Ferguson, Jaap Fabius, and others for helpful discussions, and to Michael Stuart for a careful reading of the manuscript that led to many improvements.

#### REFERENCES

- [1] C. ANTONIAK, "Mixtures of Dirichlet processes with application to some nonparametric problems," Ph.D. thesis, Mathematics Department, University of California, Los Angeles, 1969.
- [2] P. BILLINGSLEY, *Convergence of Probability Measures*, New York, Wiley, 1968.
- [3] K. A. DOKSUM, "Minimax results for IFRA scale alternatives," *Ann. Math. Statist.*, Vol. 40 (1969), pp. 1778-1783.
- [4] L. E. DUBINS and D. A. FREEDMAN, "Random distribution functions," *Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability*, Berkeley and Los Angeles, University of California Press, 1966, Vol. 2, pp. 183-214.
- [5] J. FABIUS, "Asymptotic behavior of Bayes' estimates," *Ann. Math. Statist.*, Vol. 35 (1964), pp. 846-856.
- [6] T. S. FERGUSON, "A Bayesian analysis of some nonparametric problems," 1970, submitted for publication.
- [7] D. A. FREEDMAN, "On the asymptotic behavior of Bayes' estimates in the discrete case," *Ann. Math. Statist.*, Vol. 34 (1963), pp. 1386-1403.
- [8] W. HOEFFDING, "'Optimum' nonparametric tests," *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, Berkeley and Los Angeles, University of California Press, 1951, pp. 83-92.

- [9] P. J. HUBER, "A robust version of the probability ratio test," *Ann. Math. Statist.*, Vol. 36 (1965), pp. 1753-1758.
- [10] ———, "Robust confidence limits," *Z. Wahrscheinlichkeitstheorie und Verw. Gebiete*, Vol. 10 (1968), pp. 269-278.
- [11] ———, *Théorie de l'Inférence Statistique Robuste*, Montreal, Les Presses de l'Université de Montréal, 1969.
- [12] C. H. KRAFT, "A class of distribution function processes which have derivatives," *J. Appl. Probability*, Vol. 1 (1964), pp. 385-388.
- [13] C. H. KRAFT and C. VAN EEDEN, "Bayesian bio-assay," *Ann. Math. Statist.*, Vol. 35 (1964), pp. 886-890.
- [14] E. RUIST, "Comparison of tests for non-parametric hypotheses," *Arkiv Math.*, Vol. 3 (1954), pp. 133-163.
- [15] W. SCHAAFSMA, *Hypothesis Testing Problems with the Alternative Restricted by a Number of Inequalities*, Groningen, Noordhoff, 1966.
- [16] R. J. CONNOR and J. E. MOSIMANN, "Concepts of independence for proportions with a generalization of the Dirichlet distribution," *J. Amer. Statist. Assoc.*, Vol. 64 (1969), pp. 194-206.