



A Bayesian Analysis of Some Nonparametric Problems

Thomas S. Ferguson

The Annals of Statistics, Vol. 1, No. 2 (Mar., 1973), 209-230.

Stable URL:

<http://links.jstor.org/sici?sici=0090-5364%28197303%291%3A2%3C209%3AABAOSN%3E2.0.CO%3B2-U>

Your use of the JSTOR archive indicates your acceptance of JSTOR's Terms and Conditions of Use, available at <http://www.jstor.org/about/terms.html>. JSTOR's Terms and Conditions of Use provides, in part, that unless you have obtained prior permission, you may not download an entire issue of a journal or multiple copies of articles, and you may use content in the JSTOR archive only for your personal, non-commercial use.

Each copy of any part of a JSTOR transmission must contain the same copyright notice that appears on the screen or printed page of such transmission.

The Annals of Statistics is published by Institute of Mathematical Statistics. Please contact the publisher for further permissions regarding the use of this work. Publisher contact information may be obtained at <http://www.jstor.org/journals/ims.html>.

The Annals of Statistics

©1973 Institute of Mathematical Statistics

JSTOR and the JSTOR logo are trademarks of JSTOR, and are Registered in the U.S. Patent and Trademark Office. For more information on JSTOR contact jstor-info@umich.edu.

©2003 JSTOR

A BAYESIAN ANALYSIS OF SOME NONPARAMETRIC PROBLEMS¹

BY THOMAS S. FERGUSON

University of California, Los Angeles

1. Introduction and summary. The Bayesian approach to statistical problems, though fruitful in many ways, has been rather unsuccessful in treating nonparametric problems. This is due primarily to the difficulty in finding workable prior distributions on the parameter space, which in nonparametric problems is taken to be a set of probability distributions on a given sample space. There are two desirable properties of a prior distribution for nonparametric problems.

(I) The support of the prior distribution should be large—with respect to some suitable topology on the space of probability distributions on the sample space.

(II) Posterior distributions given a sample of observations from the true probability distribution should be manageable analytically.

These properties are antagonistic in the sense that one may be obtained at the expense of the other. This paper presents a class of prior distributions, called Dirichlet process priors, broad in the sense of (I), for which (II) is realized, and for which treatment of many nonparametric statistical problems may be carried out, yielding results that are comparable to the classical theory.

In Section 2, we review the properties of the Dirichlet distribution needed for the description of the Dirichlet process given in Section 3. Briefly, this process may be described as follows. Let \mathcal{X} be a space and \mathcal{A} a σ -field of subsets, and let α be a finite non-null measure on $(\mathcal{X}, \mathcal{A})$. Then a stochastic process P indexed by elements A of \mathcal{A} , is said to be a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter α if for any measurable partition (A_1, \dots, A_k) of \mathcal{X} , the random vector $(P(A_1), \dots, P(A_k))$ has a Dirichlet distribution with parameter $(\alpha(A_1), \dots, \alpha(A_k))$. P may be considered a random probability measure on $(\mathcal{X}, \mathcal{A})$. The main theorem states that if P is a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter α , and if X_1, \dots, X_n is a sample from P , then the posterior distribution of P given X_1, \dots, X_n is also a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter $\alpha + \sum_{i=1}^n \delta_{X_i}$, where δ_x denotes the measure giving mass one to the point x .

In Section 4, an alternative definition of the Dirichlet process is given. This definition exhibits a version of the Dirichlet process that gives probability one to the set of discrete probability measures on $(\mathcal{X}, \mathcal{A})$. This is in contrast to Dubins and Freedman [2], whose methods for choosing a distribution function on the interval $[0, 1]$ lead with probability one to singular continuous distributions. Methods of choosing a distribution function on $[0, 1]$ that with probability one is absolutely continuous have been described by Kraft [7]. The

Received September 17, 1969; revised April 3, 1972.

¹ The preparation of this paper was supported in part by NSF Grant No. GP-8049.

general method of choosing a distribution function on $[0, 1]$, described in Section 2 of Kraft and van Eeden [10], can of course be used to define the Dirichlet process on $[0, 1]$.

Special mention must be made of the papers of Freedman and Fabius. Freedman [5] defines a notion of tailfree for a distribution on the set of all probability measures on a countable space \mathcal{X} . For a tailfree prior, posterior distribution given a sample from the true probability measure may be fairly easily computed. Fabius [3] extends the notion of tailfree to the case where \mathcal{X} is the unit interval $[0, 1]$, but it is clear his extension may be made to cover quite general \mathcal{X} . With such an extension, the Dirichlet process would be a special case of a tailfree distribution for which the posterior distribution has a particularly simple form.

There are disadvantages to the fact that P chosen by a Dirichlet process is discrete with probability one. These appear mainly because in sampling from a P chosen by a Dirichlet process, we expect eventually to see one observation exactly equal to another. For example, consider the goodness-of-fit problem of testing the hypothesis H_0 that a distribution on the interval $[0, 1]$ is uniform. If on the alternative hypothesis we place a Dirichlet process prior with parameter α itself a uniform measure on $[0, 1]$, and if we are given a sample of size $n \geq 2$, the only nontrivial nonrandomized Bayes rule is to reject H_0 if and only if two or more of the observations are exactly equal. This is really a test of the hypothesis that a distribution is continuous against the hypothesis that it is discrete. Thus, there is still a need for a prior that chooses a continuous distribution with probability one and yet satisfies properties (I) and (II).

Some applications in which the possible doubling up of the values of the observations plays no essential role are presented in Section 5. These include the estimation of a distribution function, of a mean, of quantiles, of a variance and of a covariance. A two-sample problem is considered in which the Mann-Whitney statistic, equivalent to the rank-sum statistic, appears naturally. A decision theoretic upper tolerance limit for a quantile is also treated. Finally, a hypothesis testing problem concerning a quantile is shown to yield the sign test.

In each of these problems, useful ways of combining prior information with the statistical observations appear.

Other applications exist. In his Ph. D. dissertation [1], Charles Antoniak finds a need to consider mixtures of Dirichlet processes. He treats several problems, including the estimation of a mixing distribution, bio-assay, empirical Bayes problems, and discrimination problems.

2. The Dirichlet distribution. The discussion of this section is well known but our definition of the Dirichlet distribution is slightly more general than the usual one. The Dirichlet distribution makes its appearance in problems involving order statistics. A discussion of these applications and of the main properties of the Dirichlet distribution may be found in the book of S. S. Wilks [11]. The

Dirichlet distribution is known to Bayesians as the conjugate prior for the parameters of a multinomial distribution. See, for example, the book by I. J. Good [7].

We denote by $\mathcal{G}(\alpha, \beta)$ the gamma distribution with shape parameter $\alpha \geq 0$, and scale parameter $\beta > 0$. For $\alpha = 0$, this distribution is degenerate at zero; for $\alpha > 0$, this distribution has density with respect to Lebesgue measure on the real line

$$(1) \quad f(z | \alpha, \beta) = \frac{1}{\Gamma(\alpha)\beta^\alpha} e^{-z/\beta} z^{\alpha-1} I_{(0,\infty)}(z)$$

where $I_S(z)$ represents the indicator function of the set S .

We define the Dirichlet distribution slightly more generally than in Wilks [11], by allowing some of the variables to be degenerate at zero.

Let Z_1, Z_2, \dots, Z_k be independent random variables with $Z_j \in \mathcal{G}(\alpha_j, 1)$, where $\alpha_j \geq 0$ for all j , and $\alpha_j > 0$ for some $j, j = 1, 2, \dots, k$.

The Dirichlet distribution with parameter $(\alpha_1, \dots, \alpha_k)$, denoted by $\mathcal{D}(\alpha_1, \dots, \alpha_k)$, is defined as the distribution of (Y_1, \dots, Y_k) , where

$$(2) \quad Y_j = Z_j / \sum_{i=1}^k Z_i \quad \text{for } j = 1, 2, \dots, k.$$

Use of the notation $\mathcal{D}(\alpha_1, \dots, \alpha_k)$ is taken to imply that $\alpha_j \geq 0$ for all j , and $\alpha_j > 0$ for some j . This distribution is always singular with respect to Lebesgue measure in k -dimensional space since $Y_1 + \dots + Y_k = 1$. In addition, if any $\alpha_j = 0$, the corresponding Y_j is degenerate at zero. However, if $\alpha_j > 0$ for all j , the $(k - 1)$ -dimensional distribution of (Y_1, \dots, Y_{k-1}) is absolutely continuous with density

$$(3) \quad \begin{aligned} f(y_1, \dots, y_{k-1} | \alpha_1, \dots, \alpha_k) \\ = \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \dots \Gamma(\alpha_k)} (\prod_{j=1}^{k-1} y_j^{\alpha_j-1}) (1 - \sum_{j=1}^{k-1} y_j)^{\alpha_k-1} I_{\mathbb{S}}(y_1, \dots, y_{k-1}) \end{aligned}$$

where \mathbb{S} is the simplex

$$\mathbb{S} = \{(y_1, \dots, y_{k-1}) : y_j \geq 0, \sum_{j=1}^{k-1} y_j \leq 1\}.$$

For $k = 2$, (3) reduces to the Beta distribution, denoted by $\mathcal{Be}(\alpha_1, \alpha_2)$.

The main property of the Dirichlet distribution as used below is

i°. If $(Y_1, \dots, Y_k) \in \mathcal{D}(\alpha_1, \dots, \alpha_k)$ and r_1, \dots, r_l are integers such that $0 < r_1 < \dots < r_l = k$, then

$$(\sum_{i=1}^{r_1} Y_i, \sum_{i=r_1+1}^{r_2} Y_i, \dots, \sum_{i=r_{l-1}+1}^{r_l} Y_i) \in \mathcal{D}(\sum_{i=1}^{r_1} \alpha_i, \sum_{i=r_1+1}^{r_2} \alpha_i, \dots, \sum_{i=r_{l-1}+1}^{r_l} \alpha_i).$$

This follows directly from the definition of the Dirichlet distribution and the additive property of the gamma distribution: If $Z_1 \in \mathcal{G}(\alpha_1, 1)$, if $Z_2 \in \mathcal{G}(\alpha_2, 1)$, and if Z_1 and Z_2 are independent, then $Z_1 + Z_2 \in \mathcal{G}(\alpha_1 + \alpha_2, 1)$.

In particular, the marginal distribution of each Y_j is Beta: $Y_j \in \mathcal{Be}(\alpha_j, (\sum_{i=1}^k \alpha_i) - \alpha_j)$.

We record for future use the first two moments of the Dirichlet distribution.

ii°. If $(Y_1, \dots, Y_k) \in \mathcal{D}(\alpha_1, \dots, \alpha_k)$, then

$$\begin{aligned} \mathcal{E}Y_i &= \alpha_i/\alpha \\ \mathcal{E}Y_i^2 &= \alpha_i(\alpha_i + 1)/(\alpha(\alpha + 1)) && \text{and} \\ \mathcal{E}Y_i Y_j &= \alpha_i \alpha_j / (\alpha(\alpha + 1)) && \text{for } i \neq j \end{aligned}$$

where $\alpha = \sum_1^k \alpha_i$.

The following Bayes property of the Dirichlet distribution is well known.

iii°. If the prior distribution of (Y_1, \dots, Y_k) is $\mathcal{D}(\alpha_1, \dots, \alpha_k)$ and if

$$\mathcal{P}\{X = j | Y_1, \dots, Y_k\} = Y_j \quad \text{a.s.} \quad \text{for } j = 1, \dots, k,$$

then the posterior distribution of (Y_1, \dots, Y_k) given $X = j$ is $\mathcal{D}(\alpha_1^{(j)}, \dots, \alpha_k^{(j)})$, where

$$\begin{aligned} \alpha_i^{(j)} &= \alpha_i && \text{if } i \neq j \\ &= \alpha_j + 1 && \text{if } i = j. \end{aligned}$$

Contained in this property is a formula that will prove useful. Let us use $D(y_1, \dots, y_k | \alpha_1, \dots, \alpha_k)$ to denote the distribution function of the Dirichlet distribution, $\mathcal{D}(\alpha_1, \dots, \alpha_k)$. Then, the equality

$$\begin{aligned} \mathcal{P}\{X = j, Y_1 \leq z_1, \dots, Y_k \leq z_k\} \\ = \mathcal{P}\{X = j\} \mathcal{P}\{Y_1 \leq z_1, \dots, Y_k \leq z_k | X = j\} \end{aligned}$$

may be expressed in terms of $D(y_1, \dots, y_k | \alpha_1, \dots, \alpha_k)$, using ii° and iii°, as

$$\begin{aligned} (4) \quad \int_0^{z_1} \dots \int_0^{z_k} y_j \, dD(y_1, \dots, y_k | \alpha_1, \dots, \alpha_k) \\ = \frac{\alpha_j}{\alpha} D(z_1, \dots, z_k | \alpha_1^{(j)}, \dots, \alpha_k^{(j)}). \end{aligned}$$

It should be noted that this formula is true even if $\alpha_j = 0$.

3. The Dirichlet process. Let \mathcal{L} be a set and let \mathcal{A} be a σ -field of subsets of \mathcal{L} . We define below a random probability, P , on $(\mathcal{L}, \mathcal{A})$ by defining the joint distribution of the random variables $(P(A_1), \dots, P(A_m))$ for every m and every sequence A_1, \dots, A_m of measurable sets ($A_i \in \mathcal{A}$ for all i). We then verify the Kolmogorov consistency conditions to show there exists a probability, \mathcal{P} , on $([0, 1]^\omega, B_{\mathcal{F}^\omega})$ yielding these distributions. Here, $[0, 1]^\omega$ represents the space of all functions from \mathcal{A} into $[0, 1]$, and $B_{\mathcal{F}^\omega}$ represents the σ -field generated by the field of cylinder sets (Kolmogorov [8]).

For our purposes, it is more convenient to define the random probability, P , by defining the joint distribution of $(P(B_1), \dots, P(B_k))$ for all k and all measurable partitions (B_1, \dots, B_k) of \mathcal{L} . (We say (B_1, \dots, B_k) is a measurable partition of \mathcal{L} if $B_i \in \mathcal{A}$ for all i , $B_i \cap B_j = \emptyset$ for $i \neq j$, and $\bigcup_{j=1}^k B_j = \mathcal{L}$.) From these distributions, the joint distribution of $(P(A_1), \dots, P(A_m))$ for arbitrary measurable sets A_1, \dots, A_m may be defined using the hoped for finite additivity of P as follows.

Given arbitrary measurable sets A_1, \dots, A_m , form the $k = 2^m$ sets obtained by taking intersections of the A_i and their complements; that is, define B_{ν_1, \dots, ν_m} for each $\nu_j = 0$ or 1 as

$$(1) \quad B_{\nu_1, \dots, \nu_m} = \bigcap_{j=1}^m A_j^{\nu_j}$$

where A_j^1 is interpreted as A_j , and A_j^0 is interpreted as A_j^c , the complement of A_j . Thus the $\{B_{\nu_1, \dots, \nu_m}\}$ form a measurable partition of \mathcal{L} . If we are given the joint distribution of

$$(2) \quad \{P(B_{\nu_1, \dots, \nu_m}); \nu_j = 0 \text{ or } 1, j = 1, \dots, m\},$$

then we may define the joint distribution of $(P(A_1), \dots, P(A_m))$ to be that obtainable from (2) upon defining for $i = 1, \dots, m$

$$(3) \quad P(A_i) = \sum_{(\nu_1, \dots, \nu_m) \ni \nu_i=1} P(B_{\nu_1, \dots, \nu_m}).$$

We note that if (A_1, \dots, A_m) was a measurable partition to start with, then this does not lead to contradictory definitions of the distribution of $(P(A_1), \dots, P(A_m))$ provided

$$(4) \quad P(\emptyset) \text{ is degenerate at } 0.$$

(We are assuming that the distributions of the random variables are defined free of their order, so that Kolmogorov's condition (2) ([8] page 29) is automatic.) Under condition (4), the distribution of $(P(A_1), \dots, P(A_m))$ for arbitrary measurable A_1, \dots, A_m is defined uniquely, once the distributions of $(P(B_1), \dots, P(B_k))$ are given for arbitrary measurable partitions (B_1, \dots, B_k) .

If we are given a system of distributions of $(P(B_1), \dots, P(B_k))$ for all k and all measurable partitions (B_1, \dots, B_k) , there is one consistency criterion we would certainly like to have satisfied; namely,

CONDITION C. If $(B'_1, \dots, B'_{k'})$ and (B_1, \dots, B_k) are measurable partitions, and if $(B'_1, \dots, B'_{k'})$ is a refinement of (B_1, \dots, B_k) with $B_1 = \bigcup_{i=1}^{r_1} B'_i$, $B_2 = \bigcup_{i=1}^{r_2} B'_i$, \dots , $B_k = \bigcup_{i=1}^{k'} B'_i$, then the distribution of

$$(\sum_{i=1}^{r_1} P(B'_i), \sum_{i=1}^{r_2} P(B'_i), \dots, \sum_{i=1}^{k'} P(B'_i))$$

as determined from the joint distribution of $(P(B'_1), \dots, P(B'_{k'}))$, is identical to the distribution of $(P(B_1), \dots, P(B_k))$.

As the following lemma shows, this condition is sufficient for the validity of the Kolmogorov consistency conditions for the distributions of $(P(A_1), \dots, P(A_m))$ defined as in (2) and (3). In fact the lemma is valid also as a description of a random finitely additive set function, with finite values (by letting $P(A_i)$ take values in the real line, \mathbb{R}). However our interest in the present paper is with random probability measures. We will say that P is a random probability measure on $(\mathcal{L}, \mathcal{A})$, if C is satisfied, if $P(A)$ takes values only in $[0, 1]$, and if $P(\mathcal{L})$ is degenerate at 1.

LEMMA 1. *If a system of joint distributions of $(P(B_1), \dots, P(B_k))$ for all k and*

measurable partitions (B_1, \dots, B_k) is defined satisfying Condition C, and if for arbitrary measurable sets A_1, \dots, A_m , the distribution of $(P(A_1), \dots, P(A_m))$ is defined as in (1), (2), and (3), then there exists a probability \mathcal{P} on $([0, 1]^\infty, \mathcal{B}_{\mathcal{F}^\infty})$ yielding these distributions.

PROOF. Since $\mathcal{L} \cup \emptyset = \mathcal{L}$, it follows from Condition C that $P(\emptyset)$ is degenerate at zero, and thus that the distribution of $(P(A_1), \dots, P(A_m))$ is well-defined by (2) and (3). To check the Kolmogorov consistency conditions, we must show that, for arbitrary m and measurable sets A_1, \dots, A_m , the marginal distribution of $(P(A_1), \dots, P(A_{m-1}))$ derived from the distribution of $(P(A_1), \dots, P(A_m))$ is identical to the defined distribution of $(P(A_1), \dots, P(A_{m-1}))$.

The marginal distribution of $(P(A_1), \dots, P(A_{m-1}))$ derived from the distribution of $(P(A_1), \dots, P(A_m))$ is identical to the distribution of

$$(5) \quad \left(\sum_{(\nu_1, \dots, \nu_m), \nu_1=1} P(B_{\nu_1, \dots, \nu_m}), \dots, \sum_{(\nu_1, \dots, \nu_m), \nu_{m-1}=1} P(B_{\nu_1, \dots, \nu_m}) \right)$$

derived from the distribution of (2). The distribution of $(P(A_1), \dots, P(A_{m-1}))$ is defined as the distribution of

$$(6) \quad \left(\sum_{(\nu_1, \dots, \nu_{m-1}), \nu_1=1} P(B_{\nu_1, \dots, \nu_{m-1}}), \dots, \sum_{(\nu_1, \dots, \nu_{m-1}), \nu_{m-1}=1} P(B_{\nu_1, \dots, \nu_{m-1}}) \right)$$

derived from the distribution of $\{P(B_{\nu_1, \dots, \nu_{m-1}}); \nu_i = 0 \text{ or } 1, i = 1, \dots, m-1\}$ where

$$B_{\nu_1, \dots, \nu_{m-1}} = \bigcup_{j=1}^{m-1} A_j^{\nu_j}.$$

Since $B_{\nu_1, \dots, \nu_{m-1}} = B_{\nu_1, \dots, \nu_{m-1}, 1} \cup B_{\nu_1, \dots, \nu_{m-1}, 0}$, Condition C implies that the distribution of $\{P(B_{\nu_1, \dots, \nu_{m-1}}); \nu_j = 0 \text{ or } 1, j = 1, \dots, m-1\}$ is identical to the distribution of $\{P(B_{\nu_1, \dots, \nu_{m-1}, 1}) + P(B_{\nu_1, \dots, \nu_{m-1}, 0}); \nu_j = 0 \text{ or } 1, j = 1, \dots, m-1\}$ as determined from the distribution of (2). Thus, the distribution of (6) can also be found from the distribution of (2) upon replacing $P(B_{\nu_1, \dots, \nu_{m-1}})$ by $P(B_{\nu_1, \dots, \nu_{m-1}, 1}) + P(B_{\nu_1, \dots, \nu_{m-1}, 0})$. With this replacement, (6) becomes formally identical to (5), proving that their distributions are identical.

DEFINITION 1. Let α be a non-null finite measure (nonnegative and finitely additive) on $(\mathcal{L}, \mathcal{A})$. We say P is a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α if for every $k = 1, 2, \dots$, and measurable partition (B_1, \dots, B_k) of \mathcal{L} , the distribution of $(P(B_1), \dots, P(B_k))$ is Dirichlet, $\mathcal{D}(\alpha(B_1), \dots, \alpha(B_k))$.

The consistency Condition C for the Dirichlet process is exactly property i° of the Dirichlet distribution. It follows from Lemma 1 that the Kolmogorov consistency conditions are satisfied so that this actually defines a random process. In addition, since $P(\mathcal{L}^c)$ is degenerate at 1, we call P a random probability measure.

The following three propositions show a close relationship between properties of the random probability measure, P , and properties of the parameter of the process, α .

PROPOSITION 1. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α ,

and let $A \in \mathcal{A}$. If $\alpha(A) = 0$, then $P(A) = 0$ with probability one. If $\alpha(A) > 0$, then $P(A) > 0$ with probability one. Furthermore, $\mathcal{E}P(A) = \alpha(A)/\alpha(\mathcal{L})$.

PROOF. By considering the partition (A, A^c) , it is seen that $P(A)$ has a Beta distribution, $\mathcal{B}e(\alpha(A), \alpha(A^c))$. The proof follows immediately.

This proposition would seem to say that α and P have the same null sets, in other words, that α and P are mutually absolutely continuous. This interpretation is false; in fact, it is shown in the next section that P is essentially a discrete distribution. Thus, α and P may be mutually singular. The point is that the null set outside of which the conclusion of Proposition 1 holds may depend upon A .

PROPOSITION 2. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α . If α is σ -additive, then so is P in the sense that for a fixed decreasing sequence of measurable sets $A_n \searrow \emptyset$, we have $P(A_n) \rightarrow 0$ with probability one.

PROOF. Since $A_n \searrow \emptyset$ and α is additive, $\alpha(A_n) \rightarrow 0$. Hence there exists a subsequence $\{n_j\}$ such that $\sum_1^\infty \alpha(A_{n_j}) < \infty$. For fixed $\varepsilon > 0$,

$$\sum_1^\infty \mathcal{P}\{P(A_{n_j}) > \varepsilon\} \leq \sum_1^\infty \varepsilon^{-1} \mathcal{E}P(A_{n_j}) = \varepsilon^{-1} \sum_1^\infty \alpha(A_{n_j})/\alpha(\mathcal{L}) < \infty .$$

Hence, from the Borel–Cantelli lemma, $\mathcal{P}\{P(A_{n_j}) > \varepsilon \text{ i.o.}\} = 0$. This proves that $P(A_{n_j}) \rightarrow 0$ with probability one. The proof is completed by noting that $P(A_n) > P(A_{n+1})$ with probability one for all n , and hence, $P(A_1) > P(A_2) > \dots$ with probability one. The converse is also true: if α is not σ -additive, then with probability one P is not σ -additive.

PROPOSITION 3. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α , and let Q be a fixed probability measure on $(\mathcal{L}, \mathcal{A})$ with $Q \ll \alpha$. Then, for any positive integer m and measurable sets A_1, \dots, A_m , and $\varepsilon > 0$,

$$\mathcal{P}\{|P(A_i) - Q(A_i)| < \varepsilon \text{ for } i = 1, \dots, m\} > 0 .$$

PROOF. Form B_{ν_1, \dots, ν_m} as in (1) and note that

$$\begin{aligned} &\mathcal{P}\{|P(A_i) - Q(A_i)| < \varepsilon \text{ for } i = 1, \dots, m\} \\ &\geq \mathcal{P}\{\sum_{(\nu_1, \dots, \nu_m) \ni \nu_i=1} |P(B_{\nu_1, \dots, \nu_m}) - Q(B_{\nu_1, \dots, \nu_m})| < \varepsilon \text{ for } i = 1, \dots, m\} . \end{aligned}$$

Therefore, it is sufficient to show

$$\mathcal{P}\{|P(B_{\nu_1, \dots, \nu_m}) - Q(B_{\nu_1, \dots, \nu_m})| < 2^{-m}\varepsilon \text{ for all } (\nu_1, \dots, \nu_m)\} > 0 .$$

If $\alpha(B_{\nu_1, \dots, \nu_m}) = 0$, then $Q(B_{\nu_1, \dots, \nu_m}) = 0$ and $P(B_{\nu_1, \dots, \nu_m}) = 0$ with probability one, so that $|P(B_{\nu_1, \dots, \nu_m}) - Q(B_{\nu_1, \dots, \nu_m})| = 0$ with probability one. For those (ν_1, \dots, ν_m) for which $\alpha(B_{\nu_1, \dots, \nu_m}) > 0$, the distribution of the corresponding $P(B_{\nu_1, \dots, \nu_m})$ gives positive weight to all open sets in the set

$$\sum_{(\nu_1, \dots, \nu_m) \ni \alpha(B_{\nu_1, \dots, \nu_m}) > 0} P(B_{\nu_1, \dots, \nu_m}) = 1$$

completing the proof.

This proposition is a version of the desirable property (I) mentioned in the introduction. To discuss the support of a random probability measure, the topology on the space of probability measures on $(\mathcal{L}, \mathcal{A})$ must be specified. If the topology is chosen to be that of pointwise convergence ($Q_n \rightarrow Q$, if for every $A \in \mathcal{A}$, $Q_n(A) \rightarrow Q(A)$), Proposition 3 states that the support of the Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α contains the set of all probability measures absolutely continuous with respect to α . It is easy to see conversely, that any measure Q not absolutely continuous with respect to α is not in the support of P .

If $(\mathcal{L}, \mathcal{A})$ is the real line with the Borel sets, we may consider the topology of convergence in distribution on the set of σ -additive probability measures on $(\mathcal{L}, \mathcal{A})$. With this topology, it may be shown that if α is σ -additive, the support of P is the set of all σ -additive probability measures whose support is contained in the support of α .

DEFINITION 2. Let P be a random probability measure on $(\mathcal{L}, \mathcal{A})$. We say that X_1, \dots, X_n is a sample of size n from P if for any $m = 1, 2, \dots$ and measurable sets $A_1, \dots, A_m, C_1, \dots, C_n$,

$$(7) \quad \mathcal{P}\{X_1 \in C_1, \dots, X_n \in C_n \mid P(A_1), \dots, P(A_m), P(C_1), \dots, P(C_n)\} \\ = \prod_{j=1}^n P(C_j) \quad \text{a.s.}$$

Roughly, X_1, \dots, X_n is a sample of size n from P , if, given $P(C_1), \dots, P(C_n)$, the events $\{X_1 \in C_1\}, \dots, \{X_n \in C_n\}$ are independent of the rest of the process, and are independent among themselves, with $\mathcal{P}\{X_j \in C_j \mid P(C_1), \dots, P(C_n)\} = P(C_j)$ a.s. for $j = 1, \dots, n$. This definition determines the joint distribution of $X_1, \dots, X_n, P(A_1), \dots, P(A_m)$, once the distribution of the process is given, since

$$(8) \quad \mathcal{P}\{X_1 \in C_1, \dots, X_n \in C_n, P(A_1) \leq y_1, \dots, P(A_m) \leq y_m\}$$

may be found by integrating (7) with respect to the joint distribution of $P(A_1), \dots, P(A_m), P(C_1), \dots, P(C_n)$ over the set $[0, y_1] \times \dots \times [0, y_m] \times [0, 1] \times \dots \times [0, 1]$. The Kolmogorov consistency conditions may easily be checked to show that (8) determines a probability \mathcal{P} over $(\mathcal{L}^n \times [0, 1]^{\mathcal{A}}, \mathcal{A}^n \times B\mathcal{F}^{\mathcal{A}})$.

PROPOSITION 4. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α and let X be a sample of size 1 from P . Then for $A \in \mathcal{A}$,

$$\mathcal{P}(X \in A) = \alpha(A) / \alpha(\mathcal{L}).$$

PROOF. Since $\mathcal{P}(X \in A \mid P(A)) = P(A)$ a.s.,

$$\begin{aligned} \mathcal{P}(X \in A) &= \mathcal{E}\mathcal{P}(X \in A \mid P(A)) \\ &= \mathcal{E}P(A) \\ &= \alpha(A) / \alpha(\mathcal{L}), \end{aligned}$$

completing the proof.

PROPOSITION 5. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α , and let X be a sample of size 1 from P . Let (B_1, \dots, B_k) be a measurable partition of

\mathcal{L} , and let $A \in \mathcal{A}$. Then,

$$(9) \quad \mathcal{P}\{X \in A, P(B_1) \leq y_1, \dots, P(B_k) \leq y_k\} \\ = \sum_{j=1}^k \frac{\alpha(B_j \cap A)}{\alpha(\mathcal{L})} D(y_1, \dots, y_k | \alpha_1^{(j)}, \dots, \alpha_k^{(j)})$$

where $D(y_1, \dots, y_k | \alpha_1, \dots, \alpha_k)$ is the distribution function of the Dirichlet distribution, $\mathcal{L}(\alpha_1, \dots, \alpha_k)$, and where

$$\alpha_i^{(j)} = \alpha(B_i) \quad \text{if } i \neq j \\ = \alpha(B_j) + 1 \quad \text{if } i = j.$$

PROOF. Define $B_{j,1} = B_j \cap A$, and $B_{j,0} = B_j \cap A^c$ for $j = 1, \dots, k$. Let $Y_{j,\nu} = P(B_{j,\nu})$ for $j = 1, \dots, k$ and $\nu = 0, 1$. Then, from (7)

$$(10) \quad \mathcal{P}\{X \in A | Y_{j,\nu}, j = 1, \dots, k, \text{ and } \nu = 0, 1\} = \sum_{j=1}^k Y_{j,1} \text{ a.s.}$$

Hence for arbitrary $y_{j,\nu} \in [0, 1]$, for $j = 1, \dots, k$ and $\nu = 0, 1$,

$$\mathcal{P}\{X \in A, Y_{j,\nu} \leq y_{j,\nu} \text{ for } j = 1, \dots, k, \text{ and } \nu = 0, 1\}$$

can be found by integrating (10) with respect to the distribution of the $Y_{j,\nu}$ over the set $\{Y_{j,\nu} \leq y_{j,\nu}, j = 1, \dots, k \text{ and } \nu = 0, 1\}$. This integration turns out to be (see (4) of Section 2)

$$\sum_{j=1}^k \frac{\alpha(B_{j,1})}{\alpha(\mathcal{L})} D(\mathbf{y} | \boldsymbol{\alpha}^{(j)})$$

where $\mathbf{y} = (y_{1,0}, \dots, y_{k,0}, y_{1,1}, \dots, y_{k,1})$ and $\boldsymbol{\alpha}^{(j)} = (\alpha_{1,0}^{(j)}, \dots, \alpha_{k,0}^{(j)}, \alpha_{1,1}^{(j)}, \dots, \alpha_{k,1}^{(j)})$, and where

$$\alpha_{i,\nu}^{(j)} = \alpha(B_{i,\nu}) \quad \text{if } i \neq j \\ = \alpha(B_{j,\nu}) + 1 \quad \text{if } i = j.$$

The conclusion of the proposition follows from this using property i° of the Dirichlet distribution, since $P(B_j) = Y_{j,0} + Y_{j,1}$ a.s., and since the process of finding marginal distributions of random variables is linear.

We are now prepared to find the conditional distribution of a Dirichlet process P , given a sample X_1, \dots, X_n from P . It turns out that this conditional distribution is also a Dirichlet process.

For $x \in \mathcal{L}$, let δ_x denote the measure on $(\mathcal{L}, \mathcal{A})$ giving mass one to the point x :

$$\delta_x(A) = 1 \quad \text{if } x \in A \\ = 0 \quad \text{if } x \notin A.$$

THEOREM 1. Let P be a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α , and let X_1, \dots, X_n be a sample of size n from P . Then the conditional distribution of P given X_1, \dots, X_n , is as a Dirichlet process with parameter $\alpha + \sum_{i=1}^n \delta_{X_i}$.

PROOF. It is sufficient to prove the theorem for $n = 1$, since the theorem would then follow by induction upon repeated application of the case $n = 1$. Let (B_1, \dots, B_k) be a measurable partition of \mathcal{L} and let $A \in \mathcal{A}$. It is easy to

check that the marginal distributions of a conditional distribution of a process are identical to the conditional distributions of the marginals. Hence, we must show that the conditional distribution of $P(B_1), \dots, P(B_k)$ given X , a sample of size one from P , has distribution function

$$(11) \quad D(y_1, \dots, y_k | \alpha(B_1) + \delta_x(B_1), \dots, \alpha(B_k) + \delta_x(B_k)) .$$

This may be done by showing that the integral of (11) with respect to the marginal distribution of X over A is equal to the probability (9). Using the marginal distribution of X as found in Proposition 4, we compute

$$\begin{aligned} \int_A D(y_1, \dots, y_k | \alpha(B_1) + \delta_x(B_1), \dots, \alpha(B_k) + \delta_x(B_k)) d\alpha(x)/\alpha(\mathcal{L}) \\ = \sum_{j=1}^k \int_{B_j \cap A} D(y_1, \dots, y_k | \alpha_1^{(j)}, \dots, \alpha_k^{(j)}) d\alpha(x)/\alpha(\mathcal{L}) \\ = \sum_{j=1}^k \frac{\alpha(B_j \cap A)}{\alpha(\mathcal{L})} D(y_1, \dots, y_k | \alpha_1^{(j)}, \dots, \alpha_k^{(j)}) , \end{aligned}$$

completing the proof.

4. An alternative definition of the Dirichlet process. In this section, we define a random probability measure which is a Dirichlet process on $(\mathcal{L}, \mathcal{A})$ with parameter α and which with probability one is a discrete probability measure on $(\mathcal{L}, \mathcal{A})$.

The basic idea is that since the Dirichlet distribution is definable, as in (2) of Section 2, as the joint distribution of a set of independent gamma variables divided by the sum, so also should the Dirichlet process be definable as a gamma process with independent ‘‘increments’’ divided by the sum. Using a representation of a process with independent increments as a sum of a countable number of jumps of random height at a countable number of random points, as found in [4], we may divide by the total heights of the jumps and obtain a discrete probability measure, which should be distributed as a Dirichlet process.

The gamma distribution, $\mathcal{G}(\alpha, 1)$, $\alpha > 0$, has characteristic function (see Gnedenko and Kolmogorov ([6] pages 86–87)),

$$(1) \quad \begin{aligned} \varphi(a) &= (1 - it)^{-\alpha} \\ &= \exp \int_0^\infty (e^{iuz} - 1) dN(x) \end{aligned}$$

where

$$(2) \quad N(x) = -\alpha \int_x^\infty e^{-y} y^{-1} dy \quad \text{for } 0 < x < \infty .$$

We define the distribution of random variables J_1, J_2, \dots as follows.

$$(3) \quad \mathcal{P}(J_1 \leq x_1) = e^{N(x_1)} \quad \text{for } x_1 > 0$$

and for $j = 2, 3, \dots$

$$(4) \quad \mathcal{P}(J_j \leq x_j | J_{j-1} = x_{j-1}, \dots, J_1 = x_1) = \exp[N(x_j) - N(x_{j-1})] \\ \text{for } 0 < x_j < x_{j-1} .$$

In other words, the distribution function of J_1 is $\exp N(x_1)$ and for $j = 2, 3, \dots$,

the distribution of J_j given J_{j-1}, \dots, J_1 , is the same as the distribution of J_1 truncated above at J_{j-1} . The following theorem is taken from the main theorem of [4].

THEOREM 1. *Let $G(t)$ be a distribution function on $[0, 1]$. Let*

$$(5) \quad Z_t = \sum_{j=1}^{\infty} J_j I_{[0, G(t)]}(U_j),$$

where (i) the distribution of J_1, J_2, \dots is given in (3) and (4), and

(ii) U_1, U_2, \dots are independent identically distributed variables, uniformly distributed on $[0, 1]$, and independent of J_1, J_2, \dots . Then, with probability one, Z_t converges for all $t \in [0, 1]$ and is a gamma process with independent increments, with $Z_t \in \mathcal{G}(\alpha G(t), 1)$.

In particular, $Z_1 = \sum_1^{\infty} J_j$ converges with probability one and $Z_1 \in \mathcal{G}(\alpha, 1)$. If we define

$$(6) \quad P_j = J_j / Z_1,$$

then $P_j \geq 0$ and $\sum_1^{\infty} P_j = 1$ with probability one.

We now define the Dirichlet process. As before, let $(\mathcal{X}, \mathcal{A})$ be a measurable space, and let $\alpha(\cdot)$ be a finite non-null measure on \mathcal{A} . Let V_1, V_2, \dots be a sequence of independent identically distributed random variables with values in \mathcal{X} , and with probability measure Q , where $Q(A) = \alpha(A) / \alpha(\mathcal{X})$. (More specifically, let $(\mathcal{X}_j, \mathcal{A}_j, Q_j)$ be identical copies of $(\mathcal{X}, \mathcal{A}, Q)$, and let V_j be the identity map from \mathcal{X}_j to \mathcal{X} . Then the V_j are extended to be defined on the infinite product space $(\prod \mathcal{X}_j, \prod \mathcal{A}_j, \prod Q_j)$ in the usual manner.)

We identify the α in formulas (1) and (2) with $\alpha(\mathcal{X})$, and define the random probability measure, P , on $(\mathcal{X}, \mathcal{A})$, as

$$(7) \quad P(A) = \sum_{j=1}^{\infty} P_j \delta_{V_j}(A).$$

THEOREM 2. *The random probability measure defined by (7) is a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter α .*

PROOF. Let (B_1, \dots, B_k) be a measurable partition of \mathcal{X} . Then

$$(P(B_1), \dots, P(B_k)) = \frac{1}{Z_1} \sum_{j=1}^{\infty} J_j (\delta_{V_j}(B_1), \dots, \delta_{V_j}(B_k)).$$

From the assumption on the distribution of V_1, V_2, \dots ,

$$\mathbf{M}_j = (\delta_{V_j}(B_1), \dots, \delta_{V_j}(B_k))$$

are independent identically distributed random vectors having a multinomial distribution with probability vector $(Q(B_1), \dots, Q(B_k))$. Hence, the distribution of $\sum_1^{\infty} J_j \mathbf{M}_j$ must be the same as the distribution of

$$(Z_{1/k}, Z_{2/k} - Z_{1/k}, \dots, Z_1 - Z_{(k-1)/k})$$

where Z_t is the gamma process defined by (5) with $G(t)$ chosen so that

$$G\left(\frac{j}{k}\right) - G\left(\frac{j-1}{k}\right) = Q(B_j) \quad j = 1, \dots, k.$$

Hence, $\sum_{j=1}^{\infty} J_j \delta_{V_j}(B_i)$ are, for $i = 1, \dots, k$, independent random variables, with $\sum_{j=1}^{\infty} J_j \delta_{V_j}(B_i) \in \mathcal{D}(\alpha(B_i), 1)$. Since Z_1 is the sum of these independent gamma variables, $(P(B_1), \dots, P(B_k)) \in \mathcal{D}(\alpha(B_1), \dots, \alpha(B_k))$ from the definition of the Dirichlet distribution. Thus, P satisfies the definition of the Dirichlet process.

THEOREM 3. *Let P be the Dirichlet process defined by (7), and let Z be a measurable real valued function defined on $(\mathcal{L}, \mathcal{A})$. If $\int |Z| d\alpha < \infty$, then $\int |Z| dP < \infty$ with probability one, and*

$$\mathcal{E} \int Z dP = \int Z d\mathcal{E}P = \alpha(\mathcal{L})^{-1} \int Z d\alpha .$$

PROOF. From (7)

$$(8) \quad \int |Z| dP = \sum_{j=1}^{\infty} |Z(V_j)| P_j$$

so that the monotone convergence theorem gives

$$\begin{aligned} \mathcal{E} \int |Z| dP &= \sum_{j=1}^{\infty} \mathcal{E} |Z(V_j)| \mathcal{E} P_j \\ &= \alpha(\mathcal{L})^{-1} \int |Z| d\alpha \sum_{j=1}^{\infty} \mathcal{E} P_j \\ &= \alpha(\mathcal{L})^{-1} \int |Z| d\alpha \end{aligned}$$

where we have used the independence of the V_j and the P_j . Therefore, $\int |Z| dP$ is finite with probability one, and hence

$$\int Z dP = \sum_{j=1}^{\infty} Z(V_j) P_j$$

is absolutely convergent with probability one. Since this series is bounded by (8), which is integrable, the bounded convergence theorem implies

$$\begin{aligned} \mathcal{E} \int Z dP &= \sum_{j=1}^{\infty} \mathcal{E} Z(V_j) \mathcal{E} P_j \\ &= \alpha(\mathcal{L})^{-1} \int Z d\alpha \end{aligned}$$

completing the proof.

This theorem emphasises the close relationship between α and the random probability measure P . It implies, in particular, that if $(\mathcal{L}, \mathcal{A})$ were the real line and the Borel sets, and if α has a finite k th moment, then with probability one P has a finite k th moment.

THEOREM 4. *Let P be the Dirichlet process defined by (7), and let Z_1 and Z_2 be measurable real valued functions defined on $(\mathcal{L}, \mathcal{A})$. If $\int |Z_1| d\alpha < \infty$, $\int |Z_2| d\alpha < \infty$ and $\int |Z_1 Z_2| d\alpha < \infty$, then*

$$\mathcal{E} \int Z_1 dP \int Z_2 dP = \frac{\sigma_{12}}{\alpha(\mathcal{L}) + 1} + \mu_1 \mu_2$$

where

$$\begin{aligned} \mu_i &= \alpha(\mathcal{L})^{-1} \int Z_i d\alpha & i = 1, 2 \text{ and} \\ \sigma_{12} &= \alpha(\mathcal{L})^{-1} \int Z_1 Z_2 d\alpha - \mu_1 \mu_2 . \end{aligned}$$

PROOF. As in Theorem 3

$$(11) \quad \begin{aligned} \int Z_1 dP \int Z_2 dP &= \sum_i Z_1(V_i) P_i \sum_j Z_2(V_j) P_j \\ &= \sum_i \sum_j Z_1(V_i) Z_2(V_j) P_i P_j \end{aligned}$$

since both series are absolutely convergent with probability one. This is bounded in absolute value by

$$(12) \quad \sum_i \sum_j |Z_1(V_i)Z_2(V_j)|P_i P_j .$$

If this is an integrable random variable, we may take an expectation of (11) inside the summation sign and obtain

$$\begin{aligned} \mathcal{E} \int Z_1 dP \int Z_2 dP &= \sum_i \sum_j \mathcal{E}(Z_1(V_i)Z_2(V_j))\mathcal{E}(P_i P_j) \\ &= \sum_i \sum_{i \neq j} \mathcal{E}Z_1(V_i)\mathcal{E}Z_2(V_j)\mathcal{E}(P_i P_j) \\ &\quad + \sum_i \mathcal{E}(Z_1(V_i)Z_2(V_i))\mathcal{E}P_i^2 \end{aligned}$$

using the independence of the P_i and the V_i , and the independence of the V_i among themselves. The equation continues

$$\begin{aligned} &= \mu_1 \mu_2 \sum_i \sum_{i \neq j} \mathcal{E}(P_i P_j) + (\sigma_{12} + \mu_1 \mu_2) \sum_i \mathcal{E}P_i^2 \\ &= \mu_1 \mu_2 + \sigma_{12} \mathcal{E} \sum_i P_i^2 . \end{aligned}$$

An analogous equation shows that (12) is integrable. The proof will be complete when we show

$$\mathcal{E} \sum_{i=1}^{\infty} P_i^2 = \frac{1}{\alpha(\mathcal{X}) + 1} .$$

This seems difficult to show directly from the definition of the P_i , so we proceed as follows. The distribution of the P_i depends on α only through the value of $\alpha(\mathcal{X})$. So choose \mathcal{X} to be the real line, α to give mass $\alpha(\mathcal{X})/2$ to -1 and mass $\alpha(\mathcal{X})/2$ to $+1$, and $Z_1(x) = Z_2(x)$ to be identically x . Then $\mu_1 = \mu_2 = 0$ and $\sigma_{12} = 1$. Hence

$$\begin{aligned} \mathcal{E} \sum_1^{\infty} P_i^2 &= \mathcal{E}(\int x dP(x))^2 = \mathcal{E}(2P(\{1\}) - 1)^2 \\ &= \frac{1}{\alpha(\mathcal{X}) + 1} \end{aligned}$$

since $P(\{1\}) \in \mathcal{B}e(\alpha(\mathcal{X})/2, \alpha(\mathcal{X})/2)$, completing the proof.

This theorem states that the covariance of the random variables $\int Z_1 dP$ and $\int Z_2 dP$ is equal to $(\alpha(\mathcal{X}) + 1)^{-1}$ times the covariance of Z_1 and Z_2 as random variables on $(\mathcal{X}, \mathcal{A}, Q)$ where $Q = \mathcal{E}P = \alpha/\alpha(\mathcal{X})$. In particular, the correlation coefficient of $\int Z_1 dP$ and $\int Z_2 dP$ is equal to the correlation coefficient of Z_1 and Z_2 as random variables on $(\mathcal{X}, \mathcal{A}, Q)$ (assuming the finiteness of $\int Z_1^2 d\alpha$ and $\int Z_2^2 d\alpha$).

5. Applications. Throughout this section, α is taken to denote a σ -additive non-null finite measure on $(\mathcal{X}, \mathcal{A})$. We write $P \in \mathcal{D}(\alpha)$ as a notation for the phrase “ P is a Dirichlet process on $(\mathcal{X}, \mathcal{A})$ with parameter α .” We let \mathbb{R} denote the real line and \mathcal{B} the σ -field of Borel sets. In most of the applications we take $(\mathcal{X}, \mathcal{A}) = (\mathbb{R}, \mathcal{B})$.

The nonparametric statistical decision problems we consider are typically

described as follows. The parameter space is the set of all probability measures P on $(\mathcal{X}, \mathcal{A})$. The statistician is to choose an action a in some space, thereby incurring a loss, $L(P, a)$. There is a sample X_1, \dots, X_n from P available to the statistician, upon which he may base his choice of action. He seeks a Bayes rule with respect to the prior distribution, $P \in \mathcal{D}(\alpha)$.

With such a prior distribution, the posterior distribution of P given the observations is $\mathcal{D}(\alpha + \sum_1^n \delta_{X_i})$, where δ_x denotes the measure giving mass one to the point x . Thus, if we can find a Bayes rule for the no-sample problem (with $n = 0$), a Bayes rule for the general problem may be found by replacing α with $\alpha + \sum_1^n \delta_{X_i}$. In the problems considered below, we first find the Bayes rule for the no-sample problem, and then state the Bayes rule for the general problem.

(a) *Estimation of a distribution function.* Let $(\mathcal{X}, \mathcal{A}) = (\mathbb{R}, \mathcal{B})$, and let the space of actions of the statistician be the space of all distribution functions on \mathbb{R} . Let the loss function be

$$L(P, \hat{F}) = \int (F(t) - \hat{F}(t))^2 dW(t)$$

where W is a given finite measure on $(\mathbb{R}, \mathcal{B})$ (a weight function), and where

$$(1) \quad F(t) = P((-\infty, t]).$$

If $P \in \mathcal{D}(\alpha)$, then $F(t) \in \mathcal{B}e(\alpha((-\infty, t]), \alpha((t, \infty)))$ for each t . The Bayes risk for the no-sample problem,

$$\mathcal{E}L(P, \hat{F}) = \int \mathcal{E}(F(t) - \hat{F}(t))^2 dW(t),$$

is minimized by choosing $\hat{F}(t)$ for each t to minimize $\mathcal{E}(F(t) - \hat{F}(t))^2$. This is achieved by choosing $\hat{F}(t)$ to be $\mathcal{E}F(t)$. Thus, the Bayes rule for the no-sample problem is

$$\hat{F}(t) = \mathcal{E}F(t) = F_0(t)$$

where

$$(2) \quad F_0(t) = \alpha((-\infty, t]) / \alpha(\mathbb{R})$$

represents our prior guess at the shape of the unknown $F(t)$.

For a sample of size n , the Bayes rule is therefore

$$(3) \quad \hat{F}_n(t | X_1, \dots, X_n) = \frac{\alpha((-\infty, t]) + \sum_1^n \delta_{X_i}((-\infty, t])}{\alpha(\mathbb{R}) + n} \\ = p_n F_0(t) + (1 - p_n) F_n(t | X_1, \dots, X_n)$$

where

$$(4) \quad p_n = \alpha(\mathbb{R}) / (\alpha(\mathbb{R}) + n)$$

and

$$F_n(t | X_1, \dots, X_n) = \frac{1}{n} \sum_{i=1}^n \delta_{X_i}((-\infty, t])$$

is the empirical distribution function of the sample.

The Bayes rule (3) is a mixture of our prior guess at F and of the empirical distribution function, with respective weights p_n and $(1 - p_n)$. If $\alpha(\mathbb{R})$ is large compared to n , little weight is given to the observations. If $\alpha(\mathbb{R})$ is small compared to n , little weight is given to the prior guess at F . One might interpret $\alpha(\mathbb{R})$ as a measure of faith in the prior guess at F measured in units of numbers of observations. As $\alpha(\mathbb{R})$ tends to zero (the “noninformative” Dirichlet prior), the Bayes estimate converges to the empirical distribution function.

It is interesting to note that whatever be the true distribution function, the Bayes estimate (3) converges to it uniformly almost surely. This follows from the Glivenko–Cantelli theorem and the observation that $p_n \rightarrow 0$ as $n \rightarrow \infty$.

The results for estimating a k -dimensional distribution function are completely analogous.

(b) *Estimation of the mean.* Again let $(\mathcal{L}, \mathcal{A}) = (\mathbb{R}, \mathcal{B})$, and suppose the statistician is to estimate the mean with squared error loss

$$L(P, \hat{\mu}) = (\mu - \hat{\mu})^2$$

where

$$(5) \quad \mu = \int x dP(x).$$

We assume $P \in \mathcal{S}(\alpha)$, where α has finite first moment. The mean of the corresponding probability measure $\alpha(\cdot)/\alpha(\mathbb{R})$ is denoted by μ_0 :

$$(6) \quad \mu_0 = \int x d\alpha(x)/\alpha(\mathbb{R}).$$

By Theorem 3, the random variable μ defined by (5) exists. The Bayes rule for the no-sample problem is the mean of μ , which, again by Theorem 3, is $\hat{\mu} = \mu_0$.

For a sample of size n , the Bayes rule is therefore

$$(7) \quad \begin{aligned} \hat{\mu}_n(X_1, \dots, X_n) &= (\alpha(\mathbb{R}) + n)^{-1} \int x d(\alpha(x) + \sum_1^n \delta_{X_i}(x)) \\ &= p_n \mu_0 + (1 - p_n) \bar{X}_n \end{aligned}$$

where p_n is given by (4) and \bar{X}_n is the sample mean,

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i.$$

The Bayes estimate is thus between the prior guess at μ , namely μ_0 , and the sample mean. As $\alpha(\mathbb{R}) \rightarrow 0$, $\hat{\mu}_n$ converges to \bar{X}_n . Also, as $n \rightarrow \infty$, $p_n \rightarrow 0$ so that, in particular, the Bayes estimate (7) is strongly consistent within the class of distributions with finite first moment.

More generally, for arbitrary $(\mathcal{L}, \mathcal{A})$, if Z is real-valued measurable defined on $(\mathcal{L}, \mathcal{A})$, and if we are to estimate

$$\theta = \int Z dP$$

with squared error loss and prior $P \in \mathcal{S}(\alpha)$, where α is such that

$$\theta_0 = \int Z d\alpha/\alpha(\mathcal{L}) < \infty,$$

then the estimate $\hat{\theta} = \theta_0$ is Bayes for the no-sample problem. For a sample of size n ,

$$\hat{\theta}_n(X_1, \dots, X_n) = p_n \theta_0 + (1 - p_n) \frac{1}{n} \sum_1^n Z(X_i)$$

is Bayes, where $p_n = \alpha(\mathcal{L})/(\alpha(\mathcal{L}) + n)$. Results for estimating a mean vector in k -dimensions are completely analogous.

(c) *Estimation of the median.* Again let $(\mathcal{X}, \mathcal{A}) = (\mathbb{R}, \mathcal{B})$ and suppose the statistician is to estimate the median

$$(8) \quad m = \text{med } P$$

of an unknown probability measure P on $(\mathbb{R}, \mathcal{B})$. If $P \in \mathcal{D}(\alpha)$, the median of P is unique with probability one. To see this, note that $F(t)$, defined by (1), increases whenever α increases, with probability one. Any multiple medians of P must occur on an interval of measure zero of α . There are only a countable number of such intervals, and the probability that any such interval is an interval of medians is just $\mathcal{P}\{F(t) = \frac{1}{2}\}$ for t any interior point of the interval. This is zero since $F(t)$ has a Beta distribution. Thus, for $P \in \mathcal{D}(\alpha)$, m defined by (8) is a random variable.

The Bayes estimate of m for the no-sample problem with Dirichlet process prior and squared error loss is the expectation of m . Unfortunately, this expectation is difficult to compute, and may, in fact, not even exist. Instead, we seek the Bayes estimate of m under absolute error loss,

$$L(P, \hat{m}) = |m - \hat{m}|.$$

As is well known, any median of the distribution of m is a Bayes estimate of m . For the Dirichlet process prior, any median of the distribution of m is a median of the expectation of P , and conversely:

$$(9) \quad \text{med}(\text{dist. med } P) = \text{med } \mathcal{E}P.$$

This may be seen as follows. A number t is a median of the distribution of m if and only if

$$(10) \quad \mathcal{P}\{m < t\} \leq \frac{1}{2} \leq \mathcal{P}\{m \leq t\}.$$

Since $\mathcal{P}\{m \leq t\} = \mathcal{P}\{F(t) > \frac{1}{2}\}$ by the definition of m , and since $F(t)$ has a $\text{Be}(\alpha((-\infty, t]), \alpha((t, \infty)))$ distribution, whose median is a non-decreasing function of t with value one-half if and only if the two parameters are equal, we see that t satisfies (10) if and only if

$$(11) \quad \frac{\alpha((-\infty, t])}{\alpha(\mathbb{R})} \leq \frac{1}{2} \leq \frac{\alpha((-\infty, t])}{\alpha(\mathbb{R})}.$$

Such t are exactly the medians of $\mathcal{E}P$, proving (9).

Thus, any number t satisfying (11) is a Bayes estimate of m for prior $\mathcal{D}(\alpha)$ and absolute error loss. For F_0 defined by (2),

$$\hat{m} = \text{median of } F_0.$$

For a sample of size n , the Bayes estimate is therefore

$$\hat{m}_n(X_1, \dots, X_n) = \text{median of } \hat{F}_n$$

where \hat{F}_n is the Bayes estimate of F given by (3).

(d) *Estimation of quantiles.* We extend the analysis of part (c) to the estimation of the q th quantile of P , denoted by t_q :

$$P((-\infty, t_q)) \leq q \leq P((-\infty, t_q]) .$$

As in the case of the median, it is easy to see that for $0 < q < 1$, the q th quantile of $P \in \mathcal{D}(\alpha)$ is unique with probability one, so that t_q is a well-defined random variable.

We consider the problem of estimating t_q with loss for some p , $0 < p < 1$,

$$\begin{aligned} L(P, \hat{t}_q) &= p(t_q - \hat{t}_q) && \text{if } t_q \geq \hat{t}_q \\ &= (1 - p)(\hat{t}_q - t_q) && \text{if } t_q < \hat{t}_q . \end{aligned}$$

As is well known, any p th quantile of the distribution of t_q is a Bayes estimate of t_q under this loss. The distribution of t_q may be found from the formula

$$\begin{aligned} (12) \quad \mathcal{S}\{t_q \leq t\} &= \mathcal{S}\{F(t) > q\} \\ &= \int_q^1 \frac{\Gamma(M)}{\Gamma(uM)\Gamma((1-u)M)} z^{uM-1}(1-z)^{(1-u)M-1} dz \end{aligned}$$

where

$$\begin{aligned} M &= \alpha(\mathbb{R}) \\ u &= \alpha((-\infty, t]) / \alpha(\mathbb{R}) = F_0(t) . \end{aligned}$$

To find the p th quantile of t_q , we set (12) equal to p and solve for t .

$$(13) \quad \int_q^1 \frac{\Gamma(M)}{\Gamma(uM)\Gamma((1-u)M)} z^{uM-1}(1-z)^{(1-u)M-1} dz = p .$$

For fixed p , q , and M , we let (13) define a function $u(p, q, M)$. The Bayes estimate of t_q for the no-sample problem is the u th quantile of F_0 ,

$$\hat{t}_q = u(p, q, \alpha(\mathbb{R}))\text{th quantile of } F_0 .$$

For a sample of size n , the Bayes estimate of t_q is therefore

$$(14) \quad \hat{t}_q(X_1, \dots, X_n) = u(p, q, \alpha(\mathbb{R}) + n)\text{th quantile of } \hat{F}_n$$

where \hat{F}_n is the Bayes estimate of F given by (3). If p and q are both $\frac{1}{2}$, this reduces to the estimate of (c), since $u(\frac{1}{2}, \frac{1}{2}, M) = \frac{1}{2}$ for all M , as is seen from (13).

If tables of the function $u(p, q, M)$ were available, it would be an easy matter to find the estimate (14). Unfortunately, it is difficult to obtain values of $u(p, q, M)$ from existing tables of the incomplete Beta function. The author has tables of u for $p = .05(.05).95$, $q = .05(.05).95$, and $M = 1(1)10$.

(e) *Estimation of a variance and a covariance.* Consider the problem of

estimating the variance of an unknown probability distribution P with squared error loss

$$L(P, \hat{\sigma}^2) = (\text{Var } P - \hat{\sigma}^2)^2.$$

If $P \in \mathcal{D}(\alpha)$, and if α has a finite second moment, then

$$\text{Var } P = \int x^2 dP(x) - (\int x dP(x))^2$$

is a random variable whose expectation (a Bayes estimate for the no-sample problem) may be computed from Theorems 3 and 4 as follows

$$\begin{aligned} \mathcal{E} \text{Var } P &= \mathcal{E} \int x^2 P(x) - \mathcal{E} (\int x dP(x))^2 \\ &= (\sigma_0^2 + \mu_0^2) - \left(\frac{\sigma_0^2}{\alpha(\mathbb{R}) + 1} + \mu_0^2 \right) \\ &= \frac{\alpha(\mathbb{R})}{\alpha(\mathbb{R}) + 1} \sigma_0^2 \end{aligned}$$

where μ_0 is given by (6) and where

$$\sigma_0^2 = \alpha(\mathbb{R})^{-1} \int x^2 d\alpha(x) - \mu_0^2$$

is the variance of F_0 , the prior guess at F .

For a sample of size n , the Bayes rule is therefore

$$\hat{\sigma}_n^2(X_1, \dots, X_n) = \frac{\alpha(\mathbb{R}) + n}{\alpha(\mathbb{R}) + n + 1} \text{Var } \hat{F}_n$$

where \hat{F}_n is given by (3). Using an easy formula for the variance of a mixture, we find

$$\begin{aligned} \hat{\sigma}_n^2(X_1, \dots, X_n) &= \frac{\alpha(\mathbb{R}) + n}{\alpha(\mathbb{R}) + n + 1} \text{Var} (p_n F_0 + (1 - p_n) F_n) \\ (15) \quad &= \frac{\alpha(\mathbb{R}) + n}{\alpha(\mathbb{R}) + n + 1} (p_n \sigma_0^2 + (1 - p_n) s_n^2 \\ &\quad + p_n(1 - p_n)(\mu_0 - \bar{X}_n)^2) \end{aligned}$$

where μ_0 , σ_0^2 , p_n , and \bar{X}_n are as before, and where s_n^2 is the sample variance

$$s_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

An alternative form of the estimate (15) expresses $\hat{\sigma}_n^2$ as a mixture of three different estimates of the variance:

$$\begin{aligned} \hat{\sigma}_n^2(X_1, \dots, X_n) \\ = \frac{\alpha(\mathbb{R}) + n}{\alpha(\mathbb{R}) + n + 1} \left(p_n \sigma_0^2 + (1 - p_n) \left(p_n \frac{1}{n} \sum_{i=1}^n (X_i - \mu_0)^2 + (1 - p_n) s_n^2 \right) \right). \end{aligned}$$

If we let our prior sample size $\alpha(\mathbb{R})$ tend to zero, keeping F_0 fixed, we find that $\hat{\sigma}_n^2$ converges to the estimate

$$\frac{1}{n + 1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

This estimate is well known as the best invariant or minimax estimate of the variance of a normal distribution under relative squared error loss (squared error divided by $(\text{Var } P)^2$). Its appearance in this problem is rather surprising.

To estimate the covariance of a distribution P in the plane,

$$\text{Cov } P = \int xy \, dP - \int x \, dP \int y \, dP,$$

a similar analysis may be carried out. Here, \mathcal{X} represents the Euclidean plane \mathbb{R}^2 , \mathcal{A} represents the Borel subsets, and α represents a finite measure thereon. With squared error loss, the Bayes estimate of $\text{Cov } P$ with respect to the prior $P \in \mathcal{D}(\alpha)$, is for the no-sample problem

$$\mathcal{E} \text{Cov } P = \frac{\alpha(\mathbb{R}^2)}{\alpha(\mathbb{R}^2) + 1} \sigma_{12}$$

where, as in Theorem 4, σ_{12} is the covariance of the distribution $\mathcal{E}P$,

$$\sigma_{12} = \alpha(\mathbb{R}^2)^{-1} \int xy \, d\alpha(x, y) - \mu_1 \mu_2$$

$$\mu_1 = \alpha(\mathbb{R}^2)^{-1} \int x \, d\alpha(x, y)$$

$$\mu_2 = \alpha(\mathbb{R}^2)^{-1} \int y \, d\alpha(x, y).$$

For a sample of size n , the Bayes estimate is

$$\begin{aligned} \hat{\sigma}_{12}(X_1, Y_1, \dots, X_n, Y_n) &= \frac{\alpha(\mathbb{R}^2) + n}{\alpha(\mathbb{R}^2) + n + 1} (p_n \sigma_{12} + (1 - p_n) s_{12}) \\ &\quad + p_n(1 - p_n)(\mu_1 - \bar{X}_n)(\mu_2 - \bar{Y}_n) \end{aligned}$$

where s_{12} is the sample covariance

$$s_{12} = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n).$$

(f) *Estimation of $\int F \, dG$ for a two-sample problem.* Let F and G be two distribution functions on the real line, and let X_1, \dots, X_m be a sample from F and Y_1, \dots, Y_n a sample from G . Consider the problem of estimating probability that $X_1 \leq Y_1$, denoted by Δ ,

$$\Delta = \int F \, dG$$

with squared error loss. As a prior distribution for (F, G) , we assume that F is the distribution function of $P_1 \in \mathcal{D}(\alpha_1)$, that G is the distribution function of $P_2 \in \mathcal{D}(\alpha_2)$ and that P_1 and P_2 are independent. A computation similar to that found in Theorems 3 and 4 shows that the Bayes rule for the no-sample problem is

$$\Delta_0 = \mathcal{E}\Delta = \int F_0 \, dG_0$$

where $F_0 = \mathcal{E}F$ and $G_0 = \mathcal{E}G$.

Given the two samples, the Bayes rule is

$$\hat{\Delta}(X_1, \dots, X_m, Y_1, \dots, Y_n) = \int \hat{F}_m \, d\hat{G}_n$$

where \hat{F}_m and \hat{G}_n are the respective Bayes estimates of F and G , as in (3). This may be written

$$\begin{aligned} \hat{\Delta}(X_1, \dots, X_m, Y_1, \dots, Y_n) &= p_{1,m} p_{2,n} \Delta_0 + p_{1,m} (1 - p_{2,n}) \frac{1}{n} \sum_1^n F_0(Y_j) \\ &\quad + (1 - p_{1,m}) p_{2,n} \frac{1}{m} \sum_1^m (1 - G_0(X_i^-)) \\ &\quad + (1 - p_{1,m})(1 - p_{2,n}) \frac{1}{mn} U \end{aligned}$$

where
$$p_{1,m} = \frac{\alpha_1(\mathbb{R})}{\alpha_1(\mathbb{R}) + m}, \quad p_{2,n} = \frac{\alpha_2(\mathbb{R})}{\alpha_2(\mathbb{R}) + n}$$

and where U , the number of pairs (X_i, Y_j) for which $X_i \leq Y_j$,

$$U = \sum_i \sum_j I_{(-\infty, Y_j]}(X_i)$$

is the Mann–Whitney statistic, a linear function of the rank-sum statistic. It is interesting to note that the estimate $\hat{\Delta}$ is a simple mixture of four separate estimates of Δ . As both $\alpha_1(\mathbb{R})$ and $\alpha_2(\mathbb{R})$ tend to zero, the estimate $\hat{\Delta}$ converges to $(mn)^{-1}U$, the usual nonparametric estimate.

(g) “Tolerance” regions. The notion of tolerance regions that we treat is not the usual one, but rather the decision theoretic analogue. Consider the problem of estimating the q th quantile t_q of an unknown distribution P on the real line by an upper “tolerance” point \underline{a} with loss function

$$L(P, \underline{a}) = pP((-\infty, \underline{a}]) + qI_{(a, \infty)}(t_q)$$

where p is a constant $0 < p < 1$. If t_q is known exactly, L is minimized by choosing $\underline{a} = t_q$. But if t_q is only vaguely known, it is best to “overestimate” t_q to keep the expectation of the second term small, provided the expectation of the first term is not too enlarged.

If $P \in \mathcal{S}(\alpha)$, the Bayes risk for the no-sample problem is

$$\begin{aligned} \mathcal{E}L(P, \underline{a}) &= p\mathcal{E}P((-\infty, \underline{a}]) + q\mathcal{P}\{t_q > \underline{a}\} \\ (16) \quad &= pF_0(\underline{a}) + q\mathcal{P}\{F(\underline{a}) \leq q\} \\ &= pu + q \int_0^q \frac{\Gamma(M)}{\Gamma(uM)\Gamma((1-u)M)} z^{uM-1}(1-z)^{(1-u)M-1} dz \end{aligned}$$

where u represents $F_0(\underline{a})$ and $M = \alpha(\mathbb{R})$. Since this Bayes risk depends on \underline{a} only through $F_0(\underline{a})$, we seek that value of u in $[0, 1]$ that minimizes (16). It is not difficult to show uniqueness of the point at which the minimum occurs, but assuming uniqueness, and letting $u = f(p, q, M)$ denote the point at which the minimum occurs, the Bayes rule for the no-sample problem is

$$\underline{a} = f(p, q, \alpha(\mathbb{R}))\text{th quantile } F_0.$$

For a sample of size n , the Bayes rule is

$$\hat{a}_n(X_1, \dots, X_n) = f(p, q, \alpha(\mathbb{R}) + n)\text{th quantile of } \hat{F}_n.$$

(h) *Tests of hypotheses involving quantiles.* Consider the problem of testing the hypothesis that the q th quantile t_q of an unknown distribution P on the real line does not exceed a given constant, taken without loss of generality to be zero. There are two actions available to the statistician, a_0 and a_1 , with loss functions

$$L(P, a_0) = w_0 I_{(0, \infty)}(t_q)$$

$$L(P, a_1) = w_1 I_{(-\infty, 0]}(t_q)$$

where w_0 and w_1 are positive constants. Suppose $P \in \mathcal{D}(\alpha)$. The Bayes rule for the no-sample problem is to select the action with smaller expected risk. This rule is: take action a_0 if

$$\mathcal{P}\{F(0) > q\} > \frac{w_0}{w_1 + w_0}$$

and take action a_1 otherwise. In terms of the function $u(p, q, M)$ defined by (13), we would take action a_0 if

$$u\left(\frac{w_0}{w_1 + w_0}, q, \alpha(\mathbb{R})\right) < F_0(0).$$

For a sample of size n , the Bayes rule is: take action a_0 if

$$u\left(\frac{w_0}{w_0 + w_1}, q, \alpha(\mathbb{R}) + n\right) < p_n F_0(0) + (1 - p_n) \frac{1}{n} W_n$$

where p_n is as in (4), and where W_n is the number of X_i less than or equal to zero. This is essentially the sign test.

Since $u(\frac{1}{2}, \frac{1}{2}, M) = \frac{1}{2}$ for all M , the above formula simplifies when testing for the median with $w_0 = w_1$. This test becomes: accept the hypothesis that the median does not exceed zero if

$$W_n > \frac{1}{2}n + \alpha(\mathbb{R})(\frac{1}{2} - F_0(0)).$$

Acknowledgments. The Dirichlet process was uncovered in the course of several long conversations with James B. MacQueen, who has an interesting way of viewing the Dirichlet process as a limit of Pólya urn schemes. The problem was suggested to me by David Blackwell, who is a great source of stimulating problems. Discussions with Charles Antoniak have helped clarify the ideas. To all three, my gratitude.

REFERENCES

[1] ANTONIAK, CHARLES (1969). Mixtures of Dirichlet processes with application to some non-parametric problems. Ph. D. Dissertation, Univ. of California, Los Angeles.
 [2] DUBINS, LESTER E. and FREEDMAN, DAVID A. (1966). Random distribution functions. *Proc. Fifth Berkeley Symp. Math. Statist. Prob.* **2** 183-214. Univ. of California Press.
 [3] FABIAN, J. (1964). Asymptotic behavior of Bayes' estimates. *Ann. Math. Statist.* **35** 846-856.
 [4] FERGUSON, THOMAS S. and KLASS, MICHAEL J. (1972). A representation of independent increment processes without Gaussian components. *Ann. Math. Statist.* **43** 1634-1643.
 [5] FREEDMAN, DAVID A. (1963). On the asymptotic behavior of Bayes' estimates in the discrete case. *Ann. Math. Statist.* **34** 1194-1216.

- [6] GNEDENKO, B. V. and KOLMOGOROV, A. N. (1949). *Limit Distributions for Sums of Independent Random Variables*, trans. K. L. Chung (1954). Addison-Wesley, Reading.
- [7] GOOD, I. J. (1965). *The Estimation of Probabilities: An Essay on Modern Bayesian Methods*. Research Monograph No. 30, M.I.T. Press.
- [8] KOLMOGOROV, A. N. (1933). *Foundations of the Theory of Probability*, 2nd ed., trans. Nathan Morrison (1956). Chelsea, New York.
- [9] KRAFT, CHARLES H. (1964). A class of distribution function processes which have derivatives. *J. Appl. Probability* **1** 385-388.
- [10] KRAFT, CHARLES H. and VAN EEDEN, CONSTANCE (1964). Bayesian bio-assay. *Ann. Math. Statist.* **35** 886-890.
- [11] WILKS, SAMUEL S. (1962). *Mathematical Statistics*. Wiley, New York.

DEPARTMENT OF MATHEMATICS
UNIVERSITY OF CALIFORNIA
LOS ANGELES, CALIFORNIA 90024