

Bayesian Methods for Function Estimation

Nidhan Choudhuri, Subhashis Ghosal and Anindya Roy

Abstract

Keywords: consistency; convergence rate; Dirichlet process; density estimation; Markov chain Monte Carlo; posterior distribution; regression function; spectral density; transition density

1. Introduction

Nonparametric and semiparametric statistical models are increasingly replacing parametric models, for the latter's lack of sufficient flexibility to address a wide variety of data. A nonparametric or semiparametric model involves at least one infinite-dimensional parameter, usually a function, and hence may also be referred to as an infinite-dimensional model. Functions of common interest, among many others, include the cumulative distribution function, density function, regression function, hazard rate, transition density of a Markov process, and spectral density of a time series. While frequentist methods for nonparametric estimation have been flourishing for many of these problems, nonparametric Bayesian estimation methods had been relatively less developed.

Besides philosophical reasons, there are some practical advantages of the Bayesian approach. On the one hand, the Bayesian approach allows one to reflect one's prior beliefs into the analysis. On the other hand, the Bayesian approach is straightforward in principle where inference is based on the posterior distribution only. Subjective elicitation of priors is relatively simple in a parametric framework, and in the absence of any concrete knowledge, there are many default mechanisms for prior specification. However, the recent popularity of Bayesian analysis comes from the availability of various Markov chain Monte Carlo (MCMC) algorithms that make the computation feasible with today's computers in almost every parametric problem. Prediction, which is sometimes the primary objective of a statistical analysis, is solved most naturally if one follows the Bayesian approach. Many non-Bayesian methods, including the maximum likelihood estimator (MLE), can have very unnatural behavior (such as staying on the boundary with high probability) when the parameter space is restricted, while

1 a Bayesian estimator does not suffer from this drawback. Besides, the optimality of a
2 parametric Bayesian procedure is often justified through large sample as well as finite
3 sample admissibility properties.

4 The difficulties for a Bayesian analysis in a nonparametric framework is threefold.
5 First, a subjective elicitation of a prior is not possible due to the vastness of the para-
6 meter space and the construction of a default prior becomes difficult mainly due to the
7 absence of the Lebesgue measure. Secondly, common MCMC techniques do not di-
8 rectly apply as the parameter space is infinite-dimensional. Sampling from the posterior
9 distribution often requires innovative MCMC algorithms that depend on the problem at
10 hand as well as the prior given on the functional parameter. Some of these techniques
11 include the introduction of latent variables, data augmentation and reparametrization of
12 the parameter space. Thus, the problem of prior elicitation cannot be separated from the
13 computational issues.

14 When a statistical method is developed, particular attention should be given to the
15 quality of the corresponding solution. Of the many different criteria, asymptotic con-
16 sistency and rate of convergence are perhaps among the least disputed. Consistency
17 may be thought of as a validation of the method used by the Bayesian. Consider an
18 imaginary experiment where an experimenter generates observations from a given sto-
19 chastic model with some value of the parameter and presents the data to a Bayesian
20 without revealing the true value of the parameter. If enough information is provided in
21 the form of a large number of observations, the Bayesian's assessment of the unknown
22 parameter should be close to the true value of it. Another reason to study consistency is
23 its relationship with robustness with respect to the choice of the prior. Due to the lack
24 of complete faith in the prior, we should require that at least eventually, the data over-
25 rides the prior opinion. Alternatively two Bayesians, with two different priors, presented
26 with the same data eventually must agree. This large sample "merging of opinions"
27 is equivalent to consistency (Blackwell and Dubins, 1962; Diaconis and Freedman,
28 1986a, 1986b; Ghosh et al., 1994). For virtually all finite-dimensional problems, the
29 posterior distribution is consistent (Ibragimov and Has'minskii, 1981; Le Cam, 1986;
30 Ghosal et al., 1995) if the prior does not rule out the true value. This is roughly a
31 consequence of the fact that the likelihood is highly peaked near the true value of the
32 parameter if the sample size is large. However, for infinite-dimensional problems, such
33 a conclusion is false (Freedman, 1963; Diaconis and Freedman, 1986a, 1986b; Doss,
34 1985a, 1985b; Kim and Lee, 2001). Thus posterior consistency must be verified before
35 using a prior.

36 In this chapter, we review Bayesian methods for some important curve estimation
37 problems. There are several good reviews available in the literature such as Hjort (1996,
38 2003), Wasserman (1998), Ghosal et al. (1999a), the monograph of Ghosh and Ra-
39 mamoorthi (2003) and several chapters in this volume. We omit many details which
40 may be found from these sources. We focus on three different aspects of the problem:
41 prior specification, computation and asymptotic properties of the posterior distribution.
42 In Section 2, we describe various priors on infinite-dimensional spaces. General results
43 on posterior consistency and rate of convergence are reviewed in Section 3. Specific
44 curve estimation problems are addressed in the subsequent sections.
45

2. Priors on infinite-dimensional spaces

A well accepted criterion for the choice of a nonparametric prior is that the prior has a large or full topological support. Intuitively, such a prior can reach every corner of the parameter space and thus can be expected to have a consistent posterior. More flexible models have higher complexity and hence the process of prior elicitation becomes more complex. Priors are usually constructed from the consideration of mathematical tractability, feasibility of computation, and good large sample behavior. The form of the prior is chosen according to some default mechanism while the key hyper-parameters are chosen to reflect any prior beliefs. A prior on a function space may be thought of as a stochastic process taking values in the given function space. Thus, a prior may be specified by describing a sampling scheme that generate random function with desired properties or they can be specified by describing the finite-dimensional laws. An advantage of the first approach is that the existence of the prior measure is automatic, while for the latter, the nontrivial proposition of existence needs to be established. Often the function space is approximated by a sequence of sieves in such a way that it is easier to put a prior on these sieves. A prior on the entire space is then described by letting the index of the sieve vary with the sample size, or by putting a further prior on the index thus leading to a hierarchical mixture prior. Here we describe some general methods of prior construction on function spaces.

2.1. Dirichlet process

Dirichlet processes were introduced by Ferguson (1973) as prior distributions on the space of probability measures on a given measurable space $(\mathcal{X}, \mathcal{B})$. Let $M > 0$ and G be a probability measure on $(\mathcal{X}, \mathcal{B})$. A Dirichlet process on $(\mathcal{X}, \mathcal{B})$ with parameters (M, G) is a random probability measure P which assigns a number $P(B)$ to every $B \in \mathcal{B}$ such that

- (i) $P(B)$ is a measurable $[0, 1]$ -valued random variable;
- (ii) each realization of P is a probability measure on $(\mathcal{X}, \mathcal{B})$;
- (iii) for each measurable finite partition $\{B_1, \dots, B_k\}$ of \mathcal{X} , the joint distribution of the vector $(P(B_1), \dots, P(B_k))$ on the k -dimensional unit simplex has Dirichlet distribution with parameters $(k; MG(B_1), \dots, MG(B_k))$.

(We follow the usual convention for the Dirichlet distribution that a component is a.s. 0 if the corresponding parameter is 0.) Using Kolmogorov's consistency theorem, Ferguson (1973) showed that a process with the stated properties exists. The argument could be made more elegant and transparent by using a countable generator of \mathcal{B} as in Blackwell (1973). The distribution of P is also uniquely defined by its specified finite-dimensional distributions in (iii) above. We shall denote the process by $\text{Dir}(M, G)$. If $(M_1, G_1) \neq (M_2, G_2)$ then the corresponding Dirichlet processes $\text{Dir}(M_1, G_1)$ and $\text{Dir}(M_2, G_2)$ are different, unless both G_1 and G_2 are degenerate at the same point. The parameter M is called the precision, G is called the center measure, and the product MG is called the base measure of the Dirichlet process. Note that

$$E(P(B)) = G(B), \quad \text{var}(P(B)) = \frac{G(B)(1 - G(B))}{1 + M}. \quad (2.1)$$

1 Therefore, if M is large, P is tightly concentrated about G justifying the terminology. 1
2 The relation (2.1) easily follows by the observation that each $P(B)$ is distributed as beta 2
3 with parameters $MG(B)$ and $M(1-G(B))$. By considering finite linear combinations of 3
4 indicator of sets and passing to the limit, it readily follows that (2.1) could be extended 4
5 to functions, that is, $E(\int \psi dP) = \int \psi dG$, and $\text{var}(\int \psi dP) = \text{var}_G(\psi)/(1+M)$. 5

6 As $P(A)$ is distributed as beta ($MG(A), MG(A^c)$), it follows that $P(A) > 0$ a.s. if 6
7 and only if $G(A) > 0$. However, this does not imply that P is a.s. mutually absolutely 7
8 continuous with G , as the null set could depend on A . As a matter of fact, the two 8
9 measures are often a.s. mutually singular. 9

10 If \mathfrak{X} is a separable metric space, the topological support of a measure on \mathfrak{X} and 10
11 the weak¹ topology on the space $\mathfrak{M}(\mathfrak{X})$ of all probability measures on \mathfrak{X} may be 11
12 defined. The support of $\text{Dir}(M, G)$ with respect to the weak topology is given by 12
13 $\{P \in \mathfrak{M}(\mathfrak{X}): \text{supp}(P) \subset \text{supp}(G)\}$. In particular, if the support of G is \mathfrak{X} , then the 13
14 support of $\text{Dir}(M, G)$ is the whole of $\mathfrak{M}(\mathfrak{X})$. Thus the Dirichlet process can be easily 14
15 chosen to be well spread over the space of probability measures. This may however 15
16 look apparently contradictory to the fact that a random P following $\text{Dir}(M, G)$ is a.s. 16
17 discrete. This important (but perhaps somewhat disappointing) property was observed 17
18 in Ferguson (1973) by using a gamma process representation of the Dirichlet process 18
19 and in Blackwell (1973) by using a Polya urn scheme representation. In the latter case, 19
20 the Dirichlet process arises as the mixing measure in de Finetti's representation in the 20
21 following continuous analogue of the Polya urn scheme: $X_1 \sim G$; for $i = 1, 2, \dots$, 21
22 $X_i = X_j$ with probability $1/(M+i-1)$ for $j = 1, \dots, i-1$ and $X_i \sim G$ with 22
23 probability $M/(M+i-1)$ independently of the other variables. This representation is 23
24 extremely crucial for MCMC sampling from a Dirichlet process. The representation also 24
25 shows that ties are expected among X_1, \dots, X_n . The expected number of distinct X 's, 25
26 as $n \rightarrow \infty$, is $M \log \frac{n}{M}$, which asymptotically much smaller than n . A simple proof of 26
27 a.s. discreteness of Dirichlet random measure, due to Savage, is given in Theorem 3.2.3 27
28 of Ghosh and Ramamoorthi (2003). 28

29 Sethuraman (1994) gave a constructive representation of the Dirichlet process. If 29
30 $\theta_1, \theta_2, \dots$ are i.i.d. G_0 , Y_1, Y_2, \dots are i.i.d. beta(1, M), $V_i = Y_i \prod_{j=1}^{i-1} (1 - Y_j)$ and 30
31

$$32 \quad P = \sum_{i=1}^{\infty} V_i \delta_{\theta_i}, \quad (2.2) \quad 32$$

33 then the above infinite series converges a.s. to a random probability measure that is dis- 33
34 tributed as $\text{Dir}(M, G)$. It may be noted that the masses V_i 's are obtained by successive 34
35 "stick-breaking" with Y_1, Y_2, \dots as the corresponding stick-breaking proportions, and 35
36 allotted to randomly chosen points $\theta_1, \theta_2, \dots$ generated from G . Sethuraman's repre- 36
37 sentation has made it possible to use the Dirichlet process in many complex problem 37
38 using some truncation and Monte Carlo algorithms. Approximations of this type are 38
39 discussed by Muliere and Tardella (1998) and Iswaran and Zarepour (2002a, 2002b). 39
40 Another consequence of the Sethuraman representation is that if $P \sim \text{Dir}(M, G)$, 40
41 $\theta \sim G$ and $Y \sim \text{beta}(1, M)$, and all of them are independent, then $Y\delta_{\theta} + (1-Y)P$ also 41
42 has $\text{Dir}(M, G)$ distribution. This property leads to important distributional equations for 42
43 43
44 44

45 ¹ What we call weak is termed as weak star in functional analysis. 45

1 functionals of the Dirichlet process, and could also be used to simulate a Markov chain 1
2 on $\mathfrak{M}(\mathfrak{X})$ with $\text{Dir}(M, G)$ as its stationary distribution. 2

3 The Dirichlet process has a very important conditioning property. If A is set with 3
4 $G(A) > 0$ (which implies that $P(A) > 0$ a.s.), then the random measure $P|_A$, the 4
5 restriction of P to A defined by $P|_A(B) = P(B|A) = P(B \cap A)/P(A)$, is distributed 5
6 as Dirichlet process with parameters $MG(A)$ and $G|_A$ and is independent of $P(A)$. The 6
7 argument can be extended to more than one set. Thus the Dirichlet process locally splits 7
8 into numerous independent Dirichlet processes. 8

9 A peculiar property of the Dirichlet process is that any two Dirichlet processes 9
10 $\text{Dir}(M_1, G_1)$ and $\text{Dir}(M_2, G_2)$ are mutually singular if G_1, G_2 are nonatomic and 10
11 $(M_1, G_1) \neq (M_2, G_2)$. 11

12 The distribution of a random mean functional $\int \psi dP$, where ψ is a measurable func- 12
13 tion, is of some interest. Although, $\int \psi dP$ has finite mean if and only if $\int |\psi| dG < \infty$, 13
14 P has a significantly shorter tail than that of G . For instance, the random P generated 14
15 by a Dirichlet process with Cauchy base measure has all moments. Distributions of 15
16 the random mean functional has been studied in many articles including Cifarelli and 16
17 Regazzini (1990) and Regazzini et al. (2002). Interestingly the distribution of $\int x dP(x)$ 17
18 is G if and only if G is Cauchy. 18

19 The behavior of the tail probabilities of a random P obtained from a Dirichlet process 19
20 is important for various purposes. Fristedt (1967) and Fristedt and Pruitt (1971) char- 20
21 acterized the growth rate of a gamma process and using their result, Doss and Sellke 21
22 (1982) obtained analogous results for the tail probabilities of P . 22

23 Weak convergence properties of the Dirichlet process are controlled by the conver- 23
24 gence of its parameters. Let G_n weakly converge to G . Then 24

- 25 (i) if $M_n \rightarrow M > 0$, then $\text{Dir}(M_n, G_n)$ converges weakly to $\text{Dir}(M, G)$; 25
26 (ii) if $M_n \rightarrow 0$, then $\text{Dir}(M_n, G_n)$ converges weakly to a measure degenerated at a 26
27 random $\theta \sim G$; 27
28 (iii) if $M_n \rightarrow \infty$, then $\text{Dir}(M_n, G_n)$ converges weakly to random measure degenerate 28
29 at G . 29
30 30

31 2.2. Processes derived from the Dirichlet process 31

32 2.2.1. Mixtures of Dirichlet processes 32

33 The mixture of Dirichlet processes was introduced by Antoniak (1974). While eliciting 33
34 the base measure using (2.1), it may be reasonable to guess that the prior mean measure 34
35 is normal, but it may be difficult to specify the values of the mean and the variance of 35
36 this normal distribution. It therefore makes sense to put a prior on the mean and the 36
37 variance. More generally, one may propose a parametric family as the base measure 37
38 and put hyper-priors on the parameters of that family. The resulting procedure has an 38
39 intuitive appeal in that if one is a weak believer in a parametric family, then instead 39
40 of using a parametric analysis, one may use the corresponding mixture of Dirichlet 40
41 to robustify the parametric procedure. More formally, we may write the hierarchical 41
42 Bayesian model $P \sim \text{Dir}(M_\theta, G_\theta)$, where the indexing parameter $\theta \sim \pi$. 42
43 43

44 In semiparametric problems, mixtures of Dirichlet priors appear if the nonparametric 44
45 part is given a Dirichlet process. In this case, the interest is usually in the posterior 45

1 distribution of the parametric part, which has a role much bigger than that of an indexing 1
2 parameter. 2

3
4 2.2.2. *Dirichlet mixtures* 4

5 Although the Dirichlet process cannot be used as a prior for estimating a density, convo- 5
6 luting it with a kernel will produce smooth densities. Such an approach was pioneered 6
7 by Ferguson (1983) and Lo (1984). Let Θ be a parameter set, typically a Euclid- 7
8 ean space. For each θ , let $\psi(x, \theta)$ be a probability density function. A nonparametric 8
9 mixture of $\psi(x, \theta)$ is obtained by considering $p_F(x) = \int \psi(x, \theta) dF(\theta)$. These mix- 9
10 tures can form a very rich family. For instance, the location and scale mixture of the 10
11 form $\sigma^{-1}k((x - \mu)/\sigma)$, for some fixed density k , may approximate any density in the 11
12 L_1 -sense if σ is allowed to approach to 0. Thus, a prior on densities may be induced by 12
13 putting a Dirichlet process prior on the mixing distribution F and a prior on σ . 13

14 The choice of an appropriate kernel depends on the underlying sample space. If the 14
15 underlying density function is defined on the entire real line, a location-scale kernel 15
16 is appropriate. If on the unit interval, beta distributions form a flexible two parameter 16
17 family. If on the positive half line, mixtures of gamma, Weibull or lognormal may be 17
18 used. The use of a uniform kernel leads to random histograms. Petrone and Veronese 18
19 (2002) motivated a canonical way of viewing the choice of a kernel through the notion 19
20 of the Feller sampling scheme, and call the resulting prior a Feller prior. 20
21

22
23 2.2.3. *Invariant Dirichlet process* 23

24 The invariant Dirichlet process was considered by Dalal (1979). Suppose that we want to 24
25 put a prior on the space of all probability measures symmetric about zero. One may let P 25
26 follow $\text{Dir}(M, G)$ and put $\bar{P}(A) = (P(A) + P(-A))/2$, where $-A = \{x: -x \in A\}$.² 26
27 More generally, one can consider a compact group \mathfrak{G} acting on the sample space \mathfrak{X} 27
28 and consider the distribution of \bar{P} as the invariant Dirichlet process where $\bar{P}(A) =$ 28
29 $\int P(gA) d\mu(g)$, μ stands for the Haar probability measure on \mathfrak{G} and P follows the 29
30 Dirichlet process. 30

31 The technique is particularly helpful for constructing priors on the error distribu- 31
32 tion F for the location problem $X = \theta + \epsilon$. The problem is not identifiable without 32
33 some restriction on F , and symmetry about zero is a reasonable condition on F ensur- 33
34 ing identifiability. The symmetrized Dirichlet process prior was used by Diaconis and 34
35 Freedman (1986a, 1986b) to present a striking example of inconsistency of the posterior 35
36 distribution. 36
37

38 2.2.4. *Pinned-down Dirichlet* 38

39 If $\{B_1, \dots, B_k\}$ is a finite partition, called control sets, then the conditional distribution 39
40 of P given $\{P(B_j) = w_j, j = 1, \dots, k\}$, where P follows $\text{Dir}(M, G)$ and $w_j \geq 0$, 40
41 $\sum_{j=1}^k w_j = 1$, is called a pinned-down Dirichlet process. By the conditioning prop- 41
42 erty of the Dirichlet process mentioned in the last subsection, it follows that the above 42
43

44 ² Another way of randomly generating symmetric probabilities is to consider a Dirichlet process P on 44
45 $[0, \infty)$ and unfold it to \tilde{P} on \mathbb{R} by $\tilde{P}(-A) = \tilde{P}(A) = \frac{1}{2}P(A)$. 45

1 process may be written as $P = \sum_{j=1}^k w_j P_j$, where each P_j is a Dirichlet process on B_j . 1
2 Consequently P is a countable mixture of Dirichlet (with orthogonal supports). 2

3 A particular case of pinned-down Dirichlet is obtained when one puts the restriction 3
4 that P has median 0. Doss (1985a, 1985b) used this idea to put a prior for the semipara- 4
5 metric location problem and showed an inconsistency result similar to Diaconis and 5
6 Freedman (1986a, 1986b) mentioned above. 6

7 2.3. Generalizations of the Dirichlet process 7

8 While the Dirichlet process is a prior with many fascinating properties, its reliance on 8
9 only two parameters may sometimes be restrictive. One drawback of Dirichlet process 9
10 is that it always produces discrete random probability measures. Another property of 10
11 Dirichlet which is sometimes embarrassing is that the correlation between the random 11
12 probabilities of two sets is always negative. Often, random probabilities of sets that 12
13 are close enough are expected to be positively related if some smoothness is present. 13
14 More flexible priors may be constructed by generalizing the way the prior probabilities 14
15 are assigned. Below we discuss some of the important generalizations of a Dirichlet 15
16 process. 16
17

18 2.3.1. Tail-free and neutral to the right process 18

19 The concept of a tail-free process was introduced by Freedman (1963) and chronolog- 19
20 ically precedes that of the Dirichlet process. A tail-free process is defined by random 20
21 allocations of probabilities to sets in a nested sequence of partitions. Let $E = \{0, 1\}$ 21
22 and E^m be the m -fold Cartesian product $E \times \dots \times E$ where $E^0 = \emptyset$. Further, set 22
23 $E^* = \bigcup_{m=0}^{\infty} E^m$. Let $\pi_0 = \{\mathcal{X}\}$ and for each $m = 1, 2, \dots$, let $\pi_m = \{B_\varepsilon: \varepsilon \in E^m\}$ 23
24 be a partition of \mathcal{X} so that sets of π_{m+1} are obtained from a binary split of the sets 24
25 of π_m and $\bigcup_{m=0}^{\infty} \pi_m$ be a generator for the Borel sigma-field on \mathbb{R} . A probabil- 25
26 ity P may then be described by specifying all the conditional probabilities $\{V_\varepsilon =$ 26
27 $P(B_{\varepsilon 0}|B_\varepsilon): \varepsilon \in E^*\}$. A prior for P may thus be defined by specifying the joint dis- 27
28 tribution of all V_ε 's. The specification may be written in a tree form. The different 28
29 hierarchy in the tree signifies prior specification of different levels. A prior for P is 29
30 said to be tail-free with respect to the sequence of partitions $\{\pi_m\}$ if the collections 30
31 $\{V_\emptyset\}, \{V_0, V_1\}, \{V_{00}, V_{01}, V_{10}, V_{11}\}, \dots$, are mutually independent. Note that, variables 31
32 within the same hierarchy need not be independent; only the variables at different lev- 32
33 els are required to be so. Partitions more general than binary partitions could be used, 33
34 although that will not lead to more general priors. 34
35

36 A Dirichlet process is tail-free with respect to any sequence of partitions. Indeed, the 36
37 Dirichlet process is the only prior that has this distinguished property; see Ferguson 37
38 (1974) and the references therein. Tail-free priors satisfy some interesting zero-one 38
39 laws, namely, the random measure generated by a tail-free process is absolutely con- 39
40 tinuous with respect to a given finite measure with probability zero or one. This follows 40
41 from the fact that the criterion of absolute continuity may be expressed as tail event 41
42 with respect to a collection of independent random variables and Kolmogorov's zero- 42
43 one law may be applied; see Ghosh and Ramamoorthi (2003) for details. Kraft (1964) 43
44 gave a very useful sufficient condition for the almost sure absolute continuity of a tail- 44
45 free process. 45

1 Neutral to the right processes, introduced by Doksum (1974), are also tail-free 1
2 processes, but the concept is applicable only to survival distribution functions. If F is a 2
3 random distribution function on the positive half line, then F is said to follow a neutral 3
4 to the right process if for every k and $0 < t_1 < \dots < t_k$, there exists independent ran- 4
5 dom variables V_1, \dots, V_k such that the joint distribution of $(1 - F(t_1), 1 - F(t_2), \dots,$ 5
6 $1 - F(t_k))$ is same as that of the successive products $(V_1, V_1 V_2, \dots, \prod_{j=1}^k V_j)$. Thus 6
7 a neutral to the right prior is obtained by stick breaking. Clearly the process is tail- 7
8 free with respect to the nested sequence $\{[0, t_1], (t_1, \infty)\}, \{[0, t_1], (t_1, t_2], (t_2, \infty)\}, \dots$ 8
9 of partitions. Note that $F(x)$ may be written as $e^{-H(x)}$, where $H(\cdot)$ is a process of 9
10 independent increments. 10

11
12 2.3.2. *Polya tree process* 12

13 A Polya tree process is a special case of a tail-free process, where besides across row 13
14 independence, the random conditional probabilities are also independent within row 14
15 and have beta distributions. To elaborate, let $\{\pi_m\}$ be a sequence of binary partition 15
16 as before and $\{\alpha_\varepsilon: \varepsilon \in E^*\}$ be a collection of nonnegative numbers. A random prob- 16
17 ability measure P on \mathbb{R} is said to possess a Polya tree distribution with parameters 17
18 $(\{\pi_m\}, \{\alpha_\varepsilon: \varepsilon \in E^*\})$, if there exist a collection $\mathcal{Y} = \{Y_\varepsilon: \varepsilon \in E^*\}$ of random variables 18
19 such that the following hold: 19

- 20 (i) The collection \mathcal{Y} consists of mutually independent random variables; 20
21 (ii) For each $\varepsilon \in E^*$, Y_ε has a beta distribution with parameters $\alpha_{\varepsilon 0}$ and $\alpha_{\varepsilon 1}$; 21
22 (iii) The random probability measure P is related to \mathcal{Y} through the relations 22
23

24
25
$$P(B_{\varepsilon_1 \dots \varepsilon_m}) = \left(\prod_{j=1; \varepsilon_j=0}^m Y_{\varepsilon_1 \dots \varepsilon_{j-1}} \right) \left(\prod_{j=1; \varepsilon_j=1}^m (1 - Y_{\varepsilon_1 \dots \varepsilon_{j-1}}) \right),$$
 25
26
27
$$m = 1, 2, \dots, \quad 27$$

28 where the factors are Y_\emptyset or $1 - Y_\emptyset$ if $j = 1$. 28
29

30 The concept of a Polya tree was originally considered by Ferguson (1974) and 30
31 Blackwell and MacQueen (1973), and later studied thoroughly by Mauldin et al. (1992) 31
32 and Lavine (1992, 1994). The prior can be seen as arising as the de Finetti measure in a 32
33 generalized Polya urn scheme; see Mauldin et al. (1992) for details. 33

34 The class of Polya trees contain all Dirichlet processes, characterized by the relation 34
35 that $\alpha_{\varepsilon 0} + \alpha_{\varepsilon 1} = \alpha_\varepsilon$ for all ε . A Polya tree can be chosen to generate only absolutely 35
36 continuous distributions. The prior expectation of the process could be easily written 36
37 down; see Lavine (1992) for details. Below we consider an important special case for 37
38 discussion, which is most relevant for statistical use. Consider \mathfrak{X} to be a subset of the 38
39 real line and let G be a probability measure. Let the partitions be obtained successively 39
40 by splitting the line at the median, the quartiles, the octiles, and in general, binary quan- 40
41 tiles of G . If $\alpha_{\varepsilon 0} = \alpha_{\varepsilon 1}$ for all $\varepsilon \in E^*$, then it follows that $E(P) = G$. Thus G will 41
42 have the role similar to that of the center measure of a Dirichlet process, and hence 42
43 will be relatively easy to elicit. Besides, the Polya tree will have infinitely many more 43
44 parameters which may be used to describe one's prior belief. Often, to avoid specifying 44
45 too many parameters, a default method is adopted, where one chooses α_ε depending 45

1 only on the length of the finite string ε . Let a_m stand for the value of α_ε when ε has
2 length m . The growth rate of a_m controls the smoothness of the Polya tree process. For
3 instance, if $a_m = c2^{-m}$, we obtain the Dirichlet process, which generate discrete prob-
4 abilities. If $\sum_{m=1}^{\infty} a_m^{-1} < \infty$ (for instance, if $a_m = cm^2$), then it follows from Kraft's
5 (1964) result that the random P is absolutely continuous with respect to G . The choice
6 $a_m = c$ leads to singular continuous distributions almost surely; see Ferguson (1974).
7 This could guide one to choose the sequence a_m . For smoothness, one should choose
8 rapidly growing a_m . One may actually like to choose according to one's prior belief
9 in the beginning of the tree deviating from the above default choice, and let a default
10 method choose the parameters at the later stages where practically no prior information
11 is available. An extreme form of this will lead to partially specified Polya trees, where
12 one chooses a_m to be infinity after a certain stage (which is equivalent to uniformly
13 spreading the mass inside a given interval).

14 Although the prior mean distribution function may have a smooth Lebesgue density,
15 the randomly sampled densities from a Polya tree are very rough, being nowhere dif-
16 ferentiable. To overcome this difficulty, mixtures of a Polya tree, where the partitioning
17 measure G involves some additional parameter θ with some prior, may be considered.
18 The additional parameter will average out jumps to yield smooth densities; see Hanson
19 and Johnson (2002). However, then the tail-freeness is lost and the resulting posterior
20 distribution could be inconsistent. Berger and Guglielmi (2001) considered a mixture
21 where the partition remains fixed and the α -parameters depend on θ , and applied the
22 resulting prior to a model selection problem.

23 2.3.3. Generalized Dirichlet process

24 The k -dimensional Dirichlet distribution may be viewed as the conditional distribution
25 of (p_1, \dots, p_k) given that $\sum_{j=1}^k p_j = 1$, where $p_j = e^{-Y_j}$ and Y_j 's are inde-
26 pendent exponential variables. In general, if Y_j 's have a joint density $h(y_1, \dots, y_k)$,
27 the conditional joint density of (p_1, \dots, p_{k-1}) is proportional to $h(-\log p_1, \dots,$
28 $-\log p_k) p_k^{-1} \cdots p_k^{-1}$, where $p_k = 1 - \sum_{j=1}^{k-1} p_j$. Hjort (1996) considered the joint
29 density of Y_j 's to be proportional to $\prod_{j=1}^k e^{-\alpha_j y_j} g_0(y_1, \dots, y_k)$, and hence the resulting
30 (conditional) density of p_1, \dots, p_{k-1} is proportional to $p_1^{\alpha_1-1} \cdots p_k^{\alpha_k-1} g(p_1, \dots, p_k)$,
31 where $g(p_1, \dots, p_k) = g_0(-\log p_1, \dots, -\log p_k)$. We may put $g(p) = e^{-\lambda \Delta(p)}$,
32 where $\Delta(p)$ is a penalty term for roughness such as $\sum_{j=1}^{k-1} (p_{j+1} - p_j)^2$, $\sum_{j=2}^{k-1} (p_{j+1} -$
33 $2p_j + p_{j-1})^2$ or $\sum_{j=1}^{k-1} (\log p_{j+1} - \log p_j)^2$. The penalty term helps maintain posi-
34 tive correlation and hence "smoothness". The tuning parameter λ controls the extent to
35 which penalty is imposed for roughness. The resulting posterior distribution is conju-
36 gate with mode equivalent to a penalized MLE. Combined with random histogram or
37 passing through the limit as the bin width goes to 0, the technique could also be applied
38 to continuous data.

39 2.3.4. Priors obtained from random series representation

40 Sethuraman's (1994) infinite series representation creates a lot of possibilities of gen-
41 eralizing the Dirichlet process by changing the distribution of the weights, the support
42 points, or even the number of terms. Consider a random probability measure given by
43
44
45

1 $P = \sum_{i=1}^N V_i \delta_{\theta_i}$, where $1 \leq N \leq \infty$, $\sum_{i=1}^N V_i = 1$ and N may be given a fur- 1
2 ther prior distribution. Note that the resulting random probability measure is almost 2
3 surely discrete. Choosing $N = \infty$, θ_i 's as i.i.d. G as in the Sethuraman representation, 3
4 $V_i = Y_i \prod_{j=1}^{i-1} (1 - Y_j)$, where Y_1, Y_2, \dots are i.i.d. beta(a, b), Hjort (2000) obtained an 4
5 interesting generalization of the Dirichlet process. The resulting process admits, as in 5
6 the case of a Dirichlet process, explicit formulae for the posterior mean and variance of 6
7 a mean functional. 7

8 From computational point of view, a prior is more tractable if N is chosen to be finite. 8
9 To be able to achieve reasonable large sample properties, either N has to depend on the 9
10 sample size n , or N must be given a prior which is infinitely supported. Given $N = k$, 10
11 the prior on (V_1, \dots, V_k) is taken to be k -dimensional Dirichlet distribution with param- 11
12 eters $(\alpha_{1,n}, \dots, \alpha_{k,n})$. The parameters θ_i 's are usually chosen as in the Sethuraman's 12
13 representation, that is i.i.d. G . Iswaran and Zarepour (2002a) studied convergence prop- 13
14 erties of these random measures. For the choice $\alpha_{j,k} = M/k$, the limiting measure is 14
15 Dir(M, G). However, the commonly advocated choice $\alpha_{j,k} = M$ leads essentially to a 15
16 parametric prior, and hence to an inconsistent posterior. 16

17 2.4. Gaussian process 17

18 Considered first by Leonard (1978), and then by Lenk (1988, 1991) in the context of 18
19 density estimation, a Gaussian process may be used in a wider generality because of 19
20 its ability to produce arbitrary shapes. The method may be applied to nonparametric 20
21 regression where only smoothness is assumed for the regression function. The mean 21
22 function reflects any prior belief while the covariance kernel may be tuned to control the 22
23 smoothness of the sample paths as well as to reflect the confidence in the prior guess. 23
24 In a generalized regression, where the function of interest has restricted range, a link 24
25 function is used to map the unrestricted range of the Gaussian process to the desired 25
26 one. A commonly used Gaussian process in the regression context is the integrated 26
27 Wiener process with some random intercept term as in Wahba (1978). Choudhuri et al. 27
28 (2004b) used a general Gaussian process prior for binary regression. 28
29 30

31 2.5. Independent increment process 31

32 Suppose that we want to put a prior on survival distribution functions, that is, dis- 32
33 tribution functions on the positive half line. Let $Z(t)$ be a process with independent 33
34 nonnegative increments such that $Z(\infty)$, the total mass of Z , is a.s. finite. Then a prior 34
35 on F may be constructed by the relation $F(t) = Z(t)/Z(\infty)$. Such a prior is necessarily 35
36 neutral to the right. When $Z(t)$ is the gamma process, that is an independent increment 36
37 process with $Z(t) \sim \text{gamma}(MG(t), 1)$, then the resulting distribution of P is Dirichlet 37
38 process Dir(M, G). 38
39

40 For estimating a survival function, it is often easier to work with the cumulative haz- 40
41 ard function, which needs only be positive. If $Z(t)$ is a process such that $Z(\infty) = \infty$ 41
42 a.s., then $F(t) = 1 - e^{-Z(t)}$ is a distribution function. The process $Z(t)$ may be 42
43 characterized in terms of its Lévy measure $N_t(\cdot)$, and is called a Lévy process. Un- 43
44 fortunately, as $Z(t)$ necessarily increases by jumps only, $Z(t)$ is not the cumulative 44
45 hazard function corresponding to $F(t)$. Instead, one may define $F(t)$ by the relation 45

1 $Z(t) = \int_0^t dF(s)/(1 - F(s-))$. The expressions of prior mean and variance, and posterior updating are relatively straightforward in terms of the Lévy measure; see Hjort
2 (1990) and Kim (1999). Particular choices of the Lévy measure lead to special priors
3 such as the Dirichlet process, completely homogeneous process (Ferguson and Phadia,
4 1979), gamma process (Lo, 1982), beta process (Hjort, 1990), beta-Stacy process
5 (Walker and Muliere, 1997) and extended beta process (Kim and Lee, 2001). Kim and
6 Lee (2001) settled the issue of consistency, and provided an interesting example of in-
7 consistency.
8

9 A disadvantage of modeling the process $Z(t)$ is that the resulting F is discrete.
10 Dykstra and Laud (1981) considered a Lévy process to model the hazard rate. How-
11 ever, this approach leads only to monotone hazard functions. Nieto-Barajas and Walker
12 (2004) replaced the independent increments process by a Markov process and obtained
13 continuous sample paths.
14

15 2.6. *Some other processes*

16 One approach to putting a prior on a function space is to decompose a function into a
17 basis expansion of the form $\sum_{j=1}^{\infty} b_j \psi_j(\cdot)$ for some fixed basis functions and then putting
18 priors on b_j 's. An orthogonal basis is very useful if the function space of interest is a
19 Hilbert space. Various popular choices of such basis include polynomials, trigonometric
20 functions, splines and wavelets among many others. If the coefficients are unrestricted,
21 independent normal priors may be used. Interestingly, when the coefficients are normally
22 distributed, the prior on the random function is a Gaussian process. Conversely, a
23 Gaussian process may be represented in this way by virtue of the Karhunen–Loévé ex-
24 pansion. When the function values are restricted, transformations should be used prior
25 to a basis expansion. For instance, for a density function, an expansion should be raised
26 to the exponential and then normalized. Barron et al. (1999) used polynomials to con-
27 struct an infinite-dimensional exponential family. Hjort (1996) discussed a prior on a
28 density induced by the Hermite polynomial expansion and a prior on the sequence of
29 cumulants.
30

31 Instead of considering an infinite series representation, one may consider a series
32 based on the first k terms, where k is deterministically increased to infinity with the
33 sample size, or is itself given a prior that has infinite support. The span of the first k
34 functions, as k tends to infinity, form approximating sieves in the sense of Grenander
35 (1981). The resulting priors are recommended as default priors in infinite-dimensional
36 spaces by Ghosal et al. (1997). In Ghosal et al. (2000), this idea was used with a spline
37 basis for density estimation. They showed that with a suitable choice of k , depending
38 on the sample size and the smoothness level of the target function, optimal convergence
39 rates could be obtained.

40 If the domain is a bounded interval then the sequence of moments uniquely deter-
41 mines the probability measure. Hence a prior on the space of probability measures could
42 be induced from that on the sequence of moments. One may control the location, scale,
43 skewness and kurtosis of the random probability by using subjective priors on the first
44 four moments. Priors for the higher-order moments are difficult to elicit, and some de-
45 fault method should be used.

1 Priors for quantiles are much easier to elicit than that for moments. One may put priors 1
2 on all dyadic quantiles honoring the order restrictions. Conceptually, this operation 2
3 is opposite to that of specifying a tree based prior such as the Polya tree or a tail- 3
4 free process. Here masses are predetermined and the partitions are chosen randomly. In 4
5 practice, one may put priors only for a finite number of quantiles, and then distribute the 5
6 remaining masses uniformly over the corresponding interval. Interestingly, if the prior 6
7 on the quantile process is induced from a Dirichlet process on the random probability, 7
8 then the posterior expectation of a quantile (in the noninformative limit $M \rightarrow 0$) is seen 8
9 to be a Bernstein polynomial smoother of the empirical quantile process. This leads to 9
10 a quantile density estimator, which, upon inversion, leads to an automatically smoothed 10
11 empirical density estimator; see Hjort (1996) for more details. 11

14 3. Consistency and rates of convergence 14

15 Let $\{(\mathcal{X}^{(n)}, \mathcal{A}^{(n)}, P_{\theta}^{(n)}): \theta \in \Theta\}$ be a sequence of statistical experiments with obser- 16
17 vations $X^{(n)}$, where the parameter set Θ is an arbitrary topological space and n is an 17
18 indexing parameter, usually the sample size. Let \mathcal{B} be the Borel sigma-field on Θ and 18
19 Π_n be a probability measure on (Θ, \mathcal{B}) , which, in general, may depend on n . The pos- 19
20 terior distribution is defined to be a version of the regular conditional probability of θ 20
21 given $X^{(n)}$, and is denoted by $\Pi_n(\cdot|X^{(n)})$. 21

22 Let $\theta_0 \in \Theta$. We say that the posterior distribution is consistent at θ_0 (with respect 22
23 to the given topology on Θ) if $\Pi_n(\cdot|X^{(n)})$ converges weakly to δ_{θ_0} as $n \rightarrow \infty$ under 23
24 $P_{\theta_0}^{(n)}$ -probability, or almost surely under the distribution induced by the parameter 24
25 value θ_0 . If the latter makes sense, it is a more appealing concept. 25
26

27 The above condition (in the almost sure sense) is equivalent to checking that ex- 27
28 cept on a θ_0 -induced null set of sample sequences, for any neighborhood U of θ_0 , 28
29 $\Pi_n(U^c|X^{(n)}) \rightarrow 0$. If the topology on Θ is countably generated (as in the case of a 29
30 separable metric space), this reduces to $\Pi_n(U^c|X^{(n)}) \rightarrow 0$ a.s. under the distribution 30
31 induced by θ_0 for every neighborhood U . An analogous conclusion holds for consis- 31
32 tency in probability. Henceforth we work with the second formulation. 32

33 Consistency may be motivated as follows. A (prior or posterior) distribution stands 33
34 for one's knowledge about the parameter. Perfect knowledge implies a degenerate prior. 34
35 Thus consistency means weak convergence of knowledge towards the perfect knowl- 35
36 edge as the amount of data increases. 36

37 Doob (1948) obtained a very general result on posterior consistency. Let the prior 37
38 Π be fixed and the observations be i.i.d. Under some mild measurability conditions 38
39 on the sample space (a standard Borel space will suffice) and model identifiability, 39
40 Doob (1948) showed that the set of all $\theta \in \Theta$ where consistency does not hold is 40
41 Π -null. This follows by the convergence of the martingale $EI(\theta \in B|X_1, \dots, X_n)$ 41
42 to $EI(\theta \in B|X_1, X_2, \dots) = I(\theta \in B)$. The condition of i.i.d. observations could 42
43 be replaced by the assumption that in the product space $\Theta \times \mathcal{X}^{\infty}$, the parameter θ is 43
44 \mathcal{A}^{∞} -measurable. Statistically speaking, the condition holds if there is a consistent esti- 44
45 mate of some bimeasurable function of θ . 45

1 The above result should not however create a false sense of satisfaction as the Π -null 1
2 set could be very large. It is important to know at which parameter values consistency 2
3 holds. Indeed, barring a countable parameter space, Doob's (1948) is of little help. On 3
4 the other hand, Doob's (1948) theorem implies that consistency holds at a parameter 4
5 point whenever there is a prior point mass there. 5

6 Freedman (1963) showed that merely having positive Π -probability in a neighbor- 6
7 hood of θ_0 does not imply consistency at that point. 7

8
9 EXAMPLE 1. Let $\Theta = \mathfrak{M}(\mathbb{Z}_+)$, the space of all discrete distribution on positive integers 9
10 with the total variation distance on Θ . Let θ_0 be the geometric distribution with 10
11 parameter 1/4. There exists a prior Π such that every neighborhood of θ_0 has positive 11
12 probability under Π , yet 12

$$13 \quad \Pi(\theta \in U | X_1, \dots, X_n) \rightarrow 1 \quad \text{a.s. } [\theta_0^\infty], \quad (3.1) \quad 13$$

14 where U is any neighborhood of θ_1 , the geometric distribution with parameter 3/4. 14
15
16

17 Indeed, the following result of Freedman (1963) shows that the above example of 17
18 inconsistency is somewhat generic in a topological sense. 18

19
20 THEOREM 1. Let $\Theta = \mathfrak{M}(\mathbb{Z}_+)$ with the total variation distance on it, and let $\mathfrak{M}(\Theta)$ 20
21 be the space of all priors on Θ with the weak topology. Put the product topology on 21
22 $\Theta \times \mathfrak{M}(\Theta)$. Then 22

$$23 \quad \left\{ (\theta, \Pi) \in \Theta \times \mathfrak{M}(\Theta) : \limsup_{n \rightarrow \infty} \Pi(\theta \in U | X_1, \dots, X_n) = 1 \right. \\
24 \quad \left. \forall U \text{ open, } U \neq \emptyset \right\} \quad (3.2) \quad 24$$

25
26
27 is the complement of a meager set.³ 27
28

29 Thus, Freedman's (1963) result tells us that except for a relatively small collection 29
30 of pairs of (θ, Π) , the posterior distribution wanders aimlessly around the parameter 30
31 space. In particular, consistency will not hold at any given θ . While this result cau- 31
32 tions us about naive uses of Bayesian methods, it does not mean that Bayesian methods 32
33 are useless. Indeed, a pragmatic Bayesian's only aim might be to just be able to find 33
34 a reasonable prior complying with one's subjective belief (if available) and obtaining 34
35 consistency at various parameter values. There could be plenty of such priors available 35
36 even though there will be many more that are not appropriate. The situation may be 36
37 compared with the role of differentiable functions among the class of all continuous 37
38 functions. Functions that are differentiable at some point form a small set in the same 38
39 sense while nowhere differentiable functions are much more abundant. 39
40

41 From a pragmatic point of view, useful sufficient conditions ensuring consistency at 41
42 a given point is the most important proposition. Freedman (1963, 1965) showed that for 42
43 estimation of a probability measure, if the prior distribution is tail-free, then (a suitable 43

44 ³ A meager set is one which can be written as a countable union of closed sets without any interior points, 44
45 and is considered to be topologically small. 45

1 version of) the posterior distribution is consistent at any point with respect to the weak 1
2 topology. The idea behind this result is reducing every weak neighborhood to a Euclid- 2
3 ean neighborhood in some finite-dimensional projection using the tail-free property. 3

4 Schwartz (1965), in a celebrated paper, obtained a general result on consistency. 4
5 Schwartz's (1965) theorem requires a testing condition and a condition on the support 5
6 of the prior. 6

7 Consider i.i.d. observations generated by a statistical model indexed by an abstract 7
8 parameter space Θ admitting a density $p(x, \theta)$ with respect to some sigma-finite mea- 8
9 sure μ . Let $K(\theta_1, \theta_2)$ denote the Kullback–Leibler divergence $\int p(x, \theta_1) \log(p(x, \theta_1)/$ 9
10 $p(x, \theta_2)) d\mu(x)$. We say that $\theta_0 \in \Theta$ is in the Kullback–Leibler support of Π , we write 10
11 $\theta_0 \in \text{KL}(\Pi)$, if for every $\varepsilon > 0$, $\Pi\{\theta: K(\theta_0, \theta) < \varepsilon\}$. As the Kullback–Leibler diver- 11
12 gence is asymmetric and not a metric, the support may not be interpreted in a topological 12
13 sense. Indeed, a prior may have empty Kullback–Leibler support even on a separable 13
14 metric space. 14

15
16 **THEOREM 2.** *Let $\theta_0 \in U \subset \Theta$. If there exists $m \geq 1$, a test function $\phi(X_1, \dots, X_m)$ 16
17 for testing $H_0: \theta = \theta_0$ against $H: \theta \in U^c$ with the property that $\inf\{E_{\theta_0}\phi(X_1, \dots, X_m):$ 17
18 $\theta \in U^c\} > E_{\theta_0}\phi(X_1, \dots, X_m)$ and $\theta_0 \in \text{KL}(\Pi)$, then $\Pi\{\theta \in U^c | X_1, \dots, X_n\} \rightarrow 0$ 18
19 a.s. $[P_{\theta_0}^\infty]$. 19*

20
21 The importance of Schwartz's theorem cannot be overemphasized. It forms the basic 21
22 foundation of Bayesian asymptotic theory for general parameter spaces. The first condi- 22
23 tion requires existence of a strictly unbiased test for testing the hypothesis $H_0: \theta = \theta_0$ 23
24 against the complement of a neighborhood U . The condition implies the existence of 24
25 a sequence of tests $\Phi_n(X_1, \dots, X_n)$ such that probabilities of both the type I error 25
26 $E_{\theta_0}\Phi_n(X_1, \dots, X_n)$ and the (maximum) type II error $\sup_{\theta \in U^c} E_{\theta}(1 - \Phi_n(X_1, \dots, X_n))$ 26
27 converges to zero exponentially fast. This existence of test is thus only a size restriction 27
28 on the model and not a condition on the prior. Writing 28

$$29 \Pi(\theta \in U^c | X_1, \dots, X_n) = \frac{\int_{U^c} \prod_{i=1}^n \frac{p(X_i, \theta)}{p(X_i, \theta_0)} d\Pi(\theta)}{\int_{\Theta} \prod_{i=1}^n \frac{p(X_i, \theta)}{p(X_i, \theta_0)} d\Pi(\theta)}, \quad (3.3)$$

30
31
32
33 this condition is used to show that for some $c > 0$, the numerator in (3.3) is smaller 33
34 than e^{-nc} for all sufficiently large n a.s. $[P_{\theta_0}^\infty]$. The condition on Kullback–Leibler sup- 34
35 port is a condition on the prior as well as the model. The condition implies that for all 35
36 $c > 0$, $e^{nc} \int_{\Theta} \prod_{i=1}^n \frac{p(X_i, \theta)}{p(X_i, \theta_0)} d\Pi(\theta) \rightarrow \infty$ a.s. $[P_{\theta_0}^\infty]$. Combining these two assertions, 36
37 we obtain the result of the theorem. The latter assertion follows by first replacing Θ by 37
38 the subset $\{\theta: K(\theta_0, \theta) < \varepsilon\}$, applying the strong law of large numbers to the integrand 38
39 and invoking Fatou's lemma. It may be noted that θ_0 needs to be in the Kullback– 39
40 Leibler support, not merely in the topological support of the prior for this argument 40
41 to go through. In practice, the condition is derived from the condition that θ_0 is in the 41
42 topological support of the prior along with some conditions on “nicety” of $p(x, \theta_0)$. 42

43 The testing condition is usually more difficult to satisfy. In finite dimension, the 43
44 condition usually holds. On the space of probability measures with the weak topology 44
45 on it, it is also not difficult to show that the required test exists; see Theorem 4.4.2 45

of Ghosh and Ramamoorthi (2003). However, in more complicated problems or for stronger topologies on densities (such as the variation or the Hellinger distance), the required tests do not exist without an additional compactness condition. Le Cam (1986) and Birgé (1983) developed an elegant theory of existence of uniformly exponentially powerful tests. However, the theory applies provided that the two hypotheses are convex. It is therefore helpful to split U^c into small balls for which required tests exist. If Θ is compact, the number of balls needed to cover U^c will be finite, and hence by taking the maximum of the resulting tests, the required test for testing $\theta = \theta_0$ against $\theta \in U^c$ may be obtained. However, the compactness condition imposes a severe restriction.

By a simple yet very useful observation, Barron (1988) concluded that it suffices that Φ_n satisfy

$$\sup_{\theta \in U^c \cap \Theta_n} E_{\theta} (1 - \Phi_n(X_1, \dots, X_n)) < a e^{-bn} \quad (3.4)$$

for some constants $a, b > 0$ and some “sieve” $\Theta_n \subset \Theta$, provided that it can be shown separately that

$$\Pi(\theta \in \Theta_n^c | X_1, \dots, X_n) \rightarrow 0 \quad \text{a.s. } [P_{\theta_0}^{\infty}]. \quad (3.5)$$

By a simple application of Fubini’s theorem, Barron (1988) concluded that (3.5) is implied by a condition only on the prior probability, namely, for some $c, d > 0$, $\Pi(\theta \in \Theta_n^c) \leq c e^{-nd}$. Now one may choose each Θ_n to be compact. However, because of dependence on n , one needs to estimate the number of balls required to cover Θ_n . From the same arguments, it follows that one needs to cover the sieve Θ_n with a maximum of e^{nc} balls, which is essentially a restriction on the covering number of the sieve Θ_n . The remaining part Θ_n^c , which may be topologically much bigger receives only a negligible prior probability by the given condition. It is interesting to note that unlike in sieve methods in non-Bayesian contexts, the sieve is merely a technical device for establishing consistency; the prior and the resulting Bayes procedure is not influenced by the choice of the sieve. Moreover, the sieve can be chosen depending on the accuracy level defined by the neighborhood U .

Barron’s (1988) useful observation made it possible to apply Schwartz’s ideas to prove posterior consistency in noncompact spaces as well. When the observations are i.i.d., one may take the parameter θ to be the density p itself. Let p_0 stand for the true density of each observation. Exploiting this idea, for a space \mathcal{P} of densities, Barron et al. (1999) gave a sufficient condition for posterior consistency in Hellinger distance $d_H(p_1, p_2) = (\int (p_1^{1/2} - p_2^{1/2})^2)^{1/2}$ in terms of a condition on bracketing Hellinger entropy⁴ a sieve $\mathcal{P}_n \subset \mathcal{P}$. Barron et al. (1999) used brackets to directly bound the likelihood ratios uniformly in the numerator of (3.4). The condition turns out to be considerably stronger than necessary in that we need to bound only an average likelihood ratio. Following Schwartz’s (1965) original approach involving test functions, Ghosal et al. (1999b) constructed the required tests using a much weaker condition on metric entropies. These authors considered the total variation distance $d_V(p_1, p_2) = \int |p_1 - p_2|$

⁴ The ε -bracketing Hellinger entropy of a set is the logarithm of the number ε -brackets with respect to the Hellinger distance needed to cover the set; see van der Vaart and Wellner (1996) for details on this and the related concepts.

(which is equivalent to d_H), constructed a test directly for a point null against a small variation ball using Hoeffding's inequality, and combined the resulting tests using the condition on the metric entropy.

For a subset S of a metric space with a metric d on it, let $N(\varepsilon, S, d)$, called the ε -covering number of S with respect to the metric d , stand for the minimum number of ε -balls needed to cover S . The logarithm of $N(\varepsilon, S, d)$ is often called the ε -entropy.

Assume that we have i.i.d. observations from a density $p \in \mathcal{P}$, a space of densities. Let p_0 stand for the true density and consider the variation distance d_V on \mathcal{P} . Let Π be a prior on \mathcal{P} .

THEOREM 3. *Suppose that $p_0 \in \text{KL}(\Pi)$. If given any $\varepsilon > 0$, there exist $\delta < \varepsilon/4$, $c_1, c_2 > 0$, $\beta < \varepsilon^2/8$ and $\mathcal{P}_n \subset \mathcal{P}$ such that $\Pi(\mathcal{P}_n^c) \leq c_1 e^{-nc_2}$ and $\log N(\delta, \mathcal{P}_n, d_V) \leq n\beta$, then $\Pi(P: d_V(P, P_0) > \varepsilon | X_1, \dots, X_n) \rightarrow 0$ a.s. [P_0^∞].*

Barron (1999) also noted that the testing condition in Schwartz's theorem is, in a sense, also necessary for posterior consistency to hold under Schwartz's condition on Kullback–Leibler support.

THEOREM 4. *Let \mathcal{P} be a space of densities, $p_0 \in \mathcal{P}$ be the true density and P_0 be the probability measure corresponding to p_0 . Let $p_0 \in \text{KL}(\Pi)$. Then the following conditions are equivalent:*

- (1) *There exists a β_0 such that $P_0\{\Pi(U^c | X_1, \dots, X_n) > e^{-n\beta_0}$ infinitely often\} = 0.*
- (2) *There exist subsets $V_n, W_n \subset \mathcal{P}$, $c_1, c_2, \beta_1, \beta_2 > 0$ and a sequence of test functions $\Phi_n(X_1, \dots, X_n)$ such that*
 - (a) $U^c \subset V_n \cup W_n$,
 - (b) $\Pi(W_n) \leq c_1 e^{-nc_2}$,
 - (c) $P_0\{\Phi_n > 0$ infinitely often\} = 0 and $\sup\{E_p(1 - \Phi_n): p \in V_n\} \leq c_2 e^{-n\beta_2}$.

In a semiparametric problem, an additional Euclidean parameter is present apart from an infinite-dimensional parameter, and the Euclidean parameter is usually of interest. Diaconis and Freedman (1986a, 1986b) demonstrated that putting a prior that gives consistent posterior separately for the nonparametric part may not lead to a consistent posterior when the Euclidean parameter is incorporated in the model. The example described below appeared to be counter-intuitive when it first appeared.

EXAMPLE 2. Consider i.i.d. observations from the location model $X = \theta + \epsilon$, where $\theta \in \mathbb{R}$, $\epsilon \sim F$ which is symmetric. Put any nonsingular prior density on θ and the symmetrized Dirichlet process prior on F with a Cauchy center measure. Then there exists a symmetric distribution F_0 such that if the X observations come from F_0 , then the posterior concentrates around two wrong values $\pm\gamma$ instead of the true value $\theta = 0$.

A similar phenomenon was observed by Doss (1985a, 1985b). The main problem in the above is that the posterior distribution for θ is close to the parametric posterior with a Cauchy density, and hence the posterior mode behaves like the M-estimator based on

1 the criterion function $m(x, \theta) = \log(1 + (x - \theta)^2)$. The lack of concavity of m leads to 1
2 undesirable solutions for some peculiar data generating distribution like F_0 . Consistency 2
3 however does obtain for the normal base measure since $m(x, \theta) = (x - \theta)^2$ is convex, 3
4 or even for the Cauchy base measure if F_0 has a strongly unimodal density. Here, 4
5 addition of the location parameter θ to the model destroys the delicate tail-free structure, 5
6 and hence Freedman's (1963, 1965) consistency result for tail-free processes cannot be 6
7 applied. Because the Dirichlet process selects only discrete distribution, it is also clear 7
8 that Schwartz's (1965) condition on Kullback–Leibler support does not hold. However, 8
9 as shown by Ghosal et al. (1999c), if we start with a prior on F that satisfies Schwartz's 9
10 (1965) condition in the nonparametric model (that is, the case of known $\theta = 0$), then the 10
11 same condition holds in the semiparametric model as well. This leads to weak consis- 11
12 tency in the semiparametric model (without any additional testing condition) and hence 12
13 consistency holds for the location parameter θ . The result extends to more general semi- 13
14 parametric problems. Therefore, unlike the tail-free property, Schwartz's condition on 14
15 Kullback–Leibler support is very robust which is not altered by symmetrization, addi- 15
16 tion of a location parameter or formation of mixtures. Thus Schwartz's theorem is the 16
17 right tool for studying consistency in semiparametric models. 17

18 Extensions of Schwartz's consistency theorem to independent, nonidentically distrib- 18
19 uted observations have been obtained by Amewou-Atisso et al. (2003) and Choudhuri 19
20 et al. (2004a). The former does not use sieves and hence is useful only when weak topol- 20
21 ogy is put on the infinite-dimensional part of the parameter. In semiparametric problems, 21
22 this topology is usually sufficient to derive posterior consistency for the Euclidean part. 22
23 However, for curve estimation problems, stronger topologies need to be considered and 23
24 sieves are essential. Consistency in probability instead of that in the almost sure sense 24
25 allows certain relaxations in the condition to be verified. Choudhuri et al. (2004a) con- 25
26 sidered such a formulation which is described below. 26

27
28 THEOREM 5. Let $Z_{i,n}$ be independently distributed with density $p_{i,n}(\cdot; \theta)$, $i =$ 28
29 $1, \dots, r_n$, with respect to a common σ -finite measure, where the parameter θ belongs 29
30 to an abstract measurable space Θ . The densities $p_{i,n}(\cdot, \theta)$ are assumed to be jointly 30
31 measurable. Let $\theta_0 \in \Theta$ and let $\bar{\Theta}_n$ and \mathcal{U}_n be two subsets of Θ . Let θ have prior 31
32 Π on Θ . Put $K_{i,n}(\theta_0, \theta) = E_{\theta_0}(\Lambda_i(\theta_0, \theta))$ and $V_{i,n}(\theta_0, \theta) = \text{var}_{\theta_0}(\Lambda_i(\theta_0, \theta))$, where 32
33 $\Lambda_i(\theta_0, \theta) = \log \frac{p_{i,n}(Z_{i,n}; \theta_0)}{p_{i,n}(Z_{i,n}; \theta)}$. 33
34

35 (A1) Prior positivity of neighborhoods. 35

36 Suppose that there exists a set B with $\Pi(B) > 0$ such that 36

37 (i) $\frac{1}{r_n^2} \sum_{i=1}^{r_n} V_{i,n}(\theta_0, \theta) \rightarrow 0$ for all $\theta \in B$, 37

38 (ii) $\liminf_{n \rightarrow \infty} \Pi \left(\left\{ \theta \in B: \frac{1}{r_n} \sum_{i=1}^{r_n} K_{i,n}(\theta_0, \theta) < \varepsilon \right\} \right) > 0$ for all $\varepsilon > 0$. 38
39

40 (A2) Existence of tests. 40

41 Suppose that there exists test functions $\{\Phi_n\}$, $\Theta_n \subset \bar{\Theta}_n$ and constants $C_1, C_2, c_1,$ 41
42 $c_2 > 0$ such that 42

43 (i) $E_{\theta_0} \Phi_n \rightarrow 0$, 43
44

45

- 1 (ii) $\sup_{\theta \in \mathcal{U}_n^c \cap \Theta_n} E_\theta(1 - \Phi_n) \leq C_1 e^{-c_1 r_n}$, 1
- 2 (iii) $\Pi(\bar{\Theta}_n \cap \Theta_n^c) \leq C_2 e^{-c_2 r_n}$. 2

3
4 Then $\Pi(\theta \in \mathcal{U}_n^c \cap \bar{\Theta}_n | Z_{1,n}, \dots, Z_{r_n,n}) \rightarrow 0$ in $P_{\theta_0}^n$ -probability. 4

5
6 Usually, the theorem will be applied to $\bar{\Theta}_n = \Theta$ for all n . If, however, condition (A2) 6
7 could be verified only on a part of Θ which may possibly depend on n , the above formu- 7
8 lation could be useful. However, the final conclusion should then be complemented by 8
9 showing that $\Pi(\bar{\Theta}_n^c | Z_1, \dots, Z_{r_n}) \rightarrow 0$ in $P_{\theta_0}^n$ -probability by some alternative method. 9

10 The first condition (A1) asserts that certain sets, which could be thought of as neigh- 10
11 borhoods of the true parameter θ_0 , have positive prior probabilities. This condition 11
12 ensures that the true value of the parameter is not excluded from the support of the prior. 12
13 The second condition (A2) asserts that the hypothesis $\theta = \theta_0$ can be tested against the 13
14 complement of a neighborhood for a topology of interest with a small probability of 14
15 type I error and a uniformly exponentially small probability of type II error on most 15
16 part of the parameter space in the sense that the prior probability of the remaining part 16
17 is exponentially small. 17

18 The above theorem is also valid for a sequence of priors Π_n provided that (A1)(i) is 18
19 strengthened to uniform convergence. 19

20 It should be remarked that Schwartz's condition on the Kullback–Leibler support is 20
21 not necessary for posterior consistency to hold. This is clearly evident in parametric 21
22 nonregular cases, where Kullback–Leibler divergence to some direction could be infin- 22
23 ity. For instance, as in Ghosal et al. (1999a), for the model $p_\theta = \text{Uniform}(0, \theta)$ density, 23
24 $0 < \theta \leq 1$, the Kullback–Leibler numbers $\int p_1 \log(p_1/p_\theta) = \infty$. However, the poster- 24
25 ior is consistent at $\theta = 1$ if the prior Π has 1 in its support. Modifying the model to 25
26 $\text{uniform}(\theta - 1, \theta + 1)$, we see that the Kullback–Leibler numbers are infinite for every 26
27 pair. Nevertheless, consistency for a general parametric family including such nonreg- 27
28 ular cases holds under continuity and positivity of the prior density at θ_0 provided that 28
29 the general conditions of Ibragimov and Has'minskii (1981) can be verified; see Ghosal 29
30 et al. (1995) for details. For infinite-dimensional models, consistency may hold with- 30
31 out Schwartz's condition on Kullback–Leibler support by exploiting special structure 31
32 of the posterior distribution as in the case of the Dirichlet or a tail-free process. For 32
33 estimation of a survival distribution using a Lévy process prior, Kim and Lee (2001) 33
34 concluded consistency from the explicit expressions for pointwise mean and variance 34
35 and monotonicity. For densities, consistency may also be shown by using some alter- 35
36 native conditions. One approach is by using the so-called Le Cam's inequality: For any 36
37 two disjoint subsets $U, V \subset \mathfrak{M}(\mathfrak{X})$, test function Φ , prior Π on $\mathfrak{M}(\mathfrak{X})$ and probability 37
38 measure P_0 on \mathfrak{X} , 38

$$\begin{aligned}
 & \int \Pi(V|x) dP_0(x) \\
 & \leq d_V(P_0, \lambda_U) + \int \Phi dP_0 + \frac{\Pi(V)}{\Pi(U)} \int (1 - \Phi) d\lambda_V, \quad (3.6)
 \end{aligned}$$

39
40
41 where $\lambda_U(B) = \int_U P(B) d\Pi(P)/\Pi(U)$, the conditional expectation of $P(B)$ with 41
42 respect to the prior Π restricted to the set U . Applying this inequality to V the comple- 42
43 ment of a neighborhood of P_0 and n i.i.d. observations, it may be shown that posterior 43
44
45

consistency in the weak sense holds provided that for any $\beta, \delta > 0$,

$$e^{n\beta} \Pi(P: d_V(P, P_0) < \delta/n) \rightarrow \infty. \quad (3.7)$$

Combining with appropriate testing conditions, stronger notions of consistency could be derived. The advantage of using this approach is that one need not control likelihood ratios now, and hence the result could be potentially used for undominated families as well, or at least can help reduce some positivity condition on the true density p_0 . On the other hand, (3.7) is a quantitative condition on the prior unlike Schwartz's, and hence is more difficult to verify in many examples.

Because the testing condition is a condition only on a model and is more difficult to verify, there have been attempts to prove some assertion on posterior convergence using Schwartz's condition on Kullback–Leibler support only. While Theorem 4 shows that the testing condition is needed, it may be still possible to show some useful results by either weakening the concept of convergence, or even by changing the definition of the posterior distribution! Barron (1999) showed that if $p_0 \in \text{KL}(\Pi)$, then

$$n^{-1} \sum_{i=1}^n E_{p_0} \left(\log \frac{p_0(X_i)}{p(X_i|X_1, \dots, X_{i-1})} \right) \rightarrow 0, \quad (3.8)$$

where $p(X_i|X_1, \dots, X_{i-1})$ is the predictive density of X_i given X_1, \dots, X_{i-1} . It may be noted that the predictive distribution is equal to the posterior mean of the density function. Hence in the Cesàro sense, the posterior mean density converges to the true density with respect to Kullback–Leibler neighborhoods, provided that the prior puts positive probabilities on Kullback–Leibler neighborhoods of p_0 . Walker (2003), using a martingale representation of the predictive density, showed that the average predictive density converges to the true density almost surely under d_H . Walker and Hjort (2001) showed that the following pseudo-posterior distribution, defined by

$$\Pi_\alpha(p \in B|X_1, \dots, X_n) = \frac{\int_B \prod_{i=1}^n p^\alpha(X_i) d\Pi(p)}{\int_B \prod_{i=1}^n p^\alpha(X_i) d\Pi(p)} \quad (3.9)$$

is consistent at any $p_0 \in \text{KL}(\Pi)$, provided that $0 < \alpha < 1$.

Walker (2004) obtained another interesting result using an idea of restricting to a subset and looking at the predictive distribution (in this case, in the posterior) somewhat similar to that in Le Cam's inequality. If V is a set such that $\liminf_{n \rightarrow \infty} d_H(\lambda_{n,V}, p_0) > 0$, where $\lambda_{n,V}(B) = (\Pi(V|X_1, \dots, X_n))^{-1} \int_V p(B) d\Pi(p|X_1, \dots, X_n)$, then $\Pi(V|X_1, \dots, X_n) \rightarrow 0$ a.s. under P_0 . A martingale property of the predictive distribution is utilized to prove the result. If V is the complement of a suitable weak neighborhood of p_0 , then $\liminf_{n \rightarrow \infty} d_H(\lambda_{n,V}, p_0) > 0$, and hence the result provides an alternative way of proving the weak consistency result without appealing to Schwartz's theorem. Walker (2004) also considered other topologies.

The following is another result of Walker (2004) proving sufficient conditions for posterior consistency in terms of a suitable countable covering.

THEOREM 6. *Let $p_0 \in \text{KL}(\Pi)$ and $V = \{p: d_H(p, p_0) > \varepsilon\}$. Let there exists $0 < \delta < \varepsilon$ and V_1, V_2, \dots a countable disjoint cover of V such that $d_H(p_1, p_2) < 2\delta$*

1 for all $p_1, p_2 \in V_j$ and for all $j = 1, 2, \dots$, and $\sum_{j=1}^{\infty} \sqrt{\Pi(V_j)} < \infty$. Then 1
2 $\Pi(V|X_1, \dots, X_n) \rightarrow 0$ a.s. [p_0^∞]. 2
3 3

4 While the lack of consistency is clearly undesirable, consistency itself is a very weak 4
5 requirement. Given a consistency result, one would like to obtain information on the 5
6 rates of convergence of the posterior distribution and see whether the obtained rate 6
7 matches with the known optimal rate for point estimators. In finite-dimensional prob- 7
8 lems, it is well known that the posterior converges at a rate of $n^{-1/2}$ in the Hellinger 8
9 distance; see Ibragimov and Has'minskii (1981) and Le Cam (1986). 9

10 Conditions for the rate of convergence given by Ghosal et al. (2000) and described 10
11 below are quantitative refinement of conditions for consistency. A similar result, but 11
12 under a much stronger condition on bracketing entropy numbers, was given by Shen 12
13 and Wasserman (2001). 13
14 14

15 THEOREM 7. Let $\varepsilon_n \rightarrow 0$, $n\varepsilon_n^2 \rightarrow \infty$ and suppose that there exist $\mathcal{P}_n \subset \mathcal{P}$, constants 15
16 $c_1, c_2, c_3, c_4 > 0$ such that 16
17 17

- 18 (i) $\log D(\varepsilon_n, \mathcal{P}_n, d) \leq c_1 n \varepsilon_n^2$, where D stands for the packing number; 18
19 (ii) $\Pi(\mathcal{P} \setminus \mathcal{P}_n) \leq c_2 e^{-(c_3+4)n\varepsilon_n^2}$; 19
20 (iii) $\Pi(p: \int p_0 \log \frac{p_0}{p} < \varepsilon_n^2, \int p_0 \log^2 \frac{p_0}{p} < \varepsilon_n^2) \geq c_4 e^{-c_3 n \varepsilon_n^2}$. 20
21 21

22 Then for some M , $\Pi(d(p, p_0) > M\varepsilon_n | X_1, X_2, \dots, X_n) \rightarrow 0$. 22
23 23

24 More generally, the entropy condition can be replaced by a testing condition, though, 24
25 in most applications, a test is constructed from entropy bounds. Some variations of the 25
26 theorem are given by Ghosal et al. (2000), Ghosal and van der Vaart (2001) and Belitser 26
27 and Ghosal (2003). 27

28 While the theorems of Ghosal et al. (2000) satisfactorily cover i.i.d. data, major 28
29 extensions are needed to cover some familiar situations such as regression with a 29
30 fixed design, dose response study, generalized linear models with an unknown link, 30
31 Whittle estimation of a spectral density and so on. Ghosal and van der Vaart (2003a) 31
32 considered the issue and showed that the basic ideas of the i.i.d. case work with 32
33 suitable modifications. Let d_n^2 be the average squared Hellinger distance defined by 33
34 $d_n^2(\theta_1, \theta_2) = n^{-1} \sum_{i=1}^n d_H^2(p_{i,\theta_1}, p_{i,\theta_2})$. Birgé (1983) showed that a test for θ_0 against 34
35 $\{\theta: d_n(\theta, \theta_1) < d_n(\theta_0, \theta_1)/18\}$ with error probabilities at most $\exp(-nd_n^2(\theta_0, \theta_1)/2)$ 35
36 may be constructed. To find the intended test for θ_0 against $\{\theta: d_n(\theta, \theta_0) > \varepsilon\}$, one 36
37 therefore needs to cover the alternative by d_n balls of radius $\varepsilon/18$. The number of such 37
38 balls is controlled by the d_n -entropy numbers. Prior concentration near θ_0 controls the 38
39 denominator as in the case of i.i.d. observations. Using these ideas, Ghosal and van der 39
40 Vaart (2003a) obtained the following theorem on convergence rates that is applicable to 40
41 independent, nonidentically distributed observations, and applied the result to various 41
42 non-i.i.d. models. 42
43 43

44 THEOREM 8. Suppose that for a sequence $\varepsilon_n \rightarrow 0$ such that $n\varepsilon_n^2$ is bounded away from 44
45 zero, some $k > 1$, every sufficiently large j and sets $\Theta_n \subset \Theta$, the following conditions 45

1 are satisfied:

$$2 \quad \sup_{\varepsilon > \varepsilon_n} \log N(\varepsilon/36, \{\theta \in \Theta_n: d_n(\theta, \theta_0) < \varepsilon\}, d_n) \leq n\varepsilon_n^2, \quad (3.10)$$

$$5 \quad \Pi_n(\Theta \setminus \Theta_n) / \Pi_n(B_n^*(\theta_0, \varepsilon_n; k)) = o(e^{-2n\varepsilon_n^2}), \quad (3.11)$$

$$8 \quad \frac{\Pi_n(\theta \in \Theta_n: j\varepsilon_n < d_n(\theta, \theta_0) \leq 2j\varepsilon_n)}{\Pi_n(B_n^*(\theta_0, \varepsilon_n; k))} \leq e^{n\varepsilon_n^2 j^2 / 4}. \quad (3.12)$$

10 Then $P_{\theta_0}^{(n)} \Pi_n(\theta: d_n(\theta, \theta_0) \geq M_n \varepsilon_n | X^{(n)}) \rightarrow 0$ for every $M_n \rightarrow \infty$.

12 Ghosal and van der Vaart (2003a) also considered some dependent cases such as
13 Markov chains, autoregressive model and signal estimation in presence of Gaussian
14 white noise.

15 When one addresses the issue of optimal rate of convergence, one considers a
16 smoothness class of the involved functions. The method of construction of the opti-
17 mal prior with the help of bracketing or spline functions, as in Ghosal et al. (2000)
18 requires the knowledge of the smoothness index. In practice, such information is not
19 available and it is desirable to construct a prior that is adaptive. In other words, we
20 wish to construct a prior that simultaneously achieves the optimal rate for every possi-
21 ble smoothness class under consideration. If only countably many models are involved,
22 a natural and elegant method would be to consider a prior that is a mixture of the opti-
23 mal priors for different smoothness classes. Belitser and Ghosal (2003) showed that
24 the strategy works for an infinite-dimensional normal. Ghosal et al. (2003) and Huang
25 (2004) obtained similar results for the density estimation problem.

26 Kleijn and van der Vaart (2002) considered the issue of misspecification, where p_0
27 may not lie in the support of the prior. In such a case, consistency at p_0 cannot hold, but
28 it is widely believed that the posterior concentrates around the Kullback–Leibler projec-
29 tion p^* of p_0 to the model; see Berk (1966) for some results for parametric exponential
30 families. Under suitable conditions which could be regarded as generalizations of the
31 conditions of Theorem 7, Kleijn and van der Vaart (2002) showed that the posterior con-
32 centrates around p^* at a rate described by a certain entropy condition and concentration
33 rate of the prior around p^* . Kleijn and van der Vaart (2002) also defined a notion of
34 covering number for testing under misspecification that turns out to be the appropriate
35 way of measuring the size of the model in the misspecified case. A weighted version
36 of the Hellinger distance happens to be the proper way of measuring distance between
37 densities that leads to a fruitful theorem on rates in the misspecified case. A useful theo-
38 rem on consistency (in the sense that the posterior distribution concentrates around p^*)
39 follows as a corollary.

40 When the posterior distribution converges at a certain rate, it is also important to
41 know whether the posterior measure, after possibly a random centering and scaling,
42 converges to a nondegenerate measure. For smooth parametric families, convergence to
43 a normal distribution holds and is popularly known as the Bernstein–von Mises theorem;
44 see Le Cam and Yang (2000) and van der Vaart (1998) for details. For a general paramet-
45 ric family which need not be smooth, a necessary and sufficient condition in terms of the

1 limiting likelihood ratio process for convergence of the posterior (to some nondegen- 1
2 erate distribution using some random centering) is given by Ghosh et al. (1994, 1995). 2
3 For infinite-dimensional cases, results are relatively rare. Some partial results were ob- 3
4 tained by Lo (1983, 1986) for Dirichlet process, Shen (2002) for certain semiparametric 4
5 models, Susarla and Van Ryzin (1978) and Kim and Lee (2004) for certain survival 5
6 models respectively with the Dirichlet process and Lévy process priors. However, it 6
7 appears from the work of Cox (1993) and Freedman (1999) that Bernstein–von Mises 7
8 theorem does not hold in most cases when the convergence rate is slower than $n^{-1/2}$. 8
9 Freedman (1999) indeed showed that for the relatively simple problem of the estima- 9
10 tion of the mean of an infinite-dimensional normal distribution with independent normal 10
11 priors, the frequentist and the Bayesian distribution of L_2 -norm of the difference of the 11
12 Bayes estimate and the parameter differ by an amount equal to the scale of interest, 12
13 and the frequentist coverage probability of a Bayesian credible set for the parameter is 13
14 asymptotically zero. However, see Ghosal (2000) for a partially positive result. 14
15
16

17 4. Estimation of cumulative probability distribution 17

18 4.1. Dirichlet process prior 18

19 One of the nicest properties of the Dirichlet distribution, making it hugely popular, is its 19
20 conjugacy for estimating a distribution function (equivalently, the probability law) with 20
21 i.i.d. observations. Consider X_1, \dots, X_n are i.i.d. samples from an unknown cumulative 21
22 distribution function (cdf) F on \mathbb{R}^d . Suppose F is given a Dirichlet process prior with 22
23 parameters (M, G) . Then the posterior distribution is again a Dirichlet process with the 23
24 two parameters updated as 24
25
26

$$27 \quad M \mapsto M + n \quad \text{and} \quad G \mapsto (MG + n\mathbb{F}_n)/(M + n), \quad (4.1) \quad 27$$

28 where \mathbb{F}_n is the empirical cdf. This may be easily shown by reducing the data to counts 28
29 of sets from a partition, using the conjugacy of the finite-dimensional Dirichlet dis- 29
30 tribution for the multinomial distribution and passing to the limit with the aid of the 30
31 martingale convergence theorem. Combining with (2.1), this implies that the posterior 31
32 expectation and variance of $F(x)$ are given by 32
33

$$34 \quad \begin{aligned} 35 \quad \tilde{\mathbb{F}}_n(x) &= \mathbb{E}(F(x)|X_1, \dots, X_n) = \frac{M}{M+n}G(x) + \frac{n}{M+n}\mathbb{F}_n(x), & 34 \\ 36 \quad \text{var}(F(x)|X_1, \dots, X_n) &= \frac{\tilde{\mathbb{F}}_n(x)(1 - \tilde{\mathbb{F}}_n(x))}{1 + M + n}. & 35 \\ 37 & & 36 \\ 38 & & 37 \end{aligned} \quad (4.2) \quad 38$$

39 Therefore the posterior mean is a convex combination of the prior mean and the empir- 39
40 ical cdf. As the sample size increases, the behavior of the posterior mean is inherited 40
41 from that of the empirical probability measure. Also M could be interpreted as the 41
42 strength in the prior or the “prior sample size”. 42

43 The above discussion may lull us to interpret the limiting case $M \rightarrow 0$ as nonin- 43
44 formative. Indeed, Rubin (1981) proposed $\text{Dir}(n, \mathbb{F}_n)$ as the Bayesian bootstrap, which 44
45 corresponds to the posterior obtained from the Dirichlet process by letting $M \rightarrow 0$. 45

1 However, some caution is needed while interpreting the case $M \rightarrow 0$ as noninformative 1
2 because of the role of M in also controlling the number of ties among samples drawn 2
3 from P , where P itself is drawn from the Dirichlet process. Sethuraman and Tiwari 3
4 (1982) pointed out that as $M \rightarrow 0$, the Dirichlet process converges weakly to the ran- 4
5 dom measure which is degenerate at some point θ distributed as G by property (ii) of 5
6 convergence of Dirichlet measures mentioned in Section 2.1. Such a prior is clearly 6
7 “very informative”, and hence is unsuitable as a noninformative prior. 7

8 To obtain posterior consistency, note that (4.1) converges a.s. to the true cdf gen- 8
9 erating data. An important consequence of the above assertions is that the posterior 9
10 distribution based on the Dirichlet process, not just the posterior mean, is consis- 10
11 tent for the weak topology. Thus, by the weak convergence property of Dirichlet 11
12 process, the posterior is consistent with respect to the weak topology. It can also be 12
13 shown that, the posterior is consistent in the Kolmogorov–Smirnov distance defined as 13
14 $d_{KS}(F_1, F_2) = \sup_x |F_1(x) - F_2(x)|$. The space of cdf’s under d_{KS} is however neither 14
15 separable nor complete. 15

16 If the posterior distribution of F is given a prior that is a mixture of Dirichlet process, 16
17 the posterior distribution is still a mixture of Dirichlet processes; see Theorem 3 of 17
18 Antoniak (1974). However, mixtures may lead to inconsistent posterior distribution, 18
19 unlike a single Dirichlet process. Nevertheless, if M_θ is bounded in θ , then posterior 19
20 consistency holds. 20
21

22 4.2. Tail-free and Polya tree priors 22

23 Tail-free priors are extremely flexible, yet have some interesting properties. If the dis- 23
24 tribution function generating the i.i.d. data is given a tail-free prior, the posterior dis- 24
25 tribution is also tail-free. Further, as mentioned in Section 3, Freedman (1963, 1965) 25
26 showed that the posterior obtained from a tail-free process prior is weakly consistent. 26
27 The tail-free property helps reduce a weak neighborhood to a neighborhood involving 27
28 only finitely many variables in the hierarchical representation, and hence the problem 28
29 reduces to a finite-dimensional multinomial distribution, where consistency holds. In- 29
30 deed Freedman’s original motivation was to avoid pitfall as in Example 1. 30
31

32 A Polya tree prior may be used if one desires some smoothness of the random cdf. 31
33 The most interesting property of a Polya tree process is its conjugacy. Conditional on 32
34 the data X_1, \dots, X_n , the posterior distribution is again a Polya tree with respect to the 33
35 same partition and α_ε updated to $\alpha_\varepsilon^* = \alpha_\varepsilon + \sum_{i=1}^n I\{X_i \in B_\varepsilon\}$. Besides, they lead to a 34
36 consistent posterior in the weak topology as Polya trees are also tail-free processes. 35
36

37 4.3. Right censored data 37

38 Let X be a random variable of interest that is right censored by another random vari- 38
39 able Y . The observation is (Z, Δ) , where $Z = \min(X, Y)$ and $\Delta = I(X > Y)$. Assume 39
40 that X and Y are independent with corresponding cdf F and H , where both F and H are 40
41 unknown. The problem is to estimate F . Susarla and Van Ryzin (1976) put a Dirichlet 41
42 process prior on F . Blum and Susarla (1977) found that the posterior distribution for 42
43 i.i.d. data can be written as a mixture of Dirichlet processes. Using this idea, Susarla 43
44 and Van Ryzin (1978) obtained that the posterior is mean square consistent with rate 44
45

1 $O(n^{-1})$, almost surely consistent with rate $O(\log n/n^{1/2})$, and that the posterior dis- 1
2 tribution of $\{F(u): 0 < u < T\}$, $T < \infty$, converges weakly to a Gaussian process 2
3 whenever F and H are continuous and that $P(X > u)P(Y > u) > 0$. The mixture 3
4 representation is however cumbersome. Ghosh and Ramamoorthi (1995) showed that 4
5 the posterior distribution can also be written as a Polya tree process (with partitions de- 5
6 pendent on the uncensored samples). They proved consistency by an elegant argument. 6
7 Doksum (1974) found that the neutral to right process for F form a conjugate family 7
8 for the right censored data. Viewed as a prior on the cumulative hazard process, the prior 8
9 can be identified with an independent increment process. An updating mechanism is 9
10 described by Kim (1999) using a counting process approach. Beta processes, introduced 10
11 by Hjort (1990), also form a conjugate family. Kim and Lee (2001) obtained sufficient 11
12 conditions for posterior consistency for a Lévy processes prior, which includes Dirichlet 12
13 processes and beta processes. Under certain conditions, the posterior also converges 13
14 at the usual $n^{-1/2}$ rate and admits a Bernstein–von Mises theorem; see Kim and Lee 14
15 (2004). 15

16 5. Density estimation 16

17 Density estimation is one of the fundamental problems of nonparametric inference 17
18 because of its applicability to various problems including cluster analysis and robust 18
19 estimation. A common approach to constructing priors on the space of probability den- 19
20 sities is to use Dirichlet mixtures where the kernels are chosen depending on the sample 20
21 space. The posterior distributions are analytically intractable and the MCMC techniques 21
22 are different for different kernels. Other priors useful for this problem are Polya tree 22
23 processes and Gaussian processes. In this section, we discuss some of the computational 23
24 issues and conditions for consistency and convergence rates of the posterior distribution. 24
25 25
26 26
27 27
28 28

29 5.1. Dirichlet mixture 29

30 Consider that the density generating the data is a mixture of densities belonging to 30
31 some parametric family, that is, $p_F(x) = \int \psi(x, \theta) dF(\theta)$. Let the mixing distribution 31
32 F be given a $\text{Dir}(M, G)$ prior. Viewing $p_F(x)$ as a linear functional of F , the prior 32
33 expectation of $p_F(x)$ is easily found to be $\int \psi(x, \theta) dG(\theta)$. To compute the posterior 33
34 expectation, the following hierarchical representation of the above prior is often conve- 34
35 nient: 35
36 36

$$37 \quad X_i \stackrel{\text{ind}}{\sim} \psi(\cdot, \theta_i), \quad \theta_i \stackrel{\text{i.i.d.}}{\sim} F, \quad F \sim \text{Dir}(M, G). \quad (5.1) \quad 37$$

38 Let $\Pi(\boldsymbol{\theta}|X_1, \dots, X_n)$ stand for the distribution of $(\theta_1, \dots, \theta_n)$ given (X_1, \dots, X_n) . 38
39 Observe that given $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)$, the posterior distribution of F is Dirichlet with 39
40 base measure $MG + n\mathbb{G}_n$, where $\mathbb{G}_n(\cdot, \boldsymbol{\theta}) = n^{-1} \sum_{i=1}^n \delta_{\theta_i}$, the empirical distribution 40
41 of $(\theta_1, \dots, \theta_n)$. Hence the posterior distribution of F may be written as a mixture of 41
42 Dirichlet processes. The posterior mean of $F(\cdot)$ may be written as 42
43 43

$$44 \quad \frac{M}{M+n}G(\cdot) + \frac{n}{M+n} \int G_n(\cdot, \boldsymbol{\theta})\Pi(d\boldsymbol{\theta}|X_1, \dots, X_n) \quad (5.2) \quad 44$$

1 and the posterior mean of the density at x becomes

$$2 \quad \frac{M}{M+n} \int \psi(x, \theta) dG(\theta) + \frac{n}{M+n} \frac{1}{n} \sum_{i=1}^n \int \psi(x, \theta_i) \Pi(d\theta | X_1, \dots, X_n). \quad 3$$

$$4 \quad (5.3) \quad 5$$

6 The Bayes estimate is thus composed of a part attributable to the prior and
7 a part due to observations. Ferguson (1983) remarks that the factor
8 $n^{-1} \sum_{i=1}^n \int \psi(x, \theta_i) \Pi(d\theta | X_1, \dots, X_n)$ in the second term of (5.3) can be viewed
9 as a partially Bayesian estimate with the influence of the prior guess reduced. The
10 evaluation of the above quantities depend on $\Pi(d\theta | X_1, \dots, X_n)$. The joint prior for
11 $(\theta_1, \theta_2, \dots, \theta_n)$ is given by the generalized Polya urn scheme

$$12 \quad G(d\theta_1) \times \frac{(MG(d\theta_2) + \delta_{\theta_1})}{M+1} \times \dots \times \frac{(MG(d\theta_n) + \sum_{i=1}^{n-1} \delta_{\theta_i})}{M+n}. \quad 13$$

$$14 \quad (5.4) \quad 15$$

16 Further, the likelihood given $(\theta_1, \theta_2, \dots, \theta_n)$ is $\prod_{i=1}^n \psi(X_i, \theta_i)$. Hence H can be writ-
17 ten down using the Bayes formula. Using the above equations and some algebra, Lo
18 (1984) obtained analytical expressions of the posterior expectation of $f(x)$. However,
19 the formula is of marginal use because the number of terms grows very fast with the
20 sample size. Computations are thus done via MCMC techniques as in the special case
21 of normal mixtures described in the next subsection; see the review article Escobar and
22 West (1998) for details.

23 5.1.1. Mixture of normal kernels

24 Suppose that the unknown density of interest is supported on the entire real line. Then
25 a natural choice of the kernel is $\phi_\sigma(x - \mu)$, the normal density with mean μ and vari-
26 ance σ^2 . The mixture distribution F is given Dirichlet process prior with some base
27 measure MG , while G is often given a normal/inverse-gamma distribution to achieve
28 conjugacy. Thus, under G , $\sigma^{-2} \sim \text{Gamma}(s, \beta)$, a gamma distribution with shape pa-
29 rameter s and scale parameter β , and $(\mu|\sigma) \sim N(m, \sigma^2)$. Let $\theta = (\mu, \sigma)$. Then the
30 hierarchical model is

$$31 \quad X_i | \theta_i \stackrel{\text{ind}}{\sim} N(\mu_i, \sigma_i^2), \quad \theta_i \stackrel{\text{i.i.d.}}{\sim} F, \quad F \sim \text{Dir}(M, G). \quad 32$$

$$33 \quad (5.5) \quad 34$$

35 Given $\theta = (\theta_1, \dots, \theta_n)$, the distribution of F may be updated analytically. Thus,
36 if one can sample from the posterior distribution of θ , Monte Carlo averages may
37 be used to find the posterior expectation of F and thus the posterior expectation
38 of $p(x) = \int \phi_\sigma(x - \mu) dF(x)$. Escobar (1994) and Escobar and West (1995) pro-
39 vided an algorithm for sampling from the posterior distribution of θ . Let $\theta_{-i} =$
40 $\{\theta_1, \dots, \theta_{i-1}, \theta_{i+1}, \dots, \theta_n\}$. Then

$$39 \quad (\theta_i | \theta_{-i}, x_1, \dots, x_n) \sim q_{i0} G_i(\theta_i) + \sum_{j=1, j \neq i}^n q_{ij} \delta_{\theta_j}(\theta_i), \quad 40$$

$$41 \quad (5.6) \quad 42$$

42 where $G_i(\theta_i)$ is the bivariate normal/inverse-gamma distribution under which

$$43 \quad \sigma_i^{-2} \sim \text{Gamma}(s + 1/2, \beta + (x_i - m)^2/2), \quad 44$$

$$44 \quad (\mu_i | \sigma_i) \sim N(m + x_i, \sigma_i^2) \quad 45$$

$$45 \quad (5.7) \quad 46$$

1 and the weights q_{ij} 's are defined by $q_{i0} \propto M\Gamma(s + 1/2)(2\beta)^s \Gamma(s)^{-1} \{2\beta + (x_i -$
2 $m)^2\}^{-(s+1/2)}$ and $q_{ij} \propto \sqrt{\pi} \phi_{\sigma_i}(x_i - \mu_i)$ for $j \neq i$. Thus a Gibbs sampler algorithm
3 is described by updating θ componentwise through the conditional distribution in (5.6).
4 The initial values of θ_i could be a sample from G_i .

5 The bandwidth parameter σ is often kept constant depending on the sample size,
6 say σ_n . This leads to only the location mixture. In that case a Gibbs sampler algorithm
7 is obtained by keeping σ_i fixed at σ_n in the earlier algorithm and updating only the
8 location components μ_i .

9 Consistency of the posterior distribution for Dirichlet mixture of normals was studied
10 by Ghosal et al. (1999b). Let p_0 stand for the true density.

11
12 THEOREM 9. If $p_0 = \int \phi_\sigma(x - \mu) dF_0(\mu, \sigma)$, where F_0 is compactly supported and
13 in the weak support of Π , then $p_0 \in \text{KL}(\Pi)$.

14 If p_0 is not a mixture of normals but is compactly supported, 0 is in the support of
15 the prior for σ , and $\lim_{\sigma \rightarrow 0} \int p_0 \log(p_0/p_0 * \phi_\sigma) = 0$, then $p_0 \in \text{KL}(\Pi)$.

16 If $p_0 \in \text{KL}(\Pi)$, the base measure G of the underlying Dirichlet process is compactly
17 supported and $\Pi(\sigma < t) \leq c_1 e^{-c_2/t}$, then the posterior is consistent at p_0 for the
18 total variation distance d_V . If the compact support G is replaced by the condition that
19 for every $\varepsilon > 0$, there exist a_n, σ_n with $a_n/\sigma_n < \varepsilon n$ satisfying $G[-a_n, a_n] < e^{-n\beta_1}$
20 and $\Pi(\sigma < \sigma_n) \leq e^{-n\beta_2}$ for $\beta_1, \beta_2 > 0$, then also consistency for d_V holds at any
21 $p_0 \in \text{KL}(\Pi)$.

22
23 The condition $p_0 \in \text{KL}(\Pi)$ implies weak consistency by Schwartz's theorem. The
24 condition for $p_0 \in \text{KL}(\Pi)$ when p_0 is neither a normal mixture nor compactly sup-
25 ported, as given by Theorem 5 of Ghosal et al. (1999b) using estimates of Dirichlet
26 tails, is complicated. However, the conditions holds under strong integrability condi-
27 tions on p_0 . The base measure for the Dirichlet could be normal and the prior on σ could
28 be a truncated inverse gamma possibly involving additional parameters. Better sufficient
29 condition for $p_0 \in \text{KL}(\Pi)$ is given by Tokdar (2003). Consider a location-scale mix-
30 ture of normal with a prior Π on the mixing measure. If p_0 is bounded, nowhere zero,
31 $\int p_0 |\log p_0| < \infty$, $\int p_0 \log(p_0/\psi) < \infty$ where $\psi(x) = \inf\{p_0(t): x - 1 \leq t \leq x + 1\}$,
32 $\int |x|^{2+\delta} p_0(x) dx < \infty$, and every compactly supported probability lies in $\text{supp}(\Pi)$,
33 then $p_0 \in \text{KL}(\Pi)$. The moment condition can be weakened to only δ -moment if Π is
34 Dirichlet. In particular, the case that p_0 is Cauchy could be covered.

35 Convergence rates of the posterior distribution were obtained by Ghosal and van
36 der Vaart (2001, 2003b) respectively the "super smooth" and the "smooth" cases. We
37 discuss below the case of location mixtures only, where the scale gets a separate inde-
38 pendent prior.

39
40 THEOREM 10. Assume that $p_0 = \phi_{\sigma_0} * F_0$, and the prior on σ has a density that is
41 compactly supported in $(0, \infty)$ but is positive and continuous at σ_0 . Suppose that F_0
42 has compact support and the base measure G has a continuous and positive density on
43 an interval containing the support of F_0 and has tails $G(|z| > t) \lesssim e^{-b|t|^\delta}$. Then the
44 posterior converges at a rate $n^{-1/2} (\log n)^{\max(\frac{2}{\delta}, \frac{1}{2}) + \frac{1}{2}}$ with respect to d_H . The condition
45

1 of compact support of F_0 could be replaced by that of sub-Gaussian tails if G is normal, 1
2 in which case the rate is $n^{-1/2}(\log n)^{3/2}$. 2

3 If instead p_0 is compactly supported, twice continuously differentiable and 3
4 $\int (p_0''/p_0)^2 p_0 < \infty$ and $\int (p_0'/p_0)^4 p_0 < \infty$, and the prior on (σ/σ_n) has a density 4
5 that is compactly supported in $(0, \infty)$, where $\sigma_n \rightarrow 0$, then the posterior converges at a 5
6 rate $\max((n\sigma_n)^{-1/2}(\log n), \sigma_n^2 \log n)$. In particular, the best rate $\varepsilon_n \sim n^{-2/5}(\log n)^{-4/5}$ 6
7 is obtained by choosing $\sigma_n \sim n^{-1/5}(\log n)^{-2/5}$. 7

8
9 The proofs are the result of some delicate estimates of the number of components a 9
10 discrete mixing distribution must have to approximate a general normal mixture. Some 10
11 further results are given by Ghosal and van der Vaart (2003b) when p_0 does not have 11
12 compact support. 12
13

14 5.1.2. Uniform scale mixtures 14

15 A nonincreasing density on $[0, \infty)$ may be written as a mixture of the form 15
16 $\int \theta^{-1} I\{0 \leq x \leq \theta\} F(d\theta)$ by a well known representation theorem of Khinchine 16
17 and Shepp. This lets us put a prior on this class from that on F . Brunner and Lo (1989) 17
18 considered this idea and put a Dirichlet prior for F . Coupled with a symmetrization 18
19 technique as in Section 2.2.3, this leads to a reasonable prior for the error distribution. 19
20 Brunner and Lo (1989) used this approach for the semiparametric location problem. 20
21 The case of asymmetric error was treated by Brunner (1992) and that of semiparametric 21
22 linear regression by Brunner (1995). 22
23

24 5.1.3. Mixtures on the half line 24

25 Dirichlet mixtures of exponential distributions may be considered as a reasonable model 25
26 for a decreasing, convex density on the positive half line. More generally, mixtures of 26
27 gamma densities, which may be motivated by Feller approximation procedure using a 27
28 Poisson sampling scheme in the sense of Petrone and Veronese (2002), may be con- 28
29 sidered to pick up arbitrary shapes. Such a prior may be chosen to have a large weak 29
30 support. Mixtures of inverse gammas may be motivated similarly by Feller approxi- 30
31 mation using a gamma sampling scheme. In general, a canonical choice of a kernel 31
32 function could be made once a Feller sampling scheme appropriate for the domain 32
33 could be specified. For a general kernel, weak consistency may be shown exploiting 33
34 Feller approximation property as in Petrone and Veronese (2002). 34
35

36 Mixtures of Weibulls or lognormals are dense in the stronger sense of total variation 36
37 distance provided that we let the shape parameter of the Weibull to approach infinity or 37
38 that of the lognormal to approach zero. To see this, observe that these two kernels form 38
39 location-scale families in the log-scale, and hence are approximate identities. Kottas and 39
40 Gelfand (2001) used these mixtures for median regression, where asymmetry is an im- 40
41 portant aspect. The mixture of Weibull is very useful to model observations of censored 41
42 data because its survival function has a simpler expression compared to that for the mix- 42
43 tures of gamma or lognormal. Ghosh and Ghosal (2003) used these mixtures to model 43
44 a proportional mean structure censored data. The posterior distribution was computed 44
45 using an MCMC algorithm for Dirichlet mixtures coupled with imputation of censored 45

1 data. Posterior consistency can be established by reducing the original model to a stan- 1
2 dard regression model with unknown error for which the results of Amewou-Atisso et 2
3 al. (2003) apply. More specifically, consistency holds if the true baseline density is in 3
4 the Kullback–Leibler support of the Dirichlet mixture prior. The last condition can be 4
5 established under reasonable conditions using the ideas of Theorem 9 and its extension 5
6 by Tokdar (2003). 6
7

8
9 *5.1.4. Bernstein polynomials*

10 On the unit interval, the family of beta distributions form a flexible two-parameter fam- 10
11 ily of densities and their mixtures form a very rich class. Indeed, mixtures of beta 11
12 densities with integer parameters are sufficient to approximate any distribution. For 12
13 a continuous probability distribution function F on $(0, 1]$, the associated Bernstein 13
14 polynomial $B(x; k, F) = \sum_{j=0}^k F(j/k) \binom{k}{j} x^j (1-x)^{k-j}$, which is a mixture of beta 14
15 distributions, converges uniformly to F as $k \rightarrow \infty$. Using an idea of Diaconis that 15
16 this approximation property may be exploited to construct priors with full topological 16
17 support, Petrone (1999a, 1999b) proposed the following hierarchical prior called the 17
18 Bernstein polynomial prior: 18

- 19
20 • $f(x) = \sum_{j=1}^k w_{j,k} \beta(x; j, k-j+1)$, 19
21 • $k \sim \rho(\cdot)$, 20
22 • $(\mathbf{w}_k = (w_{1,k}, \dots, w_{k,k})|k) \sim H_k(\cdot)$, a distribution on the k -dimensional simplex. 21

23 Petrone (1999a) showed that if for all k , $\rho(k) > 0$ and \mathbf{w}_k has full support on Δ_k , then 23
24 every distribution on $(0, 1]$ is in the weak support of the Bernstein polynomial prior, 24
25 and every continuous distribution is in the topological support of the prior defined by 25
26 the Kolmogorov–Smirnov distance. 26

27 The posterior mean, given k , is 27

28
29
30
$$E(f(x)|k, x_1, \dots, x_n) = \sum_{j=1}^k E(w_{j,k}|x_1, \dots, x_n) \beta(x; j, k-j+1), \quad (5.8)$$
 30
31

32 and the distribution of k is updated to $\rho(k|x_1, \dots, x_n)$. Petrone (1999a, 1999b) dis- 32
33 cussed MCMC algorithms to compute the posterior expectations and carried out exten- 33
34 sive simulations to show that the resulting density estimates work well. 34

35 Consistency is given by Petrone and Wasserman (2002). The corresponding results 35
36 on convergence rates are obtained by Ghosal (2001). 36
37

38
39 **THEOREM 11.** *If p_0 is continuous density on $[0, 1]$, the base measure G has support 39
40 all of $[0, 1]$ and the prior probability mass function $\rho(k)$ for k has infinite support, then 40
41 $p_0 \in \text{KL}(\Pi)$. If further $\rho(k) \lesssim e^{-\beta k}$, then the posterior is consistent for d_H .* 41

42 *If p_0 is itself a Bernstein polynomial, then the posterior converges at the rate 42
43 $n^{-1/2} \log n$ with respect to d_H .* 43

44 *If p_0 is twice continuously differentiable on $[0, 1]$ and bounded away from zero, then 44
45 the posterior converges at the rate $n^{-1/3} (\log n)^{5/6}$ with respect to d_H .* 45

1 5.1.5. Random histograms

2 Gasparini (1996) used the Dirichlet process to put a prior on histograms of different
3 bin width. The sample space is first partitioned into (possibly an infinite number of)
4 intervals of length h , where h is chosen from a prior. Mass is distributed to the intervals
5 according to a Dirichlet process, whose parameters $M = M_h$ and $G = G_h$ may depend
6 on h . Mass assigned to any interval is equally distributed over that interval. The method
7 corresponds to Dirichlet mixtures with a uniform kernel $\psi(x, \theta, h) = h^{-1}$, $x, \theta \in$
8 $[jh, (j + 1)h)$ for some j .

9 If $n_j(h)$ is the number of X_i 's in the bin $[jh, (j + 1)h)$, it is not hard to see
10 that the posterior is of the same form as the prior with $M_h G_h$ updated to $M_h G_h +$
11 $\sum_j n_j(h) I_{[jh, (j + 1)h)}$ and the prior density $\pi(h)$ of h changed to

$$\pi^*(h) = \frac{\pi(h) \prod_{j=1}^{\infty} (M_h G_h([jh, (j + 1)h)))^{(n_j(h)-1)}}{M_h + n}. \quad (5.9)$$

15 The predictive density with no observations is given by $\int f_h(x) \pi(h) dh$, where
16 $f_h(x) = h^{-1} \sum_{j=-\infty}^{\infty} G_h([jh, (j+1)h)) I_{[jh, (j+1)h)}(x)$. In view of the conjugacy prop-
17 erty, the predictive density given n observations can be easily written down. Let P_h stand
18 for the histogram of bin-width h obtained from the probability measure P . Assume that
19 $G_h(j)/G_h(j - 1) \leq K_h$. If $\int x^2 p_0(x) dx < \infty$ and $\lim_{h \rightarrow 0} \int p_0(x) \log \frac{p_0 h}{p_0} = 0$, then
20 the posterior is weakly consistent at p_0 . Gasparini (1996) also gave additional condi-
21 tions to ensure consistency of the posterior mean of p under d_H .

23 5.2. Gaussian process prior

24 For density estimation on a bounded interval I , Leonard (1978) defined a random den-
25 sity on I through $f(x) = \frac{e^{Z(x)}}{\int_I e^{Z(t)} dt}$, where $Z(x)$ is a Gaussian process with mean func-
26 tion $\mu(x)$ and covariance kernel $\sigma(x, x')$. Lenk (1988) introduces an additional parame-
27 ter ξ to obtain a conjugate family. It is convenient to introduce the intermediate lognor-
28 mal process $W(x) = e^{Z(x)}$. Denote the distribution of W by $LN(\mu, \sigma, 0)$. For each ξ de-
29 fine a positive valued random process $LN(\mu, \sigma, \xi)$ on I whose Radon–Nikodym deriva-
30 tive with respect to $LN(\mu, \sigma, 0)$ is $(\int_I W(x, \omega) dx)^\xi$. The normalization $f(x, \omega) =$
31 $\frac{W(x)}{\int_I W(t) dt}$ gives a random density and the distribution of this density under $LN(\mu, \sigma, \xi)$
32 is denoted by $LNS(\mu, \sigma, \xi)$. If X_1, \dots, X_n are i.i.d. f and $f \sim LNS(\mu, \sigma, \xi)$, then the
33 posterior is $LNS(\mu^*, \sigma, \xi^*)$, where $\mu^*(x) = \mu(x) + \sum_{i=1}^n \sigma(x_i, x)$ and $\xi^* = \xi - n$.

35 The interpretation of the parameters are somewhat unclear. Intuitively, for a station-
36 ary covariance kernel, a higher value of $\sigma(0)$ leads to more fluctuations in $Z(x)$ and
37 hence more noninformative. Local smoothness is controlled by $-\sigma''(0)$ – smaller value
38 implying a smoother curve. The parameter ξ , introduced somewhat unnaturally, is the
39 least understood. Apparently, the expression for the posterior suggests that $-\xi$ may be
40 thought of as the “prior sample size”.

42 5.3. Polya tree prior

43 A Polya tree prior satisfying $\sum_{m=1}^{\infty} a_m^{-1} < \infty$ admits densities a.s. by Kraft (1964) and
44 hence may be considered for density estimation. The posterior expected density is given
45

1 by

$$E(f(x)|X_1, \dots, X_n) = \alpha(x) \prod_{m=1}^{\infty} \frac{2a_m + 2N(B_{\epsilon(m)})}{2a_m + N(B_{\epsilon(m)}) + N(B'_{\epsilon(m)})}, \quad (5.10)$$

2
3
4
5 where $N(B_{\epsilon(m)})$ stand for the number of observations falling in $B_{\epsilon(m)}$, the set in the
6 m -level partition which contains x and $N(B'_{\epsilon(m)})$ is the number of observations falling
7 in its sibling $B'_{\epsilon(m)}$. From Theorem 3.1 of Ghosal et al. (1999c), it follows that under the
8 condition $\sum_{m=1}^{\infty} a_m^{-1/2} < \infty$, any p_0 with $\int p_0 \log(p_0/\alpha) < \infty$ satisfies $p_0 \in \text{KL}(\mathcal{I})$
9 and hence the weak consistency holds. Consistency under d_H has been obtained by
10 Barron et al. (1999) under the rather strong condition that $a_m = 8^m$. This high value
11 of 8^m appears to be needed to control the roughness of the Polya trees. Using the pseudo-
12 posterior distribution as described in Section 3, Walker and Hjort (2001) showed that the
13 posterior mean converges in d_H solely under the condition $\sum_{m=1}^{\infty} a_m^{-1/2} < \infty$. Interest-
14 ingly, they identify the posterior mean with the mean of a pseudo-posterior distribution
15 that also comes from a Polya tree prior with a different set of parameters.
16
17
18

19 **6. Regression function estimation**

20
21 Regression is one of the most important and widely used tool in statistical analysis.
22 Consider a response variable Y measured with some covariate X that may possibly
23 be multivariate. The regression function $f(x) = E(Y|X = x)$ describes the overall
24 functional dependence of Y on X and thus becomes very useful in prediction. Spatial
25 and geostatistical problems can also be formulated as regression problems. Classical
26 parametric models such as linear, polynomial and exponential regression models are in-
27 creasingly giving way to nonparametric regression model. Frequentist estimates of the
28 regression functions such as the kernel estimate, spline or orthogonal series estimators
29 have been in use for a long time and their properties have been well studied. Some non-
30 parametric Bayesian methods have also been developed recently. The Bayesian analysis
31 depends on the dependence structure of Y on X and are handled differently for different
32 regression models.
33

34 *6.1. Normal regression*

35
36 For continuous response, a commonly used regression model is $Y_i = f(X_i) + \epsilon_i$, where
37 ϵ_i are assumed to be i.i.d. mean zero Gaussian errors with unknown variance and be
38 independent of X_i 's. Leading nonparametric Bayesian techniques, among some others,
39 include those based on (i) Gaussian process prior, (ii) orthogonal basis expansion, and
40 (iii) free-knot splines.

41 Wahba (1978) considered a Gaussian process prior for f . The resulting Bayes esti-
42 mator is found to be a smoothing spline with the appropriate choice of the covariance
43 kernel of the Gaussian process. A commonly used prior for f is defined through the
44 stochastic differential equation $\frac{d^2 f(x)}{dx^2} = \tau \frac{dW(x)}{dx}$, where $W(x)$ is a Wiener process. The
45 scale parameter τ is given an inverse gamma prior while the intercept term $f(0)$ is given

1 an independent Gaussian prior. Ansley et al. (1993) described an extended state-space 1
2 representation for computing the Bayes estimate. Barry (1986) used a similar prior for 2
3 multiple covariates and provided asymptotic result for the Bayes estimator. 3

4 Another approach to putting a nonparametric prior on f is through an orthogonal 4
5 basis expansion of the form $f(x) = \sum_{j=1}^{\infty} b_j \psi_j(x)$ and then putting a prior on the 5
6 coefficients b_j 's. Smith and Kohn (1997) considered such an approach when the infi- 6
7 nite series is truncated at some predetermined finite stage k . Zhao (2000) considered a 7
8 sieve prior putting an infinitely supported prior on k . Shen and Wasserman (2001) in- 8
9 vestigated the asymptotic properties for this sieve prior and obtained a convergence rate 9
10 $n^{-q/(2q+1)}$ under some restriction on the basis function and for a Gaussian prior on the 10
11 b_j 's. Variable selection problem is considered in Shively et al. (1999) and Wood et al. 11
12 (2002a). Wood et al. (2002b) extended this approach to spatially adaptive regression, 12
13 while Smith et al. (1998) extended the idea to autocorrelated errors. 13

14 A free-knot spline approach is considered by Denison et al. (1998) and DiMatteo 14
15 et al. (2001). They modeled f as a polynomial spline of fixed order (usually cubic), 15
16 while putting a prior on the number of the knots, the location of the knots and the 16
17 coefficients of the polynomials. Since the parameter space is canonical, computations 17
18 are done through Monte Carlo averages while samples from the posterior distribution is 18
19 obtained by reversible jump MCMC algorithm of Green (1995). 19
20

21 6.2. Binary regression 21

22 In this case, $Y|X = x \sim \text{binom}(1, f(x))$ so that $f(x) = P(Y = 1|X = x) =$ 23
24 $E(Y|X = x)$. Choudhuri et al. (2004b) induced a prior on $f(x)$ by using a Gaussian 24
25 process $\eta(x)$ and mapping $\eta(x)$ into the unit interval as $f(x) = H(\eta(x))$ for some 25
26 strictly increasing continuous chosen "link function" H . The posterior distribution of 26
27 $f(x)$ is analytically intractable and the MCMC procedure depends on the choice of link 27
28 function. The most commonly used link function is the probit link in which H is the 28
29 standard normal cdf. In this case, an elegant Gibbs sampler algorithm is obtained by 29
30 introducing some latent variables following an idea of Albert and Chib (1993). 30

31 Let $\mathbf{Y} = (Y_1, \dots, Y_n)^T$ be the random binary observations measured along with the 31
32 corresponding covariate values $\mathbf{X} = (X_1, \dots, X_n)^T$. Let $\mathbf{Z} = (Z_1, \dots, Z_n)^T$ be some 32
33 unobservable latent variables such that conditional on the covariate values \mathbf{X} and the 33
34 functional parameter η , Z_i 's are independent normal random variables with mean $\eta(X_i)$ 34
35 and variance 1. Assume that the observations Y_i 's are functions of these latent variables 35
36 defined as $Y_i = I(Z_i > 0)$. Then, conditional on (η, \mathbf{X}) , Y_i 's are independent Bernoulli 36
37 random variables with success probability $\Phi(\eta(X_i))$ and thus leads to the probit link 37
38 model. Had we observed Z_i 's, the posterior distribution of η could have been obtained 38
39 analytically, which is also a Gaussian process by virtue of the conjugacy of the Gaussian 39
40 observation with a Gaussian prior for the mean. However, \mathbf{Z} is unobservable. Given the 40
41 data (\mathbf{Y}, \mathbf{X}) and the functional parameter η , Z_i 's are conditionally independent and 41
42 their distributions are truncated normal with mean $\eta(X_i)$ and variance 1, where Z_i is 42
43 right truncated at 0 if $Y_i = 0$, while Z_i is right truncated at 0 if $Y_i = 1$, then Z_i 43
44 is taken to be positive. Now, using the conditional distributions of $(\mathbf{Z}|\eta, \mathbf{Y}, \mathbf{X})$ and 44
45 $(\eta|\mathbf{Z}, \mathbf{Y}, \mathbf{X})$, a Gibbs sampler algorithm is formulated for sampling from the distribution 45

of $(\mathbf{Z}, \eta | \mathbf{Y}, \mathbf{X})$. Choudhuri et al. (2004b) also extended this Gibbs sampler algorithm to the link function that is a mixture of normal cdfs. These authors also showed that the posterior distribution is consistent under mild conditions, as stated below.

THEOREM 12. *Let the true response probability function $f_0(x)$ be continuous, $(d + 1)$ -times differentiable and bounded away from 0 and 1, and that the underlying Gaussian process has mean function and covariance kernel $(d + 1)$ -times differentiable, where d is the dimension of the covariate X . Assume that the range of X is bounded.*

If the covariate is random having a nonsingular density $q(x)$, then for any $\varepsilon > 0$, $\Pi(f: \int |f(x) - f_0(x)|q(x) dx > \varepsilon | X_1, Y_1, \dots, X_n, Y_n) \rightarrow 0$ in P_{f_0} -probability.

If the covariates are nonrandom, then for any $\varepsilon > 0$, $\Pi(f: n^{-1} \sum_{i=1}^n |f(X_i) - f_0(X_i)| > \varepsilon | Y_1, \dots, Y_n) \rightarrow 0$ in P_{f_0} -probability.

To prove the result, conditions of Theorem 3 and Theorem 5 respectively for random and nonrandom covariates, are verified. The condition on the Kullback–Leibler support is verified by approximating the function by a finite Karhunen–Loève expansion and by the nonsingularity of the multivariate normal distributions. The testing condition is verified on a sieve that is given by the maximum of f and its $(d + 1)$ derivatives bounded by some $M_n = o(n)$. The complement of the sieve has exponentially small prior probability if M_n is not of smaller order than $n^{1/2}$.

Wood and Kohn (1998) considered the integrated Wiener process prior for the probit transformation of f . The posterior is computed via Monte Carlo averages using a data augmentation technique as above. Yau et al. (2003) extended the idea to multinomial problems. Holmes and Mallick (2003) extended the free-knot spline approach to generalized multiple regression treating binary regression as a particular case.

A completely different approach to semiparametric estimation of f is to nonparametrically estimate the link function H while using a parametric form, usually linear, for $\eta(x)$. Observe that H is a nondecreasing function with range $[0, 1]$ and this is an univariate distribution function. Gelfand and Kuo (1991), and Newton et al. (1996) used a Dirichlet process prior for H . Mallick and Gelfand (1994) modeled H as a mixture of beta cdf's with a prior probability on the mixture weights, which resulted in smoother estimates. Basu and Mukhopadhyay (2000) modeled the link function as Dirichlet scale mixture of truncated normal cdf's. Posterior consistency results for these procedures were obtained by Amewou-Atisso et al. (2003).

7. Spectral density estimation

Let $\{X_t: t = 1, 2, \dots\}$ be a stationary time series with autocovariance function $\gamma(\cdot)$ and spectral density $f^*(\omega^*) = (2\pi)^{-1} \sum_{r=-\infty}^{\infty} \gamma(r) e^{-ir\omega^*}$, $-\pi < \omega^* \leq \pi$. To estimate f^* , it suffices to consider the function $f(\omega) = f^*(\pi\omega)$, $0 \leq \omega \leq 1$, by the symmetry of f^* . Because the actual likelihood of f is difficult to handle, Whittle (1957, 1962) proposed a “quasi-likelihood”

$$L_n(f | X_1, \dots, X_n) = \prod_{l=1}^v \frac{1}{f(\omega_l)} e^{-I_n(\omega_l)/f(\omega_l)}, \quad (7.1)$$

where $\omega_l = 2l/n$, v is the greatest integer less than or equal to $(n - 1)/2$, and $I_n(\omega) = |\sum_{t=1}^n X_t e^{-it\pi\omega}|^2 / (2\pi n)$ is the periodogram. A pseudo-posterior distribution may be obtained by updating the prior using this likelihood.

7.1. Bernstein polynomial prior

Normalizing f to $q = f/\tau$ with the normalizing constant $\tau = \int f$, Choudhuri et al. (2004a) induced a prior on f by first putting a Bernstein polynomial prior on q and then putting an independent prior on τ . Thus, the prior on f is described by the following hierarchical scheme:

- $f(\omega) = \tau \sum_{j=1}^k F((j - 1)/k, j/k) \beta(\omega; j, k - j + 1)$;
- $F \sim \text{Dir}(M, G)$, where G has a Lebesgue density g ;
- k has probability mass function $\rho(k) > 0$ for $k = 1, 2, \dots$;
- The distribution of τ has Lebesgue density π on $(0, \infty)$;
- F, k , and τ are a priori independent.

The pseudo-posterior distribution is analytically intractable and hence is computed by an MCMC method. Using the Sethuraman representation for F as in (2.2), (f, k, τ) may be reparameterized as $(\theta_1, \theta_2, \dots, Y_1, Y_2, \dots, k, \tau)$. Because the infinite series in (2.2) is almost surely convergent, it may be truncated at some large L . Then one may represent F as $F = \sum_{l=1}^L V_l \delta_{\theta_l} + (1 - V_1 - \dots - V_L) \delta_{\theta_0}$, where $\theta_0 \sim G$ and is independent of the other parameters. The last term is added to make F a distribution function even after the truncation. Now the problem reduces to a parametric one with finitely many parameters $(\theta_0, \theta_1, \dots, \theta_L, Y_1, \dots, Y_L, k, \tau)$. The functional parameter f may be written as a function of these univariate parameters as

$$f(\omega) = \tau \sum_{j=1}^k w_{j,k} \beta(\omega; j, k - j + 1), \tag{7.2}$$

where $w_{j,k} = \sum_{l=0}^L V_l I\{\frac{j-1}{k} < \theta_l \leq \frac{j}{k}\}$ and $V_0 = 1 - V_1 - \dots - V_L$. The posterior distribution of $(\theta_0, \theta_1, \dots, \theta_L, Y_1, \dots, Y_L, k, \tau)$ is proportional to

$$\left[\prod_{m=1}^v \frac{1}{f(2m/n)} e^{-U_m/f(2m/n)} \right] \left[\prod_{l=1}^L M(1 - y_l)^{M-1} \right] \left[\prod_{l=0}^L g(\theta_l) \right] \rho(k) \pi(\tau). \tag{7.3}$$

The discrete parameter k may be easily simulated from its posterior distribution given the other parameters. If the prior on τ is an inverse gamma distribution, then the posterior distribution of τ conditional on the other parameters is also inverse gamma. To sample from the posterior density of θ_i 's or Y_i 's conditional on the other parameters, Metropolis algorithm is within the Gibbs sampling step is used. The starting values of τ may be set to the sample variance divided by 2π , while the starting value of k may be set to some large integer K_0 . The approximate posterior mode of θ_i 's and Y_i 's given the starting values of τ and k may be considered as the starting values for the respective variables.

Let f_0^* be the true spectral density. Assume that the time series satisfies the conditions

- 1 (M1) the time series is Gaussian with $\sum_{r=0}^{\infty} r^{\alpha} \gamma(r); < \infty$ for some $\alpha > 0$; 1
2 (M2) for all ω^* , $f_0^*(\omega^*) > 0$; 2
3 3
4 and the prior satisfies 4
5
6 (P1) for all k , $0 < \rho(k) \leq C e^{-ck(\log k)^{1+\alpha'}}$ for some constants $C, c, \alpha' > 0$; 6
7 (P2) g is bounded, continuous, and bounded away from zero; 7
8 (P3) the prior on τ is degenerate at the true value $\tau_0 = \int f_0$. 8

9
10 Using the contiguity result of Choudhuri et al. (2004c), the following result was shown 10
11 by Choudhuri et al. (2004a) under the above assumptions. 11

12
13 THEOREM 13. For any $\varepsilon > 0$, $\Pi_n\{f^*: \|f^* - f_0^*\|_1 > \varepsilon\} \rightarrow 0$ in $P_{f_0^*}^n$ -probability, 13
14 where Π_n is the pseudo-posterior distribution computed using the Whittle likelihood of 14
15 and $P_{f_0^*}^n$ is the actual distribution of the data (X_1, \dots, X_n) . 15
16

17
18 REMARK 1. The conclusion of Theorem 13 still holds if the degenerated prior on τ is 18
19 replaced by a sequence of priors distribution that asymptotically bracket the true value, 19
20 that is, the prior support of τ is in $[\tau_0 - \delta_n, \tau_0 + \delta_n]$ for some $\delta_n \rightarrow 0$. A two-stage 20
21 empirical Bayes method, by using one part of the sample to consistently estimate τ and 21
22 the other part to estimate q , may be considered to construct the above asymptotically 22
23 bracketing prior. 23
24

25 7.2. Gaussian process prior 25

26
27 Since the spectral density is nonnegative valued function, a Gaussian process prior may 27
28 be assigned to $g(\omega) = \log(f(\omega))$. Because the Whittle likelihood in (7.1) arises by as- 28
29 suming that $I_n(\omega_l)$'s are approximately independent exponential random variables with 29
30 mean $f(\omega_l)$, one may obtain a regression model of the form $\log(I_n(\omega_l)) = g(\omega_l) + \epsilon_l$, 30
31 where the additive errors ϵ_l 's are approximately i.i.d. with the Gumbel distribution. 31
32

33 Carter and Kohn (1997) considered an integrated Wiener process prior for g . They 33
34 described an elegant Gibbs sampler algorithm for sampling from the posterior dis- 34
35 tribution. Approximating the distribution of ϵ_l 's as a mixture of five known normal 35
36 distribution, they introduced latent variables indicating the mixture components for the 36
37 corresponding errors. Given the latent variables, conditional posterior distribution of g 37
38 is obtained by a data augmentation technique. Given g , the conditional posterior distri- 38
39 bution of the latent variables are independent and samples are easily drawn from their 39
40 finite support. 40

41 Gangopadhyay et al. (1998) considered the free-not spline approach to modeling g . 41
42 In this case, the posterior is computed by the reversible jump algorithm of Green (1995). 42
43 Liseo et al. (2001) considered a Brownian motion process as prior on g . For sampling 43
44 from the posterior distribution, they considered the Karhunen–Lo ev e series expansion 44
45 for the Brownian motion and then truncated the infinite series to a finite sum. 45

8. Estimation of transition density

Estimation of the transition density of a discrete-time Markov process is an important problem. Let Π be a prior on the transition densities $p(y|x)$. Then the predictive density of a future observation X_{n+1} given the data X_1, \dots, X_n equals to $E(p(\cdot|X_n)|X_1, \dots, X_n)$, which is the Bayes estimate of the transition density p at X_n . The prediction problem thus directly relates to the estimation of the transition density.

Tang and Ghosal (2003) considered a mixture of normal model

$$p(y|x) = \int \phi_\sigma(y - H(x; \theta)) dF(\theta, \sigma), \tag{8.1}$$

where θ is possibly vector valued and $H(x; \theta)$ is a known function. Such models are analogous to the normal mixture models in the density estimation where the unknown probability density is modeled as $p(y) = \int \phi_\sigma(y - \mu) dF(\mu, \sigma)$. A reasonable choice for the link function H in (8.1) could be of the form $\tau + \gamma\psi(\delta + \beta x)$ for some known function ψ .

As in density estimation, this mixture model may be represented as

$$X_i \sim N(H(X_{i-1}; \theta_i), \sigma_i^2), \quad (\theta_i, \sigma_i) \stackrel{\text{i.i.d.}}{\sim} F. \tag{8.2}$$

Here, unlike a parametric model, the unknown parameters are varying along with the index of the observation, and are actually drawn as i.i.d. samples from an unknown distribution. Hence the model is “dynamic” as opposed to a “static” parametric mixture model.

Tang and Ghosal (2003) let the mixing distribution F have a Dirichlet process prior $\text{Dir}(M, G)$. As in density estimation, the hierarchical representation (8.2) helps develop Gibbs sampler algorithms for sampling from the posterior distribution. However, because of the nonstandard forms of the conditionals, special techniques, such as the “no gaps” algorithm of MacEachern and Muller (1998) need to be implemented.

To study the large sample properties of the posterior distribution, Tang and Ghosal (2003) extended Schwartz’s (1965) theorem to the context of an ergodic Markov processes. For simplicity, X_0 is assumed to be fixed below, although the conclusion extends to random X_0 also.

THEOREM 14. *Let $\{X_n, n \geq 0\}$ be an ergodic Markov process with transition density $p \in \mathcal{P}$ and stationary distribution π . Let Π be a prior on \mathcal{P} . Let $p_0 \in \mathcal{P}$ and π_0 be respectively the true values of p and π . Let U_n be a sequence of subsets of \mathcal{P} containing p_0 .*

Suppose that there exist a sequence of tests Φ_n , based on X_0, X_1, \dots, X_n for testing the pair of hypotheses $H_0: p = p_0$ against $H: p \in U_n^c$, and subsets $V_n \subset \mathcal{P}$ such that

- (i) p_0 is in the Kullback–Leibler support of Π , that is $\Pi\{p: K(p_0, p) < \varepsilon\} > 0$, where

$$K(p_0, p) = \iint \pi_0(x) p_0(y|x) \log \frac{p_0(y|x)}{p(y|x)} dy dx,$$

- 1 (ii) $\Phi_n \rightarrow 0$ a.s. $[P_{f_0}^\infty]$, 1
- 2 (iii) $\sup_{p \in U_n^c \cap V_n} E_p(1 - \Phi_n) \leq C_1 e^{-n\beta_1}$ for some constants C_1 and β_1 , 2
- 3 (iv) $\Pi(p \in V_n^c) \leq C_2 e^{-n\beta_2}$ for some constants C_2 and β_2 . 3
- 4 4

5 Then $\Pi(p \in U_n | X_0, X_1, \dots, X_n) \rightarrow 1$ a.s. $[P_0^\infty]$, where $[P_0^\infty]$ denote the distribution 5
6 of the infinite sequence (X_0, X_1, \dots) . 6

7 7

8 Assume that $p_0(y|x)$ is of the form (8.1). Let F_0 denote the true mixing distribu- 8
9 tion, and π_0 denote the corresponding invariant distribution. Let the sup- L_1 distance 9
10 on the space of transition probabilities be given by $d(p_1, p_2) = \sup_x \int |p_1(y|x) -$ 10
11 $p_2(y|x)| dy$. Let H be uniformly equicontinuous in x and the support of G be com- 11
12 pact containing the support of F_0 . Tang and Ghosal (2003) showed that (i) the test 12
13 $I\{\sum_{i=1}^k \log \frac{p_1(X_{2i}|X_{2i-1})}{p_0(X_{2i}|X_{2i-1})} > 0\}$, where $n = 2k$ or $2k + 1$, for testing p_0 against a small 13
14 ball around p_1 , has exponentially small error probabilities, (ii) the space of transition 14
15 probabilities supported by the prior is compact under the sup- L_1 distance, and (iii) the 15
16 Kullback–Leibler property holds at p_0 . By the compactness property, a single test can 16
17 be constructed for the entire alternative having exponentially small error probabilities. 17
18 It may be noted that because of the compactness of \mathcal{P} , it is not necessary to consider 18
19 sieves. Thus by Theorem 14, the posterior distribution is consistent at p_0 with respect 19
20 to the sup- L_1 distance. 20

21 The conditions assumed in the above result are somewhat stringent. For instance 21
22 $H(x, \beta, \delta, \tau) = \tau + \gamma\psi(\delta + \beta x)$, then ψ is necessarily bounded, ruling out the linear 22
23 link. If a suitable weaker topology is employed, Tang and Ghosal (2003) showed that 23
24 consistency can be obtained under weaker conditions by extending Walker’s (2004) 24
25 approach to Markov processes. More specifically, the Kullback–Leibler property holds 25
26 if H satisfies uniform equicontinuity on compact sets only. If now a topology is defined 26
27 by the neighborhood base $\{f: \int |\int g_i(y)f(y|x) dy - \int g_i(y)f_0(y|x) dy| \nu(x) dx <$ 27
28 $\varepsilon, i = 1, \dots, k\}$, where ν is a probability density, then consistency holds if σ is bounded 28
29 below and contains σ_0 in its support. If further σ is also bounded above and the θ is 29
30 supported on a compact set, then consistency also holds in the integrated- L_1 distance 30
31 integrated with respect to ν . For a linear link function $H(x, \rho, b) = \rho x + b, |\rho| < 1,$ 31
32 the compactness condition can be dropped, for instance, if the distribution of b under G 32
33 is normal. 33

34 34

35 35

36 36

37 37

38 38

39 39

40 40

41 41

42 42

43 43

44 44

45 45

9. Concluding remarks

39 In this article, we have reviewed Bayesian methods for the estimation of functions of 39
40 statistical interest such as the cumulative distribution function, density function, re- 40
41 gression function, spectral density of a time series and the transition density function 41
42 of a Markov process. Function estimation can be viewed as a problem of the estima- 42
43 tion of one or more infinite-dimensional parameter arising in a statistical model. It has 43
44 been argued that the Bayesian approach to function estimation, commonly known as 44
45 Bayesian nonparametric estimation, can provide an important, coherent alternative to 45

1 more familiar classical approaches to function estimation. We have considered the prob- 1
2 lems of construction of appropriate prior distributions on infinite-dimensional spaces. 2
3 It has been argued that, because of the lack of subjective knowledge about every de- 3
4 tail of a distribution in an infinite-dimensional space, some default mechanism of prior 4
5 specification needs to be followed. We have discussed various important priors on 5
6 infinite-dimensional spaces, and their merits and demerits. While certainly not exhaus- 6
7 tive, these priors and their various combinations provide a large catalogue of priors in a 7
8 statistician's toolbox, which may be tried and tested for various curve estimation prob- 8
9 lems including, but not restricted to, the problems we discussed. Due to the vastness 9
10 of the relevant literature and the rapid growth of the subject, it is impossible to even 10
11 attempt to mention all the problems of Bayesian curve estimation. The material pre- 11
12 sented here is mostly a reflection of the authors' interest and familiarity. Computa- 12
13 tion of the posterior distribution is an important issue. Due to the lack of useful analytical 13
14 expressions for the posterior distribution in most curve estimation problems, computa- 14
15 tion has to be done by some numerical technique, usually by the help of Markov chain 15
16 Monte Carlo methods. We described computing techniques for the curve estimation 16
17 problems considered in this chapter. The simultaneous development of innovative sam- 17
18 pling techniques and computing devices has brought tremendous computing power to 18
19 nonparametric Bayesians. Indeed, for many statistical problems, the computing power 19
20 of a Bayesian now exceeds that of a non-Bayesian. While these positive developments 20
21 are extremely encouraging, one should however be extremely cautious about naive uses 21
22 of Bayesian methods for nonparametric problems to avoid pitfalls. We argued that it 22
23 is important to validate the use of a particular prior by using some benchmark crite- 23
24 rion such as posterior consistency. We discussed several techniques of proving posterior 24
25 consistency and mentioned some examples of inconsistency. Sufficient conditions for 25
26 posterior consistency are discussed in the problems we considered. Convergence rates 26
27 of posterior distributions have also been discussed, together with the related concepts 27
28 of optimality, adaptation, misspecification and Bernstein–von Mises theorem. 28

29 The popularity of Bayesian nonparametric methods is rapidly growing among prac- 29
30 tioners as theoretical properties are increasingly better understood and the compu- 30
31 tational hurdles are being removed. Innovative Bayesian nonparametric methods for 31
32 complex models arising in biomedical, geostatistical, environmental, econometric and 32
33 many other applications are being proposed. Study of theoretical properties of non- 33
34 parametric Bayesian beyond the traditional i.i.d. set-up has started to receive attention 34
35 recently. Much more work will be needed to bridge the gap. Developing techniques of 35
36 model selection, the Bayesian equivalent of hypothesis testing, as well as the study of 36
37 their theoretical properties will be highly desirable. 37
38
39

40 References 40

- 41
42 Albert, J., Chib, S. (1993). Bayesian analysis of binary and polychotomous response data. *J. Amer. Statist.*
43 *Assoc.* **88**, 669–679. 43
44 Amewou-Atisso, M., Ghosal, S., Ghosh, J.K., Ramamoorthi, R.V. (2003). Posterior consistency for semipara- 44
45 metric regression problems. *Bernoulli* **9**, 291–312. 45

- 1 Ansley, C.F., Kohn, R., Wong, C. (1993). Nonparametric spline regression with prior information. *Bio-* 1
2 *metrika* **80**, 75–88. 2
- 3 Antoniak, C. (1974). Mixtures of Dirichlet processes with application to Bayesian non-parametric problems. 3
4 *Ann. Statist.* **2**, 1152–1174. 4
- 5 Barron, A.R. (1988). The exponential convergence of posterior probabilities with implications for Bayes 5
6 estimators of density functions. Unpublished manuscript. 5
- 6 Barron, A.R. (1999). Information-theoretic characterization of Bayes performance and the choice of priors 6
7 in parametric and nonparametric problems. In: Bernardo, J.M., et al. (Eds.), *Bayesian Statistics, vol. 6*. 7
8 Oxford University Press, New York, pp. 27–52. 8
- 8 Barron, A., Schervish, M., Wasserman, L. (1999). The consistency of posterior distributions in nonparametric 8
9 problems. *Ann. Statist.* **27**, 536–561. 9
- 10 Barry, D. (1986). Nonparametric Bayesian regression. *Ann. Statist.* **14**, 934–953. 10
- 11 Basu, S., Mukhopadhyay, S. (2000). Bayesian analysis of binary regression using symmetric and asymmetric 11
12 links. *Sankhyā, Ser. B* **62**, 372–387. 12
- 13 Belitser, E.N., Ghosal, S. (2003). Adaptive Bayesian inference on the mean of an infinite-dimensional normal 13
14 distribution. *Ann. Statist.* **31**, 536–559. 14
- 14 Berger, J.O., Guglielmi, A. (2001). Bayesian and conditional frequentist testing of a parametric model versus 14
15 nonparametric alternatives. *J. Amer. Statist. Assoc.* **96** (453), 174–184. 15
- 16 Berk, R. (1966). Limiting behavior of the posterior distribution when the model is incorrect. *Ann. Math.* 16
17 *Statist.* **37**, 51–58. 17
- 18 Birgé, L. (1983). Robust testing for independent non-identically distributed variables and Markov chains. 18
19 In: Florens, J.P., et al. (Eds.), *Specifying Statistical Models. From Parametric to Non-Parametric. Using*
20 *Bayesian or Non-Bayesian Approaches*. In: *Lecture Notes in Statistics*, vol. 16. Springer-Verlag, New
21 York, pp. 134–162. 20
- 21 Blackwell, D. (1973). Discreteness of Ferguson selection. *Ann. Statist.* **1**, 356–358. 21
- 22 Blackwell, D., Dubins, L.E. (1962). Merging of opinions with increasing information. *Ann. Math. Statist.* **33**,
23 882–886. 23
- 23 Blackwell, D., MacQueen, J.B. (1973). Ferguson distributions via Polya urn schemes. *Ann. Statist.* **1**, 353–
24 355. 24
- 25 Blum, J., Susarla, V. (1977). On the posterior distribution of a Dirichlet process given randomly right censored
26 observations. *Stochastic Process. Appl.* **5**, 207–211. 26
- 27 Brunner, L.J. (1992). Bayesian nonparametric methods for data from a unimodal density. *Statist. Probab.*
28 *Lett.* **14**, 195–199. 28
- 28 Brunner, L.J. (1995). Bayesian linear regression with error terms that have symmetric unimodal densities.
29 *J. Nonparametr. Statist.* **4**, 335–348. 29
- 30 Brunner, L.J., Lo, A.Y. (1989). Bayes methods for a symmetric unimodal density and its mode. *Ann. Sta-*
31 *tist.* **17**, 1550–1566. 31
- 32 Carter, C.K., Kohn, R. (1997). Semiparametric Bayesian inference for time series with mixed spectra. *J. Roy.*
33 *Statist. Soc., Ser. B* **59**, 255–268. 33
- 33 Choudhuri, N., Ghosal, S., Roy, A. (2004a). Bayesian estimation of the spectral density of a time series.
34 *J. Amer. Statist. Assoc.* **99**, 1050–1059. 34
- 35 Choudhuri, N., Ghosal, S., Roy, A. (2004b). Bayesian nonparametric binary regression with a Gaussian
36 process prior. Preprint. 36
- 37 Choudhuri, N., Ghosal, S., Roy, A. (2004c). Contiguity of the Whittle measure in a Gaussian time series.
38 *Biometrika* **91**, 211–218. 38
- 38 Cifarelli, D.M., Regazzini, E. (1990). Distribution functions of means of a Dirichlet process. *Ann. Statist.* **18**,
39 429–442. 39
- 40 Cox, D.D. (1993). An analysis of Bayesian inference for nonparametric regression. *Ann. Statist.* **21**, 903–923. 40
- 41 Dalal, S.R. (1979). Dirichlet invariant processes and applications to nonparametric estimation of symmetric
42 distribution functions. *Stochastic Process. Appl.* **9**, 99–107. 42
- 42 Denison, D.G.T., Mallick, B.K., Smith, A.F.M. (1998). Automatic Bayesian curve fitting. *J. Roy. Statist. Soc.,*
43 *Ser. B Stat. Methodol.* **60**, 333–350. 43
- 44 Diaconis, P., Freedman, D. (1986a). On the consistency of Bayes estimates (with discussion). *Ann. Statist.* **14**,
45 1–67. 45

- 1 Diaconis, P., Freedman, D. (1986b). On inconsistent Bayes estimates. *Ann. Statist.* **14**, 68–87. 1
- 2 DiMatteo, I., Genovese, C.R., Kass, R.E. (2001). Bayesian curve-fitting with free-knot splines. *Biometrika* **88**,
3 1055–1071. 3
- 4 Doksum, K.A. (1974). Tail free and neutral random probabilities and their posterior distributions. *Ann.*
5 *Probab.* **2**, 183–201. 4
- 6 Doob, J.L. (1948). Application of the theory of martingales. Coll. Int. du CNRS, Paris, pp. 22–28. 5
- 7 Doss, H. (1985a). Bayesian nonparametric estimation of the median. I. Computation of the estimates. *Ann.*
8 *Statist.* **13**, 1432–1444. 6
- 9 Doss, H. (1985b). Bayesian nonparametric estimation of the median. II. Asymptotic properties of the esti-
10 mates. *Ann. Statist.* **13**, 1445–1464. 7
- 11 Doss, H., Sellke, T. (1982). The tails of probabilities chosen from a Dirichlet prior. *Ann. Statist.* **10**, 1302–
12 1305. 8
- 13 Dykstra, R.L., Laud, P.W. (1981). A Bayesian nonparametric approach to reliability. *Ann. Statist.* **9**, 356–367. 9
- 14 Escobar, M. (1994). Estimating normal means with a Dirichlet process prior. *J. Amer. Statist. Assoc.* **89**,
15 268–277. 10
- 16 Escobar, M., West, M. (1995). Bayesian density estimation and inference using mixtures. *J. Amer. Statist.*
17 *Assoc.* **90**, 577–588. 11
- 18 Escobar, M., West, M. (1998). Computing nonparametric hierarchical models. In: *Practical Nonparametric*
19 *and Semiparametric Bayesian Statistics*. In: *Lecture Notes in Statistics*, vol. 133. Springer, New York,
20 pp. 1–22. 12
- 21 Ferguson, T.S. (1973). A Bayesian analysis of some nonparametric problems. *Ann. Statist.* **1**, 209–230. 13
- 22 Ferguson, T.S. (1974). Prior distribution on the spaces of probability measures. *Ann. Statist.* **2**, 615–629. 14
- 23 Ferguson, T.S. (1983). Bayesian density estimation by mixtures of Normal distributions. In: Rizvi, M.,
24 Rustagi, J., Siegmund, D. (Eds.), *Recent Advances in Statistics*, pp. 287–302. 15
- 25 Ferguson, T.S., Phadia, E.G. (1979). Bayesian nonparametric estimation based on censored data. *Ann. Sta-*
26 *tist.* **7**, 163–186. 16
- 27 Freedman, D. (1963). On the asymptotic distribution of Bayes estimates in the discrete case I. *Ann. Math.*
28 *Statist.* **34**, 1386–1403. 17
- 29 Freedman, D. (1965). On the asymptotic distribution of Bayes estimates in the discrete case II. *Ann. Math.*
30 *Statist.* **36**, 454–456. 18
- 31 Freedman, D. (1999). On the Bernstein–von Mises theorem with infinite-dimensional parameters. *Ann. Sta-*
32 *tist.* **27**, 1119–1140. 19
- 33 Fristedt, B. (1967). Sample function behavior of increasing processes with stationary independent increments.
34 *Pacific J. Math.* **21**, 21–33. 20
- 35 Fristedt, B., Pruitt, W.E. (1971). Lower functions for increasing random walks and subordinators. *Z. Wahsch.*
36 *Verw. Gebiete* **18**, 167–182. 21
- 37 Gangopadhyay, A.K., Mallick, B.K., Denison, D.G.T. (1998). Estimation of spectral density of a stationary
38 time series via an asymptotic representation of the periodogram. *J. Statist. Plann. Inference* **75**, 281–290. 22
- 39 Gasparini, M. (1996). Bayesian density estimation via Dirichlet density process. *J. Nonparametr. Statist.* **6**,
40 355–366. 23
- 41 Gelfand, A.E., Kuo, L. (1991). Nonparametric Bayesian bioassay including ordered polytomous response.
42 *Biometrika* **78**, 657–666. 24
- 43 Ghosal, S. (2000). Asymptotic normality of posterior distributions for exponential families with many para-
44 meters. *J. Multivariate Anal.* **74**, 49–69. 25
- 45 Ghosal, S. (2001). Convergence rates for density estimation with Bernstein polynomials. *Ann. Statist.* **29**,
1264–1280. 26
- Ghosal, S., van der Vaart, A.W. (2001). Entropies and rates of convergence for maximum likelihood and
Bayes estimation for mixtures of normal densities. *Ann. Statist.* **29**, 1233–1263. 27
- Ghosal, S., van der Vaart, A.W. (2003a). Convergence rates for non-i.i.d. observations. Preprint. 28
- Ghosal, S., van der Vaart, A.W. (2003b). Posterior convergence rates of Dirichlet mixtures of normal distrib-
utions for smooth densities. Preprint. 29
- Ghosal, S., Ghosh, J.K., Samanta, T. (1995). On convergence of posterior distributions. *Ann. Statist.* **23**, 2145–
2152. 30

- 1 Ghosal, S., Ghosh, J.K., Ramamoorthi, R.V. (1997). Noninformative priors via sieves and consistency. In: 1
2 Panchapakesan, S., Balakrishnan, N. (Eds.), *Advances in Statistical Decision Theory and Applications*. 2
3 Birkhäuser, Boston, pp. 119–132. 3
4 Ghosal, S., Ghosh, J.K., Ramamoorthi, R.V. (1999a). Consistency issues in Bayesian nonparametrics. In: 4
5 Ghosh, S. (Ed.), *Asymptotics, Nonparametrics and Time Series: A Tribute to Madan Lal Puri*. Marcel 5
6 Dekker, New York, pp. 639–668. 6
7 Ghosal, S., Ghosh, J.K., Ramamoorthi, R.V. (1999b). Posterior consistency of Dirichlet mixtures in density 7
8 estimation. *Ann. Statist.* **27**, 143–158. 8
9 Ghosal, S., Ghosh, J.K., Ramamoorthi, R.V. (1999c). Consistent semiparametric Bayesian inference about a 9
10 location parameter. *J. Statist. Plann. Inference* **77**, 181–193. 10
11 Ghosal, S., Ghosh, J.K., van der Vaart, A.W. (2000). Convergence rates of posterior distributions. *Ann. Sta-* 11
12 *tist.* **28**, 500–531. 12
13 Ghosal, S., Lember, Y., van der Vaart, A.W. (2003). On Bayesian adaptation. *Acta Appl. Math.* **79**, 165–175. 13
14 Ghosh, J.K., Ramamoorthi, R.V. (1995). Consistency of Bayesian inference for survival analysis with or 14
15 without censoring. In: Koul, H. (Ed.), *Analysis of Censored Data*. In: *IMS Lecture Notes Monograph* 15
16 *Series*, vol. 27. Inst. Math. Statist., Hayward, CA, pp. 95–103. 16
17 Ghosh, J.K., Ramamoorthi, R.V. (2003). *Bayesian Nonparametrics*. Springer-Verlag, New York. 17
18 Ghosh, S.K., Ghosal, S. (2003). Proportional mean regression models for censored data. Preprint. 18
19 Ghosh, J.K., Ghosal, S., Samanta, T. (1994). Stability and convergence of posterior in non-regular problems. 19
20 In: Gupta, S.S., Berger, J.O. (Eds.), *Statistical Decision Theory and Related Topics V*. Springer-Verlag, 20
21 New York, pp. 183–199. 21
22 Green, P. (1995). Reversible jump Markov chain Monte Carlo computation and Bayesian model determina- 22
23 tion. *Biometrika* **82**, 711–732. 23
24 Grenander, U. (1981). *Abstract Inference*. Wiley, New York. 24
25 Hanson, T., Johnson, W.O. (2002). Modeling regression error with a mixture of Polya trees. *J. Amer. Statist.* 25
26 *Assoc.* **97**, 1020–1033. 26
27 Hjort, N.L. (1990). Nonparametric Bayes estimators based on beta processes in models for life history data. 27
28 *Ann. Statist.* **18**, 1259–1294. 28
29 Hjort, N.L. (1996). Bayesian approaches to non- and semiparametric density estimation. In: Bernardo, J., et 29
30 al. (Eds.), *Bayesian Statistics, vol. 5*, pp. 223–253. 30
31 Hjort, N.L. (2000). Bayesian analysis for a generalized Dirichlet process prior. Preprint. 31
32 Hjort, N.L. (2003). Topics in nonparametric Bayesian statistics (with discussion). In: Green, P.J., Hjort, N., 32
33 Richardson, S. (Eds.), *Highly Structured Stochastic Systems*. Oxford University Press, pp. 455–487. 33
34 Holmes, C.C., Mallick, B.K. (2003). Generalized nonlinear modeling with multivariate free-knot regression 34
35 splines. *J. Amer. Statist. Assoc.* **98**, 352–368. 35
36 Huang, T.Z. (2004). Convergence rates for posterior distributions and adaptive estimation. *Ann. Statist.* **32**, 36
37 1556–1593. 37
38 Ibragimov, I.A., Has'minskii, R.Z. (1981). *Statistical Estimation: Asymptotic Theory*. Springer-Verlag, New 38
39 York. 39
40 Iswaran, H., Zarepour, M. (2002a). Exact and approximate sum representation for the Dirichlet process. 40
41 *Canad. J. Statist.* **26**, 269–283. 41
42 Iswaran, H., Zarepour, M. (2002b). Dirichlet prior sieves in finite normal mixture models. *Statistica Sinica*, 42
43 269–283. 43
44 Kim, Y. (1999). Nonparametric Bayesian estimators for counting processes. *Ann. Statist.* **27**, 562–588. 44
45 Kim, Y., Lee, J. (2001). On posterior consistency of survival models. *Ann. Statist.* **29**, 666–686. 45
46 Kim, Y., Lee, J. (2004). A Bernstein–von Mises theorem in the nonparametric right-censoring model. *Ann.* 46
47 *Statist.* **32**, 1492–1512. 47
48 Kleijn, B., van der Vaart, A.W. (2002). Misspecification in infinite-dimensional Bayesian statistics. Preprint. 48
49 Kottas, A., Gelfand, A.E. (2001). Bayesian semiparametric median regression modeling. *J. Amer. Statist.* 49
50 *Assoc.* **96**, 1458–1468. 50
51 Kraft, C.H. (1964). A class of distribution function processes which have derivatives. *J. Appl. Probab.* **1**, 51
52 385–388. 52
53 Lavine, M. (1992). Some aspects of Polya tree distributions for statistical modeling. *Ann. Statist.* **20**, 1222– 53
54 1235. 54

- 1 Lavine, M. (1994). More aspects of Polya tree distributions for statistical modeling. *Ann. Statist.* **22**, 1161–
2 1176. 2
- 3 Le Cam, L.M. (1986). *Asymptotic Methods in Statistical Decision Theory*. Springer-Verlag, New York. 3
- 4 Le Cam, L., Yang, G.L. (2000). *Asymptotics in Statistics*, second ed. Springer-Verlag. 4
- 5 Lenk, P.J. (1988). The logistic normal distribution for Bayesian, nonparametric, predictive densities. *J. Amer.*
6 *Statist. Assoc.* **83**, 509–516. 5
- 7 Lenk, P.J. (1991). Towards a practicable Bayesian nonparametric density estimator. *Biometrika* **78**, 531–543. 6
- 8 Leonard, T. (1978). Density estimation, stochastic processes, and prior information. *J. Roy. Statist. Soc., Ser.*
9 *B* **40**, 113–146. 7
- 10 Liseo, B., Marinucci, D., Petrella, L. (2001). Bayesian semiparametric inference on long-range dependence.
11 *Biometrika* **88**, 1089–1104. 8
- 12 Lo, A.Y. (1982). Bayesian nonparametric statistical inference for Poisson point process. *Z. Wahsch. Verw.*
13 *Gebiete* **59**, 55–66. 9
- 14 Lo, A.Y. (1983). Weak convergence for Dirichlet processes. *Sankhyā, Ser. A* **45**, 105–111. 10
- 15 Lo, A.Y. (1984). On a class of Bayesian nonparametric estimates I: Density estimates. *Ann. Statist.* **12**, 351–
16 357. 11
- 17 Lo, A.Y. (1986). A remark on the limiting posterior distribution of the multiparameter Dirichlet process.
18 *Sankhyā, Ser. A* **48**, 247–249. 12
- 19 MacEachern, S.N., Muller, P. (1998). Estimating mixture of Dirichlet process models. *J. Comput. Graph.*
20 *Statist.* **7**, 223–228. 13
- 21 Mallick, B.K., Gelfand, A.E. (1994). Generalized linear models with unknown link functions. *Biometrika* **81**,
22 237–245. 14
- 23 Mauldin, R.D., Sudderth, W.D., Williams, S.C. (1992). Polya trees and random distributions. *Ann. Statist.* **20**,
24 1203–1221. 15
- 25 Muliere, P., Tardella, L. (1998). Approximating distributions of functionals of Ferguson–Dirichlet priors.
26 *Canad. J. Statist.* **30**, 269–283. 16
- 27 Newton, M.A., Czado, C., Chappell, R. (1996). Bayesian inference for semiparametric binary regression.
28 *J. Amer. Statist. Assoc.* **91**, 142–153. 17
- 29 Nieto-Barajas, L.E., Walker, S.G. (2004). Bayesian nonparametric survival analysis via Lévy driven Markov
30 process. *Statistica Sinica* **14**, 1127–1146. 18
- 31 Petrone, S. (1999a). Random Bernstein polynomials. *Scand. J. Statist.* **26**, 373–393. 19
- 32 Petrone, S. (1999b). Bayesian density estimation using Bernstein polynomials. *Canad. J. Statist.* **26**, 373–393. 20
- 33 Petrone, S., Veronese, P. (2002). Nonparametric mixture priors based on an exponential random scheme.
34 *Statist. Methods Appl.* **11**, 1–20. 21
- 35 Petrone, S., Wasserman, L. (2002). Consistency of Bernstein polynomial posteriors. *J. Roy. Statist. Soc., Ser.*
36 *B* **64**, 79–100. 22
- 37 Regazzini, E., Guglielmi, A., Di Nunno, G. (2002). Theory and numerical analysis for exact distributions of
38 functionals of a Dirichlet process. *Ann. Statist.* **30**, 1376–1411. 23
- 39 Rubin, D. (1981). The Bayesian bootstrap. *Ann. Statist.* **9**, 130–134. 24
- 40 Schwartz, L. (1965). On Bayes procedures. *Z. Wahsch. Verw. Gebiete* **4**, 10–26. 25
- 41 Sethuraman, J. (1994). A constructive definition of Dirichlet priors. *Statistica Sinica* **4**, 639–650. 26
- 42 Sethuraman, J., Tiwari, R. (1982). Convergence of Dirichlet measures and interpretation of their parameters.
43 In: Gupta, S.S., Berger, J.O. (Eds.), *Statistical Decision Theory and Related Topics. III, vol. 2*. Academic
44 Press, New York, pp. 305–315. 27
- 45 Shen, X. (2002). Asymptotic normality of semiparametric and nonparametric posterior distributions. *J. Amer.*
Statist. Assoc. **97**, 222–235. 28
- Shen, X., Wasserman, L. (2001). Rates of convergence of posterior distributions. *Ann. Statist.* **29**, 687–714. 29
- Shively, T.S., Kohn, R., Wood, S. (1999). Variable selection and function estimation in additive nonparametric
regression using a data-based prior (with discussions). *J. Amer. Statist. Assoc.* **94**, 777–806. 30
- Smith, M., Kohn, R. (1997). A Bayesian approach to nonparametric bivariate regression. *J. Amer. Statist.*
Assoc. **92**, 1522–1535. 31
- Smith, M., Wong, C., Kohn, R. (1998). Additive nonparametric regression with autocorrelated errors. *J. Roy.*
Statist. Soc., Ser. B **60**, 311–331. 32

- 1 Susarla, V., Van Ryzin, J. (1976). Nonparametric Bayesian estimation of survival curves from incomplete observations. *J. Amer. Statist. Assoc.* **71**, 897–902. 1
- 2 2
- 3 Susarla, V., Van Ryzin, J. (1978). Large sample theory for a Bayesian nonparametric survival curve estimator based on censored samples. *Ann. Statist.* **6**, 755–768. 3
- 4 4
- 5 Tang, Y., Ghosal, S. (2003). Posterior consistency of Dirichlet mixtures for estimating a transition density. Preprint. 5
- 6 Tokdar, S.T. (2003). Posterior consistency of Dirichlet location-scale mixtures of normals in density estimation and regression. Preprint. 6
- 7 7
- 8 van der Vaart, A.W. (1998). *Asymptotic Statistics*. Cambridge University Press. 8
- 9 van der Vaart, A.W., Wellner, J.A. (1996). *Weak Convergence and Empirical Processes*. Springer-Verlag, New York. 9
- 10 Wahba, G. (1978). Improper priors, spline smoothing and the problem of guarding against model errors in regression. *J. Roy. Statist. Soc., Ser. B* **40**, 364–372. 10
- 11 11
- 12 Walker, S.G. (2003). On sufficient conditions for Bayesian consistency. *Biometrika* **90**, 482–490. 12
- 13 Walker, S.G. (2004). New approaches to Bayesian consistency. *Ann. Statist.* **32**, 2028–2043. 13
- 14 Walker, S.G., Hjort, N.L. (2001). On Bayesian consistency. *J. Roy. Statist. Soc., Ser. B* **63**, 811–821. 14
- 15 Walker, S.G., Muliere, P. (1997). Beta-Stacy processes and a generalization of the Polya-urn scheme. *Ann. Statist.* **25**, 1762–1780. 15
- 16 Wasserman, L. (1998). Asymptotic properties of nonparametric Bayesian procedures. In: Dey, D., et al. (Eds.), *Practical Nonparametric and Semiparametric Bayesian Statistics*. In: *Lecture Notes in Statistics*, vol. 133. Springer-Verlag, New York, pp. 293–304. 16
- 17 17
- 18 Whittle, P. (1957). Curve and periodogram smoothing. *J. Roy. Statist. Soc., Ser. B* **19**, 38–63. 18
- 19 Whittle, P. (1962). Gaussian estimation in stationary time series. *Bull. Int. Statist. Inst.* **39**, 105–129. 19
- 20 Wood, S., Kohn, R. (1998). A Bayesian approach to robust binary nonparametric regression. *J. Roy. Statist. Soc., Ser. B* **93**, 203–213. 20
- 21 21
- 22 Wood, S., Kohn, R., Shively, T., Jiang, W. (2002a). Model selection in spline nonparametric regression. *J. Roy. Statist. Soc., Ser. B* **64**, 119–139. 22
- 23 23
- 24 Wood, S., Jiang, W., Tanner, M. (2002b). Bayesian mixture of splines for spatially adaptive nonparametric regression. *Biometrika* **89**, 513–528. 24
- 25 Yau, P., Kohn, R., Wood, S. (2003). Bayesian variable selection and model averaging in high-dimensional multinomial nonparametric regression. *J. Comput. Graph. Statist.* **12**, 23–54. 25
- 26 26
- 27 Zhao, L.H. (2000). Bayesian aspects of some nonparametric problems. *Ann. Statist.* **28**, 532–552. 27
- 28 28
- 29 29
- 30 30
- 31 31
- 32 32
- 33 33
- 34 34
- 35 35
- 36 36
- 37 37
- 38 38
- 39 39
- 40 40
- 41 41
- 42 42
- 43 43
- 44 44
- 45 45