# Bayesian Encompassing Specification Tests of a Parametric Model against a Nonparametric Alternative

J.P. Florens[1]     J.F. Richard[2]     J.M. Rolin[3]

June 15, 2000

[1]GREMAQ and IDEI, Université des Sciences Sociales de Toulouse, France.
[2]Department of Economics, University of Pittsburg, USA.
[3]Institut de Statistique and CORE, Université Catholique de Louvain, Belgium.

## Abstract

Encompassing tests of a model $M_0$ are based upon the notion that $M_0$ ought to be able to account for results derived upon alternative models. Within a Bayesian framework, posterior distributions obtained under an alternative model $M_1$ are explicitly compared with posterior distributions obtained within $M_0$ by means of a suitable "transition" distribution for the parameters of $M_1$ conditionnally on those of $M_0$.

In our paper, we propose a generic Bayesian procedure for testing the capabilities of a parametric model $M_0$ to encompass a wide range of inferential results derived under a nonparametric alternative $M_1$. By combining the use of Dirichlet prior measures on $M_1$ with that of Monte Carlo simulation techniques, we develop exact operational and highly °exible Bayesian encompassing test procedures of $M_0$ relative to $M_1$. An application to an exponential lifetime model $M_0$ illustrates the performance of our procedure.

## Aknowledgement

1

# 1 Introduction

The objective of our paper is that of proposing a generic Bayesian procedure for testing the validity of a parametric model $M_0$. Specifically, we rely upon Monte Carlo simulation techniques to construct an exact Bayesian test of whether or not $M_0$ encompasses a general nonparametric model $M_1$.

As discussed by Florens et al (1996) Bayesian encompassing corresponds to a concept of sufficiency among models, dual to the traditional concept of sufficiency (among statistics) as defined by Blackwell (1951,1953) and Lecam (1964) and obtained from the latter by interchanging parameters and statistics. The model $M_0$ encompasses (Bayesian) inference results derived on $M_1$ if the latter can be reproduced within $M_0$ without any additional data treatment beyond that already necessitated by the analyis of $M_0$ itself. In other words, $M_1$ then becomes "inferentially redundant" relatice to $M_0$. In a wide range of scientific disciplines the capability of one's current model (or theory) to encompasses findings derived under other models is perceived to be a critical component of its validation.

Conversely, failure to encompass results obtained on $M_1$ constitutes prima facie evidence of "deficiencies" of $M_0$ which, if deemed relevant to one's objectives, would indicate a need for further refinements of $M_0$. It does not, however, warrant the conclusion that $M_1$ is to be preferred to $M_0$ since the requirement that $M_0$ ought to encompass rival models is by no means conditional upon the validity of the latter. Significant advances in sciences have resulted from the need to account for "anomalous" findings derived under invalid models and/or theories.

This latter comment illustrates a fundamental difference between an encompassing test, whereby the model $M_1$ is essentially instrumental in the construction of the test but generally does not constitute an acceptable alternative to $M_0$ itself, and a Bayesian model selection procedure whereby one chooses between two alternative models $M_0$ and $M_1$ which are generally treated symmetrically relative to one another (notwithstanding the obvious fact that loss function and prior probability assessments can be used to favor one model relative to the other). Actually, as we shall discuss further later, encompassing tests are based upon a comparison between posterior densities while model selection depends upon a ratio of predictive densities (Bayes factor, see e.g. Florens and Mouchart (1993)) This difference will turn out to be highly significant in the context discussed hereafter, whereby $M_0$ is a parametric model and $M_1$ a nonparameteric one - a situation which is aimed at providing a flexible and yet quite stringent test of the validity of $M_0$ by examining its capability to encompass a potentially broad range of results derived under a general nonparametric specification for $M_1$. In particular,

the use of a natural conjugate prior Dirichlet process for $M_1$ - which assigns probability one to discrete measure but is, nevertheless, dense within a set of continuous measures - has critical and somewhat "extreme" implications for model selection while it produces a fully operational fundamentally well-behaved and highly flexible encompassing test procedure.

In particular, flexibility originates from the fact that our encompassing test procedure allows its user to select specific functionals of interest (moment, tail probabilities,...) whose posterior densities on $M_1$ have to be encompassed by $M_0$, i.e. to decide which findings associated with $M_1$ are particularly relevant to the proprietor of $M_0$ in any given (decisional) application. In constrast, Bayesian model selection procedures are based upon a single Bayes factor and often provide limited insights as to which specific aspects of the competing models drove the final decision.

Finally, we wish to mention here that the encompassing tests we propose prolong a long tradition in sampling statistics, initiated at the turn of the century by K. Pearson with chi-squared test statistics, extended in the thirties with the Kolmogorov-Smirnov and Cramer von Mises test statistics and broadly labeled today as "goodness of fit" tests. Recent extensions are discussed e.g. by Revesz (1984). Additional related contributions are specifically discussed in the sequel of our paper after the relevant concepts and notation have been introduced.

Our paper is organized as follows : The models and notation are introduced in section 2; Posterior odds are derived in section 3; Encompassing is discussed in section 4 (principle in section 4.1. and nonparametric version in section 4.2, respectively); an application to the exponential distribution is presented in section 5 and section 6 concludes.

## 2   Models and notation

The observation $x = (x_1, ..., x_n)$ is an i.i.d. sample of a random variable $X$ with support $S$ in $\mathbb{R}^p$. Let $M_0$ denote the parametric model to be (in)validated and $M_1$ a nonparametric model to be encompassed. On $M_0$ all probability are assumed to be absolutely continuous with respect to an appropriate $\frac{3}{4}$-finite dominating measure and are represented by the corresponding density functions. On $M_1$ we are dealing with probability measures on a sampling distribution $F$. In order to unify the presentation of our results and at the cost of an abuse of notation, we shall equally represent such probabilities by a "density function" on $F$. Densities on a (finite-dimensional) parameter $\mu$, on a sample $x$ and on a distribution $F$ are denoted $\pi_0$, $p$ and $\varphi$, respectively. In addition, $p$ and $\varphi$ are subscripted with the index of the

relevant model. Conditioning is represented in the usual way. The symbol $\gg$ reads as "is distributed as".

The sampling distribution of $X$ on $M_0$ is denoted by $F^\mu$ with parameter $\mu$ in $\mathcal{E} \subseteq \mathbb{R}^k$ and with density $f(x|\mu)$. The sampling density of the sample $x$ is then given by

$$p_0(x|\mu) = \prod_{i=1}^{n} f(x_i|\mu): \qquad (1)$$

Let $\pi_0(\mu)$; $p_0(x)$ and $\pi_0(\mu|x)$ denote the prior density of $\mu$, the predictive density of $x$ and the posterior density of $\mu$, respectively. Following Bayes theorem, we have

$$\pi_0(\mu) \cdot p_0(x|\mu) = \pi_0(\mu|x) \cdot p_0(x): \qquad (2)$$

The sampling distribution of $X$ on $M_1$ is denoted by $F$. Following Ferguson (1973), a natural conjugate prior distribution for $F$ is provided by a Dirichlet measure which is parametrized by a real positive number $n_a$ and a probability $F_a$ on $\mathbb{R}^p$, say

$$F \gg Di(n_a F_a): \qquad (3)$$

In view of its importance for the sequel of the discussion, we briefly discuss here the interpretation of the Dirichlet measure (3). Consider first a fixed partition $B = (B_1; ...; B_R)$ of the support $S$ of $F$ and let $\pi^0 = (\pi_1; ...; \pi_R)$, with $\pi_r = F(B_r)$, denote the corresponding vector of sampling probabilities. A sufficient statistic for the sample $x$ is given by the vector $f^0 = (f_1; ...; f_R)$, where $f_r$ denotes the proportion of observations in $B_r$. The random variable $nf$ has multinomial distribution with parameter(s) $\pi$ (and $n$). Its density function is given by

$$g(nf|\pi) = n! \prod_{r=1}^{R} \frac{(\pi_r)^{nf_r}}{(nf_r)!}: \qquad (4)$$

A natural conjugate distribution for $\pi$ is given by a Dirichlet distribution whose support is the simplex

$$S_R = \{\pi; \pi_r > 0; \sum_{r=1}^{R} \pi_r = 1\} \qquad (5)$$

4

of dimension $R - 1$, with parameters $n_a > 0$ and $f_a \in S_R$ and whose density function is given by

$$\pi_1(\mu) = \Gamma(n_a) \prod_{r=1}^{R} \frac{(\mu_r)^{n_a f_{a,r} - 1}}{\Gamma(n_a f_{a,r})} : \tag{6}$$

The $M_1$-posterior distribution of $\mu$ is itself a Dirichlet distribution with parameters

$$n_p = n_a + n; \quad n_p f_p = n_a f_a + n f : \tag{7}$$

The Dirichlet prior measure in (3) is then de¯ned as one which assigns a Dirichlet distribution to the random vector $F(B) = (F(B_1); ::::; F(B_R))$ associated with any arbitrary partition $B$ of the support $S$ of $\mu$. That is to say that the Dirichlet prior measure in (3) assigns the following density to $F(B)$ :

$$\pi_1(F(B)) = \Gamma(n_a) \prod_{r=1}^{R} \frac{[F(B_r)]^{n_a F_a(B_r) - 1}}{\Gamma(n_a F_a(B_r))} \tag{8}$$

with support $S_R$. It follows from our discussion that the posterior distribution of $F$ on $M_1$, is a Dirichlet measure with parameters

$$n_p = n_a + n; \quad n_p F_p = n_a F_a + n F_n; \tag{9}$$

where $F_n$ is the empirical distribution of the sample

$$F_n = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}; \tag{10}$$

and $\delta_{x_i}$ denotes the Dirac measure at $x_i$ which is such that $\delta_{x_i}(A) = \mathbb{1}_A(x_i)$ for $A \subseteq \mathbb{R}^p$. Following Ferguson (1973), the predictive distribution $p_1(x)$ on $M_1$ can be represented in a variety of ways. The following sequential representation proves to be the most operational for our analysis

$$x_1 \sim F_a$$
$$\tag{11}$$
$$x_{i+1} | x_1; ::::; x_i \sim \frac{n_a}{n_a + i} F_a + \frac{i}{n_a + i} F_i; \quad \text{for } i : 1 \! \to n - 1 :$$

Within an encompassing framework it may be desirable to achieve some degree of "coherency" between the prior distributions of the two models. The concept of coherency introduced by Florens et al (1996) is too strong in that it requires that the two priors be linked together by a transition which actually depends upon sample size. An alternative approach consists of choosing the prior parameter $F_a$ in such a way that the two models be indistiguishable from one another on the basis of a single observation $x_1$ (i.e. that the Bayes factor associated with a single observation be equal to one). In view of the representation of the $M_1$ predictive in (11), this is equivalent to setting $F_a$ equal to the former that is to say

$$F_a(B) = \int_B p_0(x_1)dx_1; \quad \forall B \ 2 \ B_p \tag{12}$$

De¯nition 2.1 The prior measure $Di(n_a F_a)$ is predictive-coherent if and only if condition (12) holds.

¥

Note that, independently of whether or not condition (12) holds, $F_a(B)$ represents the prior expectation of $F(B)$ on $M_1$. As usual within a natural conjugate framework and as illustrated by (9), $n_a$ represents the prior "weight" attached to $F_a$ in a metric which is directly comparable to the actual sample size n and is, therefore, often refered to as to a "hypothetical" sample size - see e.g. Rai®a and Schlaifer (1961).

Our motivation for selecting a prior Dirichlet measure on $M_1$ is fourfold : (i) It is easy to specify since, as we just discussed, its parameters have "natural" interpretations ; (ii) It is °exible in that, following Ferguson (1973) it generates a dense subset of probabilities ; (iii) It is natural conjugate so that its two parameters are revised by means of the simple (additive) convolution rule given by (9) ; and (iv), it will prove especially operational for the construction of simulation based encompassing test procedures.

Nevertheless, two important characteristics of the Dirichlet measure deserve additional comments. Firstly, Dirichlet measures generate discrete probabilities. Actually, this property is intrinsic to inference in nonparametric models since the empirical distribution which is itself discrete constitutes a su±cient statistic and, therefore, serves as the basis of numerous goodness of ¯t tests, such as the Kolmogorov-Smirnov test. One might consider instead using prior measures whose realizations are continuous. Several approaches have been explored in the recent literature. One consists of introducing prior measures whose trajectories are the convolution of a Dirichlet realization and a of preassigned continuous probability -see e.g. Lo (1984) or

Escobar and West (1995). From our viewpoint, the main drawback of that approach lies in the complexity of the resulting posterior distributions (for example, posterior expectations are linear combinations comprising as many terms as there are partitions of the set $\{1, ..., n\}$). An alternative approach consists of selecting a logistic normal measure -see e.g. Leonard (1978) or Lenk (1988, 1991). Here again posterior distributions are complicated and, furthermore, this prior specification produces Bayesian procedures which are not convergent. See also Lavine (1992, 1994) for an alternative way of generating continuous probabilities. Preliminary investigations suggest that there are no conceptual difficulties generalizing the encompassing test procedures which are developed here to such continuous extensions and, furthermore, that actual computation are likely to prove simpler than those necessitated e.g. to produce posterior odds.

Secondly, as illustrated by (11), predictive probabilities associated with Dirichlet prior measures are intrinsically "bumpy" in the sense that they generate ties with positive probabilities. As we shall discuss further below, this property has "extreme" implications for (Bayesian) model selection procedures which essentially reduce to a test of continuity versus discreteness. It also produces infinite divergence measures in a "brute force" comparison between posterior densities. On the other hand functionals of $F$ admit posterior densities on $M_1$ under mild technical conditions and will, therefore, result in encompassing test procedures which are no longer bumpy.

# 3   Posterior odds

Let $\pm \in \{0, 1\}$ denote a binary model indicator. Let $P$ denote a probability distribution on $(\pm, X)$ which is such that

$$P(\pm = 0) = a, \quad P(\cdot | \pm = i) = P_i, \tag{13}$$

where $P_i$ denotes the predictive probability on $M_i$, for $i = 1, 2$. It is always possible to dominate $P_0$ and $P_1$ by a $\sigma$-finite measure such as, for example, $\frac{1}{2}(P_0 + P_1)$. Let, therefore, $dP_i$ denote the density of $P_i$ w.r.t. a suitable dominating measure[1]. The corresponding Bayes factor is then given by the ratio $dP_0/dP_1$ (not to be confused with the Radom-Nikodyn derivative of $P_0$ with respect to $P_1$ which does not exist in the present context). The posterior odds ratio is given by the product of the prior odds ratio $a/(1 - a)$ by the Bayes factor and the posterior probability of model $M_0$ equals

---

[1]We need to draw an explicit distinction between $dP_i$ and the predictive densities $p_i$ introduced earlier, since they are actually related to different dominating measures.

$$P(\pm = 0jx) = \frac{adP_0}{adP_0 + (1 \text{ } \text{¡} \text{ } a)dP_1}: \tag{14}$$

We consider the case where $F^\mu$ is absolutely continuous with respect to the Lebesgue measure in $\mathbb{R}^p$, with density function $f(xj\mu)$. The corresponding predictive densities on $M_0$ for a single observation $x_i$ and for the sample $x$ are given by

$$p_0(x_i) = \int_\mu f(x_ij\mu) {}^1{}_0(\mu)d\mu \tag{15}$$

and

$$p_0(x) = \int_{\text{£ } i=1}^{\text{¥}} f(x_ij\mu) {}^1{}_0(\mu)d\mu; \tag{16}$$

respectively. The posterior probability of model $M_0$ under a predictive coherent prior Dirichlet measure is then characterized in the following theorem.

**Theorem 3.1** : If $F^\mu$ is absolutely continuous w.r.t. the Lebesgue measure with density $f(xj\mu)$ and if the Dirac prior on $M_1$ is predictive-coherent in the sense of (12), then the posterior probability of $M_0$ is given by

$$P(\pm = 0jx) = \frac{ap_0(x)}{aP_0(x) + (1 \text{ } \text{¡} \text{ } a)^\circledR{}_n \prod_{i=1}^n p_0(x_i)} \mathbb{1}_{\text{f}d_n = \text{ng}}; \tag{17}$$

where $a = P(\pm = 0); p_0(x_i)$ and $p_0(x)$ are de¯ned in (15) and (16), respectively, $d_n$ denotes the number of distinct values in the sample $x$ and $^\circledR{}_n$ the predictive probability of no ties on $M_1$, which is given by

$$^\circledR{}_n = P_1(d_n = n) = \frac{\text{¡}(n_a):n_a^n}{\text{¡}(n_a + n)} \tag{18}$$

¥

Proof : The key issue is that of choosing a suitable dominating measure. Let $L$ denote the Lebesgue measure on $\mathbb{R}^p$. Clearly, $dP_0 = dL = p_0$. Since $P_1$, as de¯ned in (11), can generate ties, we decompose it as follows

$$P_i = {}^\circledR{}_n P_{11} + (1 \text{ } \text{¡} \text{ } {}^\circledR{}_n)P_{12};$$

8

where $P_{11} = P_1(\cdot|d_n = n)$ and $P_{12} = P_1(\cdot|d_n < n)$. We note that $P_{11}$ and $P_{12}$ are mutually singular since, in particular, $P_{11}(d_n = n) = 1$ and $P_{12}(d_n = n) = 0$. Note furthermore that $P_0(d_n = n) = 1$ by the absolute continuity of $F^\mu$. Let choose $P_\alpha = L + P_{12}$ as the dominating measure. It immediately follows that $dP_0 = dP_\alpha = p_0 \mathbb{1}_{\{d_n = n\}}$. Furthermore, it follows from (11) that, conditionally on $d_n = n$, the $x_i$'s are i.i.d. $F_a$ on $M_1$, whence

$$\frac{dP_1}{dP_\alpha} = \omega_n \prod_{i=1}^n p_0(x_i) \mathbb{1}_{\{d_n = n\}} + (1 - \omega_n) \mathbb{1}_{\{d_n < n\}}$$

and (17) follows. It also follows from (11) that $\omega_1 = 1$ and $\omega_{i+1} = \frac{n_a}{n_a + i}\omega_i$ which implies (18) and completes the proof. ¥

If $M_0$ consists of a simple hypothesis, i.e. if $F^\mu \equiv F_a$, then $p_0(x) = \prod_{i=1}^n p_0(x_i)$ and the following corollary obtains.

**Corollary 3.2** : Under the conditions of theorem 1 and if $F^\mu \equiv F_a$ on $M_0$, then

$$P(\pm = 0|x) = \frac{a}{a + \omega_n(1 - a)}\mathbb{1}_{\{d_n = n\}}. \tag{19}$$

¥

It is important to note that the predictive density $P_1$, as de‾ned in (11), is degenerate. It follows that Bayesian model selection procedures based upon posterior odds fundamentally amount to a test of continuity ($M_0$) versus discreteness ($M_1$), a conclusion which is fully supported by the behavior of the posterior probabilities in (17) and (19). Note, in particular, that $P(\pm = 0|x) \equiv 0$ as soon as $d_n < n$. $M_0$ is rejected as soon as a tie is observed. Conversely, if $d_n$ stays equal to $n$, then $\omega_n \to 0$ and $P(\pm = 0|x) \to 1$ as $n \to 1$, independently of the actual "validity" of $M_0$. In the absence of ties, $M_0$ ends being always accepted for large enough sample sizes. The same holds for any sample size if $n_a = 0$ (non-informative prior on $M_1$). In other words, the "implicit null hypothesis" of a posterior odds "test" procedure consists of the set of all continuous distributions, despite the fact that the Dirichlet prior on $M_1$ is dense ! This "extreme" behavior of posterior odds unequivocally suggest that they are fundamentally inadequate as an instrument for testing the "validity" of $M_0$ by confronting it to a general nonparametric alternative $M_1$.

We conclude this discussion with two additional comments. Firstly, the behavior of the posterior odds is fundamentally una®ected if we remove the

assumption of predictive-coherency. It can be proved that if $F_a$ no longer satis¯es condition (12), the posterior probabilities in (17) and (19) nevertheless remain valid, except that $p_0(x_i)$ is to be replaced by $p_1(x_i)$, the predictive density of $x_i$ on $M_1$. Secondly, it is not possible to eliminate the intrinsic "bumpy" character of the $M_1$-predictive distribution by increasing the prior weight attached to $F_a$, i.e. by increasing $n_a$ along with $n$, the actual sample size. It is well known that $d_n(n_a \ln n)^{i\ 1}$ converges almost surely to 1 on $M_1$ as n increases -see e.g. Rolin (1993). One would, therefore, be tempted to set $n_a$ equal to $n= \ln n$, a choice which preserves the consistency of the posterior distribution since $n_a=n$ still tends to zero as n tend to in¯nity. This choice, however, does not prevent $®_n$ from (rapidly) tending to zero as n tends to in¯nity. Actually, it can be shown that, in order to have $®_n$ tending toward a constant, we have to choose $n_a$'s of the order of $n^2$. Such choices are clearly unacceptable since, in particular, they destroy the consistency of posterior distributions on $M_1$.

# 4   Encompassing

We ¯rst brie°y present the encompassing principle per se, using notation introduced earlier, and then apply it to the pair of models under consideration.

## 4.1   General principle

As already discussed in the introduction our objective is not that of choosing between $M_0$ and $M_1$ but rather that of (in)validating $M_0$ by analyzing whether results derived on $M_1$ can be reproduced within $M_0$. Within a Bayesian framework, results to be encompassed typically take the form of a posterior density on a functional $_¸$ of the "parameter" F associated with $M_1$. Drawing upon the general presentation in Florens et al (1997), we brie°y discuss here how the encompassing principle applies to the models introduced in section 2.

The key step consists of prolonging $M_0$ into a model $M_¤$ which is characterized by a joint probability on $(\mu; X)$ and F. This extension is obtained by adding onto $M_0$ a conditional distribution for F, given $\mu$. The joint distribution of $(F;\mu;x)$ on $M_¤$ is then characterized by the density

$$\frac{1}{4}(F;\mu;x) = \pm(F j\mu):[p_0(xj\mu):{}^1{}_0(\mu)]:  \qquad (20)$$

Note that ¼ emboddies the key assumption that F and x are mutually independent, conditionally on $\mu$, an assumption we shall qualify further in a

10

moment. The "$M_0$-posterior density" of $F$ is then de¯ned as the posterior density of $F$ which obtains from ¼. It can be rewritten as

$$\pi_0(F\,j\,x) = \int_{\mathcal{L}} \pm(F\,j\,\mu)\,^1\pi_0(\mu\,j\,x)\,d\mu; \qquad (21)$$

from which we can obtain an $M_0$-posterior density for ¸ say, $\pi_0(\,¸\,j\,x)$. Encompassing compares the two posterior densities for ¸ : $\pi_1(\,¸\,j\,x)$, as obtained on $M_1$ and $\pi_0(\,¸\,j\,x)$ which constitutes our "reinterpretation" of $\pi_1(\,¸\,j\,x)$ within the context of the enlarged $M_0$.

There are two key reasons for imposing independence between $F$ and $x$, conditionally on $\mu$. Firstly, at a heuristic level, encompassing aims at reinterpreting $M_1$-posterior densities within $M_0$, without additional data processing beyond that already incorporated in the posterior density $^1\pi_0(\mu\,j\,x)$. Furthermore, from the viewpoint of $M_0$, it appears natural to assume that $\mu$ is a su±cient parametrization of $M_0$, i.e. that $p_0(x\,j\,\mu; F)$ ´ $p_0(x\,j\,\mu)$, an assumption which immediately validates (20).

Secondly, at a more formal level, if the transition $\pm$ were allowed to depend on $x$, then it would always be possible to produce a perfect match between $\pi_1(\,¸\,j\,x)$ and $\pi_0(\,¸\,j\,x)$, immediately voiding the encompassing principle from any meaning. As discussed in Florens et al (1997), the assumed conditional independence between $F$ and $x$ enables us to formally reinterpret encompassing as a concept of "su±ciency among models", dual to the traditional concept of su±ciency (among statistics), as de¯ned by Blackwell (1951, 1953) or Lecam (1964). That concept provides an unambiguous statistical foundation to the commonly observed scienti¯c pratice of verifying that new theories or models be capable of accounting for ¯ndings ("failures" as well as "successes") obtained from earlier models.

There remains to discuss how to formalize comparisons between $\pi_1$ and $\pi_0$. An operational procedure consists of computing a measure of "divergence" between $\pi_1$ and $\pi_0$, say, $^2{}_¸(x) = D(\pi_1(\,¸\,j\,x); \pi_0(\,¸\,j\,x))$. Such a measure, which depends upon the observed sample $x$, is refered to as to a measure of the speci¯city of $M_1$ relative to $M_0$ (in reference to ¸). As discussed e.g. in Florens et al. (1997) di®erent choices are available for the divergence $D$. In the present paper, we restrict our attention to the entropy measure

$$^2{}_¸(x) = \int_{¤} \ln \left[\frac{\pi_0(\,¸\,j\,x)}{\pi_1(\,¸\,j\,x)}\right]^{\,³} \pi_0(\,¸\,j\,x)\,d¸; \qquad (22)$$

which proves to be particularly operational for the problem under consideration.

Two additional issues need to be addressed before we close this brief presentation of encompassing. Firstly and foremost, we have to select an appropriate transition density $\pm(F|\mu)$, often refered to as to a "Bayesian Pseudo-True Value" (hereafter BPTV). In the present paper, we use a BPTV which is defined as the $M_0$-(sample) expectation of the $M_1$-posterior density of $F$, that is to say

$$\pm(F|\mu) = \int {}^{\circ}_1(F|x)p_0(x|\mu)dx. \tag{23}$$

This transition offers two keys advantages : (i) As we shall illustrate below, it can be evaluated at a high level of generality under a combination of Monte Carlo simulation and kernel smoothing; and (ii), it produces consistent encompassing test procedures since, under general technical conditions, it converges toward classical pseudo-true values (which are formally defined as plims on $M_0$ of Maximum Likelihood estimators).

The sample size $\bar{n}$ in (23) need not be set equal to n, the actual sample size. For the ease of reference, we shall explicitly distinguish between n and $\bar{n}$ in all subsequent formulae. This being said, there are good reasons for setting $\bar{n} = n$. Foremost, it is common pratice to condition BPTV's upon exogenous variables and in the context consideration, sample size truly is exogenous. Also, though we do not specifically discuss here asymptotic encompassing - see e.g. Florens and Richard (1998) for asymptotic analysis in the context of finite parameter spaces - consistency of encompassing test procedures under BPTV's does require that n and $\bar{n}$ both tend to infinity at a common rate. Finally, setting $\bar{n} = n$ facilitates calibration. We note also that, had we set $\bar{n} = 1$, the BPTV would collapse into a (degenerate) classical pseudo-true value.

A second problem is that of interpreting the actual value obtained for $^2_{,}(x)$. An obvious and fairly simple "calibration" procedure consists of evaluating (by Monte Carlo simulation) its distribution on the $M_0$-predictive distribution of x and of computing the probability of drawing a value larger than that which actually obtained (such a probability can usefully be reinterpreted as a Bayesian encompassing p-value). This "exact" calibration procedure can also usefully contribute "standardizing" specificity measures that would be evaluated under alternative divergence measures.

Finally, we ought to mention that more "decision oriented" Bayesian encompassing procedures could usefully be considered under specific circumstances. If, for example, the functionals $_{,}$ happened to be of interest to the proprietor of $M_0$ within a given decisional context (or if it was thought to be desirable that decisions reached on $M_1$ be encompassed within $M_0$), then

di®erences in posterior expected losses would provide a natural metric to evaluate the speci¯city of $M_1$ relative to $M_0$. Assume, for example, that $M_1$ were paired with a decision a 2 A and a loss function l(a; µ). One would then compute the following posterior expected losses

$$l_1^\pi(x) = \min_{d2D} \int_\Theta l(d;µ)°_1(µjx)dµ;$$

$$l_0^\pi(x) = \min_{d2D} \int_\Theta l(d;µ)°_0(µjx)dµ \qquad (24)$$

$$= \min_{d2D} \int_{\Theta£\Phi} l(d;µ) ¹_0(µj\mu)±(µj\mu)dµ\,d\mu:$$

The di®erence between $l_0^\pi(x)$ and $l_1^\pi(x)$ would provide an obvious measure of how well $M_0$ encompasses $M_1$ for the speci¯c decision problem under consideration. Here again, that di®erence could be calibrated under the $M_0$-predictive density of x.

We conclude this presentation of encompassing by mentioning here that the use of measures of divergence (entropy,...) has often been advocated as a simple alternative to a full decision-oriented evaluation. See e.g. DeGroot (1970)

## 4.2   Nonparametric encompassing

We now apply the encompassing procedure we just described to the pair of models introduced in section 2. We have already established that the $M_1$-posterior distribution of F is $Di(n_p; F_p)$ where $n_p$ and $F_p$ are de¯ned in (9). It follows that the BPTV introduced in (23) is a mixture of Dirichlet processes. Using the density notation introduced earlier we can rewrite the $M_0$-posterior density of F as follows :

$$°_0(Fjx) = \int_\Phi ±(Fj\mu) ¹_0(\mu jx)d\mu$$

$$= \int_\Phi \int_{\mathbb{R}^{np}} °_1(Fjx)p_0(xj\mu) ¹_0(\mu jx)dx d\mu$$

$$= \int_{\mathbb{R}^{np}} °_1(Fjx)p_0(xjx)dx; \qquad (25)$$

where

$$p_0(xjx) = \int_\Phi p_0(xj\mu) ¹_0(\mu jx)d\mu \qquad (26)$$

13

denotes the $M_0$-posterior predictive distribution of $x$, given $x$. It follows, therefore, that $\pi_0(F|x)$ is itself a mixture of Dirichlet measures. Actually, the only difference between $\tau(F|x)$ and $\pi_0(F|x)$ lies in the mixing distribution which is $p_0(x|\mu)$ for $\tau$ and $p_0(x|x)$ for $\pi_0$. $M_1$ and $M_0$-expectations of $F$ are given by

$$
\begin{aligned}
E_1(F) &= F_a; \\
E_1(F|x) &= \frac{n_a}{n_a + n}F_a + \frac{n}{n_a + n}F_n; \\
E_0(F|\mu) &= \frac{n_a}{n_a + \tilde{n}}F_a + \frac{\tilde{n}}{n_a + \tilde{n}}F^\mu; \\
E_0(F|x) &= \frac{n_a}{n_a + \tilde{n}}F_a + \frac{\tilde{n}}{n_a + \tilde{n}}E_0(F^\mu|x); \qquad (27)
\end{aligned}
$$

respectively. We note that, under the predictive coherent prior in (12), we have $E_0(F) = E_\mu[E_0(F|\mu)] = F_a = E_1(F)$. Furthermore, the comparison between the $M_0$-and $M_1$-posterior expectations of $F$ amounts to a comparison between $E_0(F^\mu|x)$, the $M_0$ posterior expectation of $F^\mu$ and $F_n$, the empirical distribution of $x$.

We first demonstrate that analytical characterizations of the transition $\tau(F|\mu)$ and of the $M_1$-posterior distribution $\pi_1(F|x)$ are intractable. According to Ferguson (1973), it suffices to derive the distribution of the vector $\{F(B_l)\ l : 1 ! L\}$ for any non trivial partition $\{B_l ; l : 1 ! L\}$ of $\mathbb{R}^p$, notwithstanding the fact that there might be specific partition(s) of interest to be encompassed by $M_0$ (The latter would be included in the definition of the functional "of interest" $\vartheta = \vartheta(F)$ which is introduced below).

Relative to an arbitrary partition $\{B_l ; l : 1 ! L\}$, a sufficient statistic associated with the sample $x$ is given by $\{n_l ; l : 1 ! L\}$, where $n_l$ denotes the number of $x_i$'s in $B_l$. Therefore, the $M_1$-posterior density of $\{F(B_l); l : 1 ! L\}$ is given by

$$
\pi_1(F(B_1), \dots, F(B_L)|x) = \Gamma(n_a + n) \prod_{l=1}^{\Psi} \frac{(F(B_l))^{n_a F_a(B_l)+n_l - 1}}{\Gamma(n_a F_a(B_l) + n_l)} : \qquad (28)
$$

A similar expression applies to the $M_1$-posterior distribution of $\{F(B_l); l : 1 ! L\}$ conditional on $x$, except that the $n_l$'s in (28) are replaced by $\tilde{n}_l$'s. Since, furthermore, the $x_i$'s are i.i.d. $F^\mu$ on $M_0$, the $\tilde{n}_l$'s follow a multinomial distribution with parameters $\tilde{n}$ and $\{F^\mu(B_l) ; l : 1 ! L\}$. It immediately follows that $\tau(F(B_1), \dots, F(B_L)|\mu)$ is given by

14

$$\pm(F(B_1); ::::; F(B_L)j\mu) = \overset{\mathbf{P}}{\phantom{}}_{f\mathsf{R}_l g2\mathsf{S}_{L_i 1}} \, \overset{\mathbf{h}}{\phantom{}}_i (n_a + \mathsf{R}) \overset{\mathbf{Q}_L}{\phantom{}}_{l=1} \frac{(F(B_l))^{n_a F_a(B_l)+\mathsf{R}_{li} 1}}{i (n_a F_a(B_l)+\mathsf{R}_l)} \overset{\mathbf{i}}{\phantom{}} :$$

$$\overset{\mathbf{h}}{\phantom{}}_{\mathsf{R}} \overset{\mathbf{Q}_L}{\phantom{}}_{l=1} \frac{(F^\mu(B_l))^{\mathsf{R}_l}}{\mathsf{R}_l!} \overset{\mathbf{i}}{\phantom{}} ;$$

$$(29)$$

where $\mathsf{S}_{L_i 1} = f(\mathsf{R}_1; ::::; \mathsf{R}_L); \mathsf{R}_l 2 \mathbb{N}; \overset{\mathbf{P}_L}{\phantom{}}_{l=1} \mathsf{R}_l = \mathsf{R}g$. Such summations have no (simple) analytical solutions and, furthermore, include very large numbers of terms for all but small $\mathsf{R}$ and/or $L$. It also follows that $°_1(F jx)$ is itself analytically intractable.

On the other hand, we can easily simulate trajectories of $F$ under its $M_0$- and $M_1$-posterior distributions by relying upon powerful representations of the Dirichlet process, as found, e.g., in Sethuraman (1994) or Rolin (1993). Let $F_i^x$ denote a realisation of the $M_i$-posterior distibution of $F$ ($i = 1; 2$). The following representation applies to $F_1^x$

$$F_1^x = (1_i °) \overset{\mathbf{X}}{\phantom{}}_{k=1} ®_k\pm_{»_k} + ° \overset{\mathbf{X}}{\phantom{}}_{i=1} ^-_i\pm_{x_i}; \qquad (30)$$

where

(i) $° ??f®_k; k : 1 ! 1 g??f»_k; k : 1 ! 1 g??f^-_i; i : 1 ! ng;$

(ii) $° » B(n; n_a)$, where $B(a; b)$ denotes the Beta distribution with parameters $a > 0$ and $b > 0$;

(iii) The $»_k$'s are i.i.d. $F_a$;

(iv) $®_k = v_k \overset{\mathbf{Q}_{k_i 1}}{\phantom{}}_{l=1}(1_i v_l)$ for $k : 1 ! 1$, where the $v_l$'s are i.i.d. $B(1; n_a)$;

(v) The $^-_i$'s are uniformly distributed on the simplex $S_{n_i 1}$, i.e. $(^-_1; ::::; ^-_n) » Di(n; (1:::1))$;

and where, as above, $\pm_{x_i}(\pm_{»_k})$ denotes the Dirac measure at point $x_i(»_k)$, i.e., $\pm_{x_i}(B) = \mathbb{1}_B(x_i)$.

A similar representation applies to $F_0^x$, except that $x$ in (30) is replaced by $x$ which, according to (26), is generated as follows : $\mu$ is drawn from $^1_0(\mu jx)$ and, conditionally on $\mu$, the $x_i$'s are drawn independently from one another from $F^\mu$.

Using these representations, we can generate repeated draws from the posterior distributions $°_0$ and $°_1$. The next issue to be addressed is that of

15

using there draws to evaluate a divergence measure between $\circ_1$ and $\circ_0$. It is important to immediately point out that $\circ_1$ and $\circ_0$ are mutually singular. For example, $F_1^x$ in (30) assigns a non zero probability to $x_1$, that is to say, the set of measures on $(IR^p; B_p)$ which assign positive probability to $x_1$ has measure one on $\circ_1(F jx)$. On the other hand $F_0^x$ assigns zero probability to $x_1$, since the probability that $x_1 = x_1$ on $p_0(xjx)$ is zero. It follows that conventional divergence measures will automatically be maximal, leading to automatic rejection of the corresponding encompassing test procedure (a problem which is closely related to that already discussed in section 3 for posterior odds). We should then consider metrics which can produce non degenerated comparisons of mutually singular distributions. Skorohod's distances might provide a conceptual solution in that respect but it is far from being obvious that they could lead to operational procedures.

A fully operational solution consists of restricting the encompassing comparison to the $M_0$- and $M_1$-posterior distributions of an appropriate functional of $F_1$ say

$$\text{\j} = {}'(F) \, 2 \, IR^k; \tag{31}$$

where $\text{\j}$ would represent "parameters of interest" within $M_1$ which would, therefore, be obvious targets for an encompassing test. The key advantage o®ered by this approach lies in the ¯nding that, under very mild technical condition, the $M_0$- and $M_1$-posterior distributions of $\text{\j}$ are absolutely continuous w.r.t. the Lebesgue measure -see e.g. Florens and Rolin (1994) - whence divergence comparisons do apply to $\text{\j}$. Typical choices for $\text{\j}$ are :

(i) $\text{\j} = (F(B_l); l : 1 \, ! \, l)$ for a measurable partition $(B_1; :::B_L)$ of $IR^p$ which regroups "regions of interest" in the sample space ; (ii) $\text{\j} = \int_{IR^p} h(x)F(dx)$, where $h$ denotes a Borel function from $IR^p$ to $IR^k$. For example a comparison of ¯rst and second order moments would obtain for $\text{\j} = (x; xx^0)$ and $k = p + \frac{1}{2}p(p+1)$.

The simulation method described above is then exploited in the following way to produce an operational encompassing test procedure. A set of randomly drawn trajectories $f(F_{1;r}^x; F_{0;r}^x); r : 1 \, ! \, Rg$ is transformed into a corresponding set of random drawns for $\text{\j}$, say $f\tilde{\text{\j}}_{1;r}^x; \tilde{\text{\j}}_{0;r}^x); r : 1 \, ! \, Rg$. Kernel estimates of the posterior densities $\circ_1(\text{\j}jx)$ and $\circ_0(\text{\j}jx)$ are then evaluated, which are denoted by $\mathfrak{a}_1(\text{\j}jx)$ and $\mathfrak{a}_0(\text{\j}jx)$, respectively. A measure of the $\text{\j}$-speci¯city of $M_1$ relative to $M_0$ is then given by

$$\mathfrak{X}_{\text{\j}}(x) = \frac{1}{R} \sum_{r=1}^{R} \ln \left[ \frac{\mathfrak{a}_0(\tilde{\text{\j}}_{0;r}^x jx)}{\mathfrak{a}_1(\tilde{\text{\j}}_{0;r}^x jx)} \right]; \tag{32}$$

16

Finally, in order to calibrate $\hat{z}(x)$ on $M_0$, we generate $S$ (i.i.d.) samples from $p_0(x)$, the $M_0$-predictive distribution of $x$ and compute $\hat{z}(x)$ for each such sample, a procedure which produces $S$ i.i.d. draws from $\hat{z}(x)$ on $p_0$, from which a Bayesian p-value immediately obtains as the proportion of draws which exceeds the actual sample value $\hat{z}(x)$.

# 5  An application : validation of an exponential distribution

Let X denote a positive scalar random variable, such as a lifetime. Our application is based upon a sample $x$ of size $n = 50$ which we drew from an exponential distribution with parameter $\mu = 1$. Let $M_0$ consists of an exponential distribution with parameter $\mu > 0$ under a natural conjugate gamma prior density for $\mu$ with parameters $a_0 > 0$ and $°_0 > 0$. The corresponding sample and prior density functions are given by

$$p_0(x|\mu) = \mu^n e^{i\ \mu t}; \text{ and} \tag{33}$$

$$1_0(\mu) = \frac{a_0^{°_0}}{i\,(°_0)}\mu^{°_0 i\ 1}e^{i\ a_0\mu}; \tag{34}$$

respectively, where $t = \mathbf{P}_{i=1}^{n}\ x_i$ is a su±cient statistic. Actual calculations are based upon the values $a_0 = °_0 = 2$. The $M_0$-posterior distribution is itself a gamma distribution with parameters $a_¤ = a_0 + t$ and $°_¤ = °_0 + n$. The $M_0$-predictive densities in (15) and (16) are given by

$$p_0(x_i) = °_0 a_0^{°_0}(a_0 + x_i)^{i\ (°_0+1)}; \text{ and} \tag{35}$$

$$p_0(t) = \frac{i\,(°_0 + n)}{i\,(°_0)}a_0^{°_0}(a_0 + t)^{i\ (°_0+n)}; \tag{36}$$

respectively. We note in passing that $p_0(x_i)$ represents a Pareto distribution.

On $M_1$, we use a non-informative Dirichlet prior measure on F by setting $n_a = 0$. In view of the "weakly" informative $M_0$-prior ($a_0 = °_0 = 2$), imposing the predictive-coherency condition (12) makes little di®erence while setting $n_a = 0$ simpli¯es the representation of $F_1^x$ in (30), which is now given by

$$F_1^x = \sum_{i=1}^{X} \bar{\phantom{i}}_i \pm_{x_i};\tag{37}$$

where $f^-_i; i : 1 ! \; ng \gg Di(n; (1; :::; 1))$. A similar expression applies to $F_0^x$ with x being replaced by x to be drawn from the $M_0$-predictive posterior distribution in (26).

Finally, implicitely restricting $M_1$ to distributions with ¯nite ¯rst and second order moments, we consider here that the latter are the parameters of interest to be encompassed by $M_0$, i.e.

$$\vartheta_i = \int_0^1 x^i F(dx); \text{ for } i = 1; 2:\tag{38}$$

Following the procedure described in section 4.2, we draw two samples of size R from $^\circ_1(\vartheta jx)$ and $^\circ_0(\vartheta jx)$, respectively, that is to say

$$\tilde{\vartheta}_{1;r}^x = \sum_{j=1}^{X} \tilde{\phantom{a}}_{1;j;r} \begin{pmatrix} x_j \\ x_j^2 \end{pmatrix}; \text{ and } \tilde{\vartheta}_{0;r}^x = \sum_{j=1}^{X} \tilde{\phantom{a}}_{0;j;r} \begin{pmatrix} x_{j;r} \\ x_{j;r}^2 \end{pmatrix};\tag{39}$$

where $f^\sim_{1;j;r}; j : 1 ! \; ng$ and $f^\sim_{0;j;r}; j : 1 ! \; ng$ for $r : 1 ! \; R$ are i.i.d. drawns from a $Di(n; (1; :::; 1))$ distribution, and $fx_{j;r}; j : 1 ! \; ng$ for $r : 1 ! \; R$ are i.i.d. draws from $p_0(xjx)$ (in pratice, $x_r$ is drawn from (33), conditionally on $\mu = \bar{\mu}_r$, where $\bar{\mu}_r$ itself is drawn from (34)). Next, we compute bivariate kernel estimates for the posterior distributions of $\vartheta$ say

$$^\circ_i(\vartheta jx) = \frac{1}{R} \sum_{r=1}^{R} K_2(\vartheta_i \; \tilde{\vartheta}_{i;r}^x); \text{ for } i = 1; 2;\tag{40}$$

where $K_2$ denotes the kernel

$$K_2(u_1; u_2) = \prod_{i=1}^{Y} \frac{1}{h_{i;R}} Á \begin{pmatrix} u_i \\ h_{i;R} \end{pmatrix};\tag{41}$$

Á denotes the standardized normal density and $h_{i;R} = S_i R^{i^{1=5}}$, where $S_i$ denotes the estimated standard deviation of the $\tilde{\vartheta}_{i;r}^x$ for $i = 1; 2$. The $\vartheta$-speci¯city of $M_1$ relative to $M_0$ then obtains by application of (32) and equals 0.58. In the present application we use a Monte Carlo sample size $R = 1; 000$.

We generated 1,000 auxiliary samples of size n = 50 from the $M_0$-predictive density $p_0(t)$ in (36). Actually, each such draw is generating by drawing ¯rst μ from the prior (34) and then x conditionally on μ from (33). The computations we just described are applied to each auxiliary sample in turn in order to produce 1,000 draws from the $M_0$-predictive distribution of ₂ (x). An histogram of the latter is reproduced in ¯gure 1 and a few selected "critical values" are reported in table 1. The $M_0$-predictive mean and standard deviation of ₂ (x) equal 0.8104 and 0.3374, respectively.

Finally, we also evaluated the ¸-speci¯cities of $M_1$ relative to a variety of "invalid" $M_0$'s, all constrained to have expectations equal to 1. The entropy results are reported in table 2, together with the corresponding Bayesian p-values which obtain from the histogram in ¯gure 1. In contrast with the valid exponential model, these invalid models clearly fail to encompass $M_1$ relative to ¸.

# 6   Conclusion

We have proposed a °exible and fully operational Bayesian procedure for examining whether or not a tentative model $M_0$ encompasses results derived under a general nonparametric model $M_1$. It is important to emphasize the fact that encompassing is not a model selection procedure since, in particular, $M_1$ generally is not meant to be interpreted as an alternative to $M_0$ but rather as an instrument in testing the (in)validity of $M_0$. Clearly, broad failure to encompass characteristics of interest of $M_1$ would lead to the conclusion that $M_0$ is seriously de¯cient. Whether such de¯ciencies would lead one to try to re¯ne $M_0$ further or instead to trash it can only be answered in the context of the speci¯c application under consideration after careful examination of the said de¯ciencies.

Our method o®ers considerable °exibility in this respect relative, in particular, to the selection of results to be encompassed which would typically depend upon the objectives of the proprietor of $M_0$. Encompassing also requires the selection of a transition probability and of a divergence measure. Though such choices introduce a degree of arbitrariness in our procedure, the pair consisting of the entropy measure (22) and of the BPTV (23) provides a natural Bayesian extension of concepts widely used in statistics. Furthermore, arbitrariness is tempered by the fact that all encompassing results are $M_0$ calibrated in the end. As brie°y discussed in the paper, much arbitrariness can also be removed by adopting a more decision oriented measure of encompassing di®erences.

In view of the inherent complexity of the task of constructing empirical

econometric models, we strongly believe that encompassing can provide an invaluable tool for carefully investigating potential deⁿciencies of a model and suggesting avenues for further improvements.

Figure 1 :
$M_0$-predictive distribution of $\gamma$-specificity of $M_1$.

## Table 1. Estimated critical values of the $\gamma$-specificity

| $1 - \alpha$ | 0,90 | 0,91 | 0,92 | 0,93 | 0,94 | 0,95 | 0,96 | 0,97 | 0,98 | 0,99 |
|---|---|---|---|---|---|---|---|---|---|---|
| $C_{1-\alpha}$ | 1,28 | 1,30 | 1,33 | 1,35 | 1,39 | 1,44 | 1,49 | 1,59 | 1,70 | 1,94 |

## Table 2. Bayesian p-values

| $M_0$ distribution | Variance | Entropy | Prob-value |
|---|---|---|---|
| Exponential (1) | 1,000 | 0,58 | 0,829 |
| Weibull (2) | 0,273 | 3,17 | 0,000 |
| Gamma $(\frac{1}{2}, \frac{1}{2})$ | 2,000 | 1,43 | 0,052 |
| Gamma (2,2) | 0,500 | 1,69 | 0,022 |
| Gamma (4,4) | 0,250 | 1,95 | 0,009 |
| Uniform (0,2) | 0,333 | 1,83 | 0,015 |
| Lognormal $(-\frac{1}{4}, \frac{1}{2})$ | 0,649 | 1,31 | 0,086 |
| Lognormal $(-\frac{1}{2}\sigma^2, \sigma^2)$  $\sigma^2 = 0,2231$ | 0,250 | 2,88 | 0,000 |

# References

Antoniak C.E. (1974) "Mixtures of Dirichlet processes with applications to Bayesian nonparametric problems" Annals of Statistics,2, 1152-1174.

Berk, R. (1966) "Limiting behavior of posterior distributions when the model is incorrect", The Annals of Mathematical Statistics, 37, 51-58.

Blackwell D. (1951) "Comparison of experiments", Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probabilities, Berkeley : University of California Press, 93-102.

Blackwell D. (1953) "Equivalent comparison of experiments", The Annals of Mathematical Statistics, 24, 265-272.

DeGroot M.H. (1970) "Optimal statistical decisions", McGraw-Hill, Inc, New-York.

Escobar, M.D. and M. West (1995) "Bayesian density estimation and inference using mixtures", Journal of the American Statistical Association, 90, 577-588.

Ferguson T.S. (1973) "A Bayesian analysis of some non parametric problems", Annals of Statistics, 1, 209-230.

Ferguson T.S. (1974) "Prior distributions on spaces of probability measures", Annals of Statistics, 2, 615-629.

Florens J.P. and M. Mouchart (1993) "Bayesian testing and testing Bayesians", Handbook of Statistics, vol.11, 65, 303-334, Maddala G.S., C.R. Rao and H.D. Vinod, Edts, Elsevier Science Publishers, .

Florens J.P. and J.M. Rolin (1994) "Bayes, Bootstrap, Moments" Institut de Statistique D.P. 9413, Université Catholique de Louvain, Belgium.

Florens J.P., D.P. Hendry and J.F. Richard (1996) "Encompassing and speci¯city", Econometric Theory, 12, 620-656.

Florens J.P., C. Protopopescu and J.F. Richard (1997) "Identi-¯cation and estimation of a class of game theoretic models", University of Toulouse.

Florens J.P. and J.F. Richard (1998) "Encompassing in ¯nite parametric spaces", University of Toulouse.

Huber P.J. (1967) "The behaviour of maximum likelihood estimates under non standard conditions" in Proceedings of the Fifth Berkeley Symposium on Mathematical Statistics and Probability, vol.1, 105-123, L. Lecam and J. Neyman, Edts, University of California Press, .

Lavine, M (1992) "Some aspects of polya tree distributions for statistical modelling", The Annals of Statistics, 20, 1222-1235.

Lavine, M. (1994) "More aspects of Polya tree for statistical modelling", The Annals of Statistics, 22, 1161-1176.

LeCam, L. (1964), "Su±ciency and approximate su±ciency", The Annals of Mathematical Statistics, 35, 1419-1455.

Lenk, P.J. (1988) "The logistic normal distribution for Bayesian, nonparametric, predictive densities", Journal of the American Statistical Association, 83, 509-516.

Lenk, P.J. (1991) "Towards a practicable bayesian nonparametric density estimator", Biometrika, 78, 531-543.

Leonard, T. (1978) "Density estimation, stochastic processes and prior information", The Journal of Royal Statistical Society, B 40, 113-146

Lo, A.Y. (1984) "On a class of Bayesian nonparameteric estimates I density estimates", The Annals of Statistics, 12, 351-357.

Marshall, R., M. Meurer, J.F. Richard and W. Stromquiest (1994) "Numerical analysis of asymmetric ¯rst-price auctions", Game and Economic Behavior, 7, 193-220.

Rai®a, H. and R. Schlaifer (1961), Applied statistical decision theory, Boston : Division of Research, Harvard Business School.

Révész P. (1984) "Density estimation" in Handbook of Mathematical Statistics vol. 4, 531-549, P.R. Krishnaiah and P.K. Sen, Edts, Elsevier Science Publishers, .

Rolin J.M. (1992) "Some useful properties of the Dirichlet process", Institut de Statistique D.P. 9202, Université Catholique de Louvain, Belgium.

Rolin J.M. (1993) "On the distribution of jumps of the Dirichlet process", Institut de Statistique D.P. 9302, Université Catholique de Louvain, Belgium.

Sethuraman J. (1994) "A constructive de¯nition of Dirichlet priors", Statistica Sinica, Vol 4(2), 639-650.

White H. (1982) "Maximum Likelihood estimation of misspeci‾ed models", Econometrica, 50, 1-26.