# LOCAL ROBUSTNESS OF BAYES FACTORS FOR NONPARAMETRIC ALTERNATIVES

By Cinzia Carota[1]

*University of Pavia*

In this paper we consider a particular Bayes factor B for comparing a fixed parametric model against a nonparametric alternative, and we investigate its local sensitivity to the sampling distribution. The nonparametric alternative is constructed by embedding the parametric model, characterized by a d.f. $F_0$ known up to a real parameter $\theta$, into a mixture of Dirichlet processes. More precisely, conditionally on $\theta$, $F_0$ represents the mean of a random d.f. which is assumed to be a Dirichlet Process. So, for the Bayes factor B, sensitivity to perturbations of the sampling distribution $F_0$ and sensitivity to small departures from the fixed Dirichlet process parameter are the same problem. Here we consider B as a (non ratio-linear) functional defined on a set of sampling d.f.'s and maximize its first von Mises derivative over this set. In particular, mixture and density bounded sets are considered.

**1. Introduction.** A Bayesian analysis may depend critically on the modeling assumptions which include prior, sampling distribution and loss function. Therefore, it is useful to assess the sensitivity of inferences to modest changes in the specification of the problem by means of a so called robustness analysis. On this subject there exists an extensive literature. A general discussion and comprehensive lists of references can be found in Berger (1984), (1990), (1994), Gustafson, Wasserman and Srinivasan (1994) and Wasserman (1992). Most of the literature is concerned with global sensitivity to prior specification and focusses on posterior expectations as inferences of interest. In this article we discuss *local* sensitivity to the *sampling distribution* of a particular *Bayes factor* [Carota and Parmigiani (1994)].

In general, a sensitivity analysis is performed when there is uncertainty about modeling assumptions. Such uncertainty is expressed by specifying a class of inputs (for example, a class of priors or a class of sampling distributions) instead of a single one. The local sensitivity analysis examines the rate at which the inference changes relative to small perturbations to a base input in direction of the other elements in the class. It is preferred to the global analysis when the given class of inputs contains a natural reference point, either because of a very high degree of belief or because of mathematical appeal. Sometimes a local analysis can be used to construct quite accurate global robustness bounds, when exact computations are too difficult or too time consuming.

In this paper we will study the effects of infinitesimal perturbations to a fixed sampling distribution function $F_0$, known up to a parameter $\theta$, on the Bayes factor B comparing the given parametric model against a nonparametric alternative. The alternative is constructed by embedding the parametric model in a mixture of Dirichlet processes. In particular, conditionally on $\theta$, $F_0$ represents the mean of a random d.f., $F$, which is assumed to be a Dirichlet Process. So, the sampling distribution $F_0$ is, at the same time, the baseline model, or null hypothesis in the testing terminology, and the location parameter inside the alternative model. There are at least two natural ways of measuring the local sensitivity of B to $F_0$. Here (Section 4), we consider B as a (non ratio-linear) functional defined on a set of sampling d.f.'s and maximize its first von Mises derivative over this set. In particular, mixture and density bounded sets are considered. Section 2 introduces the Bayes factor B. Section 3 discusses the use of functional derivatives in local sensitivity and in particular motivates using von Mises derivative. Section 5 contains a brief discussion.

**2. Bayes factors for nonparametric alternatives.** Let $y = (y_1, y_2, ..., y_n)$ be an observed sample from a real-valued exchangeable sequence, and let $\mathcal{F}_0 = \{F_0(\cdot|\theta), \theta \in \mathcal{R}^k, k < \infty\}$ be the parametric class of sampling distribution functions whose adequacy for $y$ we want to investigate. In this paper we will assume the absolute continuity of $F_0$ with respect to either the Lebesgue measure or the counting measure. In both cases we will denote the corresponding probability density function by $f_0(\cdot|\theta)$ and the observed likelihood by $l(\theta)$. The alternative to $\mathcal{F}_0$ is based on a random distribution function, $F$, taking values in the set of all distribution functions, $\mathcal{F}$. Conditionally on $\theta$, $F$ is assumed to be a Dirichlet process (Ferguson, 1973) with parameter $\alpha(\theta, \cdot) = \alpha(\theta, \infty)F_0(\cdot|\theta)$. So $F_0(\cdot|\theta)$ represents the conditional mean of $F$ and $\alpha(\theta, \infty)$ is the prior weight on the mean. In what follows $\alpha(\theta, \infty)$ is assumed to be independent on $\theta$ and denoted by $A$. Finally, the finite dimensional parameter $\theta$ is distributed according to the d.f. $P(\theta)$. Under these hypotheses the marginal distribution of F is a mixture of Dirichlet processes [Antoniak (1974)].

In this context, we compare the adequacy of $\mathcal{F}_0$ against the nonparametric family $\mathcal{F}$ by computing a "Bayes factor" of the form:

$$B = \frac{p(y|F \in \mathcal{F}_0)}{p(y)},$$

where $p$ denotes the probability density function of the data. In fact, B is not strictly a Bayes factor, but the Bayes factor

$$b = \frac{p(y|F \in \mathcal{F}_0)}{p(y|F \in (\mathcal{F} - \mathcal{F}_0))}$$

is an increasing function of B and use of B is conventional for testing nested models. Furthermore, $b$ reduces to $B$ when the Dirichlet prior assigns mass zero to the entire parametric family $\mathcal{F}_0$, then, for example, when one conditions on $\theta$ or when $F_0$ is absolutely continuous with respect to the Lebesgue measure. The relevance of these two cases will become clear later. Now we give an explicit form for $B$ [see Carota and Parmigiani (1994)]. To do this we introduce additional notation. Let $\tilde{y} = (\tilde{y}_1, \cdots, \tilde{y}_r)$ be the array of distinct observations in $y$, $n_i$ be the number of observations equal to $y_i$ and $r$ be the number of distinct observations in the sample. Also, let $\mu$ be the Lebesgue measure except for the points where $\alpha$ is atomic, to which $\mu$ assigns unit mass, let

$$\alpha'(\theta, y_i) = \frac{d\alpha(\theta, y_i)}{d\mu},$$

and $v(\theta, y_i) = \alpha'(\theta, y_i)$ if $y_i$ is an atom of $\alpha$ and zero otherwise.
Then

$$B = \frac{\int_{\mathcal{R}^k} A^{-n} \prod_{j=1}^{n} \alpha'(\theta, y_j) P(d\theta)}{\int_{\mathcal{R}^k} A_{(n)}^{-1} \prod_{j=1}^{r} \alpha'(\theta, \tilde{y}_j)(v(\theta, \tilde{y}_j) + 1)_{(n_j - 1)} P(d\theta)}$$

where $a_{(n)} = a(a+1) \cdots (a + n - 1), n > 0$ and $a_{(0)} = 1$.

We will refer to $B$ as to the global Bayes factor for the comparison of $\mathcal{F}_0$ and $\mathcal{F}$. Comments and criticism of this Bayes factor can be found in Carota (1994) and Carota and Parmigiani (1994). Here we are interested in local sensitivity of $B$ to the parametric family $\mathcal{F}_0$. When we perturb $\mathcal{F}_0$ the numerator of $B$ changes because we have changed the parametric family and the denominator of $B$ changes because the Dirichlet prior also depends on the parametric family.

**3. How to measure local sensitivity?** We will consider the Bayes factor described above as a functional of the sampling distribution $F_0$. Then we will measure local robustness of $B(F_0)$ by maximizing its first von Mises derivative over a given class $\mathcal{Q} \subseteq \mathcal{F}$. More precisely, $\mathcal{Q}$ is a convex neighborhood of $\mathcal{F}_0$, so that $\mathcal{F}_0 \subseteq \mathcal{Q} \subseteq \mathcal{F}$. In the following sections, distribution functions belonging to $\mathcal{Q}$ will be denoted by $Q$ and, as before, the existence of densities with respect to the Lebesgue measure or the counting measure will be assumed. These will be denoted by $q$.

DEFINITION. *Let $\rho$ be a functional on a convex set of d.f.'s and let $H$ and $G$ be two points in this convex set. Then the first von Mises derivative $\rho'_G$ of $\rho$ at $G$ is defined by*

$$\rho'_G(H - G) = \frac{d}{dt}\rho(G + t(H - G))|_{t=0}$$

*if there exists a real valued function $\phi_G$ (independent on H) such that*

$$\rho'_G(H - G) = \int \phi_G(y)d(H - G)(y).$$

One appealing feature of the von Mises derivative is that it is defined directly on $\mathcal{F}$. Often the Fréchet and Gâteaux derivatives of a posterior quantity $\rho(G)$ are used to measure its rate of change. These are defined on normed linear spaces or, more generally, on topological vector spaces. Thus, many authors [Diaconis and Freedman (1986), Srinivasan and Truszczynska (1990), Ruggeri and Wasserman (1993), etc.] first artificially extend $\rho(G)$ to the linear space of all signed measures and then apply the notion of functional derivative to quantify its local sensitivity. An alternative solution is suggested by Huber (1981), Clarke (1983) and more recently by Basu (1994). They generalize the definition of Fréchet or Gâteaux derivatives to encompass the case when $\rho$ is defined only on a convex set. It is worth noting that, when this convex set is $\mathcal{F}$, the generalized or *weak* Gâteaux derivative coincides with the von Mises derivative (von Mises 1947) [see also Fernholz (1983)]. If we equip $\mathcal{F}$ with the weak topology, that is the weakest topology for which all functionals of the form

$$\rho(G) = \int \phi(x)dG(x) \qquad G \in \mathcal{F}$$

are continuous for $\phi$ bounded and continuous, then the von Mises derivatives will be continuous on $\mathcal{F}$. Furthermore, if we embed $\mathcal{F}$ in $\mathcal{F}^*$, the space of all bounded signed measures on $\mathcal{R}$ equipped with the weak topology, and if the functional $\rho(G)$ can be extended to $\mathcal{F}^*$, then its von Mises derivative will correspond to the usual Gâteaux derivative on this space. [Because of this similarity between the two derivatives, the von Mises derivative has often been referred to as the Gâteaux derivative in the statistical literature].
An important characteristic of the von Mises derivative is that when $H$ is a degenerate d.f. $\rho'_G$ becomes the influence function [Hampel (1971)], a very important tool in classical robustness. [See also Carota (1994)].

   **4. Local sensitivity to the sampling distribution.** Let $C$ and $L$ denote the counting measure and the Lebesgue measure respectively, and let $m = \int_{\mathcal{R}^k} l(\theta)dP(\theta)$. Furthermore, define

$$\Delta_{F_0}(\theta) = A_{(n)}^{-1} \prod_{j=1}^{r} Af_0(\tilde{y}_j|\theta)(Af_0(\tilde{y}_j|\theta) + 1)_{(n_j-1)}$$

and

$$\delta_{F_0}(\theta) = A_{(n)}^{-1}[\prod_{j=1}^{r}(n_j - 1)!] \prod_{j=1}^{r} Af_0(\tilde{y}_j|\theta),$$

i.e. the likelihood of $\theta$ when $F|\theta$ is a Dirichlet process with parameter $AF_0(\cdot|\theta)$ absolutely continuous with respect to $C$ and $L$, respectively. The following result gives the von Mises derivative of the Bayes factor.

PROPOSITION 1. *The first von Mises derivative of $B$ at $F_0$ in the direction of $Q \in \mathcal{Q}$ is given by*

$$(1) \qquad B'_{F_0}(Q - F_0) = \sum_{i=1}^{n-r} \int_{\mathcal{R}^k} \left[ \frac{q(\bar{y}_i|\theta)}{f_0(\bar{y}_i|\theta)} - 1 \right] c_{i,1}(\theta) P(d\theta)$$

$$+ \sum_{i=1}^{r} \int_{\mathcal{R}^k} \left[ \frac{q(\tilde{y}_i|\theta)}{f_0(\tilde{y}_i|\theta)} - 1 \right] c_{i,2}(\theta) P(d\theta)$$

*where $\bar{y}$ denotes the array of duplicate observations in the sample, and*

*i) if $F_0 << C$ (discrete case)*

$$c_{i,1}(\theta) = \frac{l(\theta)}{\Delta_{F_0}}, \quad c_{i,2}(\theta) = \frac{l(\theta)\Delta_{F_0} - m\Delta_{F_0}(\theta) \sum_{s=0}^{n_i-1} Af_0(\tilde{y}_i|\theta)/(Af_0(\tilde{y}_i|\theta) + s)}{(\Delta_{F_0})^2}$$

*with $\Delta_{F_0} = \int_{\mathcal{R}^k} \Delta_{F_0}(\theta) P(d\theta)$;*

*ii) if $F_0 << L$ (continuous case)*

$$c_{i,1}(\theta) = \frac{l(\theta)}{\delta_{F_0}}, \qquad c_{i,2}(\theta) = \frac{l(\theta)\delta_{F_0} - m\delta_{F_0}(\theta)}{(\delta_{F_0})^2}$$

*with $\delta_{F_0} = \int_{\mathcal{R}^k} \delta_{F_0}(\theta) P(d\theta)$.*

PROOF. It is omitted for brevity.□

REMARK 1. $B'_{F_0}$ is a linear combination of terms with a different structure for distinct and duplicate observations. In particular, $c_{i,1} = c_1$ for all $i$, while $c_{i,2} = c_2$ for all $i$ only when $F_0 << L$ or when $f_0(\tilde{y}_i|\theta)$ is a constant and $n_i = n/r$ or $n_i = 1$. In the last case $B'_{F_0} = 0$ because $B$ depends on $y$ only through the sample size.

REMARK 2. It is interesting to focus on the case of a degenerate prior. Suppose that $P(\theta_0) = 1$. Then

i) if $F_0 << C$ (discrete case)

$$c_{i,1} = \frac{l(\theta_0)}{\Delta_{F_0}}, \qquad c_{i,2} = \frac{l(\theta_0)[1 - \sum_{s=0}^{n_i-1} Af_0(\tilde{y}_i|\theta_0)/(Af_0(\tilde{y}_i|\theta_0) + s)]}{\Delta_{F_0}};$$

ii) if $F_0 << L$ (continuous case)

$$c_{i,1} = \frac{l(\theta_0)}{\delta_{F_0}}, \qquad c_{i,2} = 0.$$

In general, from the linear combination (1) we can derive the influence of a given observation $y_i$, $i = 1, \cdots, n$, on $B$. Let $y_i = x$ and suppose that $y_i$ is the $(s+1) - th$ of the $n_i$ observations equal to $x$, then its influence on $B$ is:

$$\int_{\mathcal{R}^k} \left[ \frac{q(x|\theta)}{f_0(x|\theta)} - 1 \right] \left[ \frac{l(\theta)}{\Delta_{F_0}} - \frac{(m\Delta_{F_0}(\theta)Af_0(x|\theta))/(Af_0(x|\theta) + s)}{(\Delta_{F_0})^2} \right] P(d\theta).$$

A small value of the derivative $B'$ means that the Bayes factor changes little as the parametric family does, since the data is fixed. So, if the derivative is small it means that entire parametric families formed by deviations of the base parametric family $\mathcal{F}_0$ in a direction $Q$ are roughly equally good at explaining the data. If base model is poor, then nearby models are roughly equally poor and if the base model is good then nearby models are roughly equally good. [Vice versa, a large value of this derivative means that the Bayes factor changes rapidly as the parametric family changes and this implies that very near to $\mathcal{F}_0$ there are parametric families much more, or much less, effective at explaining the observed data.]

Results controlling the supremum of $B'$ as $Q$ varies over a class therefore give uniformly good control of the Bayes factor over directions of deviation. Thus, all models formed by deviating a base model infinitesimally for directions $Q$ are roughly equally good at explaining the particular data set obtained. An important problem is the calibration of the supremum of $B'$. A rough calibration can be based on dividing this by $B(F_0)$, therefore considering the relative rate of change of $B$ at $F_0$.
The following two propositions give results for mixture and density bounded classes of sampling distributions.

Define $\mathcal{P}_\Omega=$ class of all d.f.'s on $\Omega$, and let $\mathcal{Q}_M$ be a mixture class of sampling distributions

$$\mathcal{Q}_M = \{Q(\cdot|\theta) = \int_\Omega Q(\cdot|\theta, \omega) dK(\omega) \colon K \in \mathcal{P}_\Omega\}.$$

PROPOSITION 2. *Let* $\mathcal{Q} = \mathcal{Q}_M$, *then*

$$\sup_{Q \in \mathcal{Q}_M} |B'_{F_0}(Q - F_0)| = \sup_{\omega \in \Omega} |S(\omega) - c|$$

*where*

$$S(\omega) = \sum_{i=1}^{r} \int_{\mathcal{R}^k} q(\tilde{y}_i|\theta, \omega)/f_0(\tilde{y}_i|\theta)c_i(\theta)P(d\theta),$$

$$c_i(\theta) = (n_i - 1)c_{i,1}(\theta) + c_{i,2}(\theta) \quad and \quad c = \sum_{i=1}^{r} \int_{\mathcal{R}^k} c_i(\theta)P(d\theta).$$

PROOF. A convenient expression for $B'_{F_0}$ is

$$(2) \qquad B'_{F_0}(Q - F_0) = \sum_{i=1}^{r} \int_{\mathcal{R}^k} \frac{q(\tilde{y}_i|\theta)}{f_0(\tilde{y}_i|\theta)}c_i(\theta)P(d\theta) - c,$$

so that

$$\sup_{Q \in \mathcal{Q}_M} |B'_{F_0}(Q - F_0)| = \sup_{K} |\int_{\Omega} \sum_{i=1}^{r} \int_{\mathcal{R}^k} \frac{q(\tilde{y}_i|\theta, \omega)}{f_0(\tilde{y}_i|\theta)}c_i(\theta)P(d\theta)K(d\omega) - c|$$

$$= \sup_{\omega \in \Omega} |S(\omega) - c|. \square$$

REMARK 3. Let $\mathcal{Q} = \mathcal{Q}_{SU}$,

$\mathcal{Q}_{SU} = \{Q(\cdot|\theta) : Q(\cdot|\theta) \text{ is a symmetric unimodal distribution with mode } \theta \}$,

and suppose that $Q(\cdot|\theta) << L$. Then, $Q(\cdot|\theta)$ can be written as a mixture of uniform distributions of the form $U(\theta - z, \theta + z)$, and

$$\sup_{Q \in \mathcal{Q}_M} |B'_{F_0}(Q - F_0)| = \sup_{z>0} \left| \sum_{i=1}^{r} \frac{1}{2z} \int_{z-\tilde{y}}^{z+\tilde{y}} \frac{c_i(\theta)}{f_0(\tilde{y}|\theta)}P(d\theta) - c \right|.$$

Consider now the case of the density bounded class

$$\mathcal{Q}_{DB} = \{Q(\cdot|\theta) : \underline{L}(E|\theta) \le Q(E|\theta) \le \overline{U}(E|\theta) \text{ for all measurable } E\}$$

where $\underline{L}$ and $\overline{U}$ are fixed measures satisfying: $\underline{L}(E|\theta) \le \overline{U}(E|\theta)$ for all measurable $E$ and $\underline{L}(\mathcal{R}^n|\theta) \le 1 \le \overline{U}(\mathcal{R}^n|\theta)$, with densities $\underline{l}$ and $\overline{u}$, respectively.

PROPOSITION 3. *Assume that* $P(\theta_0) = 1$ *and let* $\mathcal{Q} = \mathcal{Q}_{DB}$. *Then*

$$\sup_{Q \in \mathcal{Q}_{SU}} |B'_{F_0}(Q - F_0)| = max\{|S_{\overline{U}} - c|, |S_{\underline{L}} - c|\}$$

*where*

$$S_{\overline{U}} = \sum_{i=1}^{r} \overline{u}(\tilde{y}_i|\theta_0)/f_0(\tilde{y}_i|\theta_0)c_i(\theta_0) \quad and \quad S_{\underline{L}} = \sum_{i=1}^{r} \underline{l}(\tilde{y}_i|\theta_0)/f_0(\tilde{y}_i|\theta_0)c_i(\theta_0).$$

PROOF. The rest follows straightforwardly from the fact that $c_i(\theta_0) \geq 0$ for all $i$ in expression (2) for $B'_{F_0}$. $\square$

**5. Discussion.** This paper analyzes the problem of measuring local sensitivity of the Bayes factor $B$ to small perturbations of the parametric model $\mathcal{F}_0$. The obtained results are expression for the von Mises derivative of $B$ using a direction $Q \in \mathcal{Q}$ and expressions for the supremum of this derivative as $Q$ ranges over particular classes. An alternative approach to the problem of local sensitivity of $B$ is to use standard methods of sensitivity analysis to the prior. Following Lavine (1991), we can define a class $\Gamma$ of priors on the subclass of sampling distributions $\mathcal{Q}$, $\mathcal{F}_0 \subseteq \mathcal{Q} \subseteq \mathcal{F}$, and consider $B$ as a functional defined on $\Gamma$ rather than $\mathcal{Q}$. In this case $B$ can be written as a posterior expectation and standard results are available [see, e.g., Sivaganesan (1993), Gustafson (1994) and Basu (1994)]. In general, the measure of local sensitivity corresponding to this approach is different from von Mises derivative $B'$, so we must consider carefully what kind of class, $\mathcal{Q}$ or $\Gamma$, better represents our initial uncertainty about $\mathcal{F}_0$. It could be interesting to study conditions under which the two measures are coincident.

## REFERENCES

ANTONIAK, C. E. (1974). Mixtures of Dirichlet Processes with Applications to Nonparametric Problems. *Ann. Statist.* **2** 1152-1174.

BASU, S. (1994). Local Sensitivity, Functional Derivatives and Nonlinear Posterior Quantities. Technical Report **63**, Department of Mathemetical Science. University of Arkansas, Fayetteville, AR.

BERGER, J. (1984). The robust Bayesian Viewpoint (with discussion). In *Robustness in Bayesian Statistics* (J. Kadane, ed.) 63-124. North Holland, Amsterdam.

BERGER, J. (1990). Robust Bayesian Analysis: Sensitivity to the prior. *J. Statist. Plann. Inference* **25** 303-328.

BERGER, J. (1994). An overview of Robust Bayesian Analysis (with discussion). *Test* **3**, 1 5-124.

CAROTA, C. (1994). Diagnostica di una Alternativa Nonparametrica per un Modello Statistico Discreto. In *Atti del XXXVII Convegno S.I.S.* 511-518. CISU, Roma.

CAROTA, C. AND PARMIGIANI, G. (1994). On Bayes Factors for Nonparametric Alternatives. To appear in *Bayesian Statistics V: Fifth Valencia International Meeting on Bayesian Statistics* (J. M. Bernardo, J. O. Berger, A. P. David and A. F. M. Smith, eds.).

CLARKE, B. R. (1983). Uniqueness and Fréchet differentiability of functional solutions to maximum likelihood type equations. *Ann. Statist.* **11** 1196-1205.

DIACONIS, P. AND FREEDMAN, P. (1986). On the consistency of Bayes estimates. *Ann. Statist.* **14** 1-67.

FERGUSON, T. S. (1973). A Bayesian Analysis of some Nonparametric Problems. *Ann. Statist.* **1** 209-230.

FERNHOLZ, L. T. (1983). *von Mises Calculus for Statistical Functionals.* Springer-Verlag Lecture Notes in Statistics 19, New York.

GUSTAFSON, P. A. (1994). Local Sensitivity of Posterior expectations. Ph. D. Thesis, Technical Report **596**, Department of Statistics. Carnegie Mellon University. Pittsburgh, PA.

GUSTAFSON, P. A., WASSERMAN, L. A. AND SRINIVASAN, C. (1994). Local sensitivity (whit discussion). To appear in *Bayesian Statistics V: Fifth Valencia International Meeting on Bayesian Statistics* (J. M. Bernardo, J. O. Berger, A. P. David and A. F. M. Smith, eds.).

HAMPEL, F. R. (1971). A general qualitative definition of robustness. *Ann. Math. Statist.* **42** 1886-1896.

HUBER, P. J. (1981). *Robust Statistics.* Wiley, New York.

LAVINE, M. (1991). Sensitivity in Bayesian Statistics: the Prior and the Likelihood. *J. Amer. Statist. Assoc.* **86** 396-399.

RUGGERI, F. AND WASSERMAN, L. A. (1993). Infinitesimal sensitivity of posterior distributions. *Canad. J. Statist.* **21** 195-203.

SIVAGANESAN, S. (1993). Robust Bayesian Diagnostics. *J. Statist. Plann. Inference* **35** 171-188.

SRINIVASAN ,C. AND TRUSZCZYNSKA, H. (1990). Approximation to the range of a ratio-linear posterior quantity based on Fréchet derivative. Technical Report **289**, Department of Statistics, University of Kentucky.

VON MISES, R. (1947). On the Asymptotic Distribution of Differentiable Statistical Functions. *Ann. Math. Statist.* **18** 309-348.

WASSERMAN, L. A. (1992). Recent methodological advances in robust Bayesian inference (with discussion). In *Bayesian Statistics IV: Fourth Valencia International Meeting on Bayesian Statistics* (J. M. Bernardo, J. O. Berger, A. P. David and A. F. M. Smith, eds.) 483-502. Claredon, Oxford.

ISTITUTO DI STATISTICA
UNIVERSITÀ DI PAVIA
STRADA NUOVA 65
I-27100 PAVIA
ITALY