

Conditional inference in parametric models

Michel Broniatowski⁽¹⁾, Virgile Caron⁽¹⁾

⁽¹⁾ LSTA, Université Pierre et Marie Curie, Paris, France

February 5, 2012

Abstract

This paper presents a new approach to conditional inference, based on the simulation of samples conditioned by a statistics of the data. Also an explicit expression for the approximation of the conditional likelihood of long runs of the sample given the observed statistics is provided. It is shown that when the conditioning statistics is sufficient for a given parameter, the approximating density is still invariant with respect to the parameter. A new Rao-Blackwellisation procedure is proposed and simulation shows that Lehmann Scheffé Theorem is valid for this approximation. Conditional inference for exponential families with nuisance parameter is also studied, leading to Monte carlo tests. Finally the estimation of the parameter of interest through conditional likelihood is considered. Comparison with the parametric bootstrap method is discussed.

Keywords: Conditional inference, Rao Blackwell Theorem, Lehmann Scheffé Theorem, Exponential families, Nuisance parameter, Simulation.

1 Introduction and context

This paper explores conditional inference in parametric models. A comprehensive overview on this area is the illuminating review paper by Reid (1995) [27]. Our starting point is as follows: given a model \mathcal{P} defined as a collection of continuous distributions P_θ on \mathbb{R}^d , with density p_θ where the parameter θ belongs to some subset Θ in \mathbb{R}^s and given a sample of independent copies of a random variable with distribution P_{θ_T} for some unknown value θ_T of the parameter, we intend to provide some inference about θ_T conditioning on some statistics of the data. The situations which we have in mind are of two different kinds.

The first one is the Rao-Blackwellisation of estimators, which amounts to reduce the variance of an unbiased estimator by conditioning on any statistics; it is a known fact that such method reduces its variance; when the conditioning statistics is complete and sufficient for the parameter then this procedure provides optimal reduction, as stated by Lehmann-Scheffé Theorem. This realm of questions is the motivation for the first part of this paper:

1. is it possible to provide good approximations for the density of a sample conditioned on a given statistics, and, when applied for a model where

some sufficient statistics for the parameter is known, does sufficiency w.r.t. the parameter still hold for the approximating density?

2. in the case when the first question has positive answer, is it possible to simulate samples according to the approximating density, and to propose some Rao-Blackwellised version for a given preliminary estimator? Also we would hope that the proposed method would be feasible, that the programming burden would be light, that the run time for this method be short, and that the involved techniques would keep in the range of globally known ones by the community of statisticians.

The second application of conditional inference pertains to the role of conditioning in models with nuisance parameters. There is a huge bibliography on this topic, some of which will be considered in details in the sequel. The usual frame for this field of problems is the exponential families one, for reasons related both with the importance of these models in applications and on the role of the concept of sufficiency when dealing with the notion of nuisance parameter. Conditioning on a sufficient statistics for the nuisance parameter produces a new exponential family, which gets free of this parameter, and allows for simple inference on the parameter of interest, at least in simple cases. This will also be discussed, since the reality, as known, is not that simple, and since so many complementary approaches have been developed over decades in this area. Using the approximation of the conditional density in this context and performing simulations yields Monte Carlo tests for the parameter of interest, free from the nuisance parameter. Also conditional maximum likelihood estimators will be produced. Comparison with the parametric bootstrap will also be discussed.

This paper is organized as follows. Section 2 describes a general approximation scheme for the conditional density of long runs of subsamples conditioned on a statistics, with explicit formulas. The (rather lengthy) proof of the main result of this section is presented in Broniatowski and Caron(2011)[4]. Discussion about implementation is provided. Section 3 presents two aspects of the approximating conditional scheme: we first show on examples that sufficiency is kept under the approximating scheme and, second, that this yields to an easy Rao-Blackwellisation procedure. An illustration of Lehmann-Scheffé Theorem is presented. Section 4 deals with models with nuisance parameters in the context of exponential families. We have found it useful to spend a few paragraphs on bibliographical issues. We address Monte Carlo tests based on the simulation scheme; in simple cases its performance is similar to that of parametric bootstrap; however conditional simulation based tests improve clearly over parametric bootstrap procedure when the test pertains to models for which the likelihood is multimodal with respect to the nuisance parameter; an example is provided. Finally we consider conditioned maximum likelihood based on the approximation of the conditional density; in simple cases its performance is similar to that of unconditional likelihood; however when the preliminary estimator of the nuisance is difficult to obtain, for example when it depends strongly on

some initial point for a Newton-Raphson routine (this is indeed a very common situation), then, by the very nature of sufficiency, conditional inference based on the proxy of the conditional likelihood performs better; this is illustrated with examples.

2 The approximate conditional density of the sample

Most attempts which have been proposed for the approximation of conditional densities stem from arguments developed in Lehmann (1986)[16] for inference on the parameter of interest in models with nuisance parameter; however the proposals in this direction hinge at the approximation of the distribution of the sufficient statistics for the parameter of interest given the observed value of the sufficient statistics of the nuisance parameter. We will present some of these proposals in the section devoted to exponential families. To our knowledge, no attempt has been made to approximate the conditional distribution of a sample (or of a long subsample) given some observed statistics.

However, generating samples from the conditional distribution itself (such samples are often called co-sufficient samples, following Lockhart et al.(2007) [20]) has been considered by many authors; see for example Engen and Lillegard (1997)[12], Lindqvist et al. (2003)[17] and references therein, and Lindqvist and Taraldsen (2005)[18].

In Engen and Lillegard (1997)[12], simulating exponential or normal samples under the given value of the empirical mean is proposed. For example under the exponential distribution $Exp(\theta)$, the minimal sufficient statistics for θ is the sum of the observations, say t_n ; a co-sufficient sample x^* can be created by generating an x' -sample from $Exp(1)$ and taking $x_i^* = x'_i t_n / \bar{x}'$. However, this approach may be at odd in simple cases, as for the Gamma density in the non exponential case.

Lockhart et al. (2007)[20] proposed a different framework based on the Gibbs sampler, simulating the conditioned sample one at a time through a sequential procedure. The example which is presented is for the Gamma distribution under the empirical mean, but it seems to perform well, for location parameter, when the true parameter is in some range, therefore not uniformly on the model. Their paper contains a comparative study with the parametric bootstrap procedure (introduced by Efron (1979)[11]) for similar problems. In a simple case, they argue favorably for both methods. We will turn back to parametric bootstrap in relation with conditional likelihood estimators, in the last section of this paper.

Other techniques have been developed in specific cases: for the inverse gaussian distribution see O'Reilly and Gravia-Medrano (2006)[22], Cheng (1984) [8]; for the Weibull distribution see Lockhart et Stephens (1994)[21]. No unified technique exists in the literature which would work under general models.

2.1 Approximation of conditional densities

2.1.1 Notation and hypotheses

For sake of clearness we consider the case when the model \mathcal{P} is a family of distributions on \mathbb{R} . Extension to $\mathbb{R}^d, d > 1$ can be achieved in the same way, using similar results developed in futur work.

Denote $\mathbf{X}_1^n := (\mathbf{X}_1, \dots, \mathbf{X}_n)$ a set of n independent copies of a real random variable \mathbf{X} with density $p_{\mathbf{X}, \theta_T}$ on \mathbb{R} . Let $\mathbf{x}_1^n := (\mathbf{x}_1, \dots, \mathbf{x}_n)$ denote the observed values of the data, each \mathbf{x}_i resulting from the sampling of \mathbf{X}_i . Define the r.v. $\mathbf{U} := u(\mathbf{X})$ and $\mathbf{U}_{1,n} := u(\mathbf{X}_1) + \dots + u(\mathbf{X}_n)$ where u is a real-valued measurable function on \mathbb{R} , and, accordingly, $u_{1,n} := u(\mathbf{x}_1) + \dots + u(\mathbf{x}_n)$. Denote $p_{\mathbf{U}, \theta_T}$ the density of the r.v. \mathbf{U} . We consider approximations of the density of the vector $\mathbf{X}_1^k = (\mathbf{X}_1, \dots, \mathbf{X}_k)$ on \mathbb{R} when $\mathbf{U}_{1,n} = u_{1,n}$. It will be assumed that the observed value $u_{1,n}$ is "typical", in the sense that it keeps in the range of the iterated logarithm law order of magnitude (for large n). Large deviation cases could also be handled, but conditional inference is based on the implicit assumption that such cases are excluded from the analysis. We hence assume

$$\limsup_{n \rightarrow \infty} \frac{|u_{1,n} - nE[u(\mathbf{X})]|}{\sqrt{\text{Var}(u(\mathbf{X}))} \sqrt{2n \log \log n}} = 1. \quad (\text{LIL})$$

We propose an approximation for

$$p_{u_{1,n}, \theta_T}(x_1^k) := p_{\theta_T}(x_1^k | \mathbf{U}_{1,n} = u_{1,n}) \quad (1)$$

where $\mathbf{X}_1^k := (\mathbf{X}_1, \dots, \mathbf{X}_k)$ and $k := k_n$ is an integer sequence such that

$$0 \leq \limsup_{n \rightarrow \infty} k/n \leq 1 \quad (\text{K1})$$

together with

$$\lim_{n \rightarrow \infty} n - k = \infty \quad (\text{K2})$$

which is to say that we approximate $p_{u_{1,n}, \theta_T}(x_1^k)$ on long runs. The rule which define the value of k for a given accuracy of the approximation is stated in section 3.2 of Broniatowski and Caron(2011) [4].

The hypotheses pertaining to the function u and the r.v. $\mathbf{U} = u(\mathbf{X})$ are as follows

1. u is real valued and the characteristic function of the random variable \mathbf{U} is assumed to belong to L^r for some $r \geq 1$.
2. The r.v. \mathbf{U} is supposed to fulfill the Cramer condition: its moment generating function satisfies

$$\phi_{\mathbf{U}}(t) := E \exp t\mathbf{U} < \infty$$

for t in a non void neighborhood of 0.

Define the functions $m(t), s^2(t)$ and $\mu_3(t)$ as the first, second and third derivatives of $\log \phi_{\mathbf{U}}(t)$. Denote

$$\pi_{u, \theta_T}^\alpha(x) := (x) := \frac{\exp tu(x)}{\phi_{\mathbf{U}}(t)} p_{\mathbf{X}, \theta_T}(x)$$

with $m(t) = \alpha$ and α belongs to the support of $P_{\mathbf{X}, \theta_T}$, the distribution of \mathbf{X} . The density π_{u, θ_T}^α is the *tilted* density with parameter α . Also it is assumed that this latest definition of t makes sense for all α in the support of \mathbf{X} . Conditions on $\phi_{\mathbf{U}}(t)$ which ensure this fact are referred to as *steepness properties*, and are exposed in Barndorff-Nielsen(1978)[1], p153.

We introduce a positive sequence ϵ_n which satisfies

$$\lim_{n \rightarrow \infty} \epsilon_n \sqrt{n - k} = \infty \quad (\text{E1})$$

$$\lim_{n \rightarrow \infty} \epsilon_n (\log n)^2 = 0. \quad (\text{E2})$$

2.2 The proxy of the conditional density of the sample

We recursively define the density $g_{u_{1,n}, \theta_T}(x_1^k)$ on \mathbb{R}^k , which approximates $p_{u_{1,n}, \theta_T}(x_1^k)$ sharply with relative error smaller than $\epsilon_n (\log n)^2$. The subscript θ_T will be omitted when there is no ambiguity about the value of the parameter.

Set

$$m_0 := u_{1,n}/n.$$

and

$$g_0(x_1 | x_0) := \pi_{u}^{m_0}(x_1)$$

with x_0 arbitrary, and for $1 \leq i \leq k-1$ define the density $g(x_{i+1} | x_1^i)$ recursively as follows.

Set t_i the unique solution of the equation

$$m_i := m(t_i) = \frac{u_{1,n} - u_{1,i}}{n - i} \quad (2)$$

where $u_{1,i} := u(x_1) + \dots + u(x_i)$. The tilted adaptive family of densities $\pi_{\mathbf{X}}^{m_i}$ is the basic ingredient of the derivation of approximating scheme. Let

$$s_i^2 := \frac{d^2}{dt^2} (\log E_{\pi_{u}^{m_i}} \exp tu(\mathbf{X})) (0)$$

and

$$\mu_j^i := \frac{d^j}{dt^j} (\log E_{\pi_{u}^{m_i}} \exp tu(\mathbf{X})) (0), \quad j = 3, 4$$

which are the second, third and fourth cumulants of π^{m_i} . Let

$$g(x_{i+1} | x_1^i) = C_i p_{\mathbf{X}, \theta_T}(x_{i+1}) \mathbf{n}(\alpha\beta, \beta, u(x_{i+1})) \quad (3)$$

where $\mathbf{n}(\mu, \tau, x)$ is the normal density with mean μ and variance τ at x . Here

$$\beta = s_i^2 (n - i - 1) \quad (4)$$

$$\alpha = t_i + \frac{\mu_3^i}{2s_i^4 (n - i - 1)} \quad (5)$$

and the C_i is a normalizing constant.

Define

$$g_{u_{1,n}}(x_1^k) := g_0(x_1 | x_0) \prod_{i=1}^{k-1} g(x_{i+1} | x_1^i). \quad (6)$$

It holds

Theorem 1 *Assume (K1,K2) together with (E1,E2). Then (i)*

$$p_{u_{1,n}}(x_1^k) = g_{u_{1,n}}(x_1^k)(1 + o_{P_{u_{1,n}}}(\epsilon_n (\log n)^2))$$

and (ii)

$$p_{u_{1,n}}(x_1^k) = g_{u_{1,n}}(x_1^k)(1 + o_{G_{u_{1,n}}}(\epsilon_n (\log n)^2)).$$

For the proof, see Broniatowski and Caron (2011) [4].

Statement (i) means that the conditional likelihood of any long sample path \mathbf{X}_1^k given $\mathbf{U}_{1,n} = u_{1,n}$ can be approximated by $g_{u_{1,n}}(\mathbf{X}_1^k)$ with a small relative error on typical realizations of \mathbf{X}_1^n .

The second statement states that simulating \mathbf{X}_1^k under $g_{u_{1,n}}$ produces runs which could have been sampled under the conditional density $p_{u_{1,n}}$ since $g_{u_{1,n}}$ and $p_{u_{1,n}}$ coincide sharply on larger and larger subsets of \mathbb{R}^k as n increases.

Remark 2 *Theorem 1 states that the density $g_{u_{1,n},(\theta_T, \eta_T)}$ on \mathbb{R}^k approximates $p_{u_{1,n},(\theta_T, \eta_T)}$ on the sample \mathbf{x}_1^n generated under (θ_T, η_T) . However, in some cases, the r.v.'s \mathbf{x}_i 's in Theorem 1 may at time be generated under some other parameters, say (θ_0, η_0) . Indeed, for direct applications developed in this paper, Theorem 1 have to hold when the sample is generated under an other sampling scheme. Broniatowski and Caron (2011) [4] state in Theorem 11 that the approximation scheme holds true in this case.*

Let \mathbf{Y}_1^n be i.i.d. copies of \mathbf{Z} with distribution Q and density q ; assume that Q satisfies the Cramer condition $\int (\exp tx) q(x) dx < \infty$ for t in a non void neighborhood of 0. Let $\mathbf{V}_{1,n} := u(\mathbf{Y}_1) + \dots + u(\mathbf{Y}_n)$ and define

$$q_{u_{1,n}}(y_1^k) := q(\mathbf{Y}_1^k = y_1^k | \mathbf{V}_{1,n} = u_{1,n})$$

with distribution $Q_{u_{1,n}}$. It then holds

Theorem 3 *Then, with the same hypotheses and notation as in Theorem 1,*

$$p(\mathbf{X}_1^k = Y_1^k | \mathbf{U}_{1,n} = u_{1,n}) = g_{u_{1,n}}(Y_1^k)(1 + o_{Q_{u_{1,n}}}(\epsilon_n (\log n)^2)).$$

Also the total variation distance between $Q_{u_{1,n}}$ and $P_{u_{1,n}}$ goes to 0 as n tends to infinity.

2.3 Comments on implementation

The simulation of a sample X_1^k with density $g_{u_{1,n}}$ is fast as easy. Indeed the r.v. X_{i+1} with density $g(x_{i+1}|x_1^i)$ is obtained through a standard acceptance-rejection algorithm. When θ_T is unknown, a preliminary estimator may be used. When $\mathbf{U}_{1,n}$ is sufficient for $p_{u_{1,n}}$ it is nearly sufficient for its proxy $g_{u_{1,n}}$ (see next section); indeed changing the value of this preliminary estimator does not alter the likelihood of the sample; as shown in the simulations developed here after, any value of θ can be used; call θ^* the θ chosen as initial value, using henceforth $p_{\mathbf{X},\theta^*}$ instead of $p_{\mathbf{X},\theta_T}$ in (3). In exponential families the values of the parameters which appear in the gaussian component of $g(x_{i+1}|x_1^i)$ in (3) are easily calculated; note also that due to (LIL) the parameters in $\mathbf{n}(\alpha\beta, \beta, u(x_{i+1}))$ are such that the dominating density can be chosen for all i as $p_{\mathbf{X},\theta^*}$. The constant in the acceptance rejection algorithm is then $1/\sqrt{2\pi\alpha}$. This is in contrast with the case when the conditioning value is in the range of a large deviation with respect to $p_{\mathbf{X},\theta_T}$; in this case, which appears in a natural way in Importance sampling estimation for rare event probabilities, the simulation algorithm is more complex; see [5].

3 Sufficient statistics and approximated conditional density

3.1 Keeping sufficiency under the proxy density

The density $g_{u_{1,n}}(y_1^k)$ is used in order to handle Rao-Blackellisation of estimators or statistical inference for models with nuisance parameters. The basic property is sufficiency with respect to the envolved parameter. We show on some examples that $g_{u_{1,n}}(y_1^k)$ defined in (6) inherits of the invariance with respect to a parameter when conditioning on a sufficient statistics pertaining to this parameter.

Consider the Gamma density

$$f_{\rho,\theta}(x) := \frac{\theta^{-\rho}}{\Gamma(\rho)} x^{\rho-1} \exp -x/\theta \quad \text{for } x > 0. \quad (7)$$

As r varies in \mathbb{R}^+ and θ is positive, the density runs in an exponential family $\gamma_{r,\theta}$ with parameters $r := \rho - 1$ and θ , and sufficient statistics $t(x) := \log x$ and $u(x) := x$ respectively for r and θ . Given an i.i.d. sample $X_1^n := (X_1, \dots, X_n)$ with density f_{r_T,θ_T} the resulting sufficient statistics are respectively $T_{1,n} := \log X_1 + \dots + \log X_n$ and $U_{1,n} := X_1 + \dots + X_n$. We consider two parametric models $(\gamma_{r_T,\theta}, \theta \geq 0)$ and $(\gamma_{r,\theta_T}, r > 0)$ respectively assuming r_T or θ_T known.

We first consider sufficiency of $U_{1,n}$ in the first model. The density $g_{u_{1,n}}(y_1^k)$ should be free of the current value of the true parameter θ_T of the parameter under which the data are drawn. However as appears in (6) the unknown value θ_T should be used in its very definition. We show by simulation that whatever the value of θ inserted in place of θ_T in (6) the likelihood of X_1^k under $g_{u_{1,n}}$

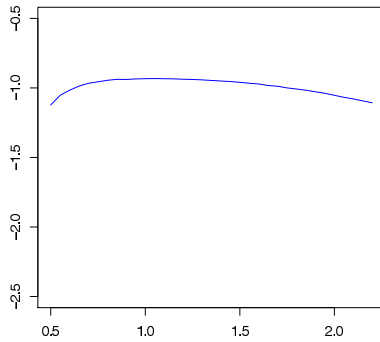


Figure 1: Proxy of the conditional likelihood of X_1^k under $g_{T_{1,n}}$ as a function of θ for $n = 100$ and $k = 80$ in the gamma case.

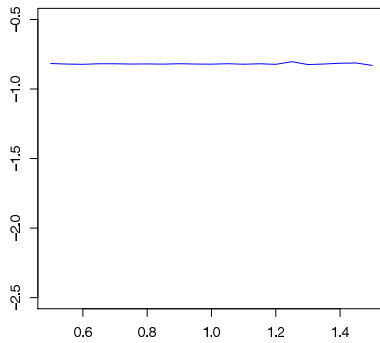


Figure 2: Proxy of the conditional likelihood of X_1^k under $g_{U_{1,n}}$ as a function of r for $n = 100$ and $k = 80$ in the gamma case.

does not depend upon θ . We thus observe that $U_{1,n}$ is sufficient for θ_T in the conditional density approximating $p_{u_{1,n}}$ as should hold as a consequence of Theorem 1 .

Similarly the same fact occurs replacing θ_T by r_T in the model $(\gamma_{r,\theta_T}, r > 0)$.

In both cases whatever the value of the parameter θ (Figure 1) or r (Figure 2), the likelihood of X_1^k remains constant.

We also consider the Inverse Gaussian distribution with density

$$f_{\lambda,\mu}(x) := \left[\frac{\lambda}{2\pi} \right]^{1/2} \exp -\frac{\lambda(x-\mu)^2}{2\mu^2 x} \quad \text{for } x > 0 \quad (8)$$

with both parameters λ and μ be positive. Given an i.i.d. sample $X_1^n := (X_1, \dots, X_n)$ with density $f_{\mu,\lambda}$, the resulting sufficient statistics are respectively $T_{1,n} := X_1 + \dots + X_n$ and $U_{1,n} := X_1^{-1} + \dots + X_n^{-1}$. Similarly as for the Gamma case we draw the likelihood of a subsample X_1^k under $g_{u_{1,n}}$ with $T_{1,n} := X_1 + \dots + X_n$, which is a sufficient statistics for μ (Figure 3), and upon $U_{1,n} :=$

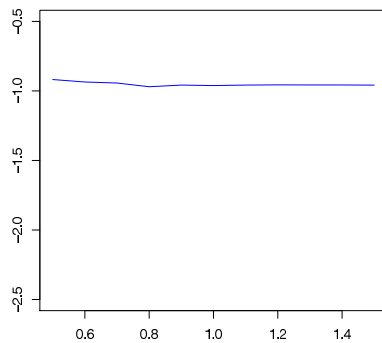


Figure 3: Conditional likelihood of X_1^k under $g_{T_{1,n}}$ as a function of μ for $n = 100$ and $k = 80$ in the Inverse Gaussian case.

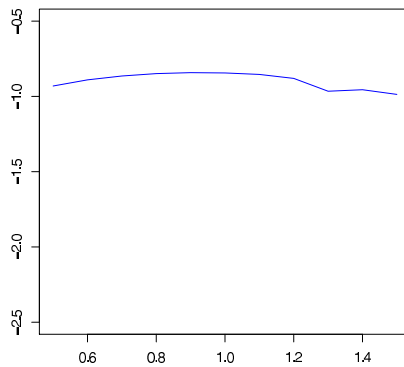


Figure 4: Conditional likelihood of X_1^k under $g_{U_{1,n}}$ as a function of λ for $n = 100$ and $k = 80$ in the Inverse Gaussian case.

$X_1^{-1} + \dots + X_n^{-1}$ which is sufficient for λ (Figure 4). In either cases the other coefficient is kept fixed at the true value of the parameter generating the sample. As for the Gamma case these curves show the invariance of the proxy of the conditional density with respect to the parameter for which the chosen statistics is sufficient.

3.2 Rao-Blackwellisation

Rao-Blackwell Theorem holds regardless of whether biased or unbiased estimators are used, since it reduces the MSE. Although its statement is rather weak, in practice, however, the improvement is often enormous. New interest in Rao-Blackwellisation procedures have risen in the recent years, conditioning on ancillary variables (see Fraser(2004) [13] for a survey on ancillaries in conditional inference); specific Rao-Blackwellisation schemes have been proposed by Casella and Robert [6], [7], Perron(1999)[26], Douc and Robert (2010)[28] and Iacobucci et al.(2010) [14]. The purpose is to improve the variance of a given statistics (for

example a tail probability) under a *known* distribution, through a simulation scheme under this distribution; the ancillary variables used in the simulation process itself are used as conditioning ones for the Rao-Blackwellisation of the statistics. The present approach is more classical in this respect, since we do not assume that the parent distribution is known; conditioning on a sufficient statistics $\mathbf{U}_{1,n}$ with respect to the parameter θ and simulating samples according to the approximating density $g_{u_{1,n}}$ will produce the improved estimator.

Since $U_{1,n}$ is sufficient for the parameter θ in $g_{u_{1,n}}$ it can be used in order to obtain improved estimators of θ_T through Rao Blackwellization. We shortly illustrate the procedure and its results on some toy cases. Consider again the Gamma family defined here-above with canonical parameters r and θ .

First the parameter to be estimated is θ_T . A first unbiased estimator is chosen as

$$\hat{\theta}_2 := \frac{X_1 + X_2}{2r_T}.$$

Given an i.i.d. sample X_1^n with density γ_{r_T, θ_T} the Rao-Blackwellised estimator of $\hat{\theta}$ is defined through

$$\theta_{RB,2} := E\left(\hat{\theta}_2 \mid U_{1,n}\right)$$

whose variance is less than $Var\hat{\theta}_2$.

Consider $k = 2$ in $g_{U_{1,n}}(y_1^k)$ and let (Y_1, Y_2) be distributed according to $g_{u_{1,n}}(y_1^2)$. Replications of (Y_1, Y_2) induce an estimator of $\theta_{RB,2}$ for fixed $u_{1,n}$. Iterating on the simulation of the runs X_1^n produces, for $n = 100$ an i.i.d. sample of $\theta_{RB,2}$'s and the $Var\theta_{RB,2}$ is estimated. The resulting variance shows a net improvement with respect to the estimated variance of $\hat{\theta}_2$. It is of some interest to confront this gain in variance as the number of terms involved in $\hat{\theta}_k$ increases together with k . As k approaches n the variance of $\hat{\theta}_k$ approaches the Cramer Rao bound. The graph below shows the decay in variance of $\hat{\theta}_k$. We note that whatever the value of k the estimated value of the variance of $\theta_{RB,k}$ is constant. This is indeed an illustration of Lehmann-Scheffé's theorem.

Remark 4 *Lockhart and O'Reilly (2005) [19] establish, under certain conditions and for fixed k , the asymptotic equivalence of the plug-in estimate for the distribution $P_{\theta_{ML}}(\mathbf{X}_1^k \in B)$ and the Rao-Blackwell estimate $P(\mathbf{X}_1^k \in B \mid \mathbf{U}_{1,n})$ where θ_{ML} is the maximum likelihood estimator of θ_T based on the whole sample \mathbf{X}_1^n (this result is known as Moore's conjecture (see Moore(1973)[23])). They also provide rates for this convergence.*

4 Exponential models with nuisance parameters

4.1 Conditional inference in exponential families

We consider the case when the parameter consists in two distinct subparameters, one of interest denoted θ and a nuisance component denoted η . As is well known, conditioning on a sufficient statistics for the nuisance parameter produces a

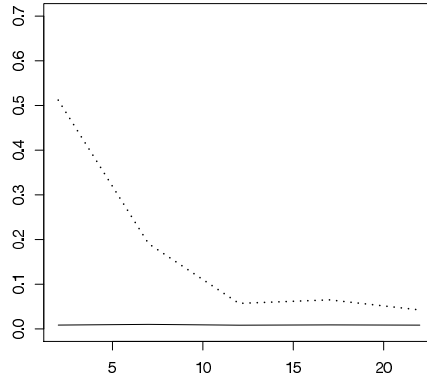


Figure 5: Variance of $\hat{\theta}_k$, the initial estimator (dotted line), along with the variance of $\theta_{RB,k}$, the Rao-Blackwellised estimator (solid line) with $n = 100$ as a function of k .

new exponential family which is free of it. Assuming the observed dataset $\mathbf{x}_1^n := (\mathbf{x}_1, \dots, \mathbf{x}_n)$ resulting from sampling of a vector $\mathbf{X}_1^n := (\mathbf{X}_1, \dots, \mathbf{X}_n)$ of i.i.d. random variables with distribution in the initial exponential model, and denoting $u(\mathbf{x}_1^n)$ a sufficient statistics for η , simulation of samples under the conditional distribution of \mathbf{X}_1^n given $u(\mathbf{X}_1^n) = u(\mathbf{x}_1^n)$ and $\theta = \theta_0$ for some θ_0 produces the basic ingredient for Monte Carlo tests with $H_0 : \theta_T = \theta_0$ where θ_T stands for the true value of the parameter of interest. Changing θ_0 for other values of the parameter of interest produces power curves as functions of the level of the test. This is the well known principle of Monte Carlo tests, and such is the goal of the present section. We consider a steep but not necessarily regular exponential family $\mathcal{P} := \{P_{\theta,\eta}, (\theta, \eta) \in \mathcal{N}\}$ defined on \mathbb{R} with canonical parametrization (θ, η) and minimal sufficient statistics (t, u) defined on \mathbb{R} through the density

$$p_{\theta,\eta}(x) := \frac{dP_{\theta,\eta}(x)}{dx} = \exp[\theta t(x) + \eta u(x) - K(\theta, \eta)] h(x). \quad (9)$$

For notational conveniency and without loss of generality both θ and η belong to \mathbb{R} . Also the model can be defined on \mathbb{R}^d , $d > 1$, at the cost of similar but more involved tools. The natural parameter space is \mathcal{N} (which is a convex set in \mathbb{R}^2) defined as the effective domain of

$$k(\theta, \eta) := \exp[K(\theta, \eta)] = \int \exp[\theta t(x) + \eta u(x)] h(x) dx. \quad (10)$$

Let $X_1^n := (X_1, \dots, X_n)$ be n i.i.d. replications of a general random variable \mathbf{X} with density (9). Denote

$$T_{1,n} := \sum_{i=1}^n t(X_i) \quad \text{and} \quad U_{1,n} := \sum_{i=1}^n u(X_i). \quad (11)$$

Basu (1977) [3] discusses ten different ways for eliminating the nuisance parameters, among which conditioning on sufficient statistics and consider UMPU tests pertaining to the parameter of interest. In most cases, the density of $T_{1,n}$ given $U_{1,n} = u_{1,n}$ is unknown. Two main ways have been developed to deal with this issue: approximating this conditional density of a statistics or simulating samples from the conditional density. We compose these two approaches in the present paper.

The classical technique is to approximate this conditional density using some expansion. Then integration produces critical values. For example, Pedersen (1979) [24] states the mixed Edgeworth-saddlepoint approximation, or the single saddlepoint approximation. However, the main issue of this technique is that the approximated density still depends on the nuisance parameter. In order to obtain the expansion, some suitable values for the parameter of interest and for the nuisance parameter have to be chosen. In the method developed here, as seen before, the conditional approximated density inherits of the invariance with respect to the nuisance parameter when conditioning on a sufficient statistics pertaining to this parameter.

Rephrasing the notation of Section 2 in the present setting it holds that the MLE (θ_{ML}, η_{ML}) satisfies

$$\left. \frac{\partial K(\theta, \eta)}{\partial \eta} \right|_{\theta_{ML}, \eta_{ML}} = u_{1,n}/n$$

and therefore $u_{1,n}/n$ converges to $\left(\frac{\partial K(\theta_T, \eta)}{\partial \eta} \right)^{-1}(\eta_T)$.

For notational clearness denote μ the expectation of $u(\mathbf{X}_1)$ and σ^2 its variance under (θ_T, η_T) , hence

$$\mu := \mu_{(\theta_T, \eta_T)} := \partial K(\theta_T, \eta_T) / \partial \eta \quad \sigma^2 := \sigma_{(\theta_T, \eta_T)}^2 := \partial^2 K(\theta_T, \eta_T) / \partial \eta^2$$

Assume at present θ_T and η_T known. It holds

$$\phi(r) := E_{(\theta_T, \eta_T)} \exp[ru(\mathbf{X})] = \exp[K(\theta_T, \eta_T + r) - K(\theta_T, \eta_T)]$$

and

$$\begin{aligned} m(r) &= \mu_{(\theta_T, \eta_T + r)} \\ s^2(r) &= \sigma_{(\theta_T, \eta_T + r)}^2 \\ \mu_3(r) &= \partial^3 K(\theta_T, \eta_T + r) / \partial \eta^3 . \end{aligned}$$

Further

$$\pi_{u, \theta_T, \eta_T}^\alpha(x) := \frac{\exp ru(x)}{\phi(r)} p_{(\theta_T, \eta_T)}(x) = p_{(\theta_T, \eta_T + r)}(x) \quad (12)$$

for any given α in the range of P_{θ_T, η_T} . In the above formula (12) the parameter r denotes the only solution of the equation

$$m(r) = \alpha.$$

For large k depending on n , using Monte Carlo tests based on runs of length k instead of n does not affect the accuracy of the results.

4.2 Application of conditional sampling to MC tests

Consider a test defined through $H_0 : \theta_T = \theta_0$ versus $H_1 : \theta_T \neq \theta_0$ Monte Carlo (MC) tests aim at obtaining p -values through simulation. where the distribution of the desired test statistics under H_0 is either unknown or very cumbersome to obtain; a comprehensive reference is Jöckel(1986), [15].

Recall the principle of those tests: denote t the observed value of the studied statistic based on the dataset and let t_2, \dots, t_L the values of the resulting test statistics obtained through the simulation of $L - 1$ samples \mathbf{X}_1^n under H_0 . If t is the M th largest value of the sample (t, t_2, \dots, t_L) , H_0 will be rejected at the $\alpha = M/L$ significance level, since the rank of t is uniformly distributed on the integer $2, \dots, L$ when H_0 holds. Calculation of power functions can be handled similarly. The above approximation of the conditional density $p_{u_{1,n}}(x_1^k)$ involves the unknown parameters θ_T and η_T in all the simulation steps. This problem is solved when simulating under $H_0 : \theta_T = \theta_0$ setting θ_0 in place of θ_T and $\hat{\eta}_{\theta_0}$ in place of η_T , where $\hat{\eta}_{\theta_0}$ is the MLE of η_T in the one parameter family $p_{\theta_0, \eta}$ defined through (9). This choice follows the commonly used one, as advocated for instance in [24] and [25]. Innumerous simulation studies support this choice in various contexts.

Consider the problem of testing the null hypothesis $H_0 : \theta_T = \theta_0$ against the alternative $H_1 : \theta_T > \theta_0$ in model (10) where η is the nuisance parameter.

When $p_{u_{1,n}, \theta_0}$ is known, the classical conditional test $H_0 : \theta_T = \theta_0$ versus $H_1 : \theta_T > \theta_0$ with level α is UMPU.

Substituting $p_{\theta_0}(\mathbf{X}_1^n = x_1^n | \mathbf{U}_{1,n} = u_{1,n})$ by $g_{u_{1,n}, \theta_0}(x_1^k)$ defined in (6), i.e. substituting the test statistics T_1^n by T_1^k and $p_{\theta_0}(\mathbf{X}_1^k = x_1^k | \mathbf{U}_{1,n} = u_{1,n})$ by $g_{u_{1,n}, \theta_0}(x_1^k)$ i.e. changing the model for a proxy while keeping the same parameter of interest θ yields the conditional test with level α

$$\psi_\alpha(x_1^k) := \begin{cases} 1 & \text{if } T_{1,k} > t_\alpha \\ \gamma & \text{if } T_{1,k} = t_\alpha \\ 0 & \text{if } T_{1,k} < t_\alpha \end{cases}$$

and

$$E_{G_{u_{1,n}}}[\psi_\alpha(X_1^k)] = \alpha$$

i.e. $\alpha := \int \mathbb{1}_{t_{1,k} > t_\alpha} g_{u_{1,n}}(x_1^k) dx_1 \dots dx_k$. Its power under a simple hypothesis $\theta_T = \theta$ is defined through

$$\beta_{\psi_\alpha}(\theta | u_n) = E_\theta[\psi_\alpha(T_{1,n}, U_{1,n}) | U_{1,n} = u_{1,n}].$$

Recall that the parametric bootstrap produces samples from a parametric model which is fitted to the data, often through maximum likelihood. In the present setting, the parameter θ is set to θ_0 and the nuisance parameter η is replaced by its estimator $\hat{\eta}_{\theta_0}$ which is the MLE of η when the parameter θ is fixed at the value θ_0 defining H_0 . Comparing their exact conditional MC tests with parametric bootstrap ones for Gamma distributions, Lockhart et al(2007)[19] conclude that no significant difference can be noticed in terms of level or in terms of power. We proceed in the same vein, comparing conditional sampling

MC tests with the parametric bootstrap ones, obtaining again similar results when the nuisance parameter is estimated accurately. However the results are somehow different when the nuisance parameter cannot be estimated accurately, which may occur in various cases.

In practice since the chosen conditioning statistics is quasi sufficient for the nuisance parameter, we plug any value for this parameter in the definition of $g_{u_{1,n}}$. This is what has been performed in all examples below.

4.3 Unimodal Likelihood: testing the coefficients of a Gamma distribution

Let \mathbf{X}_1^n be an i.i.d. sample of random variables with Gamma distribution $\Gamma(a_T, b_T)$ where a_T is the shape coefficient and b_T is the scale coefficient. As a and b vary this distribution is a two parameter exponential family. The statistics $T_{1,n} := \mathbf{X}_1 + \dots + \mathbf{X}_n$ is sufficient for the parameter a and $U_{1,n} := \log \mathbf{X}_1 + \dots + \log \mathbf{X}_n$ is sufficient for b .

MC conditional test with $H_0 : a_T = a_0$ Denote $u_{1,n} = \sum_{i=1}^n X_i$ and \hat{b}_{a_0} the MLE of b . Calculate for $l \in \{2, L\}$

$$t_l := \sum_{i=0}^k \log(Y_i(l)).$$

where the Y_i^l are a sample from $g_{u_{1,n}}^{a_0, \hat{b}_{a_0}}$.

Consider the corresponding parametric bootstrap procedure for the same test, namely simulate $Z_i(l)$, $2 \leq l \leq L$ and $0 \leq i \leq k$ with distribution $\Gamma(a_0, \hat{b}_{a_0})$; denote

$$s_l := \sum_{i=0}^k \log(Z_i(l)).$$

In this example simulation shows that for any α the M th largest value of the sample (t, t_2, \dots, t_L) is very close to the corresponding empirical M/L -quantile of s_l 's. Hence Monte Carlo tests through parametric bootstrap and conditional compete equally. Also in terms of power, irrespectively in terms of α and in terms of alternatives (close to H_0), the two methods seem to be equivalent.

MC conditional test with $H_0 : b_T = b_0$ Denote $u_{1,n} = \sum_{i=1}^n \log(X_i)$ and \hat{a}_{b_0} the MLE of a . Calculate for $l \in \{2, L\}$

$$t_l := \sum_{i=0}^k Y_i(l)$$

where the Y'_i are a sample from $g_{u_{1,n}}^{b_0, \hat{a}_{b_0}}$ and, as above define accordingly

$$s_l := \sum_{i=0}^k \log(Z_i(l))$$

where the $Z_i(l)$'s are simulated under $\Gamma(\hat{a}_{b_0}, b_0)$.

As above, parametric bootstrap and conditional sampling yield equivalent Monte Carlo tests in terms of power function under alternatives close to H_0 .

In the two cases studied above the value of k has been obtained through the rule exposed in section 3.2 of Broniatowski and Caron (2011) [4].

4.3.1 Bimodal likelihood: testing the mean of a normal distribution in dimension 2

In contrast with the above mentioned examples, the following case study shows that estimation through the unconditional likelihood may fail to provide consistent estimators when the likelihood surface has multiple critical points. This in turn yields parametric bootstrap Monte Carlo tests with unacceptable power functions.

Sundberg(2009)[30] proposes four examples that allow likelihood multimodality. Two of them can also be found in [9] and [10], and in [2], Ch 2. We consider the "Normal parabola" model which is a curved (2, 1) family (see Example 2.35 in [2], Ch 2). Two independent Gaussian variates have unknown means and known variances; their means are related by a parabolic relationship.

Let X et Y be two independent gaussian r.v.'s with same variance σ_T^2 with expectation ψ_T and ψ_T^2 . In the present example $\sigma_T^2 = 1$ and $\psi_T = 2$.

Let (X_i, Y_i) , $1 \leq i \leq n$ be an i.i.d. sample with the above distribution.

The parameter of interest is σ^2 whilst the nuisance parameters is ψ . Derivation of the likelihood function of the observed sample with respect to ψ yields the following equation

$$(U_{1,n} - \psi) + 2\psi (V_{1,n} - \psi^2) = 0$$

with $U_{1,n} := X_1 + \dots + X_n$ and $V_{1,n} := Y_1 + \dots + Y_n$. The following table shows that the likelihood function is bimodal in ψ .

Estimation of the nuisance parameter ψ is performed through the standard Newton Raphson method. The Newton-Raphson optimizer of the likelihood function converges to the true value when the initial value is larger than 1 and fails to converge to $\psi_T = 2$ otherwise. Hencefore the parametric bootstrap estimation of the likelihood function of the sample based on this preliminary estimate of the nuisance parameter may lead to erroneous estimates of the parameter of interest. Indeed according to the initial value we obtained estimators of ψ_T close to 2 or to -2 . When the estimator of the nuisance parameter is close to its true value 2 then parametric bootstrap yields Monte Carlo tests with power close to 1 for any α and any alternative close to H_0 . At the contrary when this estimate is close to the second maximizer of the likelihood (i.e. close to -2) then

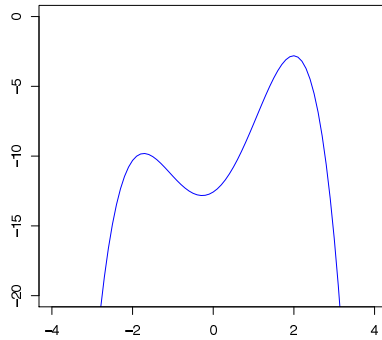


Figure 6: Bimodal likelihood in ψ .

the resulting Monte Carlo test based on parametric bootstrap has power close to 0 irrespectively of the value of α and of the alternative, when close to H_0 . In contrast with these results, Monte Carlo tests based on conditional sampling provide powers close to 1 for any α ; we have considered alternatives close to H_0 . This result is of course a consequence of quasi sufficiency of the statistics $(U_{1,n}, V_{1,n})$ for the parameter ψ of the distribution of the sample $(X_i, Y_i)_{i=1, \dots, n}$; see next paragraph for a discussion of this point.

4.4 Estimation through conditional likelihood

Considering model (10) we intend to perform an estimation of θ_T irrespectively upon the value of η_T . Denote $\hat{\eta}_\theta$ the MLE of η_T when θ holds; the model $p_{\theta, \hat{\eta}_\theta}(x)$ is a one parameter model which is fitted to the data for any peculiar choice of θ . The classic unconditional likelihood provides consistent estimators of θ_T in many cases. However, this method strongly relies on the consistency properties of $\hat{\eta}_\theta$ at any given θ .

For fixed parameter value θ of the parameter of interest, Theorem 1 means that the likelihood of the subsample \mathbf{X}_1^k with unknown distribution with parameter (θ_T, η_T) under the distribution with any parameter θ given the value of the sufficient statistics $U_{1,n}$ is approximated by $g_{u_{1,n}}(\mathbf{X}_1^k)$ when \mathbf{X}_1^k is either generated under the conditional density or under $g_{u_{1,n}}$ with parameter $\eta = \eta_T$. Substituting η_T by its estimator should yield maximal value of the approximate likelihood when (θ_T, η_T) holds, since $\hat{\eta}_\theta$ approaches η_T when $\theta = \theta_T$. In particular, this holds when \mathbf{X}_1^k is generated under θ_T, η_T which holds on the observed sample. This yields to an algorithm to estimate θ_T . For any θ calculate $\hat{\eta}_\theta$. Evaluate $g_{u_{1,n}}(\mathbf{X}_1^k)$ and optimize in θ .

In most cases, as the normal, gamma or inverse-gaussian, both estimation through the unconditional likelihood and estimation through conditional likelihood based on the proxy $g_{u_{1,n}}$ give a consistent estimator.

We consider the example of the Bimodal likelihood from the above subsection, inheriting of the notation and explore the behaviour of the proxy of the conditional likelihood of the sample (X_i, Y_i) , $1 \leq i \leq n$ when conditioning on

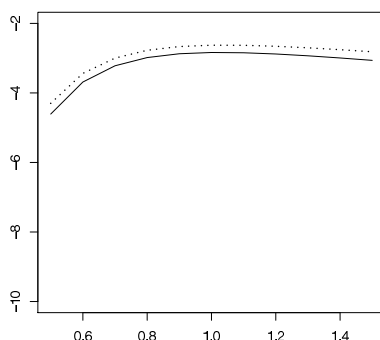


Figure 7: Proxy of the conditional likelihood (solid line) along with the empirical likelihood (dotted line) as function of σ^2 for $n = 100$ and $k = 99$ in the case where a good initial point in Newton-Raphson procedure is chosen.

$U_{1,n}$ and $V_{1,n}$, as a function of σ^2 . The likelihood writes

$$\begin{aligned} L(\sigma^2 | (X_i, Y_i), 1 \leq i \leq n, U_{1,n}, V_{1,n}) \\ = P_{\mathbf{X}_1^n}(X_1^n | U_{1,n}, \sigma^2) P_{\mathbf{Y}_1^n}(Y_1^n | V_{1,n}, \sigma^2) \end{aligned}$$

where we have used the independence of the r.v.'s X_i 's and Y_i 's.

Applying Theorem 1 to the above expression it appears that ψ cancels in the resulting density $g_{u_{1,n}}$ and $g_{v_{1,n}}$. This proves that the proxy of the conditional likelihood provides consistent estimation of σ_T^2 as shown on Figures 7 and 8 (see the solid lines).

On Figure 7, the dot line is the empirical likelihood function with consistent estimator of the nuisance parameter; the resulting maximizer in the variable σ^2 is close to $\sigma_T^2 = 1$. At the opposite in Figure 8 an inconsistent preliminary estimator of ψ_T obtained through a bad tuning of the initial point in the Newton-Raphson procedure leads to inconsistency in the estimation of σ_T^2 , the resulting likelihood function being unbounded.

References

- [1] BARNDORFF-NIELSEN, O.E. (1978). Information and exponential families in statistical theory. Wiley Series in Probability and Mathematical Statistics. Chichester: John Wiley & Sons.
- [2] BARNDORFF-NIELSEN, O.E. AND COX, R.R. (1994). Inference and Asymptotics. Chapman & Hall, London.
- [3] BASU, D. (1977). On the elimination of nuisance parameters. J. Amer. Statist. Assoc. 72 (1977), no. 358, 355–366.
- [4] BRONIATOWSKI M. AND CARON, V. (2011). Long runs under point conditioning. The real case. arXiv:1010.3616v5.

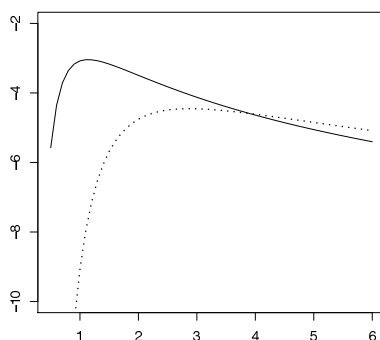


Figure 8: Proxy of the conditional likelihood (solid line) along with the empirical likelihood (dotted line) as function of σ^2 for $n = 100$ and $k = 99$ in the case where a bad initial point in Newton-Raphson procedure is chosen.

- [5] BRONIATOWSKI M. AND CARON, V. (2011). Towards zero variance estimators for rare event probabilities. arXiv:1104.1464v2.
- [6] CASELLA, G. AND ROBERT, C. P. (1996). Rao-Blackwellisation of sampling schemes. *Biometrika* 83 , no. 1, 81–94.
- [7] CASELLA, G. AND R. C. P. (1998). Post-processing accept-reject samples: recycling and rescaling. *J. Comput. Graph. Statist.* 7 , no. 2, 139–157
- [8] CHENG, R.C.H. (1984). Generation of inverse Gaussian variates with given sample mean and dispersion. *Appl. Statist.* 33, 309–16
- [9] EFRON, B. (1975). Defining the curvature of a statistical problem (with applications to second order efficiency) (with discussion). *Ann. Statist.* 3, 1189-1242.
- [10] EFRON, B. (1978). The geometry of exponential families. *Ann. Statist.* 6, 362-376.
- [11] EFRON, B. (1979), Bootstrap methods: another look at the jackknife. *Ann. Statist.* 7 , no. 1, 1–26.
- [12] ENGEN, S. AND LILLEGARD, M. (1997). Stochastic simulations conditioned on sufficient statistics. *Biometrika*, 84, 235–240.
- [13] FRASER, D. A. S. (2004). Ancillaries and conditional inference. With comments by Ronald W. Butler, Ib M. Skovgaard, Rudolf Beran and a rejoinder by the author. *Statist. Sci.* 19 , no. 2, 333–369
- [14] IACOBUCCI, A., MARIN, J.-M., ROBERT, C. (2010). On variance stabilisation in population Monte Carlo by double Rao-Blackwellisation. *Comput. Statist. Data Anal.* 54 , no. 3,

- [15] JÖCKEL, K.-H. (1986). Finite sample properties and asymptotic efficiency of Monte Carlo tests. *Ann. of Stat.*, 14, 336–347.
- [16] LEHMANN, E.L. (1986). *Testing Statistical Hypotheses*. Springer.
- [17] LINDQVIST, B.H., TARALDSEN, G., LILLEGÄRD, M. AND ENGEN, S. (2003). A counterexample to a claim about stochastic simulations. *Biometrika* 90 , no. 2, 489–490.
- [18] LINDQVIST, B.H. AND TARALDSEN, G. (2005). Monte Carlo conditioning on a sufficient statistic. *Biometrika* 92 , no. 2, 451–464.
- [19] LOCKHART, R. AND O'REILLY, F. (2005) A note on Moore's conjecture. *Statist. Probab. Lett.* 74 , no. 2, 212–220.
- [20] LOCKHART, R. A., O'REILLY, F. J. AND STEPHENS, A. (2007) Use of the Gibbs sampler to obtain conditional tests, with applications. *Biometrika* 94 , no. 4, 992–998.
- [21] LOCKHART, R.A. AND STEPHENS, M. A. (1994). Estimation and tests of fit for the three-parameter Weibull distribution. *J. Roy. Statist. Soc.. B.* 56:491–500.
- [22] O'REILLY, F. AND GRAVIA-MEDRANO, L. (2006). On the conditional distribution of goodness-of-fit tests. *Commun. Statist. A*, 35:541–9.
- [23] MOORE, DAVID S. (1973). A note on Srinivasan's goodness-of-fit test. *Biometrika* 60 , 209–211.
- [24] PEDERSEN, B.V., (1979). Approximating conditional distributions by the mixed Edgeworth-saddlepoint expansion. *Biometrika*, 66(3), 597–604.
- [25] PACE L. AND SALVAN A., (1992). A note on conditional cumulants in canonical exponential families. *Scand. J. Statist.*, 19, 185–191.
- [26] PERRON, F. (1999). Beyond accept-reject sampling. *Biometrika* 86 , no. 4, 803–813
- [27] REID, N. (1995). The roles of conditioning in inference. With comments by V. P. Godambe, Bruce G. Lindsay and Bing Li, Peter McCullagh, George Casella, Thomas J. DiCiccio and Martin T. Wells, A. P. Dawid and C. Goutis and Thomas Severini. With a rejoinder by the author. *Statist. Sci.* 10 , no. 2, 138–157, 173–189, 193–196.
- [28] DOUC, R. AND ROBERT, C.P.(2011). A vanilla Rao-Blackwellization of Metropolis-Hastings algorithms. *Ann. Statist.* 39 , no. 1, 261–277
- [29] SEVERINI T.A. (1994), On the approximate elimination of nuisance parameters by conditioning. *Biometrika*, 81(4), 649–661.

- [30] SUNDBERG, R. (2009). Flat and multimodal likelihoods and model lack of fit in curved exponential families. Research Report 2009:1, <http://www.math.su.se/matstat>.