

# WILEY



---

On the Convergence of Monte Carlo Maximum Likelihood Calculations

Author(s): Charles J. Geyer

Source: *Journal of the Royal Statistical Society. Series B (Methodological)*, Vol. 56, No. 1 (1994), pp. 261-274

Published by: [Wiley](#) for the [Royal Statistical Society](#)

Stable URL: <http://www.jstor.org/stable/2346044>

Accessed: 01/07/2014 23:36

---

Your use of the JSTOR archive indicates your acceptance of the Terms & Conditions of Use, available at <http://www.jstor.org/page/info/about/policies/terms.jsp>

JSTOR is a not-for-profit service that helps scholars, researchers, and students discover, use, and build upon a wide range of content in a trusted digital archive. We use information technology and tools to increase productivity and facilitate new forms of scholarship. For more information about JSTOR, please contact support@jstor.org.



Wiley and Royal Statistical Society are collaborating with JSTOR to digitize, preserve and extend access to *Journal of the Royal Statistical Society. Series B (Methodological)*.

<http://www.jstor.org>

# On the Convergence of Monte Carlo Maximum Likelihood Calculations

By CHARLES J. GEYER†

*University of Minnesota, Minneapolis, USA*

[Received February 1992. Revised November 1992]

## SUMMARY

Monte Carlo maximum likelihood for normalized families of distributions can be used for an extremely broad class of models. Given any family  $\{h_\theta: \theta \in \Theta\}$  of non-negative integrable functions, maximum likelihood estimates in the family obtained by normalizing the functions to integrate to 1 can be approximated by Monte Carlo simulation, the only regularity conditions being a compactification of the parameter space such that the evaluation maps  $\theta \mapsto h_\theta(x)$  remain continuous. Then with probability 1 the Monte Carlo approximant to the log-likelihood hypoconverges to the exact log-likelihood, its maximizer converges to the exact maximum likelihood estimate, approximations to profile likelihoods hypoconverge to the exact profile and level sets of the approximate likelihood (support regions) converge to the exact sets (in Painlevé-Kuratowski set convergence). The same results hold when there are missing data if a Wald-type integrability condition is satisfied. Asymptotic normality of the Monte Carlo error and convergence of the Monte Carlo approximation to the observed Fisher information are also shown.

**Keywords:** ASYMPTOTIC NORMALITY; GIBBS SAMPLER; HYPOCONVERGENCE; MARKOV CHAIN; MAXIMUM LIKELIHOOD; METROPOLIS-HASTINGS ALGORITHM; MONTE CARLO; PROFILE LIKELIHOOD

## 1. MONTE CARLO MAXIMUM LIKELIHOOD

### 1.1. *Normalized Families of Densities*

Suppose that we have a family of non-negative functions  $\{h_\theta: \theta \in \Theta\}$  on a probability space, all of which are integrable with respect to a measure  $\mu$  and none integrating to 0. Let the integrals be denoted  $c(\theta) = \int h_\theta d\mu$ . Then for each  $\theta$  in  $\Theta$  the function  $f_\theta$  defined by

$$f_\theta(x) = \frac{1}{c(\theta)} h_\theta(x)$$

is a probability density with respect to  $\mu$ . We call a family  $\{f_\theta: \theta \in \Theta\}$  of this form a *normalized family* of densities. The function  $\theta \mapsto c(\theta)$  is the *normalizer* of the family, and the functions  $h_\theta$  are the *unnormalized densities* of the family. We denote the distribution corresponding to  $\theta$  by  $P_\theta$  and expectation with respect to  $P_\theta$  by  $E_\theta$ , i.e.  $P_\theta(A) = \int_A f_\theta d\mu$  and  $E_\theta g(X) = \int g f_\theta d\mu$ .

Normalized families are interesting because they include the important special cases of exponential families and Gibbs distributions and the conditional families arising in conditional likelihood inference (Geyer and Thompson, 1992). They also have two important mathematical properties. For arbitrary functions  $h_\theta$  realizations  $X_1, X_2, \dots$  from  $P_\theta$  can be simulated without knowledge of the normalizer  $c(\theta)$  by

†Address for correspondence: School of Statistics, University of Minnesota, 270 Vincent Hall, 206 Church Street, Minneapolis, MN 55455, USA.

the Metropolis–Hastings algorithm (Metropolis *et al.*, 1953; Hastings, 1970). Moreover, maximum likelihood estimation can be carried out, again without knowledge of the normalizer or its derivatives, by using these Monte Carlo simulations (Geyer and Thompson, 1992). Somewhat surprisingly, since there is so little mathematical structure to work with, Monte Carlo maximum likelihood converges for any such family under continuity of the maps  $\theta \mapsto h_\theta(x)$ .

The log-likelihood corresponding to an observation  $x$  we take for convenience to be the log-likelihood ratio against an arbitrary fixed parameter point  $\psi$

$$l(\theta) = \log \left\{ \frac{h_\theta(x)}{h_\psi(x)} \right\} - \log \left\{ \frac{c(\theta)}{c(\psi)} \right\} = \log \left\{ \frac{h_\theta(x)}{h_\psi(x)} \right\} - \log \left\{ E_\psi \frac{h_\theta(X)}{h_\psi(X)} \right\} \quad (1)$$

since

$$E_\psi \frac{h_\theta(X)}{h_\psi(X)} = \int \frac{h_\theta(x)}{h_\psi(x)} f_\psi(x) d\mu(x) = \frac{1}{c(\psi)} \int h_\theta(x) d\mu(x) = \frac{c(\theta)}{c(\psi)}. \quad (2)$$

Although the notation suggests that  $\psi$  is a point in the parameter space of interest, this is not necessary.  $h_\psi$  can be any non-negative integrable function such that, for any  $\theta \in \Theta$ , if  $h_\psi(x) = 0$  then  $h_\theta(x) = 0$  except perhaps for  $x$  in a null set that may depend on  $\theta$ . This domination condition is necessary so that the set of points where  $h_\psi(x) = 0$  can be ignored in the integrals in equation (2). Similar domination conditions will be assumed without explicit statement throughout the paper.

Given a sample  $X_1, \dots, X_n$  from  $P_\psi$  generated by the Metropolis–Hastings algorithm, the natural Monte Carlo approximation of the log-likelihood is

$$l_n(\theta) = \log \left\{ \frac{h_\theta(x)}{h_\psi(x)} \right\} - \log \left\{ E_{n,\psi} \frac{h_\theta(X)}{h_\psi(X)} \right\} \quad (3)$$

where  $E_{n,\psi}$  denotes the ‘empirical’ expectation with respect to  $P_\psi$  defined by

$$E_{n,\psi} g(X) = \frac{1}{n} \sum_{i=1}^n g(X_i).$$

If the Markov chain  $X_1, X_2, \dots$  generated by the Metropolis–Hastings algorithm is irreducible, then  $E_{n,\psi} g(X)$  converges almost surely to  $E_\psi g(X)$  for any integrable function  $g$ . In particular,  $l_n(\theta)$  converges almost surely to  $l(\theta)$ , for any fixed  $\theta$ . The ‘almost surely’ here means for almost all sample paths of the Monte Carlo simulation; the observation  $x$  is considered fixed. Note that the null set of sample paths for which convergence fails may depend on  $\theta$ .

Let  $\hat{\theta}$  be the maximizer of  $l$  and let  $\hat{\theta}_n$  be a maximizer of  $l_n$ . Geyer and Thompson (1992) show that, if the normalized family is an exponential family, then  $\hat{\theta}_n$  converges to  $\hat{\theta}$  almost surely. They remark that an analogous result should hold outside exponential families. Section 2 gives such a theorem.

## 1.2. Missing Data

A similar but subtly different application of Monte Carlo maximum likelihood occurs with missing data (Thompson and Guo, 1991), which includes ordinary (non-Bayes) empirical Bayes methods as a special case. If  $f_\theta(x, y)$  is the joint density with  $x$  missing and  $y$  observed, then the normalizing constant for the conditional distribu-

tion of  $x$  given  $y$  is the likelihood  $f_\theta(y)$ . Again for convenience we use the likelihood ratio against the fixed parameter point  $\psi$ ; then the log-likelihood is

$$l(\theta) = \log \left\{ \frac{f_\theta(y)}{f_\psi(y)} \right\} = \log \left[ E_\psi \left\{ \frac{f_\theta(X, Y)}{f_\psi(X, Y)} \mid Y = y \right\} \right]. \quad (4)$$

The natural Monte Carlo approximation of the log-likelihood is now

$$l_n(\theta) = \log \left[ E_{n,\psi} \left\{ \frac{f_\theta(X, Y)}{f_\psi(X, Y)} \mid Y = y \right\} \right] = \log \left\{ \frac{1}{n} \sum_{i=1}^n \frac{f_\theta(X_i, y)}{f_\psi(X_i, y)} \right\} \quad (5)$$

where  $X_1, X_2, \dots$  are realizations from the conditional distribution of  $X$  given  $Y = y$ , typically simulated by using the Metropolis–Hastings algorithm when the normalizing constant  $f_\theta(y)$  is unknown.

The subtle difference between equations (3) and (5) relates not to the conditioning—in either case we need to simulate from a density known up to a constant of proportionality—but to the minus sign in equation (3). To obtain convergence results, we need to bound  $l_n$  uniformly from above on neighbourhoods, so in equation (5) the Monte Carlo average must be bounded *above*, whereas in equation (3) the average must (because of the minus sign) be bounded *below*. The former requires an integrability assumption like that imposed by Wald (1949) to obtain consistency of maximum likelihood; the latter does not.

### 1.3. Missing Data in Normalized Families

A generalization that includes both of the preceding cases has been proposed by Gelfand and Carlin (1991) for estimation in normalizing constant families with missing data. Now the unnormalized densities are  $h_\theta(x, y)$  with  $x$  missing and  $y$  observed. Then the log-likelihood, obtained by integrating over the missing data, is

$$l(\theta) = \log \left[ E_\psi \left\{ \frac{h_\theta(X, Y)}{h_\psi(X, Y)} \mid Y = y \right\} \right] - \log \left\{ E_\psi \frac{h_\theta(X, Y)}{h_\psi(X, Y)} \right\} \quad (6)$$

and its natural Monte Carlo approximation is

$$\begin{aligned} l_n(\theta) &= \log \left[ E_{n,\psi} \left\{ \frac{h_\theta(X, Y)}{h_\psi(X, Y)} \mid Y = y \right\} \right] - \log \left\{ E_{n,\psi} \frac{h_\theta(X, Y)}{h_\psi(X, Y)} \right\} \\ &= \log \left\{ \frac{1}{n} \sum_{i=1}^n \frac{h_\theta(X_i^*, y)}{h_\psi(X_i^*, y)} \right\} - \log \left\{ \frac{1}{n} \sum_{j=1}^n \frac{h_\theta(X_j, Y_j)}{h_\psi(X_j, Y_j)} \right\} \end{aligned} \quad (7)$$

where  $X_1^*, X_2^*, \dots$  are samples from the conditional distribution of  $X$  given  $Y = y$  and  $(X_1, Y_1), (X_2, Y_2), \dots$  are samples from the unconditional distribution (both for the parameter value  $\psi$ ). Gelfand and Carlin suggest maximizing equation (7) to obtain an approximation to the maximum likelihood estimate. As in the simple missing data problems of the preceding section, a Wald-type integrability condition seems to be required to assure convergence.

This double sampling is necessary only when the first term in equation (6) cannot be calculated exactly. When it can be, it is better to do so (Geyer *et al.*, 1993). Then the situation is the same as in Section 1.1. No Wald-type condition is needed for convergence.

## 2. LIKELIHOOD CONVERGENCE

### 2.1. Hypoconvergence of Monte Carlo Likelihood

Our treatment of the convergence of Monte Carlo likelihood for normalized families begins with a proof that the Monte Carlo log-likelihood (3) hypoconverges to the exact log-likelihood (1). Hypoconvergence is a type of convergence of functions that is useful in optimization theory (essentially a one-sided locally uniform convergence). The basics of the theory are given in Appendix A (or the reader may just take equations (8a) and (8b) as a definition).

**Theorem 1.** For a normalized family of densities (Section 1.1), if the parameter set  $\Theta$  is a separable metric space (e.g.  $\mathbf{R}^d$ ), if the evaluation maps  $\theta \mapsto h_\theta(x)$  are

- (a) lower semicontinuous at each  $\theta$  except for  $x$  in a  $P_\psi$  null set that may depend on  $\theta$  and
- (b) upper semicontinuous for the observed  $x$  and for  $x$  not in a  $P_\psi$  null set (that does not depend on  $\theta$ ),

and if the Metropolis–Hastings algorithm is irreducible, then the Monte Carlo log-likelihood (3) hypoconverges to the exact log-likelihood (1) with probability 1. Also the exact log-likelihood is upper semicontinuous and the normalizer of the family is lower semicontinuous.

*Proof.* What is to be shown is that  $l \leq \text{h-lim inf}_n(l_n) \leq \text{h-lim sup}_n(l_n) \leq l$  which from equations (22) in Appendix A is equivalent to

$$l(\theta) \leq \inf_{B \in \mathcal{N}(\theta)} \liminf_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\}, \quad (8a)$$

$$l(\theta) \geq \inf_{B \in \mathcal{N}(\theta)} \limsup_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\}, \quad (8b)$$

where  $\mathcal{N}(\theta)$  denotes the set of neighbourhoods of the point  $\theta$ .

By assumption there is a countable base  $\mathcal{B} = \{B_1, B_2, \dots\}$  for the topology of  $\Theta$ . For any point  $\theta$ , let  $\mathcal{N}_c(\theta) = \mathcal{B} \cap \mathcal{N}(\theta)$ . Note that the infima over the uncountable set  $\mathcal{N}(\theta)$  in inequalities (8) can be replaced by infima over the countable set  $\mathcal{N}_c(\theta)$ . Choose a countable dense subset  $\Theta_c = \{\theta_1, \theta_2, \dots\}$  as follows. For each  $n$  let  $\theta_n$  be a point of  $B_n$  satisfying

$$l(\theta_n) \geq \sup_{\phi \in B_n} \{l(\phi)\} - \frac{1}{n}.$$

We shall need

$$\lim_{n \rightarrow \infty} \left\{ E_{n,\psi} \frac{h_\theta(X)}{h_\psi(X)} \right\} = E_\psi \frac{h_\theta(X)}{h_\psi(X)} = \frac{c(\theta)}{c(\psi)} \quad (9)$$

and

$$\lim_{n \rightarrow \infty} \left[ E_{n,\psi} \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] = E_\psi \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \quad (10)$$

to hold simultaneously for all  $\theta \in \Theta_c$  and all  $B \in \mathcal{B}$ . This follows from the irreducibility assumption, since the union of a countable number of null sets (one exception

set for each limit) is still a null set. The infima in equation (10) are measurable because of assumption (b) in the theorem.

First we tackle inequality (8a). If  $B \in \mathcal{B}$  and  $\theta \in B \cap \Theta_c$

$$l(\theta) = \lim_{n \rightarrow \infty} \{l_n(\theta)\} \leq \liminf_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\}$$

by equation (9). So

$$\sup_{\phi \in B \cap \Theta_c} \{l(\phi)\} \leq \liminf_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\}$$

and

$$\inf_{B \in \mathcal{N}_c(\theta)} \sup_{\phi \in B \cap \Theta_c} \{l(\phi)\} \leq \inf_{B \in \mathcal{N}_c(\theta)} \liminf_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\}.$$

The left-hand side is equal to  $l(\theta)$  if  $l$  is upper semicontinuous by the construction of  $\Theta_c$ . Hence upper semicontinuity of  $l$  implies inequality (8a). Since  $\theta \mapsto h_\theta(x)$  is upper semicontinuous and since a sum of upper semicontinuous functions is upper semicontinuous, it remains only to be shown that  $-\log\{c(\theta)/c(\psi)\}$  is upper semicontinuous, which is true if the normalizer  $c(\theta)$  is lower semicontinuous, which follows from Fatou's lemma and the lower semicontinuity of  $\theta \mapsto h_\theta(X)$ : if  $\theta_k \rightarrow \theta$

$$c(\theta) \leq \int \liminf_{k \rightarrow \infty} \{h_{\theta_k}(x)\} d\mu(x) \leq \liminf_{k \rightarrow \infty} \left\{ \int h_{\theta_k}(x) d\mu(x) \right\} = \liminf_{k \rightarrow \infty} \{c(\theta_k)\}.$$

This establishes inequality (8a) and the assertions about upper and lower semicontinuity of the log-likelihood and the normalizer.

Now

$$\begin{aligned} \inf_{B \in \mathcal{N}_c(\theta)} \limsup_{n \rightarrow \infty} \sup_{\phi \in B} \{l_n(\phi)\} &\leq \inf_{B \in \mathcal{N}_c(\theta)} \left\{ \sup_{\phi \in B} \log \left\{ \frac{h_\phi(x)}{h_\psi(x)} \right\} \right. \\ &\quad \left. - \log \left( \liminf_{n \rightarrow \infty} \left[ E_{n,\psi} \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] \right) \right\} \\ &= \log \left\{ \frac{h_\theta(x)}{h_\psi(x)} \right\} - \log \left( \sup_{B \in \mathcal{N}_c(\theta)} \lim_{n \rightarrow \infty} \left[ E_{n,\psi} \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] \right) \\ &= \log \left\{ \frac{h_\theta(x)}{h_\psi(x)} \right\} - \log \left( \sup_{B \in \mathcal{N}_c(\theta)} \left[ E_\psi \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] \right) \end{aligned}$$

where the inequality follows from the continuity and monotonicity of the logarithm function and because of superadditivity of the supremum operation (and subadditivity of the infimum operation), and the equalities follow from the upper semicontinuity of  $\theta \mapsto h_\theta(x)$  and from equation (10). The limit will be equal to  $l(\theta)$  and establish inequality (8b) if

$$\sup_{B \in \mathcal{N}_c(\theta)} \left[ E_\psi \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] = \frac{c(\theta)}{c(\psi)}.$$

Now the integrand here satisfies

$$0 \leq \inf_{\phi \in B} \left\{ \frac{h_\phi(x)}{h_\psi(x)} \right\} \leq \frac{h_\theta(x)}{h_\psi(x)}, \quad \forall x \quad (11)$$



(since  $\theta \in B$ ). Since the right-hand side is integrable by equation (2) and the evaluation maps are assumed lower semicontinuous, dominated convergence implies

$$\sup_{B \in \mathcal{N}_c(\theta)} \left[ E_\psi \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} \right] \rightarrow E_\psi \sup_{B \in \mathcal{N}_c(\theta)} \inf_{\phi \in B} \left\{ \frac{h_\phi(X)}{h_\psi(X)} \right\} = E_\psi \frac{h_\theta(X)}{h_\psi(X)} = \frac{c(\theta)}{c(\psi)}. \quad (12)$$

This completes the proof.  $\square$

If we attempt to apply the programme of the preceding theorem to either of the missing data models (Sections 1.2 and 1.3), we find that it does not work without additional assumptions. To derive a theorem we impose a Wald-type integrability condition following Wald (1949).

**Theorem 2.** For the simple missing data problem (Section 1.2), if  $\Theta$  is a separable metric space and evaluation maps  $\theta \mapsto f_\theta(x, y)$  are

- (a) upper semicontinuous at each  $\theta$  except for  $x$  in a  $P_\psi(X|Y = y)$  null set that may depend on  $\theta$  and
- (b) lower semicontinuous except for  $x$  in a  $P_\psi(X|Y = y)$  null set (that does not depend on  $\theta$ ),

if the Metropolis–Hastings algorithm is irreducible, and if for every  $\theta \in \Theta$  there is a neighbourhood  $B$  of  $\theta$  such that

$$E_\psi \left[ \sup_{\phi \in B} \left\{ \frac{f_\phi(X, Y)}{f_\psi(X, Y)} \right\} \mid Y = y \right] < \infty \quad (13)$$

then the Monte Carlo log-likelihood (5) hypoconverges to the exact log-likelihood (4) with probability 1. Also the exact log-likelihood is continuous.

If the evaluation maps are actually continuous except for  $x$  in a  $P_\psi(X|Y = y)$  null set (that does not depend on  $\theta$ ), then the log-likelihood (5) also epiconverges to the log-likelihood (4) with probability 1.

*Proof.* The argument establishing inequality (8a) remains the same except for the invocation of Fatou's lemma. Now dominated convergence is used to prove  $l(\theta_k) \rightarrow l(\theta)$ , the dominating function being provided by inequality (13), and this gives continuity of  $l$  rather than just lower semicontinuity. The argument establishing inequality (8b) remains the same, except that infima become suprema and vice versa (because of the change in sign of the random term) and inequality (11) must be replaced by inequality (13) in justifying dominated convergence.

If the evaluation maps are almost surely continuous, then the argument in expressions (11) and (12) is still valid and proves

$$l(\theta) \leq \sup_{B \in \mathcal{N}_c(\theta)} \liminf_{n \rightarrow \infty} \inf_{\phi \in B} \{l_n(\phi)\}$$

which together with inequality (8b) implies epiconvergence.  $\square$

**Remark.** Simultaneous hypoconvergence and epiconvergence is equivalent to continuous convergence, i.e.  $\theta_n \rightarrow \theta$  implies  $l_n(\theta_n) \rightarrow l(\theta)$ . In a locally compact space (e.g.  $\mathbf{R}^d$ ) it is also equivalent to continuity of  $l$  plus convergence of  $l_n$  to  $l$  uniformly on compact sets (Rockafellar and Wets, 1993).

**Theorem 3.** For a missing data problem in a normalized family (Section 1.3) if the evaluation maps  $\theta \mapsto h_\theta(x, y)$  are

- (a) lower semicontinuous at each  $\theta$  except for  $(x, y)$  in a  $P_\psi$  null set that may depend on  $\theta$ ,
- (b) upper semicontinuous except for  $(x, y)$  in a  $P_\psi$  null set (that does not depend on  $\theta$ ) and
- (c) continuous for the observed  $y$  and for  $x$  not in a  $P_\psi(X|Y=y)$  null set (that does not depend on  $\theta$ ),

if the Metropolis–Hastings algorithm is irreducible, and if inequality (13) holds with  $f_\theta$  replaced by  $h_\theta$ , then the Monte Carlo log-likelihood (7) hypoconverges to the exact log-likelihood (6) with probability 1. Also the exact log-likelihood is upper semicontinuous.

*Proof.* This is just a combination of the two preceding proofs. The proof of theorem 1 shows that the second term in log-likelihood (7) hypoconverges, and the proof of theorem 2 shows that the first term simultaneously epiconverges and hypoconverges. The sum thus hypoconverges (see the proof of theorem 2.15 in Attouch (1984)).  $\square$

## 2.2. Convergence of Maximum Likelihood Estimate Calculation

**Theorem 4.** If

$$l_n \xrightarrow{h} l$$

with probability 1, if a sequence  $\{\hat{\theta}_n\}$  satisfies

$$l_n(\hat{\theta}_n) \geq \sup_{\theta \in \Theta} \{l_n(\theta)\} - \varepsilon_n$$

with  $\varepsilon_n \rightarrow 0$ , and if  $\{\hat{\theta}_n\}$  is contained in a compact set almost surely (or in probability), and if there is a unique maximum likelihood estimate  $\hat{\theta}$ , then  $\hat{\theta}_n \rightarrow \hat{\theta}$  and  $l_n(\hat{\theta}_n) \rightarrow l(\hat{\theta})$  almost surely (or in probability).

*Proof.* The assertion about almost sure convergence follows directly from the theorem and proposition 1 in Appendix A. If  $\{\hat{\theta}_n\}$  is contained in a compact set, then every subsequence has a convergent subsubsequence, and each such subsubsequence must converge to  $\hat{\theta}$ . Hence the whole sequence converges to  $\hat{\theta}$ . Moreover, the optimal values must converge as well.

The assertion about convergence in probability follows by almost the same argument. A sequence bounded in probability is tight; hence every subsequence has a subsubsequence which converges in distribution by Prohorov's theorem. By Skorohod representation, the convergence can be considered almost sure, in which case the only possible limit is  $\hat{\theta}$ . Hence the whole sequence and the optimal values converge in distribution to point masses at  $\hat{\theta}$  and  $l(\hat{\theta})$  (which is the same as convergence in probability).  $\square$

The theorem applies trivially when the whole parameter space  $\Theta$  is a compact set. This is the usual way in which proofs of this sort proceed, following Wald (1949), who used the one-point compactification, Kiefer and Wolfowitz (1956), who used more general compactifications, and Bahadur (1971), who gives a very general



formulation, showing that most models are compactifiable in the appropriate topology (the topology induced by vague convergence of the associated probability measures). Lacking a suitable compactification, it would be necessary to establish a uniform bound on the estimator by *ad hoc* methods.

### 2.3. Convergence of Profile Likelihoods

Suppose that  $g$  is a continuous mapping from the original parameter space  $\Theta$  to a new parameter space  $\Phi$  (both metric spaces). The *profile likelihood* is the function on  $\Phi$  defined by

$$l_p(\phi) = \sup_{\theta \in g^{-1}(\phi)} \{l(\theta)\}.$$

**Theorem 5.** If the Monte Carlo log-likelihood hypoconverges to the exact log-likelihood, and the parameter space  $\Theta$  is compact, then the Monte Carlo profile log-likelihood hypoconverges to the exact profile log-likelihood.

*Proof.* What is to be established is the analogue of inequalities (8) with  $l$  and  $l_n$  replaced by  $l_p$  and  $l_{p,n}$ . For inequality (8a) we may assume  $l_p(\phi) > -\infty$ . Then for any  $R < l_p(\phi)$  there is a  $\theta \in g^{-1}(\phi)$  such that

$$R \leq l(\theta) \leq \inf_{B \in \mathcal{N}(\theta)} \liminf_{n \rightarrow \infty} \sup_{\eta \in B} \{l_n(\eta)\} \leq \inf_{B \in \mathcal{N}(\phi)} \liminf_{n \rightarrow \infty} \sup_{\eta \in g^{-1}(B)} \{l_n(\eta)\}$$

where the second inequality is just inequality (8a) and the third inequality is true because the infimum is over a smaller set, each  $g^{-1}(B)$  being a neighbourhood of  $\theta$ . Since the right-hand side is  $\text{h-lim inf}_n(l_{p,n})$ , this establishes the analogue of inequality (8a).

For inequality (8b) we may assume that  $l_p(\phi) < +\infty$ . Hence, for every  $\varepsilon > 0$  and  $\theta \in g^{-1}(\phi)$ , there is by inequality (8b) a neighbourhood  $B_\varepsilon(\theta)$  of  $\theta$  such that

$$l(\theta) + \varepsilon \geq \limsup_{n \rightarrow \infty} \sup_{\eta \in B_\varepsilon(\theta)} \{l_n(\eta)\}.$$

By the compactness assumption there are  $\theta_1, \dots, \theta_m$  such that  $W = \bigcup_{i=1}^m B_\varepsilon(\theta_i)$  covers  $g^{-1}(\phi)$ . Also by compactness there is a neighbourhood  $B$  of  $\phi$ , such that  $g^{-1}(B) \subset W$ . Then

$$\limsup_{n \rightarrow \infty} \sup_{\eta \in g^{-1}(B)} \{l_n(\eta)\} \leq \limsup_{n \rightarrow \infty} \sup_{\eta \in W} \{l_n(\eta)\} \leq \sup_{i=1, \dots, m} \{l(\theta_i)\} + \varepsilon \leq l_p(\phi) + \varepsilon.$$

This establishes the analogue of inequality (8b). □

### 2.4. Convergence of Level Sets

Hypoconvergence also implies Painlevé–Kuratowski set convergence (Appendix A.1) for level sets of the log-likelihood  $\text{lev}_\alpha(l) = \{\theta: l(\theta) \geq \alpha\}$  which are used in forming likelihood-based interval estimates (called *support regions* in Edwards (1972)).

We may look either at a fixed level  $\alpha$  or at a fixed distance  $\gamma$  down from the maximum. The latter case makes no sense unless  $l_n(\hat{\theta}_n) \rightarrow \sup(l)$ , which need not happen, though it must under the assumptions for theorem 4.

**Theorem 6.** If

$$l_n \xrightarrow{h} l,$$

then

$$\begin{aligned}\limsup_n \text{lev}_\alpha(l_n) &\subset \text{lev}_\alpha(l), \\ \liminf_n \text{lev}_\alpha(l_n) &\supset \text{lev}_\beta(l), \quad \beta > \alpha,\end{aligned}$$

and if

$$\text{cl}\left\{\bigcup_{\beta > \alpha} \text{lev}_\beta(l)\right\} = \text{lev}_\alpha(l), \quad (14)$$

also holds, then

$$\lim_n \text{lev}_\alpha(l_n) = \text{lev}_\alpha(l). \quad (15)$$

If, in addition,  $l_n(\hat{\theta}_n) \rightarrow \sup(l)$ , then

$$\begin{aligned}\limsup_n \text{lev}_{l_n(\hat{\theta}_n) - \gamma}(l_n) &\subset \text{lev}_{\sup(l) - \gamma}(l), \\ \liminf_n \text{lev}_{l_n(\hat{\theta}_n) - \gamma}(l_n) &\supset \text{lev}_{\sup(l) - \delta}(l), \quad \delta < \gamma,\end{aligned}$$

and if equation (14) also holds for  $\alpha = \sup(l) - \gamma$ , then

$$\lim_n \text{lev}_{l_n(\hat{\theta}_n) - \gamma}(l_n) = \text{lev}_{\sup(l) - \gamma}(l). \quad (16)$$

*Proof.* The assertions about limits inferior and superior are direct consequences of theorem 3.1 in Beer *et al.* (1992), which says that  $\limsup_n \text{lev}_{\alpha_n}(l_n) \subset \text{lev}_\alpha(l)$  holds for every sequence  $\alpha_n \rightarrow \alpha$  and  $\liminf_n \text{lev}_{\alpha_n}(l_n) \supset \text{lev}_\alpha(l)$  holds for some sequence  $\alpha_n \rightarrow \alpha$ . The assertions about limits follow from the nesting of level sets and the fact that set limits are closed ( $\text{lev}_\alpha(l)$  is closed because a hypo-limit is always upper semicontinuous).  $\square$

Before leaving the subject of likelihood convergence it is worth pausing for a moment and comparing the results obtained here with the results that are obtainable for the exponential family case (Geyer, 1990; Geyer and Thompson, 1992). There the log-likelihood and its Monte Carlo approximation are concave, and this has several consequences that improve the preceding results. First, if the exact log-likelihood has a unique maximizer, the boundedness assumptions of theorem 4 can be dropped, because then a hypoconvergent sequence of concave functions is *equilevel bounded* (eventually dominated by a function with compact level sets) (Rockafellar and Wets, 1993). For the same reason the compactness assumption in theorem 5 can be dropped. Finally, equation (14) is automatically true for any level below the maximum (Rockafellar, 1970). So equations (15) and (16) hold for  $\alpha < \sup(l)$ .

### 3. ASYMPTOTIC NORMALITY

Asymptotic normality of  $n^{1/2}(\hat{\theta}_n - \hat{\theta})$  is very similar to the asymptotics of maximum likelihood.

*Theorem 7.* Suppose that the following assumptions hold.

- (a) The maximum likelihood estimate  $\hat{\theta}$  is unique and the parameter space  $\Theta$  contains an open neighbourhood of  $\hat{\theta}$  in  $\mathbf{R}^d$ .
- (b) The Monte Carlo maximum likelihood estimate  $\hat{\theta}_n$  converges in probability to  $\hat{\theta}$ .
- (c)  $c(\theta) = \int h_\theta d\mu$  can be differentiated twice under the integral sign.
- (d) 
$$n^{1/2} \nabla l_n(\hat{\theta}) \xrightarrow{\mathcal{L}} N(0, A)$$
 for some covariance matrix  $A$ .
- (e)  $B = -\nabla^2 l(\hat{\theta})$  is positive definite.
- (f)  $\nabla^3 l_n(\theta)$  is bounded in probability uniformly in a neighbourhood of  $\hat{\theta}$ .

Then

$$-\nabla^2 l_n(\hat{\theta}_n) \rightarrow B, \text{ in probability,} \quad (17)$$

and

$$n^{1/2}(\hat{\theta}_n - \hat{\theta}) \xrightarrow{\mathcal{L}} N(0, B^{-1}AB^{-1}). \quad (18)$$

A proof would be entirely classical and is omitted.

All the conditions except (d) are fairly straightforward, and one can imagine verifying them (if they hold) by standard methods. Condition (e) can be verified by using dominated convergence and ergodicity if an integrable function can be found that dominates third partial derivatives with respect to  $\theta$  of  $h_\theta/h_\psi$ .

Conclusion (17) is particularly interesting, since it gives an estimate of the observed Fisher information, which may be of interest aside from its use in expression (18). This point has also been made by Gelfand and Carlin (1991), Guo and Thompson (1992) and several discussants of Geyer and Thompson (1992).

Condition (d) is hard, if the Markov chain Monte Carlo method is being used for the simulations, because it involves a Markov chain central limit theorem. General Markov chain central limit theorems exist (Nummelin, 1984; Kipnis and Varadhan, 1986) but can be difficult to apply in practice, except when the state space is finite and the central limit theorem is automatic (Chung, 1967). The Kipnis-Varadhan theorem is the simplest for general state spaces, requiring only reversibility and summability of the autocovariances. A Metropolis-Hastings algorithm can always be arranged so that the Markov chain is reversible, a point attributed to P. Green in Besag (1986), but the summability condition is difficult. For related work in the specific context of Markov chain Monte Carlo methods see Schervish and Carlin (1992), Chan (1993), Liu *et al.* (1991), Tierney (1993) and Geyer (1992).

Assuming that (d) holds, the variance  $A$  typically cannot be calculated theoretically and must be estimated by Monte Carlo methods.

$$\nabla l_n(\theta) = \frac{\nabla h_\theta(x)}{h_\theta(x)} - \frac{E_{n,\psi} \nabla h_\theta(X)/h_\psi(X)}{E_{n,\psi} h_\theta(X)/h_\psi(X)} = \frac{E_{n,\psi} \{t_\theta(x) - t_\theta(X)\} h_\theta(X)/h_\psi(X)}{E_{n,\psi} h_\theta(X)/h_\psi(X)} \quad (19)$$

where  $t_\theta(X) = \nabla h_\theta(X)/h_\theta(X)$ . Using assumption (c) to differentiate under the integral sign

$$\nabla l(\theta) = \frac{\nabla h_\theta(x)}{h_\theta(x)} - \frac{\nabla c(\theta)}{c(\theta)}$$

$$\begin{aligned}
 &= \frac{\nabla h_\theta(x)}{h_\theta(x)} - \int \frac{\nabla h_\theta(x)}{h_\theta(x)} \frac{h_\theta(x)}{c(\theta)} d\mu(x) \\
 &= t_\theta(x) - E_\theta t_\theta(X),
 \end{aligned}$$

and this is 0 when  $\theta = \hat{\theta}$ . The denominator in equation (19) converges to  $c(\theta)/c(\psi)$ ; the expectation of the numerator with respect to  $P_\psi$  is

$$\begin{aligned}
 E_\psi \{t_\theta(x) - t_\theta(X)\} \frac{h_\theta(X)}{h_\psi(X)} &= \frac{c(\theta)}{c(\psi)} \int \{t_\theta(x) - t_\theta(y)\} f_\theta(y) d\mu(y) \\
 &= \frac{c(\theta)}{c(\psi)} \{t_\theta(x) - E_\theta t_\theta(X)\},
 \end{aligned}$$

which is also 0 when  $\theta = \hat{\theta}$ . Thus the numerator is the sample mean for a functional of the Markov chain

$$z_\theta(X) = \{t_\theta(x) - t_\theta(X)\} \frac{h_\theta(X)}{h_\psi(X)}$$

which has expectation 0 under the stationary distribution. Hence by the continuous mapping theorem

$$\frac{1}{\sqrt{n}} \sum_{i=1}^n z_\theta(X_i) \xrightarrow{\mathcal{L}} \frac{c(\theta)}{c(\psi)} N(0, A).$$

Let  $\gamma(t) = \gamma(-t)$  be the lag  $t$  autocovariance of  $z_\theta(X_i)$  at stationarity, i.e.

$$\gamma(t) = \text{cov}\{z_\theta(X_0), z_\theta(X_t)\}$$

when the starting position  $X_0$  of the Markov chain is a realization from  $P_\psi$ ; then for reversible chains (Kipnis and Varadhan, 1986)

$$A = \frac{c(\psi)^2}{c(\theta)^2} \sum_{t=-\infty}^{+\infty} \gamma(t). \quad (20)$$

Both factors in equation (20) can be estimated:  $c(\theta)/c(\psi)$  by the denominator in equation (19) and the sum by standard time series methods (for a review see Geyer (1992); see also Hastings (1970), Geweke (1992), Han (1991) and Green and Han (1992)).

#### 4. DISCUSSION

'Normalized families of densities' are an important class of statistical models. We now have two interesting properties that hold for the whole class. The Metropolis-Hastings algorithm can be used to simulate realizations from any distribution in the model, and Monte Carlo likelihood approximation can be used to do likelihood-based statistical inference. When there are no missing data, mere continuity is enough to guarantee convergence. With missing data, Wald-type integrability conditions are required. This class is extremely flexible, allowing a very wide scope for modelling and supporting the notion of a 'model liberation movement' called for by Professor A. F. M. Smith in his discussion of Geyer and Thompson (1992).

Monte Carlo likelihood may be useful even in missing data problems where the EM algorithm can be used to calculate the maximum likelihood estimate, since the Monte Carlo algorithm approximates the whole likelihood surface. The use of expression (17) to approximate the observed Fisher information may be useful in problems where analytical methods (Sundberg, 1974; Louis, 1982) are intractable. It is especially useful in conjunction with Monte Carlo EM (Tanner and Wei, 1990; Guo and Thompson, 1992) but may also be a competitor for the SEM algorithm (Meng and Rubin, 1991).

## ACKNOWLEDGEMENTS

Conversations with Elizabeth Thompson, Julian Besag and Michael Newton helped to change my focus from exponential families to the general normalized families of Section 1. The whole approach to convergence of optimization problems used in this paper comes from a course taught by Terry Rockafellar in 1990 at the University of Washington using a draft of Rockafellar and Wets (1993). Roger Wets provided the reference to Beer *et al.* (1992) and suggested the approach used in proving theorem 5. Xiaotong Shen found a mistake in my first proof of theorem 1.

## APPENDIX A

### A.1. Set Convergence

At several points the concept of Painlevé–Kuratowski set convergence (section 1.4.1 in Attouch (1984)) was needed. Given a sequence of sets  $C_n$ , the set limit superior is the set

$$\limsup_{n \rightarrow \infty} (C_n) = \bigcap_{m=1}^{\infty} \text{cl} \left( \bigcup_{n=m}^{\infty} C_n \right)$$

and the limit inferior is the set

$$\liminf_{n \rightarrow \infty} (C_n) = \text{cl} \left\{ \bigcup_{m=1}^{\infty} \bigcap_{n=m}^{\infty} \text{cl} (C_n) \right\}.$$

Note that these are topological convergence notions, different from the set theoretic notions commonly used in probability theory (defined by the same formulae without the closure operations). In a metric space the following definitions are equivalent to the preceding ones (proposition 1.34 in Attouch (1984)). The set limit superior is the set of points  $x$  such that there is a subsequence  $x_{n_k} \rightarrow x$  with  $x_{n_k} \in C_{n_k}$ , and the set limit inferior is the set of points  $x$  such that there is a sequence  $x_n \rightarrow x$  with  $x_n \in C_n$  for all  $n$  after some  $n_0$ . In short, the limit superior is the set of cluster points and the limit inferior is the set of limit points. If the set limits superior and inferior agree, then their common value is said to be the limit of the sequence.

### A.2. Epiconvergence and Hypoconvergence

Epiconvergence and hypoconvergence are types of convergence of sequences of functions that are useful in optimization problems. If a sequence of functions  $g_n$  epiconverges to a limit  $g$  (written  $g_n \xrightarrow{e} g$ ) and  $x_n$  minimizes  $g_n$  then any cluster point of the sequence  $\{x_n\}$  is a minimizer of  $g$ . Hypoconvergence is the analogous notion for maximization problems. Since  $x$  maximizes  $g$  if and only if it minimizes  $-g$ , hypoconvergence (written  $g_n \xrightarrow{h} g$ ) is defined by

$$g_n \xrightarrow{h} g$$

if and only if

$$(-g_n) \xrightarrow{e} (-g).$$

Epiconvergence is related to set convergence in the following way. The *epigraph* of an extended real-valued ( $\pm\infty$  allowed) function  $g$  with domain  $S$  is the set

$$\text{epi}(g) = \{(x, \lambda) \in S \times \mathbf{R} : g(x) \leq \lambda\}$$

of points lying on or above the graph. A sequence of functions  $g_n$  *epiconverges* to a function  $g$  if and only if the sequence of sets  $\text{epi}(g_n)$  converges to the set  $\text{epi}(g)$ .

There are several equivalent characterizations that are sometimes more useful. Given a sequence of functions  $g_n$ , the epi-limits inferior and superior are the functions (Attouch (1984), p. 26)

$$(\text{e-liminf}_n(g_n))(x) = \sup_{B \in \mathcal{N}(x)} \liminf_{n \rightarrow \infty} \inf_{y \in B} \{g_n(y)\}, \quad (21a)$$

$$(\text{e-limsup}_n(g_n))(x) = \sup_{B \in \mathcal{N}(x)} \limsup_{n \rightarrow \infty} \inf_{y \in B} \{g_n(y)\} \quad (21b)$$

where  $\mathcal{N}(x)$  denotes the set of neighbourhoods of the point  $x$ . The sequence  $g_n$  epiconverges to a function  $\text{e-lim}_n(g_n) = g$  if and only if the epi-limits inferior and superior agree and are equal to  $g$ . Similar notation with the prefix *e*- replaced by *h*- is used for hypoconvergence:

$$(\text{h-liminf}_n(g_n))(x) = \inf_{B \in \mathcal{I}(x)} \liminf_{n \rightarrow \infty} \sup_{y \in B} \{g_n(y)\}; \quad (22a)$$

$$(\text{h-limsup}_n(g_n))(x) = \inf_{B \in \mathcal{I}(x)} \limsup_{n \rightarrow \infty} \sup_{y \in B} \{g_n(y)\}. \quad (22b)$$

Another pair of conditions that are equivalent for functions on metric spaces are the following (Attouch (1984), p. 30). A sequence of functions  $g_n$  epiconverges to a function  $g$  if the following two conditions hold at every point  $x$ :

- (a)  $\liminf_n \{g_n(x_n)\} \geq g(x)$  for every sequence  $x_n \rightarrow x$ ;
- (b)  $\limsup_n \{g_n(x_n)\} \leq g(x)$  for some sequence  $x_n \rightarrow x$ .

This says that epiconvergence is a combination of one-sided locally uniform convergence (condition (a)), with something weaker than pointwise convergence from the other side (condition (b)).

The main reason for the importance of epiconvergence is the following proposition, which is theorem 1.10 in Attouch (1984).

**Proposition 1.** Suppose that  $g_n \xrightarrow{e} g$ ,  $x_n \rightarrow x$  and  $g_n(x_n) - \inf(g_n) \rightarrow 0$ ; then

$$g(x) = \inf(g) = \lim_{n \rightarrow \infty} \{g_n(x_n)\},$$

i.e., if  $x_n$  is an  $\varepsilon_n$ -minimizer of  $g_n$  with  $\varepsilon_n \rightarrow 0$ , then any convergent subsequence of  $\{x_n\}$  must converge to a point  $x$  which minimizes  $g$  and the optimal values  $g_n(x_n)$  must also converge to the asymptotic optimal value  $g(x)$ . Two points are worth comment here. First, there is no requirement that the minimizers be unique. If  $g$  has a unique minimizer  $x$ , then  $x$  is the only cluster point of the sequence  $\{x_n\}$ . Otherwise, there may be many cluster points, but all must minimize  $g$ . Second, the proposition does not rule out escape to infinity; it only describes what happens if  $x_n \rightarrow x$ . It does say that, if the sequence  $\{x_n\}$  is confined to a compact set and if  $g$  has a unique minimizer, then  $x_n$  converges to that minimizer.

## REFERENCES

- Attouch, H. (1984) *Variational Convergence of Functions and Operators*. Boston: Pitman.  
 Bahadur, R. R. (1971) *Some Limit Theorems in Statistics*. Philadelphia: Society for Industrial and Applied Mathematics.



- Beer, G., Rockafellar, R. T. and Wets, R. J.-B. (1992) A characterization of epiconvergence in terms of convergence of level sets. *Proc. Am. Math. Soc.*, **116**, 753–761.
- Besag, J. (1986) On the statistical analysis of dirty pictures (with discussion). *J. R. Statist. Soc. B*, **48**, 259–302.
- Chan, K. S. (1993) Asymptotic behavior of the Gibbs sampler. *J. Am. Statist. Ass.*, **88**, 320–326.
- Chung, K. L. (1967) *Markov Chains with Stationary Transition Probabilities*, 2nd edn, p. 99 ff. Berlin: Springer.
- Edwards, A. W. F. (1972) *Likelihood*. Cambridge: Cambridge University Press.
- Gelfand, A. E. and Carlin, B. P. (1991) Maximum likelihood estimation for constrained or missing data models. *Can. J. Statist.*, to be published.
- Geweke, J. (1992) Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments. In *Bayesian Statistics 4* (eds J. M. Bernardo, J. O. Berger, A. P. Dawid and A. F. M. Smith), pp. 169–193. Oxford: Oxford University Press.
- Geyer, C. J. (1990) Likelihood and exponential families. *PhD Dissertation*. University of Washington, Seattle.
- (1992) Practical Markov chain Monte Carlo (with discussion). *Statist. Sci.*, **7**, 473–511.
- Geyer, C. J., Ryder, O. A., Chemnick, L. G. and Thompson, E. A. (1993) Analysis of relatedness in the California condors from DNA fingerprints. *Mol. Biol. Evoln*, to be published.
- Geyer, C. J. and Thompson, E. A. (1992) Constrained Monte Carlo maximum likelihood for dependent data (with discussion). *J. R. Statist. Soc. B*, **54**, 657–699.
- Green, P. J. and Han, X.-L. (1992) Metropolis methods, gaussian proposals, and antithetic variables. *Lect. Notes Statist.*, **74**, 142–164.
- Guo, S. W. and Thompson, E. A. (1992) Monte Carlo estimation of mixed models for large complex pedigrees. *Biometrics*, to be published.
- Han, X.-L. (1991) Spectral window estimation of integrated autocorrelation time. *Research Report*. University of Bristol, Bristol.
- Hastings, W. K. (1970) Monte Carlo sampling methods using Markov chains and their applications. *Biometrika*, **57**, 97–109.
- Kiefer, J. and Wolfowitz, J. (1956) Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. *Ann. Math. Statist.*, **27**, 887–906.
- Kipnis, C. and Varadhan, S. R. S. (1986) Central limit theorem for additive functionals of reversible Markov processes and applications to simple exclusions. *Commun. Math. Phys.*, **104**, 1–19.
- Liu, J., Wong, W. H. and Kong, A. (1991) Correlation structure and convergence rate of the Gibbs sampler with various scans. *Technical Report 304*. Department of Statistics, University of Chicago, Chicago.
- Louis, T. A. (1982) Finding the observed information matrix when using the EM algorithm. *J. R. Statist. Soc. B*, **44**, 226–233.
- Meng, X.-L. and Rubin, D. B. (1991) Using EM to obtain asymptotic variance-covariance matrices: the SEM algorithm. *J. Am. Statist. Ass.*, **86**, 899–909.
- Metropolis, N., Rosenbluth, A. W., Rosenbluth, M. N., Teller, A. H. and Teller, E. (1953) Equation of state calculations by fast computing machines. *J. Chem. Phys.*, **21**, 1087–1092.
- Nummelin, E. (1984) *General Irreducible Markov Chains and Non-negative Operators*. Cambridge: Cambridge University Press.
- Rockafellar, R. T. (1970) *Convex Analysis*, theorem 7.6. Princeton: Princeton University Press.
- Rockafellar, R. T. and Wets, R. J. B. (1993) *Variational Analysis*. New York: Springer. To be published.
- Schervish, M. J. and Carlin, B. P. (1992) On the convergence rate of successive substitution sampling. *J. Comput. Graph. Statist.*, **1**, 111–127.
- Sundberg, R. (1974) Maximum likelihood theory for incomplete data from an exponential family. *Scand. J. Statist.*, **1**, 49–58.
- Tanner, M. A. and Wei, G. C. G. (1990) A Monte Carlo implementation of the EM algorithm and the poor man's data augmentation algorithms. *J. Am. Statist. Ass.*, **85**, 699–704.
- Thompson, E. A. and Guo, S. W. (1991) Evaluation of likelihood ratios for complex genetic models. *IMA J. Math. Appl. Med. Biol.*, **8**, 149–169.
- Tierney, L. (1993) Markov chains for exploring posterior distributions. *Ann. Statist.*, to be published.
- Wald, A. (1949) Note on the consistency of the maximum likelihood estimate. *Ann. Math. Statist.*, **20**, 595–601.