

Long range search for maximum likelihood in exponential families

Saisuke Okabayashi and Charles J. Geyer

School of Statistics

University of Minnesota, Minneapolis

e-mail: sai@stat.umn.edu; charlie@stat.umn.edu

Abstract: Exponential families are often used to model data sets with complex dependence. Maximum likelihood estimators (MLE) can be difficult to estimate when the likelihood is expensive to compute. Markov chain Monte Carlo (MCMC) methods based on the MCMC-MLE algorithm in [17] are guaranteed to converge in theory under certain conditions when starting from any value, but in practice such an algorithm may labor to converge when given a poor starting value. We present a simple line search algorithm to find the MLE of a regular exponential family when the MLE exists and is unique. The algorithm can be started from any initial value and avoids the trial and error experimentation associated with calibrating algorithms like stochastic approximation. Unlike many optimization algorithms, this approach utilizes first derivative information only, evaluating neither the likelihood function itself nor derivatives of higher order than first. We show convergence of the algorithm for the case where the gradient can be calculated exactly. When it cannot, it has a particularly convenient form that is easily estimable with MCMC, making the algorithm still useful to a practitioner.

Keywords and phrases: Markov chain Monte Carlo, exponential families, Potts, Ising, exponential random graph, stochastic approximation.

Received June 2010.

Contents

1	Introduction	124
1.1	Parameter estimation methods in exponential families	125
1.2	Algorithm overview	126
2	Background exponential family theory	129
3	Long range search algorithm for MLE	130
4	Refinements of algorithm	131
4.1	Search directions	131
4.2	Step size	132
4.3	MCMC approximations	132
4.4	Combining with other algorithms	134
5	Examples	134
5.1	Example: Logistic regression	135
5.2	Example: Ising model	136
6	Discussion	139

A Proofs	140
References	144

1. Introduction

Exponential families are commonly used to model phenomena with dependence structure, where the outcomes of the response variable of interest are in fact dependent on one another. For example, the Ising model [24, 37] is an exponential family model that has been used to model ferromagnetism and other spatial lattice processes [10]. A realized sample from this model is depicted in Figure 1, where neighboring pixels are more likely to have the same color. We explore this model further in Section 5.2. Other examples of phenomena with dependence structure modeled with exponential families include plant ecology [3, 4], friendship networks [18, 19, 51], protein-protein interaction networks [43], and the lifetime fitness of plants [44].



FIG 1. A realized sample from an Ising model on a 32×32 lattice with parameter $\eta = (0, \log(1 + \sqrt{2}))^T$. This value of η corresponds to the phase transition point, where the sample images are mostly one color with small but significant portions of the other color. There is no preference for the dominant color to be white or black.

The appeal of exponential families in these settings stems from their simplicity and maximum entropy property [17, 25]. By choosing statistics of interest on the data, one fully specifies a model that gives the most reasonable inference possible derived solely from those statistics. Furthermore, exponential families have been well-studied [2, 5] and utilized over the decades and have desirable properties such as a strictly concave likelihood function.

1.1. Parameter estimation methods in exponential families

Calculating the maximum likelihood estimators (MLE) for exponential families when dependence is complex, however, remains a challenging problem because the likelihood function may be computationally infeasible. In particular, the form of the likelihood is determined by the chosen statistics up to a normalizing constant, but this normalizing constant may involve a summation over an astronomical number of terms. Evaluating the likelihood function—let alone maximizing it—presents a significant challenge.

Three commonly used parameter estimation methods to circumvent this issue in exponential families are the *pseudo-likelihood* approach [4, 35, 46], which finds parameter values that maximize the pseudo-likelihood function, the *Markov chain Monte Carlo maximum likelihood estimate* (MCMC-MLE) approach [11, 17], which uses MCMC to approximate the log likelihood so that it can subsequently be maximized, and *stochastic approximation* (SA) [7, 27], which utilizes simple iterated updates of parameter estimates. The pseudo-likelihood approach is computationally expedient, but has been shown to produce unreliable results when dependence is strong [17, 49].

The MCMC-MLE approach is theoretically guaranteed to converge to the MLE if it exists and is the default algorithm in software packages such as `statnet` [21] in the R platform for network models. However, this approach has been shown in practice to be sensitive to initial parameter values when used without the trust region methodology recommended in [17], and the algorithm may require many iterations and enormous (sometimes infeasibly large) Monte Carlo sample sizes when the starting value is far from the MLE [23]. Improvement to the MCMC-MLE approach is an active area of research [22].

Variations on the Robbins-Monro stochastic approximation algorithm [39] have been applied to find the MLE similar contexts: [20, 31, 52, 53] applied MCMC stochastic approximation to spatial models and [45] to social network models (exponential random graph models). SA procedures for finding the MLE of a parameter η generate iterated estimates η_k to find the root of a gradient function $h(\eta)$:

$$\eta_{k+1} = \eta_k + \alpha_k U_k, \quad (1)$$

where α_k is a step size and is typically a member of a decreasing sequence of positive numbers, and U_k is a random variable from the distribution specified by η_k that noisily estimates the gradient function $h(\eta_k)$.

Restrictive conditions are required of α_k and U_k to establish convergence of the sequence η_k . In Robbins-Monro SA [39], the step size α_k must be a sequence

of positive constants that satisfies

$$\sum \alpha_k^2 < \infty$$

for which the choice of

$$\alpha_k = \frac{A}{B+k} \tag{2}$$

is commonly used, where A and B are constants that must be specified by the user. This specification requires experimentation and care: there can be significant variation in performance depending on choice of these constants. More recent research [27, Chapter 11] show that α_k sequences that go to 0 slower than $1/k$ can result in an improved rate of convergence, where rate of convergence is measured by the asymptotic covariance of the normalized estimates about their limit point.

The conditions on U_k are more restrictive. Popular approaches include constraining the sequence of estimators η_k to a compact set specified *a priori*, or assuming that the noise component of U_k be a martingale difference sequence. As commonly observed [1, 7, 30], these may be difficult to satisfy in practice. See [1, 30] for recent developments that impose less restrictive conditions using truncated updates.

An issue for any recursive search algorithm is the choice of starting point. It is often the case that algorithms are good at finding the MLE when the starting point is close to it, but of course the location of the MLE is unknown. Methods which rely on the Fisher information matrix may fail when the starting point for η is far from the MLE [20, 53]; for any exponential family with bounded support, Fisher information becomes singular as the natural parameter η goes to ∞ [38]. Of course, one may try different starting points until a “good” one is found, but this can be cumbersome in practice and demands patience and sophistication of the practitioner.

1.2. Algorithm overview

In this article, we propose a simple and practical line search algorithm that converges to the MLE of any regular exponential family when the MLE exists and the first derivative of the log likelihood can be calculated exactly. When it cannot, the first derivative has a particularly convenient form that is easily estimable with MCMC, making the algorithm still useful in application. The first derivative with respect to the canonical parameter vector η has the form $g(y) - E_\eta g(Y)$, where $g(Y)$ is the canonical statistic vector. Its Monte Carlo approximation is $g(y) - \frac{1}{m} \sum_{i=1}^m g(Y_i)$, where Y_1, \dots, Y_m are simulated data sets having parameter vector η . The log likelihood itself is much harder to compute [17]. The second derivative with respect to the canonical parameter vector has the form $-\text{Var}_\eta g(Y)$, and Monte Carlo estimate minus the empirical variance of the $g(Y_i)$. The second derivative is less stably estimated than the first derivative, especially when η is far from the MLE so this matrix is nearly singular. We

also show how to construct and apply confidence intervals in such a setting to increase the probability of convergence.

The appeal of this algorithm is its ease of use: no trial and error is needed. The computer can find the solution with no help from the user, thus making it suitable for use by naive users. Experimentation with multiple starting points or tuning parameters is not necessary and no unrealistic *a priori* information about the problem need be specified. It is currently used in the `aster2` contributed R package [14] as the safeguard for steepest ascent and Newton-Raphson iterations in finding the MLE for aster models.

Our algorithm generates iterated estimates η_k of the MLE $\hat{\eta}$ with the update

$$\eta_{k+1} = \eta_k + \alpha_k p_k \quad (3)$$

where α_k is a *step size* and p_k is a *search direction* and is restricted to be an ascent direction of the log likelihood. Despite the visual similarity between (1) and (3), our line search approach treats the search direction p_k in (3) as constant in the inner loop of our algorithm whereas in SA the corresponding U_k in (1) is random. Furthermore, line search algorithms have more restrictions on the step size α_k . The step size conditions in the classical gradient ascent algorithm, which is the basis for our algorithm, force a sufficiently large increase in the objective function at every step, guaranteeing convergence to a local maximum (which is the global maximum in an exponential family because of strict concavity of the log likelihood), if it exists.

Theorem 3.2 in [32] implies the global convergence of the steepest ascent algorithm for a continuously differentiable function, $\ell(\eta)$. It requires the step length α_k to satisfy the Wolfe conditions for *sufficient increase* and *curvature*:

$$\begin{aligned} \ell(\eta_k + \alpha_k p_k) &\geq \ell(\eta_k) + c_1 \alpha_k \nabla \ell(\eta_k)^T p_k \\ \nabla \ell(\eta_k + \alpha_k p_k)^T p_k &\leq c_2 \nabla \ell(\eta_k)^T p_k \end{aligned} \quad (4)$$

where ∇ is the gradient operator and $0 < c_1 < c_2 < 1$. Variations of these conditions exist in the numerical optimization literature [32, 47], but all require evaluating the objective function.

Exponential families we consider are an unusual case in optimization in that the objective function is harder to compute than its derivatives and hence not previously considered by optimization theorists. In our algorithm, we replace (4) with a single modified curvature condition:

$$0 \leq \nabla \ell(\eta_k + \alpha_k p_k)^T p_k \leq c \nabla \ell(\eta_k)^T p_k \quad (5)$$

for some $0 < c < 1$. This replacement is possible while still guaranteeing sufficient increase and convergence to the MLE (if it exists) because we have the additional property that the exponential family log likelihood function we consider is strictly concave. The restrictions on the step size α_k along a particular direction p_k and the resulting values for $\ell(\eta_{k+1})$ are depicted in Figure 2.

The desire to avoid calculation of higher order derivatives is motivated not just by computational considerations, but also by how much useful information can be extracted from them. As noted in Section 1.1, if η is far from

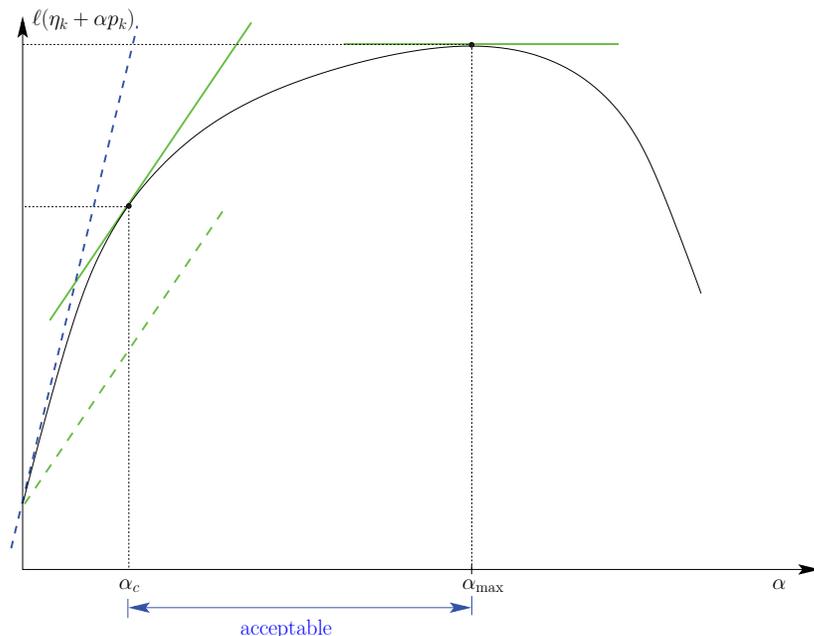


FIG 2. The acceptable region for step size α_k along a particular search direction p_k according to the modified curvature condition (5). The step sizes α_c and α_{\max} correspond to values of $\nabla \ell(\eta_k + \alpha p_k)^T p_k$ equaling $\epsilon \nabla \ell(\eta_k)^T p_k$ and 0, respectively. The condition ensures sufficient increase in the log likelihood along the search direction p_k .

the MLE, the Fisher information matrix may be near-singular and algorithms like (unsafeguarded) Newton-Raphson algorithm may fail. For this reason, the best use of our algorithm may be from “long range,” filling a gap in the MLE estimation toolbox. It may be expedient to switch to another algorithm like Newton-Raphson after significant progress is made and the Fisher information matrix becomes useful. Our line search algorithm with p_k the Newton direction provides a safeguard for Newton-Raphson that makes it safe for use from any range. The `aster2` contributed R package [14] switches p_k from the steepest ascent direction to the Newton direction after a fixed number of steps ($d/2$ where d is the dimension η) but always finds a step length α_k satisfying (5), iterating until the unsafeguarded Newton step satisfies (5).

Our algorithm can be outlined as follows. Let $\|\cdot\|$ denote the Euclidean norm function, and ϵ a small value greater than 0.

Get an initial value, η_0 .

Set $k = 0$.

Calculate $\nabla \ell(\eta_k)$, the direction of steepest ascent.

Set $p_k = \nabla \ell(\eta_k)$.

while $\|\nabla\ell(\eta_k)\| > \epsilon$

Find a step size α_k that satisfies the modified curvature condition

$$0 \leq \nabla\ell(\eta_k + \alpha_k p_k)^T p_k \leq c \nabla\ell(\eta_k)^T p_k$$

for some $0 < c < 1$.

$\eta_{k+1} = \eta_k + \alpha_k p_k$.

Calculate $\nabla\ell(\eta_{k+1})$.

Find the new search direction p_{k+1} , which must be an ascent direction.

$k = k + 1$.

end while

2. Background exponential family theory

An exponential family of distributions [2, 12] on a sample space \mathcal{Y} has log likelihood

$$\ell(\eta) = \langle g(y), \eta \rangle - c(\eta) \quad (6)$$

where $g(y)$ is a d -dimensional vector of *natural statistics* calculated from the observed data y , η a d -dimensional vector of *natural parameters*, and $\langle \cdot, \cdot \rangle$ denotes the bilinear form

$$\langle g, \eta \rangle = \sum_{i=1}^d g_i \eta_i.$$

So that the probability function integrates to 1, the *cumulant function* c must have the form

$$c(\eta) = \log \left(\int e^{\langle g(x), \eta \rangle} d\mu(x) \right), \quad (7)$$

where μ is a measure on \mathcal{Y} . Define the *natural parameter space* Ξ as the set of points $\eta = (\eta_1, \dots, \eta_d)$ that are parameter values indexing distributions in the model. An exponential family is *full* if the natural parameter space is

$$\Xi = \{\eta \in \mathbb{R}^d : c(\eta) < \infty\}, \quad (8)$$

and *regular* if, in addition, Ξ is an open set. We say an exponential family is *minimal* if $g(y)$ is not concentrated on a hyperplane. Minimality guarantees that if an MLE exists, it is unique [12].

In finite state space models with complicated dependence like an Ising model or exponential random graph model, (7) is a sum which may have no simple expression and can only be evaluated by explicitly doing the sum. When the sample space \mathcal{Y} is even moderately large, this can be prohibitively expensive. For example, the sample space \mathcal{Y} for an Ising model defined on a 32×32 square lattice where each entry takes values of 0 or 1 has $2^{1024} \approx 10^{300}$ elements. A loop with this many iterations takes too long no matter how programmed.

A useful property of all exponential families [28, p. 27] when η is in the interior of Ξ is that

$$\begin{aligned} \mathbb{E}_\eta(g(Y)) &= \nabla c(\eta) \\ \text{Var}_\eta(g(Y)) &= \nabla^2 c(\eta). \end{aligned}$$

Thus we can express first and second derivatives of the log likelihood (6) and Fisher information, $I(\eta)$, as

$$\nabla \ell(\eta) = g(y) - \mathbb{E}_\eta g(Y) \tag{9}$$

$$\nabla^2 \ell(\eta) = -\text{Var}_\eta g(Y) \tag{10}$$

$$I(\eta) = -\mathbb{E}_\eta \nabla^2 \ell(\eta) = \text{Var}_\eta g(Y) \tag{11}$$

and thereby avoid evaluation of the problematic cumulant function c so long as we make do with first and second derivatives of the log likelihood avoiding evaluation of the log likelihood itself.

3. Long range search algorithm for MLE

We now present our line search algorithm, which will converge to the MLE for any regular exponential family if the MLE exists. The theory is divided into two theorems, Theorem 3.1 and 3.2: the first presents the requirements for the algorithm and guarantees that the log likelihood gradient, when it can be calculated exactly, converges to zero. The second shows that when the MLE exists, this is equivalent to finding the MLE. Proofs are in Appendix A.

Theorem 3.1 also can be interpreted as assuring convergence to the MLE in the Barndorff-Nielsen completion [12] even when the MLE does not exist in the conventional sense. Convergence of the gradient of the log likelihood to zero is the same as convergence of the mean value parameter $\mu = \mathbb{E}_\eta g(Y)$ to $g(y)$, which is the MLE of mean value parameter in the Barndorff-Nielsen completion. This is not an efficient method of finding the MLE in the Barndorff-Nielsen completion [33, Chapters 4 and 5], but it is interesting that our algorithm behaves well even when the MLE does not exist in the conventional sense.

Theorem 3.1 (Exponential family zero gradient attainment). *Consider any line search of the form*

$$\eta_{k+1} = \eta_k + \alpha_k p_k \tag{12}$$

used to maximize the log likelihood function $\ell(\cdot)$ of a regular exponential family on a finite sample space, where the search direction p_k is a non-zero ascent direction such that the angle θ_k between the search direction p_k and steepest ascent direction $\nabla \ell(\eta_k)$ is restricted to be less than 90 degrees by

$$\cos \theta_k \geq \delta > 0 \tag{13}$$

for some fixed $\delta > 0$.

Then, unless $\nabla\ell(\eta_k) = 0$, in which case η_k is already the solution and the search is complete, it is possible to find a step length α_k that satisfies the curvature condition

$$0 \leq \nabla\ell(\eta_k + \alpha_k p_k)^T p_k \leq c \nabla\ell(\eta_k)^T p_k \quad (14)$$

for some fixed $0 < c < 1$.

Furthermore, repeated iterations of (12) satisfying (13) and (14) will produce a sequence, η_1, η_2, \dots such that

$$\lim_{k \rightarrow \infty} \|\nabla\ell(\eta_k)\| = 0.$$

Theorem 3.1 can be adapted to a more general setting to optimize any bounded, proper, upper semi-continuous, and strictly concave function assuming there are bounded level sets of this function, as detailed in [34]. However, by assuming here that the objective function is the log likelihood of an exponential family, the statement of the theorem is much simplified.

We apply Theorem 3.1 to find the MLE when it is known to exist:

Theorem 3.2. *For a regular exponential family with minimal representation where the MLE exists, the line search described in Theorem 3.1 can be applied to the log likelihood function $\ell(\eta)$ so that a search starting at any $\eta_0 \in \Xi$ will converge to the MLE of η .*

The issue of MLE existence is a problem in computational geometry, not an optimization problem, so we do not address it here. See [12, 33, 38] for further discussion of this issue.

4. Refinements of algorithm

In Theorem 3.1, we restricted our search direction p_k to be an ascent direction, so that $\nabla\ell(\eta_k)^T p_k > 0$ or, alternatively, the angle θ_k between the search direction p_k and steepest ascent direction $\nabla\ell(\eta_k)$ is less than 90 degrees. However, this still leaves many possibilities for the choice of p_k other than steepest ascent. In addition, we have specified restrictions on the step size α_k in the curvature condition (14) with $0 < c < 1$, but it would be useful to know if certain values of c are better than others.

4.1. Search directions

In our examples in Section 5, by default we use steepest ascent directions in our implementation for simplicity. Although often effective in early steps, steepest ascent directions can result in a zigzagging trajectory of the sequence η_k [47, Section 3.1]. Conjugate gradient methods address this phenomena and cover the sample space more efficiently [32, Chapter 5]. It is easy to implement a variant of the Polak-Ribière method [32, pp. 120–122] here, requiring little more in terms

of calculation or storage. The search direction p_k would update with an extra intermediate step as follows:

$$\begin{aligned}\gamma_{k+1}^{PR} &= \max\left(0, \frac{[\nabla\ell(\eta_{k+1})]^T(\nabla\ell(\eta_{k+1}) - \nabla\ell(\eta_k))}{\|\nabla\ell(\eta_k)\|^2}\right) \\ p_{k+1} &= \nabla\ell(\eta_{k+1}) + \gamma_{k+1}^{PR} p_k.\end{aligned}$$

Note that when $\gamma_{k+1}^{PR} = 0$, p_{k+1} will be just $\nabla\ell(\eta_{k+1})$, the direction of steepest ascent, and thus serves as a “reset”. The curvature condition (14) guarantees that this method always yields a ascent direction for p_{k+1} and thus Theorem 3.1 still holds.

4.2. Step size

We now turn our attention to the optimal step size α_k . Taking the derivative of $\ell(\eta_k + \alpha_k p_k)$ with respect to α_k shows that the log likelihood is maximized as a function of α_k along the direction p_k when

$$\nabla\ell(\eta_{k+1})^T p_k = 0.$$

By choosing c to be small, say 0.2, we ensure that the step taken is close to maximizing the log likelihood along the search direction. This is also apparent in Figure 2.

Making c too small, however, may make it difficult to find an α_k that meets the curvature condition (14) since this search must be done numerically. In fact, as the line search nears the MLE and $\nabla\ell(\eta_k)$ gets smaller, the rightmost term in (14) gets smaller in magnitude (it equals $c\|\nabla\ell(\eta_k)\|^2$ if using steepest ascent directions), making a numerical search for α_k more challenging.

4.3. MCMC approximations

Our algorithm requires us to be able to calculate $\nabla\ell(\eta)$ using (9). When this can be done exactly, our algorithm is straightforward to apply, as illustrated in the logistic regression example in Section 5.1. However, for many applications, we will need to approximate $E_\eta g(Y)$ using MCMC. That is,

$$\nabla\ell(\eta) = g(y) - E_\eta g(Y) \approx g(y) - \frac{1}{m} \sum_{i=1}^m g(Y_i), \quad (15)$$

where Y_1, \dots, Y_m are MCMC draws from the distribution with parameter η . There are many MCMC algorithms such as Metropolis-Hastings [15] or Swensen-Wang [48, 50], used for the Ising model example in Section 5.2.

The accuracy of the approximation in (15) increases with Monte Carlo sample size m . When the current estimate is far away from the MLE, we can use

smaller m to save time and work with a fairly noisy approximation of the gradient. However, when the current estimate approaches the MLE, larger m are necessary.

Our algorithm relies on the computed values of $\nabla\ell(\eta)$ in the curvature condition (14), as well as the stop condition for the algorithm, $\|\nabla\ell(\eta_k)\| < \epsilon$. Given that we may only have approximations of $\nabla\ell(\eta)$, we cannot know for certain if either of these conditions are truly met. We can ameliorate this by constructing confidence intervals for each of the inequalities.

For the inequalities in (14), we can estimate asymptotic standard errors of $\nabla\ell(\eta_k + \alpha_k p_k)^T p_k$ and $c\nabla\ell(\eta_k)^T p_k - \nabla\ell(\eta_k + \alpha_k p_k)^T p_k$ by appealing to the Markov chain Central limit theorem [6, 26, 40, 41]. The `initseq` function from the R package `mcmc` [13] can be used to estimate asymptotic standard errors for univariate functionals of reversible Markov chains: given an MCMC sample for a univariate quantity, `initseq` returns a value (divided by sample size) that is an estimate of the asymptotic variance in the Markov chain central limit theorem. Both of the quantities in (14) are univariate. In the second expression, $c\nabla\ell(\eta_k)^T p_k - \nabla\ell(\eta_k + \alpha_k p_k)^T p_k$, the MCMC sample generated for $\nabla\ell(\eta_k + \alpha_k p_k)^T p_k$ is independent of the sample generated for $c\nabla\ell(\eta_k)^T p_k$. Thus `initseq` can be applied to each sample separately and the results summed for an estimated variance. We can then be approximately 95% confident (non-simultaneously) that α_k satisfies (14) if

$$\begin{aligned} \nabla\ell(\eta_k + \alpha_k p_k)^T p_k - 1.645 \cdot \text{se}_1 &> 0 \\ c\nabla\ell(\eta_k)^T p_k - \nabla\ell(\eta_k + \alpha_k p_k)^T p_k - 1.645 \cdot \text{se}_2 &> 0 \end{aligned}$$

where se_1 and se_2 are the asymptotic standard errors for $\nabla\ell(\eta_k + \alpha_k p_k)^T p_k$ and $c\nabla\ell(\eta_k)^T p_k - \nabla\ell(\eta_k + \alpha_k p_k)^T p_k$, respectively, calculated as described.

The delta method can be applied to estimate a standard error for $\|\nabla\ell(\eta_k)\|$. The asymptotic variance is calculated by

$$V(\|\nabla\ell(\eta_k)\|) = \frac{1}{\|\nabla\ell(\eta_k)\|^2} \nabla\ell(\eta_k)^T \Sigma \nabla\ell(\eta_k),$$

where Σ is the variance matrix of $\nabla\ell(\eta_k)$ and can be estimated by the sample variance matrix of the batch mean vectors of $g(Y_1), \dots, g(Y_n)$ divided by the number of batches (the `initseq` function requires a univariate vector and so cannot be used here). We can be approximately 95% confident that $\|\nabla\ell(\eta_k)\| > \epsilon$ if

$$\|\nabla\ell(\eta_k)\| - 1.645\sqrt{V(\|\nabla\ell(\eta_k)\|)} > \epsilon.$$

In practice, however, use of confidence intervals does not appear necessary with Monte Carlo sample sizes that are set large enough so that these standard errors are initially small relative to the point estimates. The ratio of point estimate to standard error of course decreases as the algorithm progresses and the estimate of the parameter nears the MLE, reflected in $\nabla\ell(\eta_k)$ nearing 0. Thus these confidence intervals are most useful as a guide for when to increase the MCMC sample size, or when to switch methods, or when to terminate the algorithm.

4.4. Combining with other algorithms

We believe the best use of this algorithm is in combination with other faster methods like MCMC-MLE or Newton-Raphson safeguarded by our line search algorithm. Our algorithm with steepest ascent or conjugate gradient search direction should be used initially from long range, when one has no good intuition for an initial value. It is well known that when the objective function is quadratic, the conjugate gradient method with exact arithmetic converges to the solution in at most d steps, where d is the dimension of the problem [32, Chapter 5]. As a rule of thumb, we think using our algorithm for $2d$ steps before switching seems reasonable when using conjugate gradient directions. Determining a more precise criteria for when we are inside the “radius of convergence” for algorithms like Newton-Raphson or MCMC-MLE is an area for further research.

In the case of MCMC-MLE, [22] show that getting a distribution so that $g(y)$ is within the convex hull of the MCMC samples is a necessary condition for this algorithm to converge (in fact, this appears to be the impetus for their steplength MCMC-MLE approach). However, this is not sufficient. An effective approach may be to examine the importance sampling weights used in the log likelihood approximation in the previous iteration, and look for these to stabilize. This “rearview” approach should inform us if the previous value for η_k was close enough to apply MCMC-MLE; applying MCMC-MLE to the current distribution should then converge.

5. Examples

We illustrate the application of our algorithm in two settings for regular exponential families, where the MLE is known to exist:

1. Logistic regression. The gradient of the log likelihood, $\nabla\ell(\eta)$, can be calculated exactly. In this setting, Theorem 3.2 guarantees convergence to the MLE from any starting point.
Of course, logistic regression is done more efficiently through iterated reweighted least-squares, as in the `glm` function in the R platform. Our purpose here is to show the application of our algorithm in a familiar setting. We choose a poor starting point, where an algorithm like Newton-Raphson would fail. Other optimization algorithms such as those based on (4) would also work here.
2. Ising model. The gradient of the log likelihood, $\nabla\ell(\eta)$, can only be approximated via MCMC. In this setting, Theorem 3.2 does not guarantee convergence, but our algorithm is still effective in practice, even for poor starting values. In this case, optimization algorithms based on (4) cannot be applied since they depend on evaluating the log likelihood.

In both cases, we compare our algorithm with stochastic approximation to clarify the distinctions made in Section 1.1 between the two approaches.

5.1. Example: Logistic regression

We apply our algorithm to the case of a logistic regression with a starting point far from the solution. In such a case, the Hessian matrix is often near-singular and algorithms such as Newton-Raphson which rely on it will fail. For classical SA with step size $1/k$, the magnitudes of the updates diminish too quickly for the parameter estimates to approach the MLE in a reasonable amount of time.

The response vector Y has components that are Bernoulli trials with mean vector p . The natural parameter is $\theta_i = \log\left(\frac{p_i}{1-p_i}\right)$, which is modeled component-wise as a linear function of the predictors $1, x_1, \dots, x_q$, so that

$$\theta_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_q x_{qi} = \beta^T x_i \quad i = 1, \dots, n$$

where $\beta = (\beta_0, \dots, \beta_q)^T$ and $x_i = (1, x_{1i}, \dots, x_{qi})^T$.

Defining the model matrix M to be the $n \times (q+1)$ matrix with the x_i as rows, we can express $\theta = M\beta$. This in turn allows us to reparameterize the exponential family as one with β as the natural parameter vector and $M^T y$ the vector of statistics with log likelihood

$$\ell(\beta) = \beta^T (M^T y) - c(\beta),$$

where y is the vector of observed Bernoulli responses. By (9), the gradient is

$$\nabla \ell(\beta) = M^T y - E_\beta(M^T Y) = M^T (y - E_\beta(Y)),$$

where $E_\beta(Y) = \frac{1}{1 + \exp(-M\beta)}$ can be calculated exactly. This allows us to directly apply Theorem 3.2.

We specified our true 100-dimensional parameter value to be

$$\beta = (-0.748, 0.357, 0.727, 0.296, -0.904, 0.960, 0.262, -0.353, \dots, -0.162)^T$$

with component values generated from a Uniform($-1, 1$) distribution. Then, using 1000 independent draws from a correlated multivariate normal distribution centered at 0 as our predictors, we generated data for this model (the data vector has length 1000 like the predictor vectors). Fitting these data using the R function `glm`, we find the MLE of β to be

$$\hat{\beta}_{\text{MLE}} = (-1.051, 0.862, 0.908, 0.229, -1.192, 1.187, 0.294, -0.655, \dots, -0.377)^T,$$

where the disparity to the true value of β results from a relatively small sample size of $n = 1000$. We then use

$$\beta_0 = (5, -5, 2, 0, 3, 4, 3, 0, \dots, 1)^T$$

for the starting point for our line search algorithm, a point for which Newton-Raphson fails due to a nearly singular Hessian matrix.

We measure the performance of our algorithm in terms of the total number of iterations used, where each iteration requires evaluation of the gradient, $\nabla \ell(\beta_k +$

TABLE 1
 Comparison of MLEs of β for Example 1: MLE = *glm*, Steep = line search using steepest ascent, CG = line search using conjugate gradient, and SA = SA with step size = $1/k$ terminated at 10,000 iterations, n = number of iterations. Only first 8 components out of 100 shown. Our proposed algorithm arrives at nearly identical MLE estimates to *glm*

	n	$\beta[1]$	$\beta[2]$	$\beta[3]$	$\beta[4]$	$\beta[5]$	$\beta[6]$	$\beta[7]$	$\beta[8]$
True β		-0.748	0.357	0.727	0.296	-0.904	0.960	0.262	-0.353
$\hat{\beta}_{\text{MLE}}$		-1.051	0.862	0.908	0.229	-1.192	1.187	0.294	-0.655
$\hat{\beta}_{\text{Steep}}$	538	-1.052	0.862	0.909	0.229	-1.192	1.188	0.294	-0.655
$\hat{\beta}_{\text{CG}}$	292	-1.051	0.862	0.908	0.229	-1.192	1.187	0.294	-0.655
$\hat{\beta}_{\text{SA}}$	10^5	-132.33	90.92	120.18	28.40	-128.58	99.011	24.39	-56.02

$\alpha_k p_k$). Typically, several iterations take place in an inner loop to find a step size α_k that meets the curvature condition (14), a process that grows increasingly difficult as the estimates near the MLE since the rightmost term in (14) gets smaller in magnitude. Once an acceptable step size is found, the parameter estimate β_k is updated and a new search direction is determined, requiring another evaluation of the gradient.

Our algorithm took 538 iterations over 204 different search directions to get $\|\nabla\ell(\beta_k)\| < 0.01$ and arrive at an estimate for the MLE that differs from the *glm* result by 3.309 in Euclidean distance (See Table 1). Using the Polak-Ribière conjugate gradient method described in the previous section resulted in comparably sharp MLE estimates (see Table 1) in fewer iterations—292 over 116 search directions—a noticeable improvement.

We also applied SA with step size $1/k$ (setting $A = 1$, $B = 0$ in (2)) from the same starting point β_0 . The choice of constants A and B in the step size is of course not likely to be optimal; however, we want to apply SA without trial and error experimentation. After 100,000 iterations, the parameter estimates look nothing at all like the MLE (See Table 1). The initial step sizes are far too large, then diminish rapidly so that the algorithm does not converge in a reasonable amount of time. Table 2 shows the first 20 step sizes used by SA and our line search. Our line search continues to use step sizes of relatively stable magnitude even well into the process. It should be noted that these 20 step sizes correspond to the first 20 iterations of SA but 50 iterations of our line search algorithm using steepest ascent, since we spends several iterations finding an acceptable step size for each update, and 45 iterations using conjugate gradient directions.

5.2. Example: Ising model

In this example, we apply our gradient-based line search algorithm to an Ising model [24] on a toroidal square lattice. Here the gradient of the log likelihood cannot be calculated exactly as in the logistic example and so Theorem 3.2 cannot be applied directly. However, as discussed in Section 4.3, the gradient can be approximated using MCMC, allowing our algorithm to still be effective in finding the MLE.

Ising models are exponential families where each entry in the square lattice takes the value of either zero or one. A realized sample is shown in Figure 1. The

TABLE 2

The first 20 step sizes used by SA (with step size $1/k$) and our algorithm for Example 1.
The step sizes used by our algorithm do not diminish like $1/k$

k	$\alpha_{SA} = 1/k$	α_{Steep}	α_{CG}
1	1.000	0.002	0.002
2	0.500	0.055	0.055
3	0.333	0.007	0.012
4	0.250	0.071	0.073
5	0.200	0.002	0.003
6	0.167	0.029	0.073
7	0.143	0.003	0.015
8	0.125	0.078	0.013
9	0.111	0.004	0.008
10	0.100	0.003	0.022
11	0.091	0.005	0.005
12	0.083	0.003	0.030
13	0.077	0.005	0.015
14	0.071	0.003	0.006
15	0.067	0.007	0.009
16	0.063	0.002	0.006
17	0.059	0.011	0.009
18	0.056	0.002	0.004
19	0.053	0.024	0.010
20	0.050	0.002	0.004

sufficient statistic vector is two-dimensional, comprising the number of entries with value one and the number of “neighbor” entries with the same value. Entries are considered “neighbors” if they are adjacent to one another horizontally or vertically (but not diagonally).

Here we describe the toroidal square lattice as an $n \times n$ matrix Y and each entry as Y_{ij} , where i and j take values in $1, \dots, n$ considered as a cyclical set (addition is done modulo n). The sufficient statistic, $g(y)$, has components:

$$g_1(y) = \sum_{i=1}^n \sum_{j=1}^n I(Y_{ij} = 1),$$

$$g_2(y) = \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n [I(Y_{ij} = Y_{i-1,j}) + I(Y_{ij} = Y_{i,j-1})$$

$$+ I(Y_{ij} = Y_{i+1,j}) + I(Y_{ij} = Y_{i,j+1})],$$

where $I(\cdot)$ denotes the indicator function taking logical expressions to the numbers zero and one, false expressions to zero and true expressions to one.

Because of the dependence of neighboring entries in the lattice, there is no closed form expressing $\nabla \ell(\eta)$. Instead, we need to approximate $\nabla \ell(\eta)$ using MCMC as described by (15). The MCMC draws are performed here using the Swendsen-Wang algorithm [48, 50], available in the contributed R package `potts` [16].

We choose $\eta = (0, \log(1 + \sqrt{2}))^T$ to generate a 32×32 lattice, which we use as our observed data (Figure 1). This value for η is of particular interest

TABLE 3

Comparison of MLEs for η for Example 2: MLE = Newton-Raphson starting from the true η , Steep = line search using steepest ascent, CG = line search using conjugate gradient, and SA = SA with step size = $1/k$. All algorithms converged

	MC Samples (thousands)	$\eta[1]$	$\eta[2]$
True η		0.000	0.881
$\hat{\eta}_{\text{MLE}}$		-0.007	0.896
$\hat{\eta}_{\text{Steep}}$	620	-0.011	0.895
$\hat{\eta}_{\text{CG}}$	450	-0.008	0.895
$\hat{\eta}_{\text{SA}}$	1368	-0.010	0.895

TABLE 4

The first 17 step sizes used by SA (with step size $1/k$) and our algorithm for Example 2. The step sizes used by our algorithm are initially much smaller than $1/k$

k	$\alpha_{\text{SA}} = 1/k$	α_{Steep}	α_{CG}
1	1.0000	0.0029	0.0029
2	0.5000	0.0005	0.0005
3	0.3333	0.0017	0.0017
4	0.2500	0.0013	0.0045
5	0.2000	0.0017	0.0007
6	0.1667	0.0011	0.0002
7	0.1429	0.0021	0.0015
8	0.1250	0.0009	
9	0.1111	0.0020	
10	0.1000	0.0007	
11	0.0909	0.0018	
12	0.0833	0.0006	
13	0.0769	0.0013	
14	0.0714	0.0006	
15	0.0667	0.0007	
16	0.0625	0.0003	
17	0.0588	0.0013	

because it corresponds to the phase transition point [37] and has been shown to be difficult to estimate [9]. In order to get a good estimate of the MLE to which we can compare our algorithm's results, we use 10 iterations of MCMC Newton-Raphson [36] starting at the true value of η so that it will converge.

We apply our line search algorithm to this data using a far off initial value of $\eta^{(0)} = (2, 0.001)$ and a fixed MCMC sample size of 10,000. Our algorithm used 62 iterations (gradient evaluations) over 17 search directions to get $\|\nabla\ell(\eta_k)\| < 0.005$ and arrive at an estimate of the MLE that differs from Newton-Raphson by 0.0037 (see Table 3). Using the Polak-Ribière conjugate gradient method resulted in comparably sharp MLE estimates using 45 iterations over 7 search directions. The total MCMC sample sizes used were $62 \times 10,000 = 620,000$ and $45 \times 10,000 = 450,000$, respectively.

We also applied MCMC SA, again with step size $1/k$ from the same starting point $\eta^{(0)}$, and used a MCMC sample size of 1,000 for gradient calculation. Here SA converged in 1368 iterations or 1,368,000 MC samples, comparable to our algorithm (see Table 3). Table 4 shows the first 17 step sizes used by SA and our line search. The step sizes used by our line search are initially very small

compared to $1/k$, but stay in a range of about $1/300$ to $1/3000$. So, the $1/k$ step size used by SA in fact occasionally satisfies our curvature condition when k is large.

6. Discussion

We have presented a simple line search algorithm for finding the MLE of a regular exponential family when the MLE exists (or the MLE in the Barndorff-Nielsen completion when the MLE does not exist in the conventional sense). The algorithm avoids any trial-and-error experimentation involving tuning parameters or starting points commonly associated with optimization routines not invented by optimization specialists. Our algorithm is modeled after algorithms discussed in optimization textbooks [8, 32, 47], all of which are safeguarded to ensure rapid automatic convergence.

Convergence is guaranteed when the gradient can be calculated exactly. Even when the gradient cannot be calculated exactly and is only estimable via MCMC, the algorithm is still useful in practice, as demonstrated by the Ising model example. We have also described a way to construct and use confidence intervals to make convergence highly probable.

The algorithm can be computationally demanding. When the current iteration approaches the solution, the curvature condition for step size becomes more difficult to satisfy and the method may require several iterations of MCMC sampling and perhaps an increase in MCMC sample size. Eventual increase in MCMC sample size is unavoidable, because the achievable accuracy is inversely proportional to the square root of the MCMC sample size, as in all Monte Carlo. Thus we believe the best use of this algorithm is in combination with other faster methods like MCMC-MLE or Newton-Raphson safeguarded by our line search algorithm. Our algorithm should be used from “long range”, when one has no good intuition for an initial value and is concerned about picking one that is far from the MLE. The switch between types of search direction (steepest ascent, conjugate gradient, or Newton) within our algorithm or the switch to another algorithm (such as MCMC-MLE) need not require manual intervention. When used in combination, we do not think the confidence intervals are necessary as the curvature condition is quite easily satisfied when the current iteration is far from the MLE.

One way to improve performance is to use conjugate gradient search directions rather than steepest ascent. In our examples, this reduced the number of iterations by over 25%. However, in other problems we tried with different dimensionality, this performance varied significantly and it appears that no guarantee can be made about quantity of improvement in performance, though in all cases we examined, it never did worse. This is no surprise, because the necessity of “preconditioning” for good performance of the conjugate gradient algorithm is well known (but no good “preconditioner” is available for maximum likelihood in exponential families).

There are several outstanding issues. Most notably, we have not showed convergence of the algorithm when the gradient is approximated via MCMC. This

is a more difficult theoretical problem and is the motivation for stochastic approximation research. Further work is necessary to determine if one can adapt our restrictive curvature condition (14) to the approach of [1, 30] in MCMC stochastic approximation.

Another remaining issue is the stopping criterion: what value should be chosen for ϵ in the exit condition $\|\nabla\ell(\eta_k)\| < \epsilon$? Because the value of $\|\nabla\ell(\eta_k)\|$ can only be approximated via MCMC, one cannot be certain if this condition is actually satisfied. Here again, the switch to another methodology may be appropriate, though at least in our Ising model example, our use of 10,000 for the MCMC sample size and 0.005 for ϵ were successful in obtaining a reasonable parameter estimate.

A final remaining issue is estimation of Monte Carlo error of the estimates. Here too we recommend switching to another algorithm at the end. The MCMC-MLE procedure gives accurate error estimates [11]. For very small steps these are essentially the same as the Monte Carlo error of a single unsafeguarded Newton-Raphson step, so the method in [11] can be used for either.

Appendix A: Proofs

Proof of Theorem 3.1. Let $f(\cdot)$ represent the negative log likelihood $-\ell(\cdot)$, the objective function to be minimized. We proceed from the perspective of a minimization of a function $f(\cdot)$ since this is the convention in the optimization literature [32, 42].

The negative log likelihood function $-\ell(\cdot)$ is strictly convex by (10), and continuous since it is infinitely differentiable by Theorem 5.8 in [28]. It is bounded below by the negative log likelihood of the limiting conditional model of this exponential family described in Theorem 6 of [12], which is guaranteed to have a global minimum.

Then, unless $\nabla f(x_k) = 0$ in which case x_k is already the solution, for each k , we can uniquely define α_{c_k} as follows:

$$\nabla f(x_k + \alpha_{c_k} p_k)^T p_k = c \nabla f(x_k)^T p_k \quad (16)$$

The point α_{c_k} is uniquely defined because it is the minimizer of the function $\alpha \mapsto f(x_k + \alpha p_k) - \alpha c \nabla f(x_k)^T p_k$. We may also define α_{\min_k} as follows:

$$\alpha_{\min_k} = \begin{cases} \alpha \text{ s.t. } \nabla f(x_k + \alpha p_k)^T p_k = 0 & \text{if such an } \alpha \text{ exists} \\ +\infty & \text{otherwise.} \end{cases} \quad (17)$$

These values appear on the α -axis in Figure 3 for the case where a minimizer exists for $\alpha \mapsto f(x_k + \alpha p_k)$.

By the strict convexity of f and Theorem 2.14(b) in [42],

$$f(x_k + \alpha_{c_k} p_k) < f(x_k) + [\nabla f(x_k + \alpha_{c_k} p_k)]^T \alpha_{c_k} p_k.$$

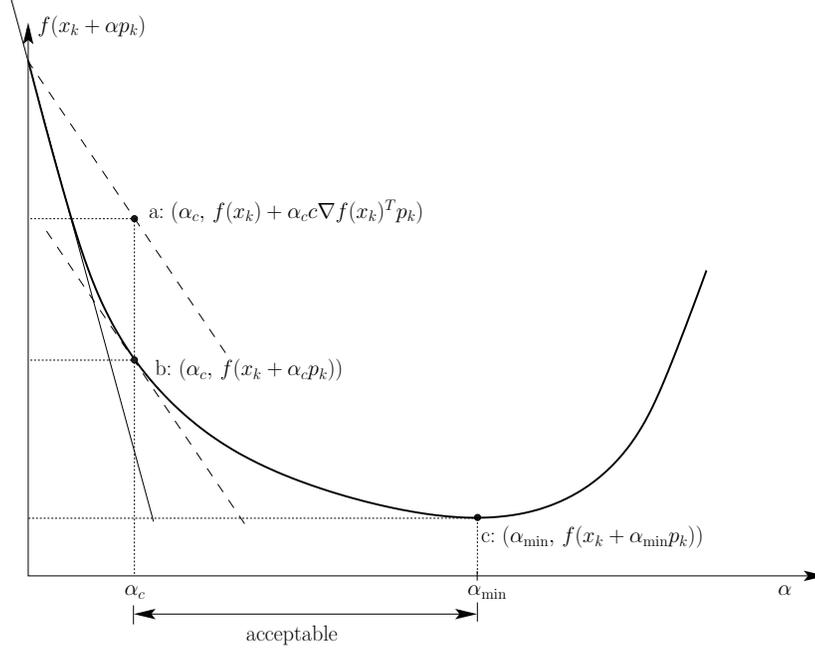


FIG 3. Acceptable region for α according to curvature condition (14) when restricting to direction p_k .

Applying (16) to the right hand side of the above gives

$$f(x_k + \alpha_{c_k} p_k) < f(x_k) + \alpha_{c_k} c \nabla f(x_k)^T p_k. \quad (18)$$

(See points a and b in Figure 3).

The subproblem $\alpha \mapsto f(x_k + \alpha p_k)$ is strictly convex and hence monotonically decreasing at α_k such that $\alpha_{c_k} \leq \alpha_k < \alpha_{\min_k}$ (in Figure 3, see points b and c). That is,

$$f(x_k + \alpha_{\min} p_k) \leq f(x_k + \alpha_k p_k) \leq f(x_k + \alpha_{c_k} p_k), \quad (19)$$

where the left-hand side denotes $\inf_{\alpha \in \mathbb{R}} f(x_k + \alpha p_k)$ when $\alpha_{\min} = \infty$.

Combining the second inequality of (19) with (18), we have

$$f(x_k + \alpha_k p_k) < f(x_k) + \alpha_{c_k} c \nabla f(x_k)^T p_k, \quad (20)$$

We now turn our attention to (16). Define $x_{c_k} = x_k + \alpha_{c_k} p_k$. Then

$$\nabla f(x_{c_k})^T p_k = c \nabla f(x_k)^T p_k.$$

Subtracting $\nabla f(x_k)^T p_k$ from both sides gives

$$(\nabla f(x_{c_k}) - \nabla f(x_k))^T p_k = (c - 1) \nabla f(x_k)^T p_k. \quad (21)$$

By (10), $\nabla^2 \ell(\eta)$ is bounded for finite state space $g(\mathcal{Y})$, which is true by assumption. Thus $|\nabla^2 f(x)| \leq K$ for some constant K for all x . Then by Theorems 9.2 and 9.7 in [42], $\nabla f(x)$ is Lipschitz continuous relative to the convex set \mathbb{R}^d .

Thus there exists a constant $L < \infty$ such that

$$\|\nabla f(x) - \nabla f(\tilde{x})\| \leq L\|x - \tilde{x}\| \quad \text{for all } x, \tilde{x} \in \mathbb{R}^d.$$

Applying this relation to x_{c_k} and x_k , we have

$$\|\nabla f(x_{c_k}) - \nabla f(x_k)\| \leq L\|x_{c_k} - x_k\| = L\|\alpha_{c_k} p_k\|. \quad (22)$$

The rest of this proof is nearly identical to the proof for Theorem 3.2 in [32, pp. 43–44]. Multiplying both sides of (22) by $\|p_k\|$ and then applying Cauchy-Schwartz gives

$$(\nabla f(x_{c_k}) - \nabla f(x_k))^T p_k \leq \alpha_{c_k} L \|p_k\|^2. \quad (23)$$

Substituting (21) into the left-hand side of (23) gives

$$-\alpha_{c_k} \leq \frac{(1-c)}{L} \frac{\nabla f(x_k)^T p_k}{\|p_k\|^2}. \quad (24)$$

Substituting (24) into (20), we obtain

$$f(x_{k+1}) < f(x_k) - c \frac{(1-c)}{L} \frac{(\nabla f(x_k)^T p_k)^2}{\|p_k\|^2}. \quad (25)$$

The angle θ_j between the search direction p_k and steepest descent direction $-\nabla f(x_k)$ can be expressed by $\cos \theta_j = \frac{-\nabla f(x_j)^T p_j}{\|\nabla f(x_j)\| \cdot \|p_j\|}$. Substituting this relation into (25) gives

$$f(x_{k+1}) < f(x_k) - c \frac{(1-c)}{L} \|\nabla f(x_k)\|^2 \cos^2 \theta_k.$$

By summing this expression over all indices less than or equal to k ,

$$f(x_{k+1}) < f(x_0) - c \frac{(1-c)}{L} \sum_{j=0}^k \|\nabla f(x_j)\|^2 \cos^2 \theta_j.$$

Because $f(x)$ is bounded below, there exists some $M < \infty$ such that $f(x_0) - f(x_{k+1}) < M$ for all k , so that

$$\frac{c(1-c)}{L} \sum_{j=0}^k \|\nabla f(x_j)\|^2 \cos^2 \theta_j < M < \infty.$$

Taking $k \rightarrow \infty$ while noting that $0 < c < 1$ gives

$$\sum_{j=0}^{\infty} \|\nabla f(x_j)\|^2 \cos^2 \theta_j < \infty. \quad (26)$$

With the additional restriction on the search direction p_k such that $\cos \theta_k \geq \delta > 0$ for some choice of δ , for all choices of k , the convergent series in (26) implies that

$$\lim_{k \rightarrow \infty} \|\nabla f(x_k)\| = 0.$$

□

The inequality (26) has been referred to as *Zoutendijk's condition* [32], though we arrive at this result via different assumptions.

Theorem 3.1 shows that the gradient of the objective function converges to 0. The proof for Theorem 3.2 is concerned with the conditions for mapping this convergence to the convergence of the iterated parameter estimates η_k to the unique MLE. In particular, the mapping from η_k to the gradient must be globally invertible.

Proof of Theorem 3.2. The Fisher information for an exponential family with minimal representation is non-singular by (11) and thus invertible. If we consider the map defined by

$$h(\eta) = \nabla c(\eta)$$

where c is the cumulant function (7), its first derivative matrix is

$$\nabla h(\eta) = \nabla^2 c(\eta) = I(\eta) \tag{27}$$

which is again non-singular. Since this is true for any η , by the inverse function theorem, h is everywhere locally invertible.

In fact, h is globally invertible. For any μ in the range of h , consider the function

$$q(\eta) = \mu^T \eta - c(\eta).$$

Since $\nabla^2 q(\eta) = -I(\eta)$ by (27), q is strictly concave. There exists an η such that

$$\nabla q(\eta) = \mu - h(\eta) = 0$$

because of the assumption that μ is in the range of h . Thus this η is a stationary point of q , which is a global maximizer of q by concavity and the unique global maximizer by strict concavity. This says that for every μ in the range of h , there exists a unique η such that $\mu = h(\eta)$, which is global invertibility of h .

Since c is infinitely differentiable by Theorem 2.7.1 in [29], so is h , and by the inverse function theorem, so is h^{-1} (even if we do not know the form of h^{-1}). The first derivative of h^{-1} can be expressed as

$$\nabla h^{-1}(\mu) = [\nabla h(\eta)]^{-1} = [I(\eta)]^{-1}, \quad \text{when } \mu = h(\eta),$$

and is thus non-singular everywhere, including at the MLE of η , $\hat{\eta}_{\text{MLE}}$.

Thus our algorithm, which concludes that $\|\nabla\ell(\eta_k)\| = \|g(y) - h(\eta_k)\| \rightarrow 0$, implies that

$$\mu_k = h(\eta_k) \rightarrow g(y),$$

or

$$h^{-1}(\mu_k) \rightarrow h^{-1}(g(y)),$$

or

$$\eta_k \rightarrow \hat{\eta}_{\text{MLE}}$$

because if the MLE exists, then the gradient of the log likelihood is zero at the MLE which is $g(y) = h(\hat{\eta}_{\text{MLE}})$. \square

References

- [1] ANDRIEU, C., MOULINES, E. and PRIOURET, P. (2005). Stability of Stochastic Approximation under Verifiable Conditions. *SIAM Journal on Control and Optimization* **44** 283–312. [MR2177157](#)
- [2] BARNDORFF-NIELSEN, O. (1978). *Information and Exponential Families in Statistical Theory*. John Wiley & Sons. [MR0489333](#)
- [3] BESAG, J. (1974). Spatial Interaction and the Statistical Analysis of Lattice Systems. *Journal of the Royal Statistical Society, Series B* **36** 192-236. [MR0373208](#)
- [4] BESAG, J. (1975). Statistical Analysis of Non-lattice Data. *The Statistician* **24** 179-195.
- [5] BROWN, L. D. (1986). *Fundamentals of Statistical Exponential Families: with Applications in Statistical Decision Theory*. Institute of Mathematical Statistics, Hayward, CA. [MR0882001](#)
- [6] CHAN, K. S. and GEYER, C. J. (1994). Discussion of the Paper by Tierney. *Annals of Statistics* **22** 1747–1758.
- [7] CHEN, H.-F. (2002). *Stochastic Approximation and Its Applications*. Kluwer Academic Publishers, Dordrecht. [MR1942427](#)
- [8] FLETCHER, R. (1987). *Practical Methods of Optimization*, Second ed. John Wiley & Sons. [MR0955799](#)
- [9] GEYER, C. J. (1990). Likelihood and Exponential Families PhD thesis, University of Washington. [MR2685353](#)
- [10] GEYER, C. J. (1991). Markov chain Monte Carlo Maximum Likelihood. In *Computing Science and Statistics: Proc. 23rd Symp. Interface* (E. KERAMIDAS, ed.) 156–163. Interface Foundation.
- [11] GEYER, C. J. (1994). On the Convergence of Monte Carlo Maximum Likelihood Calculations. *Journal of the Royal Statistical Society, Series B* **56** 261-274. [MR1257812](#)
- [12] GEYER, C. J. (2009a). Likelihood Inference in Exponential Families and Directions of Recession. *Electronic Journal of Statistics* **3** 259–289. [MR2495839](#)

- [13] GEYER, C. J. (2009b). `mcmc`: Markov chain Monte Carlo. R package version 0.7-3. [MR2495839](#)
- [14] GEYER, C. J. (2010). `aster2`: Aster models. R package version 0.1.
- [15] GEYER, C. J. (2011). Introduction to MCMC. In *Handbook of Markov Chain Monte Carlo* (S. P. Brooks, A. E. Gelman, G. L. Jones and X. L. Meng, eds.) Chapman & Hall/CRC, Boca Raton, FL.
- [16] GEYER, C. J. and JOHNSON, L. T. (2010). `potts`: Markov chain Monte Carlo for Potts Models. R package version 0.4.
- [17] GEYER, C. J. and THOMPSON, E. A. (1992). Constrained Monte Carlo Maximum Likelihood for Dependent Data. *Journal of the Royal Statistical Society, Series B* **54** 657-699. [MR1185217](#)
- [18] GOODREAU, S. M. (2007). Advances in Exponential Random Graph (p*) Models Applied to a Large Social Network. *Social Networks* **29** 231-248.
- [19] GOODREAU, S. M., KITTS, J. A. and MORRIS, M. (2009). Birds of a Feather, or Friend of a Friend? Using Exponential Random Graph Models to Investigate Adolescent Social Networks. *Demography* **46** 103-125.
- [20] GU, M. G. and ZHU, H.-T. (2001). Maximum Likelihood Estimation for Spatial Models by Markov Chain Monte Carlo Stochastic Approximation. *Journal of the Royal Statistical Society, Series B* **63** 339-355. [MR1841419](#)
- [21] HANDCOCK, M. S., HUNTER, D. R., BUTTS, C. T., GOODREAU, S. M. and MORRIS, M. (2003). `statnet`: Software Tools for the Statistical Modeling of Network Data. Version 2.0. Project home page at <http://statnetproject.org>.
- [22] HUMMEL, R., HUNTER, D. R. and HANDCOCK, M. S. (2010). A Steplength Algorithm for Fitting ERGMs Technical Report No. 10-03, Pennsylvania State University.
- [23] HUNTER, D. R., HANDCOCK, M. S., BUTTS, C. T., GOODREAU, S. M. and MORRIS, M. (2008). `ergm`: A Package to Fit, Simulate and Diagnose Exponential-Family Models for Networks. *Journal of Statistical Software* **24**.
- [24] ISING, E. (1925). Beitrag zur Theorie des Ferromagnetismus. *Zeitschrift für Physik A Hadrons and Nuclei* **31** 253-258.
- [25] JAYNES, E. T. (1978). Where Do We Stand on Maximum Entropy? In *The Maximum Entropy Formalism* (R. D. Levine and M. Tribus, eds.) Cambridge: Massachusetts Institute of Technology Press.
- [26] JONES, G. L. (2004). On the Markov Chain Central Limit Theorem. *Probability Surveys* **1** 299-320. [MR2068475](#)
- [27] KUSHNER, H. J. and YIN, G. G. (1997). *Stochastic Approximation Algorithms and Applications*. Springer, New York. [MR1453116](#)
- [28] LEHMANN, E. L. and CASELLA, G. (1998). *Theory of Point Estimation*, Second ed. Springer. [MR1639875](#)
- [29] LEHMANN, E. L. and ROMANO, J. P. (2005). *Testing Statistical Hypotheses*, 3rd ed. Springer. [MR2135927](#)
- [30] LIANG, F. (2010). Trajectory Averaging for Stochastic Approximation MCMC Algorithms. *The Annals of Applied Statistics* **38** 2823-2856. [MR2722457](#)

- [31] MOYEED, R. A. and BADDELEY, A. J. (1991). Stochastic Approximation of the MLE for a Spatial Point Pattern. *Scandinavian Journal of Statistics* **18** 39–50. [MR1115181](#)
- [32] NOCEDAL, J. and WRIGHT, S. J. (1999). *Numerical Optimization*, First ed. Springer. [MR1713114](#)
- [33] OKABAYASHI, S. (2011a). Parameter Estimation in Social Network Models PhD thesis, University of Minnesota.
- [34] OKABAYASHI, S. (2011b). Supporting Theory and Data Analysis for “Long range search for maximum likelihood in exponential families” Technical Report No. 686, University of Minnesota.
- [35] OKABAYASHI, S., JOHNSON, L. and GEYER, C. J. (2011). Extending Pseudo-likelihood for Potts Models. *Statistica Sinica* **21** 331–347. [MR2796865](#)
- [36] PENTTINEN, A. (1984). Modelling Interactions in Spatial Point Patterns: Parameter Estimation by the Maximum Likelihood Method. *Jyväskylä Studies in Computer Science, Economics and Statistics* **7**.
- [37] POTTS, R. B. (1952). Some Generalized Order-Disorder Transformations. *Proceedings of the Cambridge Philosophical Society* **48** 106–109. [MR0047571](#)
- [38] RINALDO, A., FIENBERG, S. E. and ZHOU, Y. (2009). On the Geometry of Discrete Exponential Families with Application to Exponential Random Graph Models. *Electronic Journal of Statistics* **3** 446–484. [MR2507456](#)
- [39] ROBBINS, H. and MONRO, S. (1951). A Stochastic Approximation Method. *Annals of Mathematical Statistics* **22** 400–407. [MR0042668](#)
- [40] ROBERTS, G. O. and ROSENTHAL, J. S. (1997). Geometric Ergodicity and Hybrid Markov Chains. *Electronic Communications in Probability* **2** 13–25. [MR1448322](#)
- [41] ROBERTS, G. O. and ROSENTHAL, J. S. (2004). General State Space Markov Chains and MCMC Algorithms. *Probability Surveys* **1** 20–71. [MR2095565](#)
- [42] ROCKAFELLAR, R. T. and WETS, R. J.-B. (2004). *Variational Analysis. corrected second printing*. Springer-Verlag, Berlin. [MR1491362](#)
- [43] SAUL, Z. M. and FILKOV, V. (2007). Exploring Biological Network Structure using Exponential Random Graph Models. *Bioinformatics* **23** 2604–02611.
- [44] SHAW, R. G., GEYER, C. J., WAGENIUS, S., HANGELBROEK, H. H. and ETTERSON, J. R. (2008). Unifying Life-History Analyses for Inference of Fitness and Population Growth. *The American Naturalist* **172** E35–E47.
- [45] SNIJDERS, T. A. B. (2002). Markov Chain Monte Carlo Estimation of Exponential Random Graph Models. *Journal of Social Structure* **3**.
- [46] STRAUSS, D. and IKEDA, M. (1990). Pseudolikelihood Estimation for Social Networks. *Journal of the American Statistical Association* **85** 204–212. [MR1137368](#)
- [47] SUN, W. and YUAN, Y.-X. (2006). *Optimization Theory and Methods: Nonlinear Programming*. Springer. [MR2232297](#)
- [48] SWENDSEN, R. H. and WANG, J.-S. (1987). Nonuniversal Critical Dynamics in Monte Carlo Simulations. *Physics Review Letters* **58** 86–88.

- [49] VAN DUIJN, M. A. J., GILE, K. J. and HANDCOCK, M. S. (2009). A Framework for the Comparison of Maximum Pseudo-likelihood and Maximum Likelihood Estimation of Exponential Family Random Graph Models. *Social Networks* **31** 52-62.
- [50] WANG, J. S. and SWENDSEN, R. H. (1990). Cluster Monte Carlo Algorithms. *Physics A* **167** 565–579. [MR1075564](#)
- [51] WASSERMAN, S. and PATTISON, P. (1996). Logit Models and Logistic Regression for Social Networks: I. An Introduction to Markov Graphs and p^* . *Psychometrika* **61** 401-425. [MR1424909](#)
- [52] YOUNES, L. (1988). Estimation and Annealing for Gibbsian Fields. *Ann. Inst. Henri Poincare* **24** 269–294. [MR0953120](#)
- [53] YOUNES, L. (1989). Parametric Inference for Imperfectly Observed Gibbsian Fields. *Probability Theory and Related Fields* **82** 625–645. [MR1002904](#)