

# Convergence of stochastic approximation algorithms under irregular conditions

Jian Zhang\*

*Institute of Mathematics, Statistics and Actuarial Science, University of Kent, Canterbury, Kent CT2 7NF, UK*

Faming Liang†

*Department of Statistics, Texas A&M University, College Station, TX 77843-3143, USA*

We consider a class of stochastic approximation (SA) algorithms for solving a system of estimating equations. The standard condition for the convergence of the SA algorithms is that the estimating functions are locally Lipschitz continuous. Here, we show that this condition can be relaxed to the extent that the estimating functions are bounded and continuous almost everywhere. As a consequence, the use of the SA algorithm can be extended to some problems with irregular estimating functions. Our theoretical results are illustrated by solving an estimation problem for exponential power mixture models.

*Keywords and Phrases:* stochastic approximation algorithm, M-estimator, exponential power mixture models.

## 1 Introduction

Suppose that we have  $p$ -dimensional observations  $(X_1, \dots, X_N)$  from a parametric or semiparametric model, which depends on an unknown parameter  $\theta \in \Theta \subset \mathbb{R}^d$ . Inference about  $\theta$ , e.g. finding its M-estimate and empirical likelihood confidence intervals, often involves a search for roots of the simultaneous equations

$$\frac{1}{N} \sum_{i=1}^N H_j(\theta, X_i) = 0, \quad j = 1, \dots, d, \quad (1)$$

for given estimating functions  $H_j(\cdot, \cdot): \mathbb{R}^d \times \mathbb{R}^p \mapsto \mathbb{R}$ ,  $j = 1, \dots, d$  (See HUBER 1981; NEWBY, 1993; TIAN *et al.*, 2004). As the partial derivative of an objective function for optimization,  $H = (H_j)_{1 \leq j \leq d}^T$  is usually non-additive with respect to  $X_i$ . When  $H$  is smooth and  $\partial H / \partial \theta$  is of full rank, the equations can be solved directly by the standard Newton–Raphson iteration. However, when  $H$  is not smooth or  $\partial H / \partial \theta$  is

---

\*J.Zhang@kent.ac.uk

†fliang@stat.tamu.edu

not of full rank, the stochastic approximation (SA) algorithm is found to be more appropriate for solving (1).

The subject of stochastic approximation was founded by ROBBINS and MONRO (1951). After five decades of continual development, this method has been widely used in adaptive control, system identification and optimization problems. See LAI (2003) for a review. The standard convergence analysis of the SA algorithm is based on the condition that  $H$  is locally Lipschitz continuous, although the weaker assumptions are found as well (TADIĆ, 1997; CHEN, 2002; KUSHNER and YIN, 2003). Among those conditions, KUSHNER (1981) required that  $E[H(x, Z_n)]$  or  $E[H(x, Z_n) | Z_{n-1}, \dots]$  is a smooth function of  $x$ . HE, FU and MARCUS (2003) required a convexity condition for the objective function. These conditions turn out to be either invalid or hard to verify. In this note, we show that for the convergence of the SA algorithm, the condition of local Lipschitz continuity can be relaxed to the extent that  $H$  is bound and continuous almost everywhere. As a consequence, the use of the SA algorithm can be extended to some problems with irregular estimating functions.

The remainder of this paper is organized as follows. In section 2, a general form for the SA algorithm with non-additive noise is introduced. In section 3, two sets of conditions are presented for proving the convergence of the SA algorithm. The application to estimating exponential power mixture models is described in section 4. A real data application is presented in section 5. The conclusions are made in section 6. Some basic lemmas are given in the last section.

## 2 A stochastic approximation algorithm with non-additive noise

Let  $H(\cdot, \cdot)$  be an  $\mathbb{R}^d \times \mathbb{R}^p \mapsto \mathbb{R}^d$  measurable function as defined before, and  $\mathcal{P}$  be a probability measure for a random variable  $Z$  drawn from  $\mathbb{R}^p$ . The problem is to seek the roots of the following system of equations,

$$h_i(\theta) = \int H_i(\theta, z) \mathcal{P}(dz) = 0, \quad i = 1, \dots, d. \quad (2)$$

In what follows, we denote  $H(\theta, Z) = (H_1(\theta, Z), \dots, H_d(\theta, Z))^{\tau}$  and  $h(\theta) = (h_1(\theta), \dots, h_d(\theta))^{\tau}$ . Assuming that solutions to (2) exist, we investigate when these roots can be found by the iterative algorithm:

$$\theta_{n+1} = \theta_n + \gamma_n H(\theta_n, Z_{n+1}), \quad (3)$$

where  $\theta_n$  denotes the estimate obtained at iteration  $n$ ,  $\gamma_n$  is the gain factor,  $H(\theta_n, Z_{n+1})$  is a non-additive function of  $Z_{n+1}$  and  $\{Z_n\}$  is a random process defined on  $(\Omega, \mathcal{F}, \mathcal{P})$  and is referred to as noise. Here,  $H(\theta_n, Z_{n+1})$  is non-additive if for any two deterministic functions  $H_1$  and  $H_2$ , it cannot be written in the form of

$$H(\theta_n, Z_{n+1}) = H_1(\theta_n) + H_2(Z_{n+1}).$$

As  $H(\theta_n, Z_{n+1})$  is a non-additive function of  $Z_{n+1}$ , algorithm (3) is referred to as a SA algorithm with non-additive noise. The SA algorithm circumvents the

differentiability requirement of the Newton–Raphson method by taking advantages of the bootstrap and the SA (EFRON AND TIBSHINARI, 1993).

### 3 A new convergence theorem for the SA algorithm

In some situations, the objective function (or the associated Lyapunov function)  $v(\theta)$  of problem (2) may not be differentiable but have the bounded directional partial derivatives at some points. Even when  $v(\theta)$  is differentiable, its derivative  $\nabla v(\theta)$  may not be Hölder continuous. The question is whether the SA algorithm converges under such irregular settings. The answer is partially positive as shown in the remainder of this section.

To facilitate presentation, we define a few more notations as follows. For any  $u_1 \in \mathbb{R}^d$ , letting

$$A = \left\{ \lim_{u_i \rightarrow u} \frac{\partial G(u_i)}{\partial u} : \frac{\partial G(u_i)}{\partial u} \text{ exists and is finite} \right\},$$

we define the generalized gradient  $\partial G(u) = \text{col } A$  as a convex closure of  $A$  (CLARK, 1983). When  $G(u)$  is differentiable at  $u$ , we simply set  $\partial G(u) = \{\partial G(u)/\partial u\}$  as a singleton. Similarly, the generalized directional derivative is defined as

$$G^0(u; v) = \limsup_{y \rightarrow u, t \rightarrow 0^+} \frac{G(y + tv) - G(y)}{t}.$$

$G$  is called regular if for all  $u$ : (i)  $G$  has the usual one-sided directional derivative; (ii) this derivative is equal to the corresponding generalized directional derivative.

Let  $Z_k$  be drawn from the finite set  $\{X_1, \dots, X_N\}$  randomly with replacement,  $\{Z_1, \dots, Z_n\}$  can be regarded as a random sample of  $P_N$ , where  $P_N$  is an empirical distribution defined on the set  $\{X_1, \dots, X_N\}$ . Hence,  $h(\theta)$  can be written as

$$h(\theta) = \frac{1}{N} \sum_{k=1}^N H(\theta, X_k). \tag{4}$$

In the following, we let  $\mathbb{R}^+$  denote the set of non-negative reals,  $\|\cdot\|$  denote the Euclidian norm in  $\mathbb{R}^d$  and the Frobenius norm in  $\mathbb{R}^{d \times d}$  and  $d(\cdot, \cdot)$  denote the distance induced by the Euclidian norm. Assume that there is an objective function  $v(\theta)$  with the set of stationary points  $E_* = \{\theta \in \Theta : 0 = h(\theta)\}$ . Set

$$\Theta_0 = \{\theta : h(\cdot) \text{ is not continuous at } \theta, \text{ or } v(\theta) \text{ is not differentiable at } \theta\}.$$

$$\Xi = \{\{\theta_n, n \geq 0\} : \text{there is a closed subset } D \text{ of } \Theta - \Theta_0 \text{ such that } \theta_n \in D, n \geq 0\}.$$

Let  $b_{n+1} = H(\theta_n, Z_{n+1}) - h(\theta_n)$  and  $\theta_0$  be any point in  $\Theta$ . Then the SA algorithm has the form

$$\theta_{n+1} = \theta_n + \gamma_n(h(\theta_n) + b_{n+1}).$$

We analyze the SA algorithm under the following conditions:

- B1.  $0 < \lim_{n \rightarrow \infty} \gamma_n = 0, \quad \sum_{n=1}^{\infty} \gamma_n = \infty, \quad \sum_{n=1}^{\infty} \gamma_n^2 < \infty, \quad \sup_{n \geq 0} E \|b_{n+1}\|^2 < +\infty.$
- B2.  $v(\theta)$  is continuous with finite generalized gradient  $\partial v(\theta)$ . For  $\theta \in \Theta, -h(\theta) \in \partial v(\theta)$ . For any compact subset  $Q \subset \Theta - E_*$ ,

$$\sup_Q -\|h(\theta)\| < 0.$$

THEOREM 1. Assume that the conditions B1–B2 hold. Then, as  $n \rightarrow \infty$ , we have

$$\lim_{n \rightarrow \infty} d(\theta_n, E_*) = 0 \text{ almost surely on } \Xi,$$

where  $d(\theta_n, E_*) = \inf_{\theta \in E_*} \|\theta_n - \theta\|$ .

Proof. Let  $t_0 = 0, t_k = \sum_{i=0}^{k-1} \gamma_i$  and  $m(t)$  be the  $k$  such that  $t \in [t_k, t_{k+1})$ . For  $t \in \mathbb{R}^+$ , define

$$\theta^k(t) = \theta_k + \sum_{i=k}^{m(t+t_k)-1} \gamma_i h(\theta_i) + \sum_{i=k}^{m(t+t_k)-1} \gamma_i b_{i+1}.$$

Then

$$\theta^k(t) = \theta_k + \int_0^t h(\theta^k(s)) ds + B^k(t) + \rho^k(t)$$

with

$$B^k(t) = \sum_{i=k}^{m(t+t_k)-1} \gamma_i b_i, \quad \rho^k(t) = \int_{t_{m(t+t_k)}}^t h(\theta^k(s)) ds.$$

By conditions B1 and B2, as  $k \rightarrow \infty$ ,

$$\sup_{0 \leq t < \infty} |\rho^k(t)| \leq \sup_{0 \leq t < \infty} \gamma_{m(t+t_k)} \|h(\theta_{m(t+t_k)})\| \rightarrow 0.$$

Note that  $\{\gamma_i b_{i+1}, \mathcal{F}_i : i \geq k\}$  is a martingale sequence and

$$\begin{aligned} \sup_{m \geq 1} E \left\| \sum_{i=1}^m \gamma_i b_{i+1} \right\| &\leq \sup_{m \geq 1} \left[ E \sum_{i=1}^m \gamma_i^2 b_{i+1}^2 \right]^{1/2} \\ &\leq \left[ \sum_{i=1}^m \gamma_i^2 \right]^{1/2} \sup_{n \geq 0} (E \|b_{n+1}\|^2)^{1/2} < \infty. \end{aligned}$$

It follows from Doob’s theorem (SHIRYAYEV, 1984, Theorem 1, p. 476) that

$$\lim_{k \rightarrow \infty} \sup_{m \geq k} \left\| \sum_{i=k}^m \gamma_i b_{i+1} \right\| = 0, \text{ almost surely,}$$

which yields

$$\sup_{0 \leq t < \infty} \|B^k(t)\| \leq \sup_{m \geq k} \left\| \sum_{i=k}^m \gamma_i b_{i+1} \right\| \rightarrow 0$$

almost surely as  $k \rightarrow \infty$ .

If we define  $V^k(t) = \int_0^t h(\theta^k(s)) ds$ , then

$$\sup_{0 \leq t < \infty} \|\theta^k(t) - \theta^k(0) - V^k(t)\| = o(1), \tag{5}$$

and

$$\|V^k(t) - V^k(s)\| \leq \sup_{t \leq u \leq s} \|h(\theta^k(u))\| |t - s|,$$

which implies that  $\theta^k(t)$  is also equicontinuous on any finite interval in  $\mathbb{R}^+$  provided  $\{\theta_n\}$  is contained in a compact subset  $D$  of  $\Theta$ . By the Arzelà–Ascoli theorem, there exist a subsequence, say  $\{k_j\}$ , and a  $\mathbb{R}^d$ -valued function  $\theta(t)$  on  $\mathbb{R}^+$  such that

$$\theta^{k_j}(\cdot) \rightarrow \theta(\cdot),$$

where  $\theta(\cdot)$  is also equicontinuous on any finite interval in  $\mathbb{R}^+$ .

By conditions B2,  $h(\cdot)$  is bounded on any closed subset of  $\Theta - \Theta_0$  and  $h(\theta^k(s)) \rightarrow h(\theta(s))$  when  $\theta(s) \notin \Theta_0$ . And the set  $\{0 \leq s \leq t : \theta(s) \notin \Theta_0\}$  has the Lebesgue measure tending to 0. Using the dominated convergence theorem, we have

$$V^k(t) = \int_0^t h(\theta^k(s)) ds \rightarrow \int_0^t h(\theta(s)) ds.$$

This together with (5) yields

$$\theta(t) = \theta(0) + \int_0^t h(\theta(s)) ds.$$

Thus, the limit  $\theta(s)$  of any convergent subsequence of  $\{\theta_n : n \geq 0\}$  in  $\Xi$  satisfies the differential equation:

$$\dot{\theta}(s) = h(\theta(s)), \quad s \in \mathbb{R}^+.$$

Invoking the chain rule of derivatives and condition B2, we have

$$\dot{v}(\theta(s)) = -h(\theta(s))^{\tau} \dot{\theta}(s) = -\|h(\theta(s))\|^2.$$

Therefore,

$$v(\theta(t)) - v(\theta(0)) = -\int_0^t \|h(\theta(s))\|^2 ds.$$

Note that  $v(\theta)$  is continuous. On  $\Xi$ , for any convergent subsequence  $\{\theta(t_j) : j = 1, 2, \dots\}$  of  $\{\theta(t), t \in \mathbb{R}^+\}$ , with  $\theta(t_j) \rightarrow \theta_*$ , we have that  $v(\theta(t_j)) \rightarrow v(\theta_*)$  and

$$v(\theta_*) - v(\theta(0)) = - \int_0^\infty \|h(\theta(s))\|^2 ds > -\infty,$$

which implies that for any  $\delta > 0$ , as  $t_j \rightarrow \infty$ ,

$$\int_{t_j - \delta/2}^{t_j + \delta/2} \|h(\theta(s))\|^2 ds \rightarrow 0. \tag{6}$$

Note that  $h(\theta)$  is continuous at  $\theta_*$  and  $\theta(t)$  is equicontinuous. Thus, for any  $\epsilon > 0$ , there exists a  $\delta > 0$  such that

$$\|h(\theta(t)) - h(\theta_*)\| < \epsilon,$$

when  $|t - t_j| < \delta$ . This together with (6) yields that

$$(\|h(\theta_*)\| - \epsilon)\delta \leq 0.$$

Therefore  $h(\theta_*) = 0$ , meaning that  $\theta_* \in \{\theta : 0 = h(\theta)\}$ . Now using the same technique in CHEN (2002, p. 15), we can easily prove the desired result.  $\square$

#### 4 Exponential power mixture model estimation

Suppose that  $(X_1, \dots, X_N)$  is a sample from a mixture of exponential power distributions,

$$f(x | \theta) = \sum_{k=1}^m \omega_k \varphi(x | \mu_k, \Sigma_k, \beta_k), \tag{7}$$

where  $\theta = (\mu_1, \Sigma_1, \beta_1; \dots; \mu_m, \Sigma_m, \beta_m; \omega_1, \dots, \omega_{m-1})$  contains all parameters of the model;  $m$  is the total number of components;  $\omega_k$  is the weight of the  $k$ th component,  $0 < \omega_k \leq 1$ , and  $\sum_{k=1}^m \omega_k = 1$ ; and  $\varphi(x | \mu_k, \Sigma_k, \beta_k)$  is the Kotz-type distribution, which has the density

$$\begin{aligned} \varphi(x | \mu_k, \Sigma_k, \beta_k) &= \frac{\beta_k \Gamma(p/2)}{2\pi^{p/2} |\Sigma|^{1/2} \Gamma(p/\beta_k)} \exp\{-[(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)]^{\beta_k/2}\}, \\ \mu_k &\in \mathcal{R}^p, \quad \Sigma_k \in \mathcal{S}, \quad |\beta_k| \leq c_0, \end{aligned}$$

where  $\mathcal{S}$  is the set of all  $p \times p$  positive definite matrices and  $c_0$  is a sufficiently large constant (see FANG, KOTZ and NG, 1990, p. 76). The location parameter  $\mu_k$  stands for the centre of the distribution, the matrix  $\Sigma_k$  for the dispersion and the shape parameter  $\beta_k$  for the rate of exponential decay. For  $\beta_k = 2$ , the distribution is  $N(\mu_k, \Sigma_k/2)$ , a multivariate normal distribution; for  $0 < \beta_k < 1$ , the distributed is heavy tailed; and for  $\beta_k > 1$ , the distribution is light tailed. To reduce the technical burden, we focus on the case where  $\Sigma_k = \lambda_k^2 I_p$  with  $\lambda_k > 0$  and  $I_p$  being a  $p \times p$  unit matrix. Then

$$\varphi(x | \mu_k, \Sigma_k, \beta_k) = \frac{\beta_k \Gamma(p/2)}{2\pi^{p/2} \lambda_k^p \Gamma(p/\beta_k)} \exp\left\{-\left(\frac{\|x - \mu_k\|}{\lambda_k}\right)^{\beta_k}\right\}.$$

Here  $\lambda_k$  is used to show the volume of the  $k$ th component and  $\beta_k$  is employed to describe the shape of the  $k$ th component. We re-parameterize them by  $\lambda_k^* = \log(\lambda_k)$  and  $\beta_k^* = \log(\beta_k)$ . In addition, we re-parameterize  $\omega_1, \dots, \omega_{m-1}$  by  $\omega_1^*, \dots, \omega_{m-1}^*$  by setting  $\omega_m^* = 0$  and

$$\omega_k = \frac{e^{\omega_k^*}}{\sum_{j=1}^m e^{\omega_j^*}}, \quad k = 1, \dots, m - 1.$$

For uniformity of notations, we define  $\mu_k^* \equiv \mu_i$ . With the above notations, model (7) can be reparametrized by  $\theta = (\mu_1^*, \lambda_1^*, \beta_1^*, \dots; \mu_m^*, \lambda_m^*, \beta_m^*, \omega_1^*, \dots, \omega_{m-1}^*)$ . In what follows, for simplicity, we still use  $f(x | \theta)$  to denote the mixture density with  $\theta \in \mathbb{R}^d$  and  $d = (p + 3)m - 1$ . Note that when  $\beta_k^* < 0$ ,  $\log f(z | \theta)$  is not differential at  $\mu_k = X_k$  and  $\|\partial \log f(X_k | \theta) / \partial \mu_k\|$  tends to infinity as  $\mu_k \rightarrow X_k$ . Removing these singular points from  $\mathbb{R}^d$ , we define

$$\begin{aligned} R_z &= \{\theta \in \mathbb{R}^d : \text{there exist } i, k \text{ such that } \mu_i \in \{X_1, \dots, X_N\} \text{ and } \beta_k^* < 0\}, \\ \Theta &= \mathbb{R}^d \setminus R_z. \end{aligned}$$

If there exists  $k$  such that  $\beta_k^* = 0$ , then the corresponding  $\log(f(X_k | \theta))$  is not differentiable at  $\mu_k = X_k$  but has a generalized gradient. For other  $\theta \in \Theta$ , the corresponding  $\log f(x | \theta)$  is continuously differentiable.

For a given value of  $m$ , model (7) can be estimated by maximizing the function

$$l_N(\theta) = \sum_{k=1}^N \log f(X_k | \theta),$$

which is asymptotically equivalent to minimize the Kullback–Leibler distance between  $f(\cdot | \theta)$  and the underlying density of  $X$ . Differentiating  $l_N(\theta)$  with respect to  $\theta$ , we have

$$\frac{\partial l_N(\theta)}{\partial \mu_k^*} = \sum_{i=1}^N P(k | X_i) \frac{\beta_k}{\lambda_k^{\beta_k}} \| |X_i - \mu_k| |^{\beta_k - 2} (X_i - \mu_k), \tag{8}$$

$$\frac{\partial l_N(\theta)}{\partial \lambda_k^*} = \sum_{i=1}^N P(k | X_i) \left[ p - \frac{\beta_k}{\lambda_k^{\beta_k}} \| |X_i - \mu_k| |^{\beta_k} \right], \tag{9}$$

$$\begin{aligned} \frac{\partial l_N(\theta)}{\partial \beta_k^*} &= \sum_{i=1}^N P(k | X_i) \\ &\times \left\{ 1 + \frac{p}{\beta_k} \frac{\Gamma'(p/\beta_k)}{\Gamma(p/\beta_k)} - \beta_k \left( \frac{\| |X_i - \mu_k| |}{\lambda_k} \right)^{\beta_k} \log(\| |X_i - \mu_k| |) \right\}, \end{aligned} \tag{10}$$

$$\frac{\partial l_N(\theta)}{\partial \omega_k^*} = \sum_{i=1}^N [P(k | X_i) - \omega_k]. \tag{11}$$

where the index  $k$  ranges from 1 to  $m$  in (8)–(10) and ranges from 1 to  $m - 1$  in (11);

$$P(k | x) = \omega_k \varphi(x | \mu_k, \lambda_k^2 I_p, \beta_k) / f(x | \theta)$$

is the posterior probability that  $x$  is from the  $k$ -component; and

$$\Gamma'(x) / \Gamma(x) = -1/x - v + \sum_{k=1}^{\infty} [1/k - 1/(k+x)].$$

Here,  $v \approx 0.577216$  is called Euler’s constant. Note that if  $\beta_k^* = 0$  (i.e.  $\beta_k = 1$ ) and  $\mu_k^* = X_j, \mu_k^* \neq X_i, i \neq j$ , then  $\partial l_N(\theta) / \partial \mu_k^*$  in (8) does not exist. However, under this setting,  $l_N(\theta)$  is regular and has the generalized gradient

$$\partial l_N(\theta) = \sum_{i \neq j}^N P(k | X_i) \frac{1}{\lambda_k} \frac{(X_i - \mu_k)}{\|X_i - \mu_k\|} + P(k | X_j) \frac{1}{\lambda_k} S_p,$$

where  $S_p = \{a \in \mathbb{R}^p : \|a\| \leq 1\}$ . Define

$$H_{\mu_k^*}(\theta, x) = -P(k | x) \frac{\beta_k}{\lambda_k^{\beta_k}} \|x - \mu_k\|^{\beta_k - 2} (x - \mu_k), \quad k = 1, \dots, m, \tag{12}$$

$$H_{\lambda_k^*}(\theta, x) = -P(k | x) \left[ p - \frac{\beta_k}{\lambda_k^{\beta_k}} \|x - \mu_k\|^{\beta_k} \right], \quad k = 1, \dots, m, \tag{13}$$

$$H_{\beta_k^*}(\theta, x) = -P(k | x) \times \left\{ 1 + \frac{p \Gamma'(p/\beta_k)}{\beta_k \Gamma(p/\beta_k)} - \beta_k \left( \frac{\|x - \mu_k\|}{\lambda_k} \right)^{\beta_k} \log(\|x - \mu_k\|) \right\}, \quad k = 1, \dots, m \tag{14}$$

$$H_{\omega_k^*}(\theta, x) = P(i | x) - \omega_k, \quad k = 1, \dots, m - 1, \tag{15}$$

where we set

$$\frac{(x - \mu_k)}{\|x - \mu_k\|} = a \quad \text{with } \|a\| = 1 \quad \text{when } x = \mu_k.$$

Applying the stochastic approximation method defined in section 2 to solve the system of equations,

$$\begin{cases} \int H_{\mu_k^*}(\theta, z) dF_N(z) = 0, & k = 1, \dots, m, \\ \int H_{\lambda_k^*}(\theta, z) dF_N(z) = 0, & k = 1, \dots, m, \\ \int H_{\beta_k^*}(\theta, z) dF_N(z) = 0, & k = 1, \dots, m, \\ \int H_{\omega_k^*}(\theta, z) dF_N(z) = 0, & k = 1, \dots, m - 1, \end{cases} \tag{16}$$

where  $F_N$  is the empirical distribution introduced by  $(X_1, \dots, X_N)$ . Let

$$\theta_t = (\mu_1^{*(t)}, \lambda_1^{*(t)}, \beta_1^{*(t)}, \dots, \mu_m^{*(t)}, \lambda_m^{*(t)}, \beta_m^{*(t)}; \omega_1^{*(t)}, \dots, \omega_{m-1}^{*(t)})$$



denote the estimate of  $\theta$  obtained at iteration  $t$  and  $Z_t$  denotes a sample randomly drawn from the set  $\{X_1, \dots, X_N\}$  at iteration  $t$ . In practice, we set the gain factor sequence

$$\gamma_t = \frac{\gamma_0 t_0}{\max(t_0, t)}, \quad t = 1, 2, \dots, \tag{17}$$

for some values of  $t_0 > 1$  and  $\gamma_0 > 0$ . Our default setting for them is  $t_0 = 10,000$  and  $\gamma_0 = 0.02$ . If other values were used, the values would be specified in the context. We also set a default value for the total number of iterations. It is  $500 \times N$ , where  $N$  is the sample size. See LIANG and ZHANG (2008) for details.

To show the convergence of the algorithm, we first verify conditions B1–B2 under the restriction  $\beta_{*k} > 0, 1 \leq k \leq m$ . This can be performed as follows.

By construction of the gain factor sequence, condition B1 is satisfied. Define  $v(\theta) = -I_N(\theta)/N$ . Clearly,  $v(\theta)$  is bounded below, because the exponential power density function and thus its mixture are bounded above. According to the order of the components of  $\theta$ , we define  $H_j(\theta, x)$  to be the respective functions  $H_{\mu_k^*}(\theta, x), H_{\alpha_k^*}(\theta, x), H_{\beta_k^*}(\theta, x)$  or  $H_{\omega_k^*}(\theta, x)$  defined in (16). By definition,

$$\nabla v(\theta) = - \int H(\theta, z) d\mathcal{P}(z) = -h(\theta) \quad \text{and} \quad \dot{v}(\theta) = - \|h(\theta)\|^2 \leq 0.$$

When  $\beta_k^* \geq 0, 1 \leq k \leq m$ , condition B<sub>2</sub> is satisfied by defining  $E_* = \{\theta : h(\theta) = \mathbf{0}\}$ .

### 5 Applications

In microarray data analysis, one often needs to test a large number of hypotheses simultaneously. To control the probability of erroneously rejecting true null hypotheses, we need to control the false discovery rate (FDR), the expected proportion of false-positive findings among all the rejected hypotheses. The main difficulty in calculating the FDR comes from modelling the observed test scores (EFRON, 2004).

To set up notations, let  $H_1, \dots, H_N$  denote the collection of  $N$  null hypotheses,  $P_1, \dots, P_N$  denote the corresponding  $P$ -values of the  $N$  tests, and  $Z_i = \Phi^{-1}(P_i)$  or  $Z_i = \Phi^{-1}(1 - P_i)$  denote the corresponding test scores, where  $\Phi$  is the cumulative distribution function (CDF) of the standard normal distribution. LIANG and ZHANG (2008) proposed the following model for these test scores:

$$f(z) = \pi_0 f_0(z) + (1 - \pi_0) f_1(z), \tag{18}$$

where  $\pi_0$  is the prior probability that a null hypothesis is true,  $f_0$  is the empirical null distribution and  $f_1$  is the alternative one.  $f_0$  and  $f_1$  are modelled by two exponential power mixture models. The unknown parameters in  $f_0$  and  $f_1$  are then estimated by the SA algorithm introduced in the previous section.

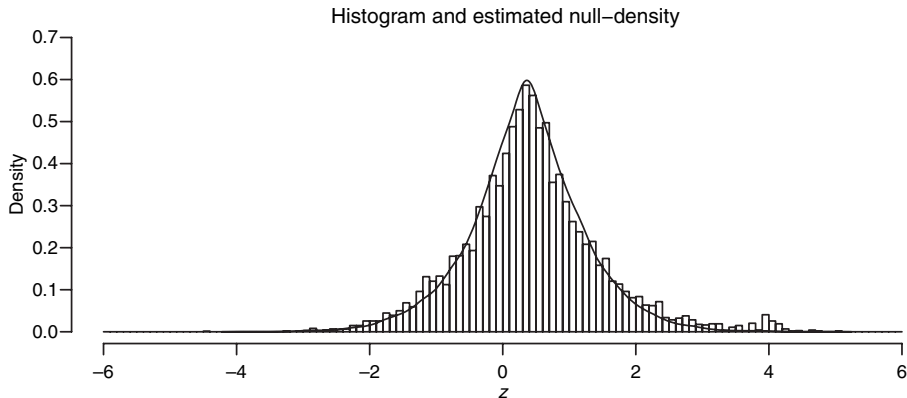


Fig. 1.  $z$ -scores of avian pineal gland gene expression data under the DD condition. The solid line is the estimated  $f_0$  by the SA.

### 5.1 A real data analysis

The avian pineal gland contains both circadian oscillators and photoreceptors to produce rhythms in biosynthesis of the hormone melatonin *in vivo* and *in vitro*. It is of great interest to understand the genetic mechanism driving the rhythms. For this purpose, Dr V. Cassone's laboratory at Texas A&M University measured the expression levels of pineal gland genes under light–dark (LD) and constant darkness (DD) conditions. Under LD, the birds were killed at 2, 6, 10, 14, 18, 22 hours Zeitgeber time (ZT) to obtain mRNA to produce adequate cDNA libraries. Four microarray chips per time point were produced, and there were two replicates for each gene in each chip. The experiment was then repeated under DD. Each chip produced gene expressions for 7400+ genes. Throughout the experiment, samples from LD ZT18 were used as controls. Relative gene expression levels to the controls were recorded and processed. Our goal was to identify genes that are differentially expressed at different time points. Mixed effect analysis with the fixed effect being the different time points and the random effects corresponding to chips and biological batches were applied to the relative gene expression levels in log scale. Normalization procedures have been adopted but will not be listed here as they are not the focus of the paper. Under both LD and DD conditions,  $P$ -values,  $P_i$ s, for testing the existence of different time effects were produced and transformed to test scores using  $\Phi^{-1}(1 - P_i)$ . Figure 1 shows the histogram of the test scores under the DD condition and the corresponding estimated density curve of  $f_0$ . It suggests that  $f_0$  can be well estimated by the SA method.

## 6 Discussion

In this paper, we have introduced a stochastic approximation method for finding the roots of multiple estimating equations. One important problem is whether the

proposed algorithm converges. The classical convergence analysis is often based on the regular condition that the estimating functions are locally Lipschitz continuous. Here, we have shown that these conditions can be relaxed to the extent that the estimating functions are bounded and continuous almost everywhere.

The SA algorithm has been found a number of applications such as model-based clustering and semiparametric estimating. Here, we applied the SA algorithm to estimate the null density of the test scores in multiple testing problem. A real data application has been presented.

### Acknowledgements

This paper is devoted to EURANDOM for celebrating its 10th anniversary. The first author would like to thank Professor Alessandro Di Bucchianico for his kind invitation to the jubilee conference.

### References

- CHEN, H. F. (2002), *Stochastic approximation and its applications*, Kluwer Academic Publishers, London.
- CLARK, F. (1983), *Optimization and nonsmooth analysis*, Wiley, New York.
- EFRON, B. (2004), Large-scale simultaneous hypothesis testing: the choice of a null hypothesis, *Journal of the American Statistical Association* **99**, 96–104.
- EFRON, B. and R. J. TIBSHIRANI (1993), *An introduction to the bootstrap*, Chapman & Hall, New York.
- FANG, K., S. KOTZ and K. NG (1990), *Symmetric multivariate and related distributions*, Chapman & Hall, New York.
- HE, Y., M. C. FU and S. I. MARCUS (2003), Convergence of simultaneous perturbation stochastic approximation for nondifferentiable optimization, *IEEE Transactions on Automatic Control* **48**, 1459–1463.
- HUBER, P. (1981), *Robust statistics*, Wiley, New York.
- KUSHNER, H. J. (1981), Stochastic approximation with discontinuous dynamics and state dependent noise: w.p.1 and weak convergence, *Journal of Mathematical Analysis and Applications* **82**, 527–542.
- KUSHNER, H. J. and G. YIN (2003), *Stochastic approximation and recursive algorithms and applications*, Springer, New York.
- LAI, T. L. (2003), Stochastic approximation, *The Annals of Statistics* **31**, 391–406.
- LIANG, F. and J. ZHANG (2008), Estimating FDR using the stochastic approximation algorithm, *Biometrika*, in press.
- NEWBY, W. K. (1993), Efficient estimation of models with conditional moment restrictions, in: G. S. MADDALA, C. R. RAO and H. D. VINOD (eds), *Econometrics, Handbook of Statistics 11*, North-Holland, Amsterdam, pp. 419–453.
- ROBBINS, H. and S. MONRO (1951), A stochastic approximation method, *Annals of Mathematical Statistics* **22**, 400–407.
- SHIRYAYEV, A. N. (1984), *Probability*, Springer, New York.
- TADIĆ, V. (1997), On the convergence of stochastic iterative algorithms and their applications to machine learning. A short version of this paper was published in *Proceedings of the 36th IEEE Conference on Decision & Control*, San Diego, CA, USA, pp. 2281–2286.
- TIAN, L., J. S. LIU, Y. ZHAO and L. J. WEI (2004), Statistical inferences based on non-smooth estimating functions, *Biometrika* **91**, 943–954.

Received: March 2008, Revised: April 2008.