

# Adaptive Stochastic Approximation by the Simultaneous Perturbation Method

Nifei Lin

October 2021

- SPSA
- 2SPSA
  - Strong Convergence
  - Asymptotic Normality

- Fabian, V. (1968). On asymptotic normality in stochastic approximation. *The Annals of Mathematical Statistics*, 1327-1332.
- Spall, J. C. (1992). Multivariate stochastic approximation using a simultaneous perturbation gradient approximation. *IEEE transactions on automatic control*, 37(3), 332-341.
- Spall, J. C. (1997, December). Accelerated second-order stochastic optimization using only function measurements. In *Proceedings of the 36th IEEE Conference on Decision and Control (Vol. 2, pp. 1417-1424)*. IEEE.
- Spall, J. C. (1998). Implementation of the simultaneous perturbation algorithm for stochastic optimization. *IEEE Transactions on aerospace and electronic systems*, 34(3), 817-823.

Spall, J. C. (1998). An overview of the simultaneous perturbation method for efficient optimization. Johns Hopkins apl technical digest, 19(4), 482-492.

Spall, J. C., & Cristion, J. A. (1998). Model-free control of nonlinear stochastic systems with discrete-time measurements. IEEE transactions on automatic control, 43(9), 1198-1210.

Spall, J. C. (2000). Adaptive stochastic approximation by the simultaneous perturbation method. IEEE transactions on automatic control, 45(10), 1839-1853.

- Problem: finding a root  $\theta^*$  of the gradient equation

$$g(\theta) \equiv \frac{\partial L(\theta)}{\partial \theta} = 0$$

- SA standard form

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k \hat{g}(\hat{\theta}_k)$$

- The central FD estimator of  $\hat{g}$  is

$$\hat{g}(\hat{\theta}_k) = \frac{1}{2c} \begin{pmatrix} y(\hat{\theta}_k + c\mathbf{e}_1) - y(\hat{\theta}_k - c\mathbf{e}_1) \\ y(\hat{\theta}_k + c\mathbf{e}_2) - y(\hat{\theta}_k - c\mathbf{e}_2) \\ \vdots \\ y(\hat{\theta}_k + c\mathbf{e}_d) - y(\hat{\theta}_k - c\mathbf{e}_d) \end{pmatrix}$$

Let  $e_i$  denote the  $i$ th column of a  $d \times d$  identity matrix.

- Let  $\Delta_k \in R^p$  be a vector of  $p$  mutually independent mean-zero random variables  $\{\Delta_{k1}, \Delta_{k2}, \dots, \Delta_{kp}\}$
- Let  $\{\Delta_k\}$  be a mutually independent sequence with  $\Delta_k$  independent of  $\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k$ .
- We have available noisy measurements of  $L(\cdot)$ :

$$y_k^{(+)} = L(\hat{\theta}_k + c_k \Delta_k) + \epsilon_k^{(+)}$$

$$y_k^{(-)} = L(\hat{\theta}_k - c_k \Delta_k) + \epsilon_k^{(-)}$$

where  $\epsilon_k^{(+)}, \epsilon_k^{(-)}$  represent measurement noise terms that satisfy

$$E(\epsilon_k^{(+)} - \epsilon_k^{(-)} | \mathcal{F}_k, \Delta_k) = 0 \text{ a.s. } \forall k, \mathcal{F}_k \equiv \{\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k\}$$

- SPSA estimator of  $g(\cdot)$  at the  $k$ th iteration is

$$\hat{g}_k(\hat{\theta}_k) = \begin{bmatrix} \frac{y_k^{(+)} - y_k^{(-)}}{2c_k \Delta_{k1}} \\ \vdots \\ \frac{y_k^{(+)} - y_k^{(-)}}{2c_k \Delta_{kp}} \end{bmatrix}$$

## Lemma [Spall, 1992]

Consider all  $k \geq K$  for some  $K < \infty$ . Suppose that for each such  $k$  the  $\{\Delta_{ki}\}$  are i.i.d. ( $i = 1, 2, \dots, p$ ) and symmetrically distributed about 0 with  $|\Delta_{ki}| \leq \alpha_0$  a.s. and  $E|\Delta_{ki}^{-1}| \leq \alpha_1$ . For almost all  $\hat{\theta}_k$  (at each  $k \geq K$ ) suppose that  $\forall \theta$  in an open neighborhood of  $\hat{\theta}_k$  that is not a function of  $k$  or  $\omega$ ,  $L^{(3)}(\theta) \equiv \partial^3 L / \partial \theta^T \partial \theta^T \partial \theta^T$  exists continuously with individual elements satisfying  $\left| L_{i_1 i_2 i_3}^{(3)}(\theta) \right| \leq \alpha_2$ . Then for almost all  $\omega \in \Omega$

$$\begin{aligned} b_k(\hat{\theta}_k) &\equiv E\left(\hat{g}_k(\hat{\theta}_k) - g(\hat{\theta}_k) \mid \hat{\theta}_k\right) \\ &= E\left(\hat{g}_k(\hat{\theta}_k) - g(\hat{\theta}_k) \mid \mathcal{F}_k\right) \\ &= O(c_k^2) \quad (c_k \rightarrow 0) \end{aligned}$$



**Proof:** Consider any  $l \in \{1, 2, \dots, p\}$  (let  $\bar{\Delta}_{kl} = c_k \Delta_k$ )

First, note that  $E \left[ (\epsilon_k^{(+)} - \epsilon_k^{(-)}) / 2\bar{\Delta}_{kl} \mid \hat{\theta}_k \right] = 0$  a.s.

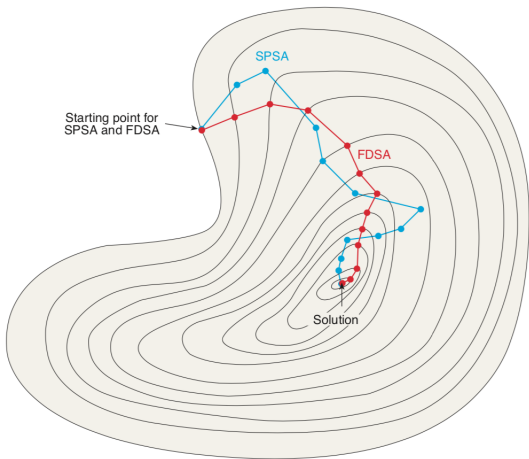
Then by the continuity of  $L^{(3)}$  near  $\hat{\theta}_k$  and uniform boundedness of  $|\Delta_{ki}|$  for all  $k$  sufficiently large, we have by Taylor's theorem for all such  $k$

$$b_{kl}(\hat{\theta}_k) = \frac{1}{12} E \left\{ \bar{\Delta}_{kl}^{-1} \left[ L^{(3)}(\bar{\theta}_k^+) + L^{(3)}(\bar{\theta}_k^-) \right] \bar{\Delta}_k \otimes \bar{\Delta}_k \otimes \bar{\Delta}_k \mid \hat{\theta}_k \right\}$$

where  $\bar{\theta}_k^+, \bar{\theta}_k^-$  are on the line segment between  $\hat{\theta}_k$  and  $\hat{\theta}_k \pm \bar{\Delta}_k$ , respectively, and  $b_{kl}$  denotes the  $l$ th term of the bias  $b_k$ .

By the mean value theorem, the term on the r.h.s. , is bounded in magnitude by

$$\begin{aligned} & \frac{\alpha_2 c_k^2}{6} \sum_{i_1} \sum_{i_2} \sum_{i_3} E \left| \frac{\Delta_{ki_1} \Delta_{ki_2} \Delta_{ki_3}}{\Delta_{kl}} \right| \\ & \leq \frac{\alpha_2 c_k^2}{6} \cdot \left\{ [p^3 - (p-1)^3] \alpha_0^2 + (p-1)^3 \alpha_1 \alpha_0^3 \right\} \end{aligned}$$



**Figure 5.** Example of relative search paths for SPSA and FDSA in  $p = 2$  problem. Deviations of SPSA from FDSA average out in reaching a solution in the same number of iterations; FDSA nearly follows the gradient descent path (perpendicular to level curves) in the low-noise setting.

Basic form of algorithm (composed of two parallel recursions: one for  $\theta$  and one for the Hessian of  $L(\theta)$ )

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k \bar{H}_k^{-1} G_k(\hat{\theta}_k), \quad \bar{H}_k = f_k(\bar{H}_k)$$

$$\bar{H}_k = \frac{k}{k+1} \bar{H}_{k-1} + \frac{1}{k+1} \hat{H}_k, \quad k = 0, 1, 2, \dots$$

- a stochastic analog of the well-known Newton-Raphson algorithm of deterministic search and optimization.
- recursive calculation of the sample mean of the per-iteration Hessian estimates

Notations:

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k \bar{\bar{H}}_k^{-1} G_k(\hat{\theta}_k), \quad \bar{\bar{H}}_k = f_k(\bar{H}_k)$$

$$\bar{H}_k = \frac{k}{k+1} \bar{H}_{k-1} + \frac{1}{k+1} \hat{H}_k, \quad k = 0, 1, 2, \dots$$

$a_k$ : a nonnegative scalar gain coefficient

$G_k(\hat{\theta}_k)$ : the input information related to  $g(\hat{\theta}_k)$  (i.e., the gradient approximation from  $y(\cdot)$  measurements in the gradient-free case or the direct observation as in the Robbins-Monro gradient-based case)

$f_k: \mathbf{R}^{p \times p} \rightarrow \{ \text{positive definite } p \times p \text{ matrices} \}$ : a mapping designed to cope with possible nonpositive definiteness of  $\bar{H}_k$

$\hat{H}_k$ : a per-iteration estimate of the Hessian

The formula for estimating the Hessian at each iteration is:

$$\hat{H}_k = \frac{1}{2} \left[ \frac{\delta G_k^T}{2c_k \Delta_k} + \left( \frac{\delta G_k^T}{2c_k \Delta_k} \right)^T \right]$$

where

$$\delta G_k = G_k^{(1)}(\hat{\theta}_k + c_k \Delta_k) - G_k^{(1)}(\hat{\theta}_k - c_k \Delta_k)$$

- for 2SG, usually  $G_k^{(1)}(\cdot) = G_k(\cdot)$ . We will suppose that  $G_k(\cdot) = G_{k^k}^{(1)}(\cdot)$  is an unbiased direct measurement of  $g(\cdot)$  (i.e.,  $G_k(\cdot) = G_k^{(1)}(\cdot) = g(\cdot) + \text{mean-zero noise}$ ).
- for 2 SPSA:  $G_k^{(1)}$  is a one-sided gradient approximation (in order to reduce the total number of function evaluations versus the two-sided form usually recommended for  $G_k(\cdot)$ )

$$G_k^{(1)}(\hat{\theta}_k \pm c_k \Delta_k) = \frac{y(\hat{\theta}_k \pm c_k \Delta_k + \tilde{c}_k \tilde{\Delta}_k) - y(\hat{\theta}_k \pm c_k \Delta_k)}{\tilde{c}_k} \begin{bmatrix} \tilde{\Delta}_{k1}^{-1} \\ \tilde{\Delta}_{k2}^{-1} \\ \vdots \\ \tilde{\Delta}_{kp}^{-1} \end{bmatrix}$$

with  $\tilde{\Delta}_k = (\tilde{\Delta}_{k1}, \tilde{\Delta}_{k2}, \dots, \tilde{\Delta}_{kp})^T$  generated in the same statistical manner as  $\Delta_k$ , but independently of  $\Delta_k$  and  $\tilde{c}_k$  satisfying conditions similar to  $c_k$

Two specific implementations of the ASP approach :

- 2SPSA (second-order SPSA) for applications in the gradient-free case (four function measurements  $y(\cdot)$  are needed at each iteration)
- 2SG (second-order stochastic gradient) for applications in the Robbins-Monro gradient-based case. (three gradient measurements  $g(\cdot)$  are needed at each iteration)

## Theorem 1a

Consider the SPSA estimate for  $G(\cdot)$  with  $G(\cdot)^{(1)}$ . Let conditions C.0-C.7 hold. Then  $\hat{\theta}_k - \theta^* \rightarrow 0$  a.s.



## Conditions

- C.0:  $E \left( \varepsilon_k^{(+)} - \varepsilon_k^{(-)} \mid \hat{\theta}_k; \Delta_k; \bar{H}_k \right) = 0$  a.s.  $\forall k$ , where  $\varepsilon_k^{(\pm)}$  is the effective SA measurement noise, i.e.,

$$\varepsilon_k^{(\pm)n} \equiv y \left( \hat{\theta}_k \pm c_k \Delta_k \right) - L \left( \hat{\theta}_k \pm c_k \Delta_k \right)$$

- C.1:  $a_k, c_k > 0 \forall k$ ;  $a_k \rightarrow 0, c_k \rightarrow 0$  as  $k \rightarrow \infty$ ;  $\sum_{k=0}^{\infty} a_k = \infty$ ,  $\sum_{k=0}^{\infty} (a_k/c_k)^2 < \infty$
- C.2: For some  $\delta, \rho > 0$  and  $\forall k, \ell$ ,  $E \left( \left| y \left( \hat{\theta}_k \pm c_k \Delta_k \right) / \Delta_{k\ell} \right|^{2+\delta} \right) \leq \rho$ ,  $|\Delta_{k\ell}| \leq \rho$ ,  $\Delta_{k\ell}$  is symmetrically distributed about 0, and  $\{\Delta_{k\ell}\}$  are mutually independent.

- C.3: For some  $\rho > 0$  and almost all  $\hat{\theta}_k$ , the function  $g(\cdot)$  is continuously twice differentiable with a uniformly (in  $k$ ) bounded second derivative for all  $\theta$  such that  $\|\hat{\theta}_k - \theta\| \leq \rho$
- C.4: For each  $k \geq 1$  and all  $\theta$ , there exists a  $\rho > 0$  not dependent on  $k$  and  $\theta$ , such that  $(\theta - \theta^*)^T \bar{g}_k(\theta) \geq \rho \|\theta - \theta^*\|$ .
- C.5: For each  $i = 1, 2, \dots, p$  and any  $\rho > 0$ ,  $P\left(\left\{\bar{g}_{ki}(\hat{\theta}_k) \geq 0 \text{ i.o.}\right\} \cap \left\{\bar{g}_{ki}(\hat{\theta}_k) < 0 \text{ i.o.}\right\} \mid \left\{\left|\hat{\theta}_{ki} - (\theta^*)_i\right| \geq \rho \quad \forall k\right\}\right) = 0$

- C.6:  $\overline{\overline{H}}_k^{-1}$  exists a.s.  $\forall k, c_k^2 \overline{\overline{H}}_k^{-1} \rightarrow 0$  a.s., and for some  $\delta, \rho > 0$ ,  $E(\|\overline{\overline{H}}_k^{-1}\|^{2+\delta}) \leq \rho$
- C.7: For any  $\tau > 0$  and nonempty  $S \subseteq \{1, 2, \dots, p\}$ , there exists a  $\rho'(\tau, S) > \tau$  such that

$$\limsup_{k \rightarrow \infty} \left| \frac{\sum_{i \notin S} (\theta - \theta^*)_i \bar{g}_{ki}(\theta)}{\sum_{i \in S} (\theta - \theta^*)_i \bar{g}_{ki}(\theta)} \right| < 1 \quad \text{a.s.}$$

for all  $|(\theta - \theta^*)_i| < \tau$  when  $i \notin S$  and  $|(\theta - \theta^*)_i| \geq \rho'(\tau, S)$  when  $i \in S$

We define  $\bar{g}(\hat{\theta}_k) = \overline{\overline{H}}^{-1} g(\hat{\theta}_k)$

**Proof** The proof will proceed in 3 parts.

- 1.  $\tilde{\theta}_k \equiv \hat{\theta}_k - \theta^*$  does not diverge in magnitude to  $\infty$
- 2.  $\tilde{\theta}_k$  converges a.s. to some random vector;
- 3. this random vector is the constant 0

**Part 1:** Letting  $M_j = a_j \bar{H}_j^{-1} E(G_j(\hat{\theta}_j) | \hat{\theta}_j) = a_j \bar{H}_j^{-1} (g_j(\hat{\theta}_j) + b_j)$  and  $M'_j = a_j \bar{H}_j^{-1} (\hat{g}_j(\hat{\theta}_j) - \bar{g}_j(\hat{\theta}_j))$ , we can write

$$\tilde{\theta}_{k+1} + \sum_{j=0}^k M_j = \tilde{\theta}_0 - \sum_{j=0}^k M'_j$$

$\left\{ \sum_{j=1}^k M'_j \right\}$  is a martingale sequence (in  $k$ )

$$E \left\| \sum_{j=0}^k M'_j \right\|^2 \leq 2 \sum_{j=0}^k E \|M'_j\|^2 < \infty$$

Then by the martingale convergence theorem

$$\tilde{\theta}_{k+1} + \sum_{j=0}^k M_j \xrightarrow{\text{a.s.}} X$$

where  $X$  is some integrable random vector.

Let us now show that  $P\left(\limsup_{k \rightarrow \infty} \|\tilde{\theta}_k\| = \infty\right) = 0$ . Since the arguments below apply along any subsequence, represented as

$$\begin{aligned} & \bigcup_S \left\{ \tilde{\theta}_{ki} \rightarrow \infty \quad \forall i \in S \right\} \\ & \subseteq \bigcup_{\tau > 0, S} \left\{ \left\{ \left\{ \tilde{\theta}_{ki} \geq \rho'(\tau, S) \quad \forall i \in S, \quad \tilde{\theta}_{ki} \leq \tau \quad \forall i \notin S \right. \right. \right. \\ & \left. \left. \left. k \geq K(\tau, S) \right\} \cap \limsup_{k \rightarrow \infty} \left\{ M_{ki} < 0 \quad \forall i \in S \right\} \right\} \\ & \bigcup \left\{ \left\{ \tilde{\theta}_{ki} \rightarrow \infty \quad \forall i \in S \right\} \cap \liminf_{k \rightarrow \infty} \left\{ M_{ki} < 0 \quad \forall i \in S \right\}^c \right\} \end{aligned}$$

### For the first event:

Assuming there exists a subsequence  $\{k_0, k_1, k_2, \dots\}$ ,  $k_0 \geq K(\tau, S)$  such that  $\{\tilde{\theta}_{k_j i} \geq \rho'(\tau, S) \forall i \in S\} \cap \{M_{k_j i} < 0 \forall i \in S\}$  is true. Then, from C.6

and  $M_j = a_j \overline{H_j}^{-1} E(G_j(\hat{\theta}_j) | \hat{\theta}_j) = a_j(\bar{g}(\hat{\theta}_j) + \overline{H_j}^{-1} o(c_j^2))$ ,

$$\sum_{i \in S} \tilde{\theta}_{k_j i} \bar{g}_{k_j i}(\hat{\theta}_{k_j}) < 0 \quad \text{a.s.}$$

for all  $k_j$ .

By C.4,  $\tilde{\theta}_{k_j}^T \bar{g}_{k_j}(\hat{\theta}_{k_j}) \geq \rho \|\tilde{\theta}_{k_j}\|$  a.s. which, by C.7,  $\rho'(\tau, S) \geq \tau$  and  $\dim(S) \geq 1$  implies, for all  $j$  sufficiently large,

$$\sum_{i \in S} \tilde{\theta}_{k_j i} \bar{g}_{k_j i}(\hat{\theta}_{k_j}) \geq \frac{\rho}{2} \|\tilde{\theta}_{k_j}\| \geq \left(\frac{\rho}{2}\right) \dim(S) \rho'(\tau, S) \geq \frac{\rho\tau}{2} \quad \text{a.s.}$$

That's a contradiction. So the first event has probability 0.

## For the second event

$\{\tilde{\theta}_{ki} \rightarrow \infty \quad \forall i \in S\} \cap \liminf_{k \rightarrow \infty} \{M_{ki} < 0 \quad \forall i \in S\}^c$ , from

$$\tilde{\theta}_{k+1} + \sum_{j=0}^k M_j \xrightarrow{\text{a.s.}} X$$

for almost all sample points,  $\sum_{k=0}^{\infty} M_{ki} \rightarrow -\infty \quad \forall i \in S$

However, at each  $k$ , the event  $\{M_{ki} < 0 \forall i \in S\}^c$  is composed of the union of  $2^{\dim(S)} - 1$  events, each of which has  $M_{ki} \geq 0$  for at least one  $i \in S$ .

Here is a contradiction. Hence, the probability of the second event is 0.

So  $\limsup_{k \rightarrow \infty} \|\tilde{\theta}_k\| < \infty$

**Part 2:** To show that  $\tilde{\theta}_k$  converges a.s. to a unique (finite) limit, we show that

$$P \left( \liminf_{k \rightarrow \infty} \tilde{\theta}_{ki} < a < b < \limsup_{k \rightarrow \infty} \tilde{\theta}_{ki} \right) = 0 \quad \forall i$$

for any  $a < b$ .

There exist two subsequences, one with convergence to a point  $< a$  and one with convergence to a point  $> b$ .

From  $\tilde{\theta}_{k+1} + \sum_{j=0}^k M_j \xrightarrow{\text{a.s.}} X$  and the conclusion of Part 1, each of these subsequences has a sub-subsequence  $\{k_{j_l}\}$  such that

$$\limsup_{l \rightarrow \infty} \left| \sum_{k=1}^{k_{j_l}} M_{ki} \right| < \infty \text{ a.s.}$$

Supposing that the event within the probability statement is true, we know that for any  $\rho > 0$  and corresponding sample point we can choose  $m > n$  sufficiently large so that for each  $i$  and combined sub-subsequence (from both sub-subsequences mentioned above)



$$\left| \sum_{k=k_{j_n}}^{k_{j_{m-1}}} M_{ki} \right| \leq \rho$$

$$\left| \tilde{\theta}_{k_{j_m}i} - \tilde{\theta}_{k_{j_n}i} + \sum_{k=k_{j_n}}^{k_{j_{m-1}}} M_{ki} \right| \leq \frac{b-a}{3}$$

$$\tilde{\theta}_{k_{j_n}i} < a < b < \tilde{\theta}_{k_{j_m}i}$$

Picking  $\rho < (b-a)/3$  implies that

$$\left| \tilde{\theta}_{k_{j_n}i} - \tilde{\theta}_{k_{j_m}i} \right| \leq 2(b-a)/3$$

it requires that

$$\tilde{\theta}_{k_{j_m}i} - \tilde{\theta}_{k_{j_n}i} > b-a$$

which is a contradiction. So  $\tilde{\theta}_k$  converges a.s. to a unique limit.

**Part 3:** Let us now show that the unique finite limit from Part 2 is 0 . we have  $\limsup_{n \rightarrow \infty} |\sum_{k=0}^{\infty} M_{ki}| < \infty$  a.s.  $\forall i$ .

Then the result to be shown follows if

$$P \left( \lim_{k \rightarrow \infty} \tilde{\theta}_k \neq 0, \left\| \sum_{k=0}^{\infty} M_k \right\| < \infty \right) = 0$$

Suppose that the event is true, and let  $I \subseteq \{1, 2, \dots, p\}$  represent those indexes  $i$  such that  $\tilde{\theta}_{ki} \not\rightarrow 0$  as  $k \rightarrow \infty$ . Then, there exists some  $0 < a' < b' < \infty$  and  $K(a', b') < \infty$  such that

$\forall k \geq K, 0 < a' \leq |\tilde{\theta}_{ki}| \leq b' < \infty$  when  $i \in I (I \neq \emptyset)$  and  $|\tilde{\theta}_{ki}| < a'$  when  $i \in I^c$ . From C.4 and the conditions above, it follows that

$$\sum_{k=K+1}^n a_k \sum_{i \in I} \tilde{\theta}_{ki} \bar{g}_{ki} (\hat{\theta}_k) \geq a' \rho \sum_{k=K+1}^n a_k$$

At least one  $i \in I$

$$\limsup_{n \rightarrow \infty} \left| \frac{\rho a' \sum_{k=K+1}^n a_k}{\sum_{k=K+1}^n a_k \bar{g}_{ki}(\hat{\theta}_k)} \right| < \infty$$

Recall that  $a_k \bar{g}_k(\hat{\theta}_k) = M_k - a_k \bar{H}_k^{-1} b_k$  and  $b_k = O(c_k^2)$  a.s.

Then,  $|\sum_{k=K+1}^{\infty} M_{ki}| = \infty$ . Hence,  $|\sum_{k=0}^{\infty} M_{ki}| = \infty$  with probability  $> 0$  for at least one  $i$ . However, this is inconsistent with the event  $\|\sum_{k=0}^{\infty} M_k\| < \infty$ , showing that the event does, in fact, have probability 0. This completes Part 3, which completes the proof.

Q.E.D.

# Strong Convergence For Hessian Estimator

## Theorem 2a

Let conditions C.0, C.1'', C.2, C.3, and C.4-C.9 hold. Then  $\bar{H}_k \rightarrow H(\theta^*)$  a.s.

C.1'': The conditions of C.1 hold plus  $\sum_{k=0}^{\infty} (k+1)^{-2} (c_k \tilde{c}_k)^{-2} < \infty$  with  $\tilde{c}_k = O(c_k)$

C.3': Change "thrice differentiable" in C.3 to "four-times differentiable" with all else unchanged.

C.9:  $\tilde{\Delta}_k$  satisfies the assumptions for  $\Delta_k$  in C.2 (i.e.,  $\forall k, \ell, |\tilde{\Delta}_{k\ell}| \leq \rho$  and  $\tilde{\Delta}_{k\ell}$  is symmetrically distributed about 0;  $\{\tilde{\Delta}_{k\ell}\}$  are mutually independent);  $\Delta_k$  and  $\tilde{\Delta}_k$  are independent;  $E(\Delta_{k\ell}^{-2}) \leq \rho, E(\tilde{\Delta}_{k\ell}^{-2}) \leq \rho \forall k, \ell$  and some  $\rho > 0$ .

C.8: For some  $\rho > 0$  and all  $k, \ell, m$ ,

$$E \left[ y \left( \hat{\theta}_k \pm c_k \Delta_k + \tilde{c}_k \tilde{\Delta}_k \right)^2 / \left( \Delta_{k\ell} \tilde{\Delta}_{km} \right)^2 \right] \leq \rho$$

and

$$E \left[ y \left( \hat{\theta}_k \pm c_k \Delta_k \right)^2 / \left( \Delta_{k\ell} \Delta_{km} \right)^2 \right] \leq \rho$$

$$E \left[ \tilde{\varepsilon}_k^{(\pm)} - \varepsilon_k^{(\pm)} \mid \hat{\theta}_k; \tilde{\Delta}_k; \bar{H}_k \right] = 0$$

and

$$E \left[ \left( \tilde{\varepsilon}_k^{(\pm)} - \varepsilon_k^{(\pm)} \right)^2 / \left( \Delta_{k\ell} \tilde{\Delta}_{km} \right)^2 \right] \leq \rho$$

where  $\tilde{\varepsilon}_k^{(\pm)} = y \left( \hat{\theta}_k \pm c_k \Delta_k + \tilde{c}_k \tilde{\Delta}_k \right) - L \left( \hat{\theta}_k \pm c_k \Delta_k + \tilde{c}_k \tilde{\Delta}_k \right)$

## Proof:

$$\begin{aligned} & E \left[ G_{kl}^{(1)} \left( \hat{\theta}_k \pm c_k \Delta_k \right) \mid \hat{\theta}_k, \Delta_k \right] \\ &= E \left[ \frac{1}{\tilde{c} \tilde{\Delta}_{kl}} [\tilde{c}_k g \left( \hat{\theta}_k \pm c_k \Delta_k \right)^T \tilde{\Delta}_k + \frac{\tilde{c}_k^2}{2} \tilde{\Delta}_k^T H \left( \hat{\theta}_k \pm c_k \Delta_k \right) \tilde{\Delta}_k \right. \\ &\quad \left. + \frac{\tilde{c}_k^3}{6} \sum_{h,i,j} L_{hij}^{(3)} \left( \bar{\theta}_k^\pm \right) \tilde{\Delta}_{kh} \tilde{\Delta}_{ki} \tilde{\Delta}_{kj}] \mid \hat{\theta}_k, \Delta_k \right] \\ &= g_l \left( \hat{\theta}_k \pm c_k \Delta_k \right) + \frac{1}{6} \tilde{c}_k^2 E \left[ \tilde{\Delta}_{kl}^{-1} \sum_{h,i,j} L_{hij}^{(3)} \left( \bar{\theta}_k^\pm \right) \tilde{\Delta}_{kh} \tilde{\Delta}_{ki} \tilde{\Delta}_{kj} \mid \hat{\theta}_k, \Delta_k \right] \end{aligned}$$

$\tilde{\theta}_k^\pm$  are points on the line segments between  $\hat{\theta}_k \pm c_k \Delta_k + \tilde{c}_k \tilde{\Delta}_k$  and  $\hat{\theta}_k \pm c_k \Delta_k$ ;

Let

$$B_{kl} = \frac{1}{6} E \left[ \tilde{\Delta}_{kl}^{-1} \sum_{h,i,j} \left( L_{hij}^{(3)}(\bar{\theta}_k^+) - L_{hij}^{(3)}(\bar{\theta}_k^-) \right) \cdot \tilde{\Delta}_{kh} \tilde{\Delta}_{ki} \tilde{\Delta}_{kj} \mid \hat{\theta}_k, \Delta_k \right]$$

we have  $B_{kl} \sim o(c_k)$  for all  $k$  sufficiently large. Hence,

$$\begin{aligned} & E \left( \hat{H}_{k,\ell m} \mid \hat{\theta}_k \right) \\ &= E \left( \frac{G_{kl}^{(1)}(\hat{\theta}_k + c_k \Delta_k) - G_{kl}^{(1)}(\hat{\theta}_k - c_k \Delta_k)}{2c_k \Delta_{km}} \mid \hat{\theta}_k \right) \\ &= E \left( \frac{g_\ell(\hat{\theta}_k + c_k \Delta_k) - g_\ell(\hat{\theta}_k - c_k \Delta_k) + \tilde{c}_k^2 B_{kl}}{2c_k \Delta_{km}} \mid \hat{\theta}_k \right) \\ &= E \left( \frac{2c_k [\partial g_\ell / \partial \theta^T]_{\theta=\hat{\theta}_k} \Delta_k + O(c_k^3)}{2c_k \Delta_{km}} \mid \hat{\theta}_k \right) \\ &= H_{\ell m}(\hat{\theta}_k) + O(c_k^2) \end{aligned}$$

Since

$$\frac{1}{n+1} \sum_{k=0}^n [\hat{H}_k - E(\hat{H}_k | \hat{\theta}_k)] \rightarrow 0 \text{ a.s.}$$

Then, by the continuity of  $H$  near  $\hat{\theta}_k$ , and the fact that  $\hat{\theta}_k \rightarrow \theta^*$  a.s. (Theorem 1a)

$$\begin{aligned} & \frac{1}{n+1} \sum_{k=0}^n E(\hat{H}_k | \hat{\theta}_k) \\ &= \frac{1}{n+1} \sum_{k=0}^n \left( H(\hat{\theta}_k) + O(c_k^2) \right) \rightarrow H(\theta^*) \text{ a.s.} \end{aligned}$$

Given that  $\bar{H}_k = (n+1)^{-1} \sum_{k=0}^{n+1} \hat{H}_k$  Q.E.D.



## Theorem 3a

Suppose that C.0, C. 1'', C. 2, C. 3', and C.4-C.9 hold (implying convergence of  $\hat{\theta}_k$  and  $\bar{H}_k$ ). Then, if C.10 and C.11 hold and  $H(\theta^*)^{-1}$  exists,

$$k^{\beta/2} (\hat{\theta}_k - \theta^*) \xrightarrow{\text{dist}} N(\mu, \Omega)$$

where  $\mu = \{0 \text{ if } 3\gamma - \alpha/2 > 0; H(\theta^*)^{-1} T / (a - \beta_+/2) \text{ if } 3\gamma - \alpha/2 = 0\}$ , the  $j$  th element of  $T$  is

$$-\frac{1}{6}ac^2\xi^2 \left[ L_{jjj}^{(3)}(\theta^*) + 3 \sum_{\substack{i=1 \\ i \neq j}}^p L_{ijj}^{(3)}(\theta^*) \right]$$

$\Omega = a^2c^{-2}\sigma^2\rho^2H(\theta^*)^{-2} / (8a - 4\beta_+)$ , and  $\beta_+ = \beta$  if  $\alpha = 1$  and  $\beta_+ = 0$  if  $\alpha < 1$ .

C.10:  $E \left( \varepsilon_k^{(+)} - \varepsilon_k^{(-)} \right)^2 \mid \hat{\theta}_k, \bar{H}_k \rightarrow \sigma^2$  a.s. for some  $\sigma^2 > 0$  For almost all  $\hat{\theta}_k, \left\{ E \left( \left( \varepsilon_k^{(+)} - \varepsilon_k^{(-)} \right)^2 \mid \hat{\theta}_k, c_k \Delta_k = \eta \right) \right\}$  is an equicontinuous sequence at  $\eta = 0$ , and is continuous in  $\eta$  on some compact, connected set containing the actual (observed) value of  $c_k \Delta_k$  a.s.

C.11: In addition to implicit conditions an  $\alpha$  and  $\gamma$  via C.1",  $3\gamma - \alpha/2 \geq 0$  and  $\beta > 0$ . Further, when  $\alpha = 1, a > \beta/2$ . Let  $f_k(\cdot)$  in (2.1a) be chosen such that  $\bar{\bar{H}}_k - \bar{H}_k \rightarrow 0$  a.s.

**Proof:**

Beginning with the expansion  $E \left( G_k \left( \hat{\theta}_k \right) \mid \hat{\theta}_k \right) = H \left( \bar{\theta}_k \right) \left( \hat{\theta}_k - \theta^* \right) + b_k$ , where  $\bar{\theta}_k$  is on the line segment between  $\hat{\theta}_k$  and  $\theta^*$  and  $b_k$  is the bias the estimation error can be represented as

$$\hat{\theta}_{k+1} - \theta^* = \left( I - k^{-\alpha} \Gamma_k \right) \left( \hat{\theta}_k - \theta^* \right) k^{-(\alpha+\beta)/2} \Phi_k V_k + k^{\alpha-\beta/2} \bar{\bar{H}}_k^{-1} T_k$$

where

$$\Gamma_k = a \bar{\bar{H}}_k^{-1} H \left( \bar{\theta}_k \right)$$

$$\Phi_k = -a \bar{\bar{H}}_k^{-1}$$

$$V_k = k^{-\gamma} \left[ G_k \left( \hat{\theta}_k \right) - E \left( G_k \left( \hat{\theta}_k \right) \mid \hat{\theta}_k \right) \right]$$

$$T_k = -a k^{\beta/2} b_k$$

The result will be shown if conditions (2.2.1) (2.2.2) and (2.2.3) of Fabian(1968) hold.

2.2. THEOREM. Suppose  $k$  is a positive integer,  $\mathcal{F}_n$  a non-decreasing sequence of  $\sigma$ -fields,  $\mathcal{F}_n \subset \mathcal{S}$ ; suppose  $U_n, V_n, T_n \in \mathbb{R}^k, T \in \mathbb{R}^k, \Gamma_n, \Phi_n \in \mathbb{R}^{k \times k}, \Sigma, \Gamma, \Phi, P \in \mathbb{R}^{k \times k}$ ,  $\Gamma$  is positive definite,  $P$  is orthogonal and  $P' \Gamma P = \Lambda$  diagonal. Suppose  $\Gamma_n, \Phi_{n-1}, V_{n-1}$  are  $\mathcal{F}_n$ -measurable,  $C, \alpha, \beta \in \mathbb{R}$  and

$$(2.2.1) \quad \Gamma_n \rightarrow \Gamma, \quad \Phi_n \rightarrow \Phi, \quad T_n \rightarrow T \quad \text{or} \quad E \|T_n - T\| \rightarrow 0,$$

$$(2.2.2) \quad E_{\mathcal{F}_n} V_n = 0, \quad C > \|E_{\mathcal{F}_n} V_n V_n' - \Sigma\| \rightarrow 0,$$

and, with  $\sigma_{j,r}^2 = E \chi \{ \|V_j\|^2 \geq r j^\alpha \} \|V_j\|^2$ , let either

$$(2.2.3) \quad \lim_{j \rightarrow \infty} \sigma_{j,r}^2 = 0 \quad \text{for every } r > 0,$$

or

$$(2.2.4) \quad \alpha = 1, \quad \lim_{n \rightarrow \infty} n^{-1} \sum_{j=1}^n \sigma_{j,r}^2 = 0 \quad \text{for every } r > 0.$$

Suppose that, with  $\lambda = \min_i \Lambda^{(ii)}, \beta_+ = \beta$  if  $\alpha = 1, \beta_+ = 0$  if  $\alpha \neq 1$ ,

$$(2.2.5) \quad 0 < \alpha \leq 1, \quad 0 \leq \beta, \quad \beta_+ < 2\lambda$$

and

$$(2.2.6) \quad U_{n+1} = (I - n^{-\alpha} \Gamma_n) U_n + n^{-(\alpha+\beta)/2} \Phi_n V_n + n^{-\alpha-\beta/2} T_n.$$

Then the asymptotic distribution of  $n^{\beta/2} U_n$  is normal with mean  $(\Gamma - (\beta_+/2)I)^{-1} T$  and covariance matrix  $PMP'$  where

$$(2.2.7) \quad M^{(ij)} = (P' \Phi \Sigma \Phi' P)^{(ij)} (\Lambda^{(ii)} + \Lambda^{(jj)} - \beta_+)^{-1}.$$

If  $3\gamma - \alpha/2 > 0$ ,  $T_k \rightarrow 0$  a.s. by the fact that  $b_k(\hat{\theta}_k) = O(k^{-2\gamma})$  a.s.

If  $3\gamma - \alpha/2 = 0$ ,

$$b_{kl}(\hat{\theta}_k) - \frac{1}{6} \frac{c^2}{k^{2\gamma}} L^{(3)}(\theta^*) E[\Delta_{kl}^{-1}(\Delta_k \otimes \Delta_k \otimes \Delta_k)] \rightarrow 0 \text{ a.s.}$$

$$T_{kl} \rightarrow -\frac{1}{6} ac^2 \xi^2 \left\{ L_{lll}^{(3)}(\theta^*) + \sum_{\substack{i=1 \\ i \neq l}}^p [L_{lil}^{(3)}(\theta^*) + L_{lil}^{(3)}(\theta^*) + L_{lil}^{(3)}(\theta^*)] \right\} \text{ a.s.}$$

We have thus shown that  $T_k$  converges for  $3\gamma - \alpha/2 \geq 0$ .

$$E(V_k V_k^T | \mathcal{F}_k) = k^{-2\gamma} E \left\{ \left[ \frac{L(\hat{\theta}_k + \bar{\Delta}_k) - L(\hat{\theta}_k - \bar{\Delta}_k)}{2ck^{-\gamma}\Delta_k} \right]^2 \mid \hat{\theta}_k \right\}$$

$$\begin{aligned}
& + k^{-2\gamma} E \left\{ \Delta_k^{-1} (\Delta_k^{-1})^T \left[ \frac{\epsilon_k^{(+)} - \epsilon_k^{(-)}}{2ck^{-\gamma}} \right] \cdot \left[ \frac{L(\hat{\theta}_k + \bar{\Delta}_k) - L(\hat{\theta}_k - \bar{\Delta}_k)}{2ck^{-\gamma}} \right] \mid \mathcal{F}_k \right\} \\
& + k^{-2\gamma} E \left\{ \Delta_k^{-1} (\Delta_k^{-1})^T \left[ \frac{\epsilon_k^{(+)} - \epsilon_k^{(-)}}{2ck^{-\gamma}} \right]^2 \mid \mathcal{F}_k \right\} \\
& - k^{-2\gamma} \left[ g(\hat{\theta}_k) + b_k(\hat{\theta}_k) \right] \left[ g(\hat{\theta}_k) + b_k(\hat{\theta}_k) \right]^T
\end{aligned}$$

For the third term

$$\begin{aligned}
& E \left[ \Delta_k^{-1} (\Delta_k^{-1})^T \left( \epsilon_k^{(+)} - \epsilon_k^{(-)} \right)^2 \mid \mathcal{F}_k \right] \\
& = \int_{\Omega_\Delta} \Delta_k^{-1} (\Delta_k^{-1})^T E \left[ \left( \epsilon_k^{(+)} - \epsilon_k^{(-)} \right)^2 \mid \mathcal{F}_k, \bar{\Delta}_k \right] dP_\Delta \\
& E \left( V_k V_k^T \mid \mathcal{F}_k \right) \rightarrow \frac{1}{4} c^{-2} \sigma^2 \rho^2 I \quad \text{a.s.}
\end{aligned}$$

This completes the proof of Fabian's conditions (2.2.1) and (2.2.2)

We now show that condition (2.2.3) holds, which is

$$\lim_{k \rightarrow \infty} E \left( \mathcal{I}_{\{\|V_k\|^2 \geq rk^\alpha\}} \|V_k\|^2 \right) = 0 \quad \forall r > 0$$

where  $\mathcal{I}_{\{\cdot\}}$  denotes the indicator for  $\{\cdot\}$ . By Holder's inequality and for any  $0 < \delta' < \delta/2$ , the above limit is bounded above by

$$\begin{aligned} & \lim_{k \rightarrow \infty} \sup P \left( \|V_k\|^2 \geq rk^\alpha \right)^{\delta'/(1+\delta')} \left( E \|V_k\|^{2(1+\delta')} \right)^{1/(1+\delta')} \\ & \leq \limsup_{k \rightarrow \infty} \left( \frac{E \|V_k\|^2}{rk^\alpha} \right)^{\delta'/(1+\delta')} \left( E \|V_k\|^{2(1+\delta')} \right)^{1/(1+\delta')} \end{aligned}$$

Note that

$$\begin{aligned} \|V_k\|^{2(1+\delta')} & \leq 2^{2(1+\delta')} k^{-2(1+\delta')\gamma} \left[ \left\| \hat{g} \left( \hat{\theta}_k \right) \right\|^{2(1+\delta')} \right. \\ & \quad \left. + \left\| g \left( \hat{\theta}_k \right) \right\|^{2(1+\delta')} + \left\| b_k \left( \hat{\theta}_k \right) \right\|^{2(1+\delta')} \right] \end{aligned}$$

# Thanks