

## AVERAGING AND DERIVATIVE ESTIMATION WITHIN STOCHASTIC APPROXIMATION ALGORITHMS

Fatemeh Sadat Hashemi

Raghu Pasupathy

Virginia Tech  
Blacksburg, VA 24061, USA

Virginia Tech  
Blacksburg, VA 24061, USA

### ABSTRACT

Stochastic Approximation (SA) is arguably the most investigated amongst algorithms for solving local continuous simulation optimization problems. Despite its enduring popularity, the prevailing opinion is that the finite-time performance of SA-type algorithms is still not robust to SA's sequence of algorithm parameters. In the last two decades, two major advances have been proposed toward alleviating this issue: (i) Polyak-Ruppert averaging where SA is executed in multiple time scales to allow for the algorithm iterates to use large (initial) step sizes for better finite time performance, without sacrificing the asymptotic convergence rate; and (ii) efficient derivative estimation to allow for better searching within the solution space. Interestingly, however, all existing literature on SA seems to treat each of these advances separately. In this article, we present two results which characterize SA's convergence rates when both (i) and (ii) are applied simultaneously. Our results should be seen as simply providing a theoretical basis for applying ideas that seem reasonable in practice.

### 1 INTRODUCTION

The broad setting of this paper is stochastic approximation (SA), the famous iteration originally introduced by Robbins and Monro (1951) as a method to identify the zero of a function. The modern version of Robbins and Monro's SA iteration usually takes the form

$$Z_n = Z_{n-1} - \gamma_n \bar{H}_n^{-1} \tilde{h}(Z_{n-1}), n = 1, 2, \dots \quad (1)$$

where  $\tilde{h}(\cdot)$  is an unbiased estimator of a vector valued function  $h : \mathbf{R}^q \rightarrow \mathbf{R}^q$ ,  $\{\gamma_n\}$  is a positive sequence converging to 0, and  $\bar{H}_n$  is an estimator of the Jacobian matrix of the function  $h$  at the point  $Z_{n-1}$ . The iteration has been widely used in both the optimization and root-finding contexts. When used in the root-finding context, the objective of the iteration in (1) is identifying a zero of the function  $h(\cdot)$ , while the estimator  $\tilde{h}(\cdot)$  provides noisy observations of the function  $h(\cdot)$ . When used in the optimization context, the objective of the iteration in (1) is identifying a stationary point of a real-valued function  $g : \mathbf{R}^q \rightarrow \mathbf{R}$  that is "observed" using an unbiased estimator  $G(\cdot)$ . In such a case, the quantity  $\tilde{h}(Z_{n-1})$  appearing in (1) estimates the gradient of the function  $g(\cdot)$  at the point  $Z_{n-1}$ , and is usually constructed using forward or central differencing (Spall 2003). In the optimization setting,  $\bar{H}_n$  then estimates the Hessian (matrix of second derivatives) of the function  $g(\cdot)$  at the point  $Z_{n-1}$ , and is again calculated using some form of differencing (Spall 2003). The SA iteration as stated in (1) is for the unconstrained context. Extending it to tackle problems with a (deterministically) constrained feasible region is usually done by performing an appropriate projection operation back into the feasible region whenever the iterates drift outside the feasible region.

SA is arguably the most popular current method of solving continuous local optimization and root-finding problems when the functions involved can only be estimated (and the constraints are known and deterministic). Owing to its simplicity, its interpretation as the natural stochastic analogue of Newton's method, and its attractive asymptotic properties, SA has seen a tremendous amount of application (Kushner

and Yin 2003). A lot has been written on the topic, and the finite-time and infinite-time behavior of the recursion in (1) is well-understood. (There are several books that will serve as good entry points into this literature, e.g., (Kushner and Yin 2003; Borkar 2008; Spall 2003; Wasan 1969).)

Despite SA's enduring popularity and the six decades of research supporting its advance, the prevailing opinion is that choosing the gain sequence  $\{\gamma_n\}$  to ensure robust and efficient SA performance is challenging (Spall 1998; Spall 2006; Broadie, Cicek, and Zeevi 2010; Broadie, Cicek, and Zeevi 2009; Pasupathy and Kim 2011; Pasupathy and Schmeiser 2010). In other words, while it is possible to tune the gain sequence and "make" SA perform well for a given problem, or even a class of problems, formulating rules that *automatically* tune the gain sequence to achieve good finite-time performance is still an open problem (albeit loosely defined). This opinion is also supported by continuing efforts to devise rules that either dynamically choose the gain sequence (Broadie, Cicek, and Zeevi 2010; Broadie, Cicek, and Zeevi 2009; Yousefian, Nedić, and Shanbhag 2011) based on the observed history of algorithm evolution, or by mitigating the effect of the gain sequence (Nemirovski, Juditsky, Lan, and Shapiro 2009; Pasupathy and Schmeiser 2010).

Of particular interest in this paper are two advances that have been crucial milestones in SA's history. The first is what is popularly called "Polyak Averaging" (Polyak and Juditsky 1992) which involves the simple idea of averaging SA's iterates. To elaborate, various authors (Chung 1954; Derman 1956; Fabian 1968) prior to 1997 had shown that the best possible convergence rate of SA's iterates (to the correct solution) is  $O(1/\sqrt{n})$ , achieved when the gain sequence  $\gamma_n = O(1/n)$ . (Rigorously, this implies that when  $\gamma_n = K/n$  and  $K$  is larger than half the inverse of the smallest eigen value of the function  $g$ 's Hessian at the solution, the iterates can be shown to satisfy  $\sqrt{n}(Z_n - x^*) \xrightarrow{D} N(0, V)$  where  $x^*$  is a solution to the problem,  $V$  is a covariance matrix, and  $\xrightarrow{D}$  denotes convergence in distribution.) While this result is useful, finite-time performance considerations suggested using step sizes that were larger, i.e., converged slower, than the  $O(1/n)$  suggested by asymptotic performance considerations. The dilemma was that choosing a slowly converging gain sequence, e.g.,  $\gamma_n = O(1/n^\alpha)$ ,  $\alpha \in (0, 1)$ , while often producing better finite-time performance, degraded SA's asymptotic convergence rate. Polyak and Juditsky (1992), and simultaneously Frees and Ruppert (1990), provided an elegant solution for this dilemma. Polyak and Juditsky (1992) showed that SA can be executed on two timescales to enjoy good finite-time performance while not sacrificing asymptotic performance. Specifically, he suggested executing SA on the "fast timescale"  $Z_n = Z_{n-1} - \gamma_n \tilde{h}(Z_{n-1})$ ,  $\gamma_n = O(n^{-\alpha})$ ,  $\alpha \in (0, 1)$  and then averaging the iterates  $Z_n, n = 1, 2, \dots$  offline to get  $Y_n = n^{-1} \sum_{i=1}^n Z_i$ . He demonstrated the remarkable result that, under certain conditions, such a two timescale averaging produced the averaged iterates  $\{Y_n\}$  having the best possible convergence rate  $O(1/\sqrt{n})$ . (He also showed that the iterates  $Z_n$  attain the degraded convergence rate of  $O(\gamma_n)$ .) Polyak's paper was written within the context of root-finding. This was extended to the optimization context by Dippon and Renz (1997).

The second milestone of interest in this paper is the efficient use of derivatives within SA. It is clear from the corresponding literature in the deterministic context that knowledge of the Jacobian matrix of derivatives  $H : \mathbb{R}^q \rightarrow \mathbb{R}^q$  of the function  $h(\cdot)$  can help immensely with efficient searching within the solution space. This prompted modifying the original SA iteration to incorporate derivative estimates to obtain the modified iteration  $Z_n = Z_{n-1} - \gamma_n \bar{H}_n^{-1} \tilde{h}(Z_{n-1})$ . While this sounds reasonable, the main issue with this approach turns out to be the computation involved in estimating the derivative estimate  $\bar{H}_n$ . For instance, if one indulged in estimating every entry of  $\bar{H}_n$  using a method such as forward differences, this would involve  $O(q^2)$  simulations just to obtain the estimated derivative at the incumbent point. This led Spall (2003) to investigate more efficient methods to obtain a derivative estimate. While a lot has been written on this particular issue, the crux of Spall's work is that with just  $O(q)$  simulations, it is possible to obtain derivative estimates that do not degrade the asymptotic convergence rate of the SA iteration. In other words, using just a crude calculation of the derivative  $\bar{H}_n$ , one could potentially enjoy the benefits of better finite-time performance without sacrificing the asymptotic convergence rate of  $O(1/\sqrt{n})$ .

We see both of the above milestones (Polyak averaging and incorporation of derivatives) as measures that were undertaken explicitly to improve the finite-time performance of SA while retaining the fastest possible  $O(1/\sqrt{n})$  convergence. Interestingly, however, we have found no evidence of any analysis in the literature that incorporates both of these ideas. (Even the most recent literature on this topic (Mokkadem and Pelletier 2011; Nemirovski, Juditsky, Lan, and Shapiro 2009; Yousefian, Nedić, and Shanbhag 2011) do not incorporate estimated Hessians into the SA iteration, most likely due to computational considerations.) Towards addressing this gap in the SA literature, we ask the following two questions.

- Q.1 When Polyak averaging and derivative estimates are included within the SA iteration, what conditions ensure that the averaged iterates retain the  $O(1/\sqrt{n})$  convergence?  
 Q.2 Can anything be said about the convergence characteristics of SA's faster timescale iterates?

We start by answering Q.2. We demonstrate that, amongst other conditions, if the sequence  $\{\gamma_n \bar{H}_n^{-1}\}$  satisfies a certain stochastic-matrix analogue of regularly varying sequences, and the Hessian estimator  $\bar{H}_n$  is consistent in a certain precise sense, the faster timescale iterates  $\{Z_n\}$  converge in mean square at the rate  $O(\gamma_n \log n)$ . This rate is slightly slower than that obtained without the Hessian estimator. The condition we impose on  $\{\gamma_n \bar{H}_n^{-1}\}$  is not entirely new and is closely related to conditions established in Polyak and Juditsky (1992) and Mokkadem and Pelletier (2011).

In answering Q.1, we show that conditions similar to that used in answering Q.2 ensure that the slower timescale sequence  $\{Y_n\}$  retains the  $O(1/\sqrt{n})$  convergence rate. This should come as no surprise to the reader and should be seen simply as theoretical confirmation of what seems intuitively clear.

The remainder of the paper is organized as follows. In Section 2.1, we outline the sufficient conditions to retrieve the rate at which the fast timescale sequence converges to the root which is established in part 2.2. In section 2.3 we establish the almost-sure convergence of the averaged iterates within the SA iteration together with a simple extension of the SA iteration for relaxing the need to pre-specify the gain sequence. Concluding remarks are made finally in section 3.

## 2 MAIN RESULTS

In everything that follows, the SA iteration of interest is the two timescale recursion given by

$$\begin{aligned} Z_n &= Z_{n-1} - \Lambda_n \tilde{h}_n; \\ Y_n &= \left(1 - \frac{1}{n+1}\right) Y_{n-1} + \frac{1}{n+1} Z_n; \end{aligned} \tag{2}$$

where  $\Lambda_n = \gamma_n \bar{H}_n^{-1}$ ,  $\bar{H}_n$  is a consistent estimator of the derivative of  $h(\cdot)$  at  $Z_{n-1}$ , and  $\{\gamma_n\}$  is a positive sequence converging to 0. Also, as is common SA settings, we assume that  $\tilde{h}_n = h(Z_{n-1}) + \varepsilon_n$  is the noisy observation of the function  $h$  at the point  $Z_{n-1}$ , where  $\varepsilon_n$  is a random disturbance. We emphasize that for purposes of this paper, our interest will be limited to the context of root-finding within the unconstrained context.

### 2.1 Assumptions and Notation

- C.1 Suppose that the solution to the vector equation  $h(Z) = 0$  is  $z^*$ . We assume the existence of  $\eta > 1$  and a neighborhood  $\mathcal{N}(z^*)$  of  $z^*$  such that  $h(z) = H(z - z^*) + O(\|z - z^*\|^\eta)$  for  $z \in \mathcal{N}(z^*)$ , where the matrix  $-H$  is Hurwitz.  
 C.2 There exists  $\rho_1 > 0$  for which  $\mathbb{E}\|\tilde{h}(z)\|^2 \leq \rho_1(1 + \|z - z^*\|^2)$ .  
 C.3 For the consistent estimator  $\bar{H}_n$ , we assume the boundedness of moments, i.e., the existence of a positive  $\rho$  such that  $E(\|\bar{H}_n^{-1}\|^2) \leq \rho$ .  
 C.4 For each  $n \geq 1$  and all  $z$  there exists  $\rho_2 > 0$  not dependent on  $n$  and  $z$  such that  $(z - z^*)^T \tilde{h}_n(Z) \geq \rho_2 \|z - z^*\|^2$ , where  $\tilde{h}_n(Z) = \bar{H}_n^{-1} \tilde{h}(Z_n)$ .

- C.5 (a)  $\lim_{n \rightarrow \infty} n(I - \Lambda_n^{-1} \Lambda_{n+1}) \xrightarrow{P} \alpha I$ , where  $\xrightarrow{P}$  denotes convergence in probability,  $1/2 < \alpha < 1$ ,  $\Lambda_n = \gamma_n \bar{H}_n^{-1}$ , and  $\gamma_n \rightarrow 0$ ,  $\gamma_n > 0$ .  
 (b)  $\lim_{n \rightarrow \infty} \log n / n \gamma_n \rightarrow 0$  and  $\sum_{n=1}^{\infty} \frac{(\gamma_n \log n)^{\eta}}{\sqrt{n}} < \infty$  for  $\eta > 1$ .
- C.6  $E(\varepsilon_{n+1} | \mathcal{F}_n) = 0$  where  $\mathcal{F}_n = \{\varepsilon_0, \varepsilon_1, \dots, \varepsilon_{n-1}, Z_n, \bar{H}_n\}$  and there exists a non-random, positive definite matrix  $\Gamma$  such that  $\lim_{n \rightarrow \infty} E(\varepsilon_{n+1} \varepsilon_{n+1}^T | \mathcal{F}_n) = \Gamma$  almost surely.

Assumptions similar to C.1 and C.2 are common within the literature on Polyak averaging. For instance, Polyak and Juditsky (1992), Brodie, Cicek, and Zeevi (2010) and Mokkadem and Pelletier (2011) impose similar conditions. The assumptions C.3 and C.4 are prevalent in the SA literature that uses an estimated derivative within the recursion. For example, Spall (2000) makes this assumption. The main condition of interest is C.5(a) and it should be seen as the matrix-analogue of the original condition introduced by Polyak and Juditsky (1992). Assuming consistency of  $\bar{H}_n$ , this assumption implies that  $\gamma_n$  is a regularly varying sequence as originally introduced by Galombos and Seneta (2008). Interestingly, C.5 implies that  $\sum_{k=1}^n \gamma_k \rightarrow \infty$  and  $\sum_{k=1}^n \gamma_k^2 < \infty$  as  $n \rightarrow \infty$ , conditions that are routinely assumed (directly) within the SA literature.

## 2.2 Behavior of the Fast Timescale Iterates

In this section, we present a result that characterizes the rate at which the fast timescale sequence  $\{Z_n\}$  converges to the root  $z^*$ .

**Theorem 1** Let assumptions C.1 – C.5 hold, and let  $\rho_2^2 < \rho \rho_1$ . Then the mean squared error  $\text{mse}(Z_n, z^*)$  of  $Z_n$  with respect to  $z^*$  satisfies  $\text{mse}(Z_n, z^*) = O(\gamma_n \log n)$ .

*Proof.* Let  $A_{n+1} = \|Z_{n+1} - z^*\|^2$ . Then

$$A_{n+1} = \|Z_n - z^* - \Lambda_n \tilde{h}(Z_n)\|^2 = \|Z_n - z^*\|^2 - 2\gamma_n (Z_n - z^*)^T \bar{h}(Z_n) + \gamma_n^2 \|\bar{h}(Z_n)\|^2. \quad (3)$$

By assumptions C.2 and C.3, we get

$$E[\|\bar{h}(Z_n)\|^2 | Z_n] \leq \rho \rho_1 (1 + A_n);$$

and in view of assumption C.4 we have

$$E[-(Z_n - z^*)^T \bar{h}(Z_n) | Z_n] \leq -\rho_2 A_n.$$

Taking expectations on both sides in (3) after conditioning on  $Z_n$  we get

$$E[A_{n+1} | Z_n] \leq A_n (1 - 2\rho_2 \gamma_n + \rho \rho_1 \gamma_n^2) + \rho \rho_1 \gamma_n^2. \quad (4)$$

If we now let  $b_n := E[A_{n+1}]$ , we get

$$b_n \leq b_1 \prod_{i=1}^n p_i + \sum_{i=2}^{n-1} \prod_{j=i+1}^n q_i p_j + q_n := u_n,$$

where  $p_i = (1 - 2\rho_2 \gamma_i + \rho \rho_1 \gamma_i^2)$ ,  $q_i = \rho \rho_1 \gamma_i^2$ . Since we have chosen  $\rho, \rho_1, \rho_2$  in such a way that  $\rho \rho_1 > \rho_2^2 > 0$  for all  $i$ ,  $p_i$  and  $q_i$  are positive.

Define  $n_0 := \sup\{n \geq 1 : \rho_2 < 2\rho \rho_1 \gamma_n, \rho_2 < 2\rho \rho_1 \alpha \frac{\gamma_n}{n}, n \gamma_n < \frac{2\alpha}{\rho_2}, \text{ and } \log n - 1 < \frac{2\rho \rho_1}{\rho_2}\} + 1$  and choose  $c$  large enough to satisfy the following

$$\frac{u_{n_0+1}}{\gamma_{n_0} \log n_0} \leq c.$$

Then one can see by induction that for all  $n \geq 1$ ,  $b_{n+1} \leq c\gamma_n \log n$ , where

$$c = \max\left\{1, \max_{1 \leq n \leq n_0} \left\{ \frac{u_{n+1}}{\gamma_n \log n} \right\}\right\}.$$

□

Theorem 1 asserts that the fast timescale iterates converge at the rate  $O(\gamma_n \log n)$ . Since the sequence  $\{\gamma_n\}$  converges slower than  $O(1/n)$ , this points to a degraded rate of convergence for the fast timescale iterates.

### 2.3 Weak Convergence of the Slow Timescale Iterates

In this section, we present the main result of the paper. Theorem 2 demonstrates that the averaged iterates within the SA iteration in (2) attain the best possible convergence rate in a weak sense.

**Theorem 2** Under assumptions C.1–C.6, we have

- (i)  $\sqrt{n}(Y_n - Z^*) \xrightarrow{D} N(0, H^{-1}\Gamma[H^{-1}]^T)$ , where  $H$  represents the Jacobian of  $h(z)$  at  $z = z^*$ ;
- (ii)  $Y_n - Z^* \rightarrow 0$  almost surely.

*Proof of (i).* (Assume that the underlying function  $h$  is linear)

Let  $\Delta_n = Z_n - z^*$  and  $\bar{\Delta}_n = Y_n - z^*$ . Then

$$\begin{aligned} Z_n &= Z_{n-1} - \Lambda_n(H\Delta_{n-1} + \varepsilon_n); \\ \Delta_{n-1} &= H^{-1}\Lambda_n^{-1}(Z_{n-1} - Z_n) - H^{-1}\varepsilon_n; \\ \Delta_n &= H^{-1}\Lambda_{n+1}^{-1}(Z_n - Z_{n+1}) - H^{-1}\varepsilon_{n+1}. \end{aligned}$$

On the other hand

$$\begin{aligned} Y_n &= \left(1 - \frac{1}{n+1}\right)Y_{n-1} + \frac{1}{n+1}Z_n; \\ Y_n - z^* &= \left(1 - \frac{1}{n+1}\right)(Y_{n-1} - z^*) + \frac{1}{n+1}(Z_n - z^*); \\ \bar{\Delta}_n &= \left(1 - \frac{1}{n+1}\right)\bar{\Delta}_{n-1} + \frac{1}{n+1}\Delta_n. \end{aligned}$$

Set  $\prod_{j=n+1}^n (1 - \frac{1}{j+1}) = 1$ . So we get

$$\begin{aligned} \bar{\Delta}_n &= \prod_{j=1}^n \left(1 - \frac{1}{j+1}\right) \Delta_0 \\ &\quad + \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1}\Lambda_{k+1}^{-1}(Z_k - Z_{k+1}) \\ &\quad - \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1}\varepsilon_{k+1}. \end{aligned}$$

Let

$$\begin{aligned} R_{n+1}^1 &= \prod_{j=1}^n \left(1 - \frac{1}{j+1}\right) \Delta_0; \\ R_{n+1}^2 &= \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} \Lambda_{k+1}^{-1} (Z_k - Z_{k+1}); \\ R_{n+1}^3 &= \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} \epsilon_{k+1}. \end{aligned}$$

Note that as  $n \rightarrow \infty$ ,  $\sqrt{n}R_{n+1}^1 = \frac{\sqrt{n}}{n+1} \Delta_0 \rightarrow 0$  a.s.. Also by Polyak and Juditsky (1992),  $\sqrt{n}R_{n+1}^3 \xrightarrow{D} N(0, H^{-1} \Gamma [H^{-1}]^T)$ . So we just need to prove that  $\sqrt{n}R_{n+1}^2 \rightarrow 0$  in probability:  $\sqrt{n}R_{n+1}^2 \xrightarrow{P} 0$ .

$$\begin{aligned} R_{n+1}^2 &= \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} \Lambda_{k+1}^{-1} (Z_k - Z_{k+1}) \\ &= \frac{1}{n+1} \sum_{k=1}^n H^{-1} \Lambda_{k+1}^{-1} [(Z_k - z^*) - (Z_{k+1} - z^*)] \\ &= \frac{1}{n+1} H^{-1} \Lambda_2^{-1} (Z_1 - z^*) + \frac{1}{n+1} \sum_{k=2}^n H^{-1} \Lambda_{k+1}^{-1} (Z_k - z^*) \\ &\quad - \frac{1}{n+1} \sum_{k=1}^{n-1} H^{-1} \Lambda_{k+1}^{-1} (Z_{k+1} - z^*) + \frac{1}{n+1} H^{-1} \Lambda_{n+1}^{-1} (Z_{n+1} - z^*). \end{aligned}$$

Since we have

$$\frac{1}{n+1} \sum_{k=1}^{n-1} H^{-1} \Lambda_{k+1}^{-1} (Z_{k+1} - z^*) = \frac{1}{n+1} \sum_{k=2}^n H^{-1} \Lambda_k^{-1} (Z_k - z^*),$$

we can write

$$\begin{aligned} R_{n+1}^2 &= \frac{1}{n+1} \sum_{k=2}^n [H^{-1} \Lambda_{k+1}^{-1} (I - \Lambda_{k+1} \Lambda_k^{-1}) (Z_k - z^*)] \\ &\quad - \frac{1}{n+1} H^{-1} \Lambda_{n+1}^{-1} (Z_{n+1} - z^*) + \frac{1}{n+1} H^{-1} \Lambda_2^{-1} (Z_1 - z^*). \end{aligned}$$

In view of Theorem 1 and assumption C.4, we then get:

$$\begin{aligned} R_{n+1}^2 &= \frac{1}{n+1} \sum_{k=2}^n [H^{-1} \Lambda_{k+1}^{-1} o_p\left(\frac{1}{k+1}\right) o(\sqrt{\gamma_k \log k})] - \frac{1}{n+1} H^{-1} \Lambda_{n+1}^{-1} o(\sqrt{\gamma_{n+1} \log n + 1}) + o\left(\frac{1}{\sqrt{n+1}}\right) \\ &= \frac{1}{n+1} \sum_{k=2}^n o_p\left[\frac{1}{k+1} \sqrt{\frac{\gamma_k \log k}{\gamma_{k+1}^2}}\right] - o_p\left[\frac{1}{n+1} \sqrt{\frac{\gamma_{n+1} \log n + 1}{\gamma_{n+1}^2}}\right] + o\left(\frac{1}{\sqrt{n+1}}\right). \end{aligned}$$

Hence

$$\begin{aligned} \sqrt{n}R_{n+1}^2 &= \frac{1}{\sqrt{n+1}} \sum_{k=2}^n o_p\left(\frac{1}{\sqrt{k+1}}\right) - o_p(1) + o(1) \\ &= o_p(1). \end{aligned}$$

*Proof of part (i).* (Assume that the underlying function  $h$  is nonlinear)

By assumption C.1 we get:

$$\begin{aligned}\bar{\Delta}_n &= \prod_{j=1}^n \left(1 - \frac{1}{j+1}\right) \Delta_0 \\ &+ \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} \Lambda_{k+1}^{-1} (Z_k - Z_{k+1}) \\ &+ \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} o(\|Z_k - z^*\|^\eta) \\ &- \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} \varepsilon_{k+1}.\end{aligned}$$

Let

$$\tilde{R}_{n+1} = \sum_{k=1}^n \prod_{j=k+1}^n \left(1 - \frac{1}{j+1}\right) \frac{1}{k+1} H^{-1} o(\|Z_k - z^*\|^\eta).$$

Thus, we are only yet to prove that  $\sqrt{n}\tilde{R}_{n+1} \rightarrow 0$  as  $n \rightarrow \infty$ . By Theorem 1, we have

$$\begin{aligned}\sqrt{n}\tilde{R}_{n+1} &= \frac{1}{\sqrt{n}} \sum_{k=1}^n o((\gamma_k \log k)^{\frac{\eta}{2}}) \\ &= \frac{1}{\sqrt{n}} \sum_{k=1}^n \sqrt{k} o\left(\frac{(\gamma_k \log k)^{\frac{\eta}{2}}}{\sqrt{k}}\right).\end{aligned}$$

The claim then follows by assumption C.5(b) and Kronecker's lemma.

*Proof of part (ii).*

$$R_{n+1}^1 = \prod_{j=1}^n \left(1 - \frac{1}{j+1}\right) \Delta_0 = \frac{\Delta_0}{n+1}$$

and so  $R_{n+1}^1 \rightarrow 0$  as  $n \rightarrow \infty$ .

$$R_{n+1}^2 = \frac{1}{n+1} \sum_{k=2}^n o_p\left[\frac{1}{k} \sqrt{\frac{\gamma_k \log k}{\gamma_{k+1}^2}}\right] - \frac{1}{n+1} \gamma_{n+1}^{-1} H^{-1} \bar{H}_{n+1} o(\sqrt{\gamma_{n+1} \log n + 1}) + o\left(\frac{1}{\sqrt{n}}\right),$$

and by Cesaro summability (Billingsley 1995) and C.5(a),  $R_{n+1}^2 \rightarrow 0$  as  $n \rightarrow \infty$ . Finally,

$$R_{n+1}^3 = \frac{1}{n+1} \sum_{k=1}^n H^{-1} \varepsilon_{k+1}$$

and so by the strong law of large numbers (Billingsley 1995) we get  $R_{n+1}^3 \rightarrow 0$  as  $n \rightarrow \infty$ . □

In conclusion and as a further step in the direction of completely relaxing the need to pre-specify the gain sequence, we now propose a simple extension of the SA iteration considered thus far.

$$Z_{j+1} = Z_j - \Lambda_{t_j} \tilde{h}(Z_j), j = 1, 2, \dots \quad (5)$$

where  $t_j := \text{Min}\{t : N_t \geq j\}$  for  $j = 1, 2, \dots$ ,  $\Lambda_{t_j} = \gamma_{t_j} \bar{H}_j^{-1}$ , and  $\bar{H}_j, \tilde{h}(Z_j)$  are as defined in (2). It can be seen that the iteration in (5) is constructed to facilitate designing heuristics that dynamically change the step sizes based on observed history of the SA iteration.

The following theorem establishes the asymptotic efficiency of (5) under suitable conditions.

**Theorem 3** Let  $(N_t)_{t \geq 0}$  be an increasing sequence of random variables with  $N_0 = 0$  and let  $\Delta_t = N_t - N_{t-1}$ .

- (I) Suppose assumption C.2 – C.4 and C.6 hold true. Further suppose that  $\Delta_t$  is uniformly bounded for all  $t$ . If the gain sequence  $\{\gamma_t\}$  satisfies  $\sum_{t=1}^{\infty} \gamma_t = \infty$  and  $\sum_{t=1}^{\infty} \gamma_t^2 < \infty$ , the iteration in (5) converges to  $z^*$  a.s.
- (II) Moreover, consider the following two time scale SA algorithm:

$$\begin{aligned} Z_{n+1} &= Z_n - \Lambda_{t_n} \tilde{h}(Z_n); \\ \bar{Z}_{n+1} &= \left(1 - \frac{1}{n+2}\right) \bar{Z}_n + \frac{1}{n+2} Z_{n+1}; \end{aligned} \tag{6}$$

Suppose conditions C.1 – C.4 and C.5 – C.6 hold, and for a  $\zeta > 0$ ,  $\frac{t_n}{n} \rightarrow \zeta$  a.s. as  $n \rightarrow \infty$ . Then we have

- (i)  $\sqrt{n}(\bar{Z}_n - Z^*) \xrightarrow{D} N(0, H^{-1} \Gamma [H^{-1}]^T)$ ;
- (ii)  $\bar{Z}_n - Z^* \rightarrow 0$  almost surely.

### 3 CONCLUDING REMARKS

While Polyak averaging and the use of derivative estimates within SA are well-studied individually, there exist no results (to our knowledge) on the behavior of the SA’s iterates when both of these strategies are used simultaneously. We have presented two simple results that together characterize the behavior of SA’s iterates under this joint scenario. The results present no surprises and assert that the averaged iterates retain the best possible convergence rate under mild stipulations on the quality of the derivative estimates and the gain sequence. The results should be seen as providing a strong theoretical basis to simultaneously apply two strategies that make sense in practice. Our treatment in this paper was limited to the context of root-finding, but extensions to the optimization context seem evident.

### REFERENCES

Billingsley, P. 1995. *Probability and Measure*. New York, NY: Wiley.

Borkar, V. S. 2008. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Cambridge, UK: Cambridge University Press.

Broadie, M., D. M. Cicek, and A. Zeevi. 2009, December. “An adaptive multidimensional version of the Kiefer-Wolfowitz Stochastic Approximation Algorithm”. In *Proceedings of the 2009 Winter Simulation Conference*, edited by M. D. Rossetti, R. R. Hill, B. Johansson, A. Dunkin, and R. G. Ingalls, 601–612. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.

Broadie, M., D. M. Cicek, and A. Zeevi. 2010. “General Bounds and Finite-Time Improvement for the Kiefer-Wolfowitz Stochastic Approximation Algorithm”. *Operations Research* 59:1211–1224. To appear.

Chung, K. L. 1954. “On a Stochastic Approximation Method”. *Annals of Mathematical Statistics* 25:463–483.

Derman, C. 1956. “An application of Chung’s lemma to the Kiefer-Wolfowitz stochastic approximation procedure”. *Annals of Mathematical Statistics* 27:532–536.

Dippon, J., and J. Renz. 1997. “Weighted means in stochastic approximation of minima”. *SIAM Journal on Control and Optimization* 35:1811–1827.

Fabian, V. 1968. “On Asymptotic Normality in Stochastic Approximation”. *Annals of Mathematical Statistics* 39:1327–1332.

Frees, E. W., and D. Ruppert. 1990. “Estimation Following a Robbins-Monro Designed Experiment”. *Journal of American Statistical Association* 85:1123–1129.

Galombos, J., and E. Seneta. 2008. “Regularly varying sequences”. *Proceedings of the American Mathematical Society* 41 (4): 110–116.



- Kushner, H. J., and G. G. Yin. 2003. *Stochastic Approximation and Recursive Algorithms and Applications*. New York, NY: Springer-Verlag.
- Mokkadem, A., and M. Pelletier. 2011. “A Generalization of the Averaging Procedure: The Use of Two-Time-Scale Algorithms”. *SIAM Journal on Control and Optimization* 49:1523.
- Nemirovski, A., A. Juditsky, G. Lan, and A. Shapiro. 2009. “Robust Stochastic Approximation Approach to Stochastic Programming”. *SIAM Journal on Optimization* 19 (4): 1574–1609.
- Pasupathy, R., and S. Kim. 2011. “The stochastic root-finding problem: overview, solutions, and open questions”. *ACM TOMACS* 21 (3): 19.
- Pasupathy, R., and B. W. Schmeiser. 2010, December. “DARTS — Dynamic Adaptive Random Target Shooting”. In *Proceedings of the 2010 Winter Simulation Conference*, edited by B. Johansson, S. Jain, J. Montoya-Torres, J. Hukan, and E. Yücesan, 12551262. Piscataway, New Jersey: Institute of Electrical and Electronics Engineers, Inc.
- Polyak, B. T., and A. B. Juditsky. 1992. “Acceleration of Stochastic Approximation by Averaging”. *SIAM Journal on Control and Optimization* 30 (4): 838–855.
- Robbins, H., and S. Monro. 1951. “A stochastic approximation method”. *Annals of Mathematical Statistics* 22:400–407.
- Spall, J. C. 1998. “Implementation of the simultaneous perturbation algorithm for stochastic optimization”. *IEEE Transactions on Aerospace and Electronic Systems* 34:817–823.
- Spall, J. C. 2000. “Adaptive stochastic approximation by the simultaneous perturbation method”. *IEEE Transactions on Automatic Control* 45:1839–1853.
- Spall, J. C. 2003. *Introduction to Stochastic Search and Optimization*. Hoboken, NJ: John Wiley & Sons, Inc.
- Spall, J. C. 2006, June. “Feedback and weighting mechanisms for improving Jacobian (Hessian) estimates in the adaptive simultaneous perturbation algorithm”. In *Proceedings of the American Control Conference*, 35–40.
- Wasan, M. T. 1969. *Stochastic Approximation*. Cambridge, UK: Cambridge University Press.
- Yousefian, F., A. Nedić, and U. Shanbhag. 2011. “On stochastic gradient and subgradient methods with adaptive steplength sequences”. *Automatica* 45:56–67.

## AUTHOR BIOGRAPHIES

**FATEMEH SADAT HASHEMI** is a PhD student in the Grado Department of Industrial and Systems Engineering at Virginia Tech. Her research interests include Monte Carlo methods and stochastic optimization. Her email address is [fatemeh@vt.edu](mailto:fatemeh@vt.edu).

**RAGHU PASUPATHY** is an associate professor in the Grado Department of Industrial and Systems Engineering at Virginia Tech. His research interests lie broadly in Monte Carlo methods with a specific focus on simulation optimization. He is a member of INFORMS, IIE, and ASA, and serves as an associate editor for *Operations Research*, *ACM TOMACS* and *INFORMS Journal on Computing*. His email address is [pasupath@vt.edu](mailto:pasupath@vt.edu) and his web page is <https://filebox.vt.edu/users/pasupath/pasupath.htm>.