

## Global Random Optimization by Simultaneous Perturbation Stochastic Approximation

John L. Maryak and Daniel C. Chin

The Johns Hopkins University, Applied Physics Laboratory  
11100 Johns Hopkins Road, Laurel, Maryland 20723-6099

john.maryak@jhuapl.edu

Phone (240) 228-4959

FAX (240) 228-1093

### ABSTRACT

A desire with iterative optimization techniques is that the algorithm reach the global optimum rather than get stranded at a local optimum value. In this paper, we examine the theoretical and numerical global convergence properties of a certain “gradient free” stochastic approximation algorithm called “SPSA,” that has performed well in complex optimization problems. We establish two theorems on the global convergence of SPSA. The first provides conditions under which SPSA will converge in probability to a global optimum using the well-known method of injected noise. The injected noise prevents the algorithm from converging prematurely to a local optimum point. In the second theorem, we show that, under different conditions, “basic” SPSA *without injected noise* can achieve convergence in probability to a global optimum. This occurs because of the noise *effectively* (and automatically) introduced into the algorithm by the special form of the SPSA gradient approximation. This global convergence without injected noise can have important benefits in the setup (tuning) and performance (rate of convergence) of the algorithm. The discussion is supported by numerical studies showing favorable comparisons of SPSA to simulated annealing and genetic algorithms.

**KEYWORDS:** *Stochastic Optimization, Global Convergence, Stochastic Approximation, Simultaneous Perturbation Stochastic Approximation (SPSA), Recursive Annealing*

### 1. INTRODUCTION

A problem of great practical importance is the problem of stochastic optimization, which may be stated as the problem of finding a minimum point,  $\theta^* \in R^P$ , of a real-valued function  $L(\theta)$ , called the “loss function,” that is observed in the presence of noise. Many approaches have been devised for numerous applications over the long history of this problem. A common desire in many applications is that the algorithm reach the global minimum rather than get stranded at a local minimum value. In this paper, we consider the popular stochastic optimization technique of stochastic approximation (SA), in particular, the form that may be called “gradient-free” SA. This refers to the case where the gradient,  $g(\theta) = \partial L(\theta) / \partial \theta$ , of the loss function is not readily available or not directly measured (even with noise). This is a common occurrence, for example, in complex systems where the exact functional relationship between the loss function value and the parameters,  $\theta$ , is not known and the loss function is evaluated by measurements on the system (or by other means, such as simulation). In such cases, one uses instead an approximation

to  $g(\theta)$  (the well-known form of SA called the Kiefer-Wolfowitz type is an example).

The usual form of this type of SA recursion is:

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

where  $\hat{g}_k(\theta)$  is an approximation (at the  $k^{\text{th}}$  step of the recursion) of the gradient  $g(\theta)$ , and  $\{a_k\}$  is a sequence of positive scalars that decreases to zero (in the standard implementation) and satisfies other properties. This form of SA has been extensively studied, and is known to converge to a local minimum of the loss function under various conditions.

Several authors (e.g., Chin (1994), Gelfand and Mitter (1991), Kushner (1987), and Styblinski and Tang (1990)) have examined the problem of *global* optimization using various forms of gradient-free SA. The usual version of this algorithm is based on using the standard “finite difference” gradient approximation for  $\hat{g}_k(\theta)$ . It is known that carefully injecting noise into the recursion based on this standard gradient can result in an algorithm that converges (in some sense) to the global minimum. For a discussion of the conditions, results, and proofs, see, e.g., Fang et al. (1997), Gelfand and Mitter (1991), and Kushner (1987). The amplitude of the injected noise is decreased over time (a process called “annealing”), so that the algorithm can finally converge when it reaches the neighborhood of the global minimum point.

A somewhat different version of SA is obtained by using a “simultaneous perturbation” gradient approximation, as described in Spall (1992) for multivariable ( $p > 1$ ) problems. The gradient approximation in simultaneous-perturbation SA (SPSA) is much faster to compute than the finite-difference approximation in multivariable problems. More significantly, using SPSA often results in a recursion that is much more economical, in terms of loss-function evaluations, than the standard version of SA. The loss function evaluations can be the most expensive part of an optimization, especially if computing the loss function requires making measurements on the physical system. Several studies (e.g., Spall (1992), Chin (1997)) have shown SPSA to be very effective in complex optimization problems. A considerable body of theory has been developed for SPSA (Spall (1992), Chin (1997), Dippon and Renz (1997), Spall (2000), and the references therein), but, because of the special form of its gradient approximation, existing theory on global convergence of standard SA algorithms is not directly applicable to SPSA. In Section 2 of this paper, we present a theorem showing that SPSA can achieve global convergence (in probability) by the technique of injecting noise. The “convergence in probability” results of our Theorem 1 (Section 2) and Theorem 2 (Section 3) are standard types of global convergence results. Several authors have shown or discussed

global convergence in probability or in distribution (Chiang *et al.* (1987), Gelfand and Mitter (1991), Gelfand and Mitter (1993), Geman and Geman (1984), Fang *et al.* (1997), Hajek (1988), Kushner (1987), Yakowitz *et al.* (2000), and Yin (1999)). Stronger “almost sure” global convergence results seem only to be available by using generally infeasible exhaustive search (Dippon and Fabian (1994)) or random search methods (Yakowitz (1993)), or for cases of optimization in a discrete ( $\theta$ -) space (Alrefaei and Andradottir (1999)).

The method of injection of noise into the recursions has proven useful, but naturally results in a relative slowing of the rate of convergence of the algorithm (e.g., Yin (1999)) due to the continued injection of noise when the recursion is near a global solution. In addition, the implementation of the extra noise terms adds to the complexity of setting up the algorithm. In Section 3, we present a theorem showing that, under different (more demanding) conditions, the basic version of SPSA can perform as a global optimizer *without* the need for injected noise. Section 4 contains numerical studies demonstrating SPSA’s performance compared to two other popular strategies for global optimization, namely, simulated annealing and genetic algorithms; and Section 5 is a summary. The Appendix provides some technical details.

## 2. SPSA WITH INJECTED NOISE AS A GLOBAL OPTIMIZER

Our first theorem applies to the following algorithm, which is the basic SPSA recursion indicated in equation (1), modified by the addition of extra noise terms:

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k \hat{g}_k(\hat{\theta}_k) + q_k \omega_k, \quad (2)$$

where  $\omega_k \in R^p$  is i.i.d.  $N(0, I)$  injected noise,  $a_k = a/k$ ,

$q_k^2 = q/k \log \log k$ ,  $a > 0$ ,  $q > 0$ , and  $\hat{g}_k(\bullet)$  is the “simultaneous perturbation” gradient defined as follows: (3)

$$\hat{g}_k(\theta) \equiv (2c_k \Delta_k)^{-1} [L(\theta + c_k \Delta_k) - L(\theta - c_k \Delta_k) + \varepsilon_k^{(+)} - \varepsilon_k^{(-)}],$$

where  $c_k, \varepsilon_k^{(\pm)}$  are scalars,  $\Delta_k \in R^p$ , and the inverse of a vector is defined to be the vector of inverses. This gradient definition follows that given in Spall (1992). The  $\varepsilon_k$  terms represent (unknown) additive noise that may contaminate the loss function observation,  $c_k$  and  $\Delta_{kl}$  are parameters of the algorithm, the  $c_k$  sequence decreases to zero, and the  $\Delta_{kl}$  components are chosen randomly according to the conditions in Spall (1992), usually (but not necessarily) from the Bernoulli ( $\pm 1$ ) distribution. (Uniformly or normally distributed perturbations are *not* allowed by the regularity conditions.)

In this Section, we will refer to Gelfand and Mitter (1991) as GM91. Our theorem on global convergence of SPSA using injected noise is based on a result in GM91. In order to state the theorem, we need to develop some notation, starting with the definition of a key probability measure,  $\pi^\eta$ , used in hypothesis H8 below. Define for any  $\eta > 0$ :

$$d\pi^\eta(\theta)/d\theta = \exp(-2L(\theta)/\eta^2)/Z^\eta, \text{ where}$$

$$Z^\eta = \int_{R^p} \exp(-2L(\theta)/\eta^2) d\theta. \text{ Next, define an important}$$

constant,  $C_0$ , for convergence theory as follows (GM91). For  $t \in R$  and  $v_1, v_2 \in R^p$ , let

$$I(t, v_1, v_2) = \inf_{\phi} \frac{1}{2} \int_0^t |d\phi(s)/ds + g(\phi(s))|^2 ds,$$

where the *inf* is taken over all absolutely continuous functions  $\phi: R \rightarrow R^p$  such that  $\phi(0) = v_1$  and  $\phi(t) = v_2$ , and  $\|\bullet\|$  is the Euclidean norm. Let  $V(v_1, v_2) = \lim_{t \rightarrow \infty} I(t, v_1, v_2)$ , and

$S_0 = \{\theta \mid g(\theta) = 0\}$ . Then

$$C_0 \equiv \frac{3}{2} \sup_{v_1, v_2 \in S_0} (V(v_1, v_2) - 2L(v_2)).$$

We will also need the following definition of *tightness*. If  $K$  is a compact subset of  $R^p$  and  $\{X_k\}$  is a sequence of random  $p$ -dimensional vectors, then  $\{X_k\}$  is tight in  $K$  if  $X_0 \in K$  and for any  $\varepsilon > 0$ , there exists a compact subset  $K_\varepsilon \subset R^p$  such that

$P(X_k \in K_\varepsilon) > 1 - \varepsilon, \forall k > 0$ . Finally, let  $\zeta_k^* \equiv \hat{g}_k(\hat{\theta}_k) - g(\hat{\theta}_k)$  and let superscript prime ( $'$ ) denote transpose.

The following are the hypotheses used in Theorem 1.

**H1.** Let  $\Delta_k \in R^p$  be a vector of  $p$  mutually independent mean-zero random variables  $\{\Delta_{k1}, \Delta_{k2}, \dots, \Delta_{kp}\}'$  such that  $\{\Delta_k\}$  is a mutually independent sequence that is also independent of the sequences  $\{\hat{\theta}_1, \dots, \hat{\theta}_{k-1}\}$ ,  $\{\varepsilon_1^{(\pm)}, \dots, \varepsilon_{k-1}^{(\pm)}\}$ , and  $\{\omega_1, \dots, \omega_{k-1}\}$ , and such that  $\Delta_{ki}$  is symmetrically distributed about zero,

$|\Delta_{ki}| \leq \alpha_1 < \infty$  a.s. and  $E|\Delta_{ki}^{-2}| \leq \alpha_2 < \infty$ , a.s.  $\forall i, k$ .

**H2.** Let  $\varepsilon_k^{(+)}$  and  $\varepsilon_k^{(-)}$  represent random measurement noise terms that satisfy  $E_k(\varepsilon_k^{(+)} - \varepsilon_k^{(-)}) = 0$  a.s.  $\forall k$ , where  $E_k$  denotes the conditional expectation given  $\mathfrak{S}_k \equiv$  the sigma algebra induced by  $\{\hat{\theta}_0, \omega_1, \dots, \omega_{k-1}, \zeta_1^*, \dots, \zeta_{k-1}^*\}$ . The  $\{\varepsilon_k^{(\pm)}\}$  sequences are not assumed independent. Assume that

$$E_k[(\varepsilon_k^{(\pm)})^2] \leq \alpha_3 < \infty \text{ a.s. } \forall k.$$

**H3.**  $L(\theta)$  is a thrice continuously differentiable map from  $R^p$  into  $R^1$ ;  $L(\theta)$  attains the minimum value of zero; as  $|\theta| \rightarrow \infty$ , we have  $L(\theta) \rightarrow \infty$  and  $|g(\theta)| \rightarrow \infty$ ;  
 $\inf(|g(\theta)|^2 - Lap(L(\theta))) > -\infty$  ( $Lap$  here is the Laplacian, i.e., the sum of the second derivatives of  $L(\theta)$  with respect to each of its components);  $L^{(3)}(\theta) \equiv \partial^3 L(\theta)/\partial \theta' \partial \theta' \partial \theta'$  exists continuously with individual elements satisfying

$$|L_{i_1 i_2 i_3}^{(3)}(\theta)| \leq \alpha_5 < \infty.$$

**H4.** The algorithm parameters have the form

$$a_k = a/k, \quad c_k = c/k^\gamma, \text{ for } k = 1, 2, \dots, \text{ where } a, c > 0, \quad q/a > C_0, \text{ and } \gamma \in [1/6, 1/2].$$

**H5.**  $[(4p-4)/(4p-3)]^{1/2} < \liminf_{|\theta| \rightarrow \infty} (g(\theta)' \theta / (|g(\theta)| \|\theta\|))$ .

**H6.**  $E_k(L(\hat{\theta}_k \pm c_k \Delta_k))^2 \leq \alpha_4 < \infty$  a.s.  $\forall k$ .

**H7.** Let  $\omega_k$  be an i.i.d.  $N(0, I)$  sequence, independent of the sequences  $\{\hat{\theta}_1, \dots, \hat{\theta}_{k-1}\}$ ,  $\{\varepsilon_1^{(\pm)}, \dots, \varepsilon_{k-1}^{(\pm)}\}$ , and  $\{\Delta_1, \dots, \Delta_{k-1}\}$ .

**H8.** For any  $\eta > 0$ ,  $Z^\eta < \infty$ ;  $\pi^\eta$  has a unique weak limit  $\pi$  as  $\eta \rightarrow 0$ .

**H9.** There exists a compact subset  $K$  of  $R^p$  such that  $\{\hat{\theta}_k\}$  is tight in  $K$ .

**Comments:**

(a) Assumptions H3, H5, and H8 correspond to assumptions (A1) through (A3) of GM91; assumptions H4 and H9 supply the hypotheses stated in GM91's Theorem 2; and the definitions of  $a_k$  and  $g_k$  given in equation (2) correspond to those used in GM91. Since we will show that assumption (A4) of GM91 is satisfied by our algorithm, this allows us to use the conclusion of their Theorem 2.

(b) The domain of  $\gamma$  given in H4 is one commonly assumed for convergence results (e.g., Spall (1992)).

We can now state our first theorem as follows:

**Theorem 1:** Under hypotheses H1 through H9,  $\hat{\theta}_k$  converges in probability to the set of global minima of  $L(\theta)$ .

**Proof:** See Maryak and Chin (1999), and the remark on convergence in probability in GM91, p. 1003.

**3. SPSSA WITHOUT INJECTED NOISE AS A GLOBAL OPTIMIZER**

As indicated in the introduction above, the injection of noise into an algorithm, while providing for global optimization, introduces some difficulties such as the need for more "tuning" of the extra terms and retarded convergence in the vicinity of the solution, due to the continued addition of noise. This effect on the rate of convergence of an algorithm using injected noise is technically subtle, but may have an important influence on the algorithm's performance. In particular, Yin (1999) shows that an algorithm of the form (2) converges at a rate proportional to

$\sqrt{\log \log(k + \text{const})}$ , while the nominal local convergence rate for

an algorithm *without* injected noise is  $k^{1/3}$ , i.e.,  $k^{1/3}(\hat{\theta}_k - \theta^*)$  converges in distribution (Spall (1992)). These rates indicate a significant difference in performance between the two algorithms.

A certain characteristic of the SPSSA gradient approximation led us to question whether SPSSA needed to use injected noise for global convergence. Although this gradient approximation tends to work very well in an SA recursion, the SPSSA gradient, evaluated at any single point in  $\theta$ -space, tends to be less accurate than the standard finite-difference gradient approximation evaluated at  $\theta$ . So, one is led to consider whether the *effective* noise introduced (automatically) into the recursion by this inaccuracy is sufficient to provide for global convergence *without* a further injection of additive noise. It turns out that *basic* SPSSA (i.e., *without* injected noise) does indeed achieve the same type of global convergence as in Theorem 1, but under a different, and more difficult to check, set of conditions.

In this Section, we designate Kushner (1987) as K87, and Kushner and Yin (1997) as KY97. Here we are working with the basic SPSSA algorithm having the same form as equation (1):

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k \hat{g}_k(\hat{\theta}_k), \tag{4}$$

where  $\hat{g}_k(\bullet)$  is the simultaneous-perturbation approximate gradient defined in Section 2, and now (obviously) no extra noise is injected into the algorithm. For use in the subsequent discussion, it will be convenient to define

$$b_k(\hat{\theta}_k) \equiv E(\hat{g}_k(\hat{\theta}_k) - g(\hat{\theta}_k) | \mathfrak{N}_k), \text{ and}$$

$$e_k(\hat{\theta}_k) \equiv \hat{g}_k(\hat{\theta}_k) - E(\hat{g}_k(\hat{\theta}_k) | \mathfrak{N}_k),$$

where  $\mathfrak{N}_k$  denotes the  $\sigma$ -algebra generated by  $\{\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k\}$ , which allows us to write equation (4) as

$$\hat{\theta}_{k+1} = \hat{\theta}_k - a_k [g(\hat{\theta}_k) + e_k(\hat{\theta}_k) + b_k(\hat{\theta}_k)]. \tag{5}$$

Another key element in the subsequent discussion is the ordinary differential equation (ODE):

$$\dot{\theta} = g(\theta), \tag{6}$$

which, in Lemma 1 of the Appendix is shown to be the "limit mean" ODE for algorithm (4).

Now we can state our assumptions for Theorem 2, as follows:

**J1.** Let  $\Delta_k \in R^p$  be a vector of  $p$  mutually independent mean-zero random variables  $\{\Delta_{k1}, \Delta_{k2}, \dots, \Delta_{kp}\}'$  such that  $\{\Delta_k\}$  is a mutually independent sequence and  $\Delta_k$  is independent of the sequence  $\{\hat{\theta}_1, \dots, \hat{\theta}_{k-1}\}$ , and such that  $\Delta_{ki}$  is  $\forall i, k$  symmetrically distributed about zero,  $|\Delta_{ki}| \leq \alpha_1 < \infty$  a.s. and  $E|\Delta_{ki}^{-2}| \leq \alpha_2 < \infty$ .

**J2.** Let  $\varepsilon_k^{(+)}$  and  $\varepsilon_k^{(-)}$  represent random measurement noise terms that satisfy  $E((\varepsilon_k^{(+)} - \varepsilon_k^{(-)}) | \mathfrak{N}_k) = 0$  a.s.  $\forall k$ . The  $\{\varepsilon_k^{(\pm)}\}$  sequences need not be assumed independent. Assume that  $E((\varepsilon_k^{(\pm)})^2 | \mathfrak{N}_k) \leq \alpha_3 < \infty$  a.s.  $\forall k$ .

**J3(a).**  $L(\theta)$  is thrice continuously differentiable and the individual elements of the third derivative satisfy

$$|L_{i_1 i_2 i_3}^{(3)}(\theta)| \leq \alpha_5 < \infty.$$

(b).  $|L(\theta)| \rightarrow \infty$  as  $|\theta| \rightarrow \infty$ .

**J4.** The algorithm parameters satisfy the following: the gains

$$a_k > 0, a_k \rightarrow 0 \text{ as } k \rightarrow \infty, \text{ and } \sum_{k=1}^{\infty} a_k = \infty.$$

The sequence  $\{c_k\}$  is of form  $c_k = c/k^\gamma$ , where  $c > 0$  and  $\gamma \in [1/6, 1/2)$ , and

$$\sum_{k=0}^{\infty} (a_k / c_k)^2 < \infty.$$

**J5.** The gradient  $g(\theta)$  is bounded and Lipschitz continuous.

**J6.** The ODE (6) has a unique solution for each initial condition.

**J7.** For the ODE (6), suppose that there exists a finite collection of disjoint compact stable invariant sets (see K87)  $K_1, K_2, \dots, K_m$ ,

such that  $\bigcup_i K_i$  contains all the limit sets for (6). These sets are interpreted as closed sets containing all local (including global) minima of the loss function.

**J8.** For any  $\eta > 0$ ,  $Z^\eta < \infty$ ;  $\pi^\eta$  has a unique weak limit  $\pi$  as

$$\eta \rightarrow 0 \text{ (} Z^\eta \text{ and } \pi^\eta \text{ are defined in Section 2)}.$$

**J9.**  $E |\sum_{i=1}^k e_i(\hat{\theta}_i)| < \infty \forall k$ .

**J10.** For any asymptotically stable (in the sense of Liapunov) point,  $\bar{\theta}$ , of the ODE (6), there exists a neighborhood of the origin in  $R^P$  such that the closure,  $Q_2$ , of that neighborhood satisfies

$\bar{\theta} + Q_2 \equiv \{\bar{\theta} + y : y \in Q_2\} \subset \Theta$ , where  $\Theta \subset R^P$  denotes the allowable  $\theta$ -region. There is a neighborhood,  $Q_1$ , of the origin in

$R^P$  and a real-valued function  $H_1(\psi_1, \psi_2)$ , continuous in  $Q_1 \times Q_2$ , whose  $\psi_1$ -derivative is continuous on  $Q_1$  for each fixed  $\psi_2 \in Q_2$ , and such that the following limit holds. For any  $\chi, \Delta > 0$ , with  $\chi$  being an integral multiple of  $\Delta$ , and any functions  $\psi_1(\bullet), \psi_2(\bullet)$  taking values in  $Q_1 \times Q_2$  and being constant on the intervals  $[i\Delta, i\Delta + \Delta)$ ,  $i\Delta < \chi$ , we have

$$\int_0^\chi H_1(\psi_1(s), \psi_2(s)) ds = \quad (7)$$

$$\limsup_{m,n} \frac{\Delta}{m} \log E \exp \left[ \sum_{i=0}^{(\chi/\Delta)-1} \psi_1'(i\Delta) \sum_{j=im}^{im+m-1} b_{n+j}(\hat{\theta}_{n+j}) \right].$$

Also, there is a function  $H_2(\psi_3)$  that is continuous and differentiable in a small neighborhood of the origin, and such that

$$\int_0^\chi H_2(\psi_1(s)) ds = \quad (8)$$

$$\limsup_{m,n} \frac{\Delta}{m} \log E \exp \left[ \sum_{i=0}^{(\chi/\Delta)-1} \psi_1'(i\Delta) \sum_{j=im}^{im+m-1} e_{n+j}(\hat{\theta}_{n+j}) \right].$$

A bit more notation is needed. Let  $T > 0$  be interpreted such that  $[0, T]$  is the total time period under consideration in ODE (6). Let

$$\bar{H}(\psi_1, \psi_2) = 0.5[H_1(2\psi_1, \psi_2) + H_2(2\psi_1)],$$

$$\bar{L}(\beta, \psi_2) = \sup_{\psi_1} [\psi_1'(\beta - g(\psi_2)) - \bar{H}(\psi_1, \psi_2)],$$

and, for  $\phi(0) = x \in R^1$ , define the function

$$S(T, \phi) = \int_0^T \bar{L}(\dot{\phi}(s), \phi(s)) ds,$$

if  $\phi(\bullet)$  is a real-valued absolutely-continuous function on  $[0, T]$  and to take the value  $\infty$  otherwise.  $S(T, \phi)$  is the usual action functional of the theory of large deviations (adapted to our context). Define  $t_n \equiv \sum_{i=0}^{n-1} a_i$ , and  $t_k^n \equiv \sum_{i=0}^{k-1} a_{n+i}$ . Define

$\{\hat{\theta}_k^n\}$  and  $\theta^n(\bullet)$  by

$$\hat{\theta}_0^n = x \in \Theta, \quad \hat{\theta}_{k+1}^n = \hat{\theta}_k^n - a_{n+k} \hat{g}_{n+k}(\hat{\theta}_k^n), \text{ and}$$

$$\theta^n(t) = \hat{\theta}_k^n \text{ for } t \in [t_k^n, t_{k+1}^n).$$

Now we can state the last two assumptions for Theorem 2:

**J11.** For each  $\delta > 0$  and  $i = 1, 2, \dots, m$ , there is a  $\rho$ -neighborhood of  $K_i$ , denoted  $N_\rho(K_i)$ , and  $\delta_\rho > 0, T_\rho < \infty$  such that, for each  $x, y \in N_\rho(K_i)$ , there is a path,  $\phi(\bullet)$ , with  $\phi(0) = x, \phi(T_y) = y$ , where  $T_y \leq T_\rho$  and  $S(T_\rho, \phi) \leq \delta$ .

**J12.** There is a sphere,  $D_1$ , such that  $D_1$  contains  $\bigcup_i K_i$  in its

interior, and the trajectories of  $\theta^n(\bullet)$  stay in  $D_1$ . All paths of ODE (6) starting in  $\bar{D}_1$  stay in  $D_1$ .

**Note 1.** Assumptions J1, J2, and J3(a) are from Spall (1992), and are used here to characterize the noise terms  $b_k(\hat{\theta}_k)$  and  $e_k(\hat{\theta}_k)$ . Assumption J3(b) is used on page 178 of K87. Assumption J4 expresses standard conditions on the algorithm parameters (see Spall (1992)), and implies hypothesis (A10.2) in KY97, p. 174. Assumptions J5 and J6 correspond to hypothesis (A10.1) in KY97, p. 174. Assumption J7 is from K87, p. 175. Assumption J8 concerns the limiting distribution of  $\hat{\theta}_k$ . Assumption J9 is used to establish the ‘‘mean’’ criterion for the martingale sequence in Lemma 2. Assumptions J11 and J12 are the ‘‘controllability’’ hypothesis A4.1 and the hypothesis A4.2, respectively, of K87, p. 176.

**Note 2.** Assumption J10 corresponds to hypotheses (A10.5) and (A10.6) in KY97, pp. 179-181. Although these hypotheses are standard forms for this type of large deviation analysis, it is important to justify their reasonableness. The first part (equation (7), involving noise terms  $b_k(\hat{\theta}_k)$ ) of J10 is justified by the discussion in KY97, p. 174, which notes that the results of their subsection 6.10 are valid if the noise terms (that they denote  $\xi_n$ ) are bounded. This discussion is applicable to our algorithm since the  $b_k(\hat{\theta}_k)$  noise terms were shown by Spall (1992) to be  $O(c_k^2)$  ( $c_k \rightarrow 0$ ) a.s. The second part (equation (8), involving noise terms  $e_k(\hat{\theta}_k)$ ) is justified by the discussion in KY97, p. 174, which notes that the results in their subsection 6.10 are valid if the noise terms they denote  $\delta M_n$  (corresponding to our noise terms  $e_k(\hat{\theta}_k)$ ) satisfy the martingale difference property that we have established in Lemma 2 of the Appendix.

Now we can state our main theorem:

**Theorem 2.** Under assumptions J1 through J12,  $\hat{\theta}_k$  converges in probability to the set of global minima of  $L(\theta)$ .

The idea of the proof is as follows (see the Appendix for the details). This theorem follows from results (in a different context) in K87 for an algorithm  $\hat{\theta}_{k+1} = \hat{\theta}_k - a_k [g(\hat{\theta}_k) + \zeta_k]$ , where  $\zeta_k$  is i.i.d. Gaussian (injected) noise. In order to prove our Theorem 2, we start by writing the SPSA recursion as  $\hat{\theta}_{k+1} = \hat{\theta}_k - a_k [g(\hat{\theta}_k) + \zeta_k^*]$ , where  $\zeta_k^* \equiv \hat{g}_k(\hat{\theta}_k) - g(\hat{\theta}_k)$  is the ‘‘effective noise’’ introduced by the inaccuracy of the SPSA gradient approximation. So, our algorithm has the same form as that in K87. However, since  $\zeta_k^*$  is not i.i.d. Gaussian, we cannot use K87’s result directly. Instead, we use material in Kushner and Yin (1997) to establish a key ‘‘large deviation’’ result related to our algorithm (4), which allows the result in K87 to be used with  $\zeta_k^*$  replacing the  $\zeta_k$  in his algorithm.

## 4. NUMERICAL STUDIES: SPSA WITHOUT INJECTED NOISE

### 4.1. Two-Dimensional Problem

A study was done to compare the performance of SPSA to a recently published application of the popular genetic algorithm (GA). The loss function is the well-known Griewank function (see Haataja (1999)) defined for a two-dimensional  $\theta = (t_1, t_2)'$ , by:

$$L(\theta) = \cos(t_1 - 100) \cos[(t_2 - 100) / \sqrt{2}] - [(t_1 - 100)^2 + (t_2 - 100)^2] / 4000 - 1,$$

which has thousands of local minima in the vicinity of a single global minimum at  $\theta = (100, 100)'$  at which  $L(\theta) = 0$ . Haataja (1999) describes the application of a GA to this function (actually, to find the *maximum* of  $-L(\theta)$ ) based on noise-free evaluations of  $L(\theta)$  (i.e.,  $\varepsilon_k = 0$ ). This study achieved a success rate of 66% (see Haataja's Table 1.3, p.16) in 50 independent trials of the GA, using 300 generations and 9000  $L(\theta)$  evaluations in each run of the GA. Haataja's definition of a successful solution is a reported solution where the norm of the solution minus the correct value,  $\theta^*$ , is less than 0.2, and the value of the loss function at the reported solution is within 0.01 of the correct value of zero. We examined the performance of basic SPSA (without adding injected noise) on this problem, using  $a_k = a / (A + k)^\alpha$ , with  $A = 60$ ,  $a = 100$  and  $\alpha = .602$ , a slowly decreasing gain sequence of a form that has been used in many applications (see Spall (1998)). For the gradient approximation (equation (3)), we chose each component of  $\Delta_k$  to be an independent sample from a

Bernoulli ( $\pm 1$ ) distribution, and  $c_k = c / k^\gamma$ , with  $c = 10$  and  $\gamma = .101$ . Since we used the exact loss function, the  $\varepsilon_k$  noise terms were zero. We ran SPSA, allowing 3000 function evaluations in each of 50 runs, and starting the algorithm (each time) at a point randomly chosen in the domain  $[-200, 400] \times [-200, 400]$ . Haataja's  $\theta$ -domain was also constrained to lie in a box, but the dimensions of the box were not specified. Hence we chose a domain that is a cube centered at the global minimum, in which there are many local minima of  $L(\theta)$  (as seen in Haataja's (1999) Figure 1.1). SPSA successfully located the global minimum in all 50 runs (100% success rate).

### 4.2. Ten-Dimensional Problem

For a more ambitious test of the global performance of SPSA, we applied SPSA to a loss function given in Example 6 of Styblinski and Tang (1990), which we will designate for convenience as ST90. The loss function is:

$$L(\theta) = (2p)^{-1} \sum_{i=1}^p t_i^2 - 4p \prod_{i=1}^p \cos(t_i),$$

where  $p = 10$  and  $\theta = (t_1, \dots, t_p)'$ . This function has the global minimum value of  $-40$  at the origin, and a large number of local minima. As in the two-dimensional study above, we used the exact loss function. Our goal is to compare the performance of SPSA without injected noise with simulated annealing and with a GA.

For the simulated annealing algorithm, we use the results reported in ST90. They used an advanced form of simulated

annealing called fast simulated annealing (FSA). According to ST90, FSA has proven to be much more efficient than classical simulated annealing due to using Cauchy (rather than Gaussian) sampling and using a fast (inversely linear in time) cooling scheme. For more details on FSA, see ST90. The results of their application of FSA to the above  $L(\theta)$  are given in Table 1 below (FSA values taken from Table 10 of ST90). Table 1 shows the results of 10 independent runs of each algorithm. In each case (each run of each algorithm), the best value of  $L(\theta)$  found by the algorithm is shown. In their study, although FSA was allowed to use 50,000 function evaluations for each of the runs, the algorithm showed very limited success in locating the global minimum. It should be noted that the main purpose of the ST90 paper was to examine a relatively new algorithm, stochastic approximation combined with convolution smoothing. This algorithm, which they call SAS, was much more effective than FSA, yielding results between those shown in Table 1 for GA and SPSA.

For the genetic algorithm (GA), we implemented a GA using the popular features of elitism (elite members of the old population pass unchanged into the new population), tournament selection (tournament size = 2), and real-number encoding (see Mitchell (1996), pp. 168, 170, and 157, respectively). After considerable experimentation, we found the following settings for the GA algorithm to provide the best performance on this problem. The population size was 100, the number of elite members (those carried forward unchanged) in each generation was 10, the crossover rate was 0.8, and mutation was accomplished by adding a Gaussian random variable with mean zero and standard deviation 0.01 to each component of the offspring. The original population of 100 (10-dimensional)  $\theta$ -vectors was created by uniformly randomly generating points in the 10-dimensional hypercube centered at the origin, with edges of length 6 (so, all components had absolute value less than or equal to 3 radians). We constrained all component values in subsequent generations to be less than or equal to 4.5 in absolute value. This worked a bit better than constraining them to be less than 3, since, with the tighter constraints, the GA got stuck at the constraint boundary and could not reach local minima that were just over the boundary. All runs of the GA algorithm reported here used 50,000 evaluations of the loss function. The results of the 10 independent runs of GA are shown in Table 1. Although the algorithm did reasonably well in getting close to the minimum loss value of  $-40$ , it only found the global minimum in one of the 10 runs (run #8). In the other nine cases, a few (typically two or four) of the components were trapped in a local minimum (around  $\pm \pi$  radians), while the rest of the components (approximately) achieved the correct value of zero. Note that the nature of the loss function is such that the value of  $L(\theta)$  is very close to an integer (e.g.,  $-39.0$  or  $-38.0$ ) when an even number (e.g., 2 or 4) of components of  $\theta$  are near  $\pm \pi$  radians.

We examined the performance of basic SPSA (without adding injected noise), using the algorithm parameters defined in Subsection 4.1 with  $A = 60$ ,  $a = 1$ ,  $\alpha = .602$ ,  $c = 2$ , and  $\gamma = .101$ . We started  $\theta$  at  $t_i = 3$  radians,  $i = 1, \dots, p$ , resulting in an initial loss function value of  $-31$ . This choice of starting point was at the outer boundary of the domain in which we chose initial values for the GA algorithm, and we did not constrain the search space for SPSA as we did for GA (the initialization and search

space for FSA were not reported in ST90). We ran 10 Monte Carlo trials (i.e., randomly varying the choices of  $\Delta_k$ ). The results are tabulated in Table 1. The results of these numerical studies show a strong performance of the basic SPSA algorithm in difficult global optimization problems.

**Table 1. Best Loss Function Value in Each of 10 Independent Runs of Three Algorithms**

Run	SPSA	GA	FSA
1	-40.0	-38.0	-24.9
2	-40.0	-39.0	-15.5
3	-40.0	-39.0	-29.0
4	-40.0	-38.0	-32.1
5	-40.0	-37.0	-30.2
6	-40.0	-39.0	-30.1
7	-40.0	-38.0	-27.9
8	-40.0	-40.0	-20.9
9	-40.0	-38.0	-28.5
10	-40.0	-39.0	-34.6
Average Value	-40.0	-38.5	-27.4
Number of Function Evaluations	2,500	50,000	50,000

## 5. SUMMARY

SPSA is an efficient gradient-free SA algorithm that has performed well on a variety of complex optimization problems. We showed in Section 2 that, as with some standard SA algorithms, adding injected noise to the basic SPSA algorithm can result in a global optimizer. More significantly, in Section 3, we showed that, under certain conditions, the basic SPSA recursion can achieve global convergence *without the need for injected noise*. The use of basic SPSA as a global optimizer can ease the implementation of the global optimizer (no need to tune the injected noise) and result in a significantly faster rate of convergence (no extra noise corrupting the algorithm in the vicinity of the solution). In the numerical studies, we found significantly better performance of SPSA as a global optimizer than for the popular simulated annealing and genetic algorithm methods, which are often recommended for global optimization. In particular, in the case of a 10-dimensional optimization parameter ( $\theta$ ), the fast simulated annealing and genetic algorithms generally failed to find the global solution.

## APPENDIX (LEMMA'S RELATED TO THEOREM 2 AND PROOF OF THEOREM 2)

In this Appendix, we designate Kushner (1987) as K87, and Kushner and Yin (1997) as KY97. Here we are working with the basic SPSA algorithm as defined in equation (4):

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

We first establish an important preliminary result that is needed in order to apply the results from K87 and KY97 in the proof of Theorem 2.

**Lemma 1.** The ordinary differential equation (eq. (6) above),

$$\dot{\theta} = g(\theta),$$

is the “limit mean ODE” for algorithm (4).

**Proof:** Examining the definition of limit mean ODE given in KY97, pp. 174 & 138, it is clear that we need to prove that  $\frac{1}{m} \sum_{k=n}^{m+n-1} [g(\theta) + e_k(\hat{\theta}_k) + b_k(\hat{\theta}_k)] \rightarrow g(\theta)$  w.p. 1 as  $m, n \rightarrow \infty$ . Since Spall (1992) has shown that  $b_k(\hat{\theta}_k) \rightarrow 0$  w.p. 1, we can conclude using Cesaro summability that the contribution of the  $b_k(\hat{\theta}_k)$  terms to the limit is zero w.p. 1. For the  $e_k(\hat{\theta}_k)$  terms, we have by definition that  $E[e_k(\hat{\theta}_k)] = 0$ ; hence, by the law of large numbers, the contribution of the  $e_k(\hat{\theta}_k)$  terms to the limit is also zero. Q.E.D.

Our next Lemma relates to Note 2 in Section 3.

**Lemma 2.** Under assumptions J1, J3(a), and J9, the sequence  $\{e_k(\hat{\theta}_k)\}$  is a  $\mathfrak{N}_k$ -martingale difference.

**Proof:** It is sufficient to show that  $M_k \equiv \sum_{i=1}^k e_k(\hat{\theta}_k)$  is a  $\mathfrak{N}_k$ -martingale. Assumption J9 satisfies the first requirement (see KY97, p.68) of the martingale definition, that  $E[M_k] < \infty$ . For the main requirement, we have for any  $k$ :

$$\begin{aligned} E[M_{k+1} | M_k, \dots, M_1] &= E[e_{k+1}(\hat{\theta}_{k+1}) + M_k | M_k, \dots, M_1] \\ &= M_k + E[e_{k+1}(\hat{\theta}_{k+1}) | M_k, \dots, M_1] \\ &= M_k + E\{[\hat{g}_{k+1}(\hat{\theta}_{k+1}) - E(\hat{g}_{k+1}(\hat{\theta}_{k+1}) | \hat{\theta}_{k+1})] | M_k\} = \\ &= M_k + E_{\hat{\theta}_{k+1}} E\{[\hat{g}_{k+1}(\hat{\theta}_{k+1}) - E(\hat{g}_{k+1}(\hat{\theta}_{k+1}) | \hat{\theta}_{k+1})] | M_k, \hat{\theta}_{k+1}\} \\ &= M_k + E_{\hat{\theta}_{k+1}} E(\hat{g}_{k+1}(\hat{\theta}_{k+1}) | M_k, \hat{\theta}_{k+1}) - \\ &E_{\hat{\theta}_{k+1}} E[E(\hat{g}_{k+1}(\hat{\theta}_{k+1}) | M_k, \hat{\theta}_{k+1})] = M_k, \end{aligned}$$

where  $E_{\hat{\theta}_{k+1}}$  denotes expectation conditional on  $\hat{\theta}_{k+1}$ , and all equalities concerning conditional expectations are w.p. 1. Q.E.D.

A key step in the proof of our main result (Theorem 2 below) is establishing the following “large deviation” result (Lemma 3). Let  $B_x$  be a set of continuous functions on  $[0, T]$  taking values in  $\Theta$  and with initial value  $x$ . Let  $B_x^0$  denote the interior of  $B_x$ , and  $\bar{B}_x$  denote the closure.

**Lemma 3.** Under assumptions J4, J5, J6, and J10, we have

$$\begin{aligned} - \inf_{\phi \in B_x^0} S(T, \phi) &\leq \liminf_n \log P_x^n \{\theta^n(\bullet) \in B_x\} \\ &\leq \limsup_n \log P_x^n \{\theta^n(\bullet) \in B_x\} \leq - \inf_{\phi \in \bar{B}_x} S(T, \phi), \quad (9) \end{aligned}$$

where  $P_x^n$  denotes the probability under the condition that

$$\theta^n(0) = x.$$

**Proof:** This result is adapted from Theorem 10.4 in KY97, p. 181. Note that our assumption J10 is a modified form of their assumptions (A10.5) and (A10.6), using “equals” signs rather than inequalities. The two-sided inequality in (9) follows from J10 by an argument analogous to the proof of KY87’s Theorem 10.1 (p. 178), which uses an “equality” assumption ((A10.4), p. 174) to arrive at a two-sided large deviation result analogous to (9) above. Q.E.D.

We restate our main theorem:

**Theorem 2:** Under hypotheses J1 through J12,  $\hat{\theta}_k$  converges in probability to the set of global minima of  $L(\theta)$ .

**Proof:** This result follows from a discussion in K87. Theorem 2 of K87, (p.177) describes probabilities involving expected times for the SA algorithm (system (1.1) of K87) to transition from one  $K_i$  to another. The SA algorithm he uses can be written in our notation as  $\hat{\theta}_{k+1} = \hat{\theta}_k - a_k [g(\hat{\theta}_k) + \zeta_k]$ , where  $\zeta_k$  is i.i.d. Gaussian (injected) noise. The K87 Theorem 2 uses the i.i.d. Gaussian assumption only to arrive at a large deviation result exactly analogous to our Lemma 3. The subsequent results in K87 are based on this large deviation result. Recall that the SPSA algorithm without injected noise can be written in the form  $\hat{\theta}_{k+1} = \hat{\theta}_k - a_k [g(\hat{\theta}_k) + \zeta_k^*]$ . Since we have established Lemma 3 for SPSA, the results of K87 hold for the SPSA algorithm with its “effective” noise  $\{\zeta_k^*\}$  replacing the  $\{\zeta_k\}$  sequence used in K87. In particular, K87’s discussion (pp. 178, 179) of his Theorem 2 is applicable to our Theorem 2 context (SPSA without injected noise), which corresponds to K87’s “potential case.” Note that our formulation corresponds to the K87 setup where  $b(x, \xi) = \bar{b}(x)$  in his notation, which, by the comment in K87, p. 179, means that his discussion is applicable to his system (1.1) and hence to our setup. In his discussion on p. 179, K87 indicates that the difference between the measure of  $X_n$  (which corresponds to our  $\hat{\theta}_k$ ) and the invariant measure (which we have denoted  $\pi^n$ ) converges asymptotically ( $n, k \rightarrow \infty, \eta \rightarrow 0$ ) to the zero measure weakly. This means that, in the limit as  $k \rightarrow \infty$ ,  $\hat{\theta}_k$  is equivalent to  $\pi$  in the same sense as in Theorem 2 of Gelfand and Mitter (1991), and the desired convergence in probability follows as in Theorem 1 above. Q.E.D.

#### ACKNOWLEDGEMENT

This work was partially supported by the JHU/APL Independent Research and Development Program and U.S. Navy contract N00024-98-D-8124.

#### REFERENCES

- Alrefaie, M.H. and Andradottir, S. (1999), “A Simulated Annealing Algorithm with Constant Temperature for Discrete Stochastic Optimization,” *Management Science*, **45**, pp. 748-764.
- Chaing T-S., Hwang, C-R., and Sheu, S-J. (1987), “Diffusion for Global Optimization in  $R^n$ ,” *SIAM J. Control Optim.*, **25**, pp. 737-753.
- Chin, D.C. (1997), “Comparative Study of Stochastic Algorithms for System Optimization Based on Gradient Approximations,” *IEEE Trans. Systems, Man, and Cybernetics – Part B: Cybernetics*, **27**, pp. 244-249.
- Chin, D.C. (1994), “A More Efficient Global Optimization Algorithm Based on Styblinski and Tang,” *Neural Networks*, **7**, pp. 573-574.
- Dippon, J. and Fabian, V. (1994), “Stochastic Approximation of Global Minimum Points,” *J. Statistical Planning and Inference*, **41**, pp. 327-347.
- Dippon, J. and Renz, J. (1997), “Weighted Means in Stochastic Approximation of Minima,” *SIAM J. Control Optim.*, **35**, pp. 1811-1827.
- Fang, H., Gong, G., and Qian, M. (1997), “Annealing of Iterative Stochastic Schemes,” *SIAM J. Control Optim.*, **35**, pp.1886-1907.
- Gelfand, S.B. and Mitter, S.K. (1991), “Recursive Stochastic Algorithms for Global Optimization in  $R^d$ ,” *SIAM J. Control Optim.*, **29**, pp. 999-1018.
- Gelfand, S.B. and Mitter, S.K. (1993), “Metropolis-Type Annealing Algorithms for Global Optimization in  $R^d$ ,” *SIAM J. Control Optim.*, **31**, pp. 110-131.
- Geman S. and Geman, D. (1984), “Stochastic Relaxation, Gibbs Distributions, and the Bayesian Restoration of Images,” *IEEE Trans. Pattern Analysis and Machine Intelligence*, **PAMI-6**, pp. 721-741.
- Haataja, J. (1999), “Using Genetic Algorithms for Optimization: Technology Transfer in Action,” Chapter 1, pp. 3 – 22, in *Evolutionary Algorithms in Engineering and Computer Science*, Edited by K. Miettinen, M.M. Makela, P. Neittaanmaki, and J. Periaux, Wiley, Chichester.
- Hajek, (1988), “Cooling Schedules for Optimal Annealing,” *Mathematics of Operations Research*, **13**, pp. 311-329.
- Kushner, H.J. and Yin, G.G. (1997), *Stochastic Approximation Algorithms and Applications*, Springer, New York.
- Kushner, H.J. (1987), “Asymptotic Global Behavior for Stochastic Approximation and Diffusions with Slowly Decreasing Noise Effects: Global Minimization Via Monte Carlo,” *SIAM J. Appl. Math.*, **47**, pp. 169-185.
- Maryak, J.L. and Chin, D.C. (1999), “Efficient Global Optimization Using SPSA,” *Proc. Amer. Control Conf.*, San Diego, June 2-4, pp. 890-894.
- Mitchell, M. (1996), *An Introduction to Genetic Algorithms*, MIT Press, Cambridge, Mass.
- Spall, J.C. (2000), “Adaptive Stochastic Approximation by the Simultaneous Perturbation Method,” *IEEE Trans. Automat. Control*, **45**, pp. 1839-1853.
- Spall, J.C. (1998), “Implementation of the Simultaneous Perturbation Algorithm for Stochastic Optimization,” *IEEE Trans. Aerospace and Electronic Systems*, **34**, pp. 817-823.
- Spall, J.C. (1992), “Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation,” *IEEE Trans. Automat. Control*, **37**, pp. 332-341.
- Styblinski, M.A. and Tang, T.-S. (1990), “Experiments in Nonconvex Optimization: Stochastic Approximation with Function Smoothing and Simulated Annealing,” *Neural Networks*, **3**, pp. 467-483.
- Yakowitz, S. (1993), “A Globally Convergent Stochastic Approximation,” *SIAM J. Control Optim.* **31** pp. 30-40.
- Yakowitz, S., L’Ecuyer, P., and Vazquez-Abad, F. (2000), “Global Stochastic Optimization with Low-Dispersion Point Sets,” *Operations Research*, **48**, pp. 939-950.
- Yin, G. (1999), “Rates of Convergence for a Class of Global Stochastic Optimization Algorithms,” *SIAM J. Optim.*, **10**, pp. 99-120.