

Discrete Simultaneous Perturbation Stochastic Approximation on Loss Function with Noisy Measurements

Qi Wang and James C. Spall

Abstract— Consider the stochastic optimization of a loss function defined on p -dimensional grid of points in Euclidean space. We introduce the middle point discrete simultaneous perturbation stochastic approximation (DSPSA) algorithm for such discrete problems and show that convergence to the minimum is achieved. Consistent with other stochastic approximation methods, this method formally accommodates noisy measurements of the loss function.

Keywords— Stochastic optimization; recursive estimation; DSPSA; noisy data; discrete optimization.

I. INTRODUCTION

THE optimization of real-world stochastic systems typically involves the use of a mathematical algorithm that iteratively seeks out the solution. It is often the case that the domain of optimization is discrete. Resource allocation, for instance, involves the distribution of discrete amount of some resource to a finite number of users in the face of uncertainty; other problems of interest within this framework include weapons assignment, plant location, network resource and experimental design. This paper introduces a method for stochastic discrete optimization that is based on stochastic approximation techniques customarily used in continuous optimization problems.

Many methods have been proposed to deal with discrete optimization problems. These methods include random search [2], simulated annealing [1], stochastic comparison [7], ordinal optimization [11], nested partitions [18]. Recently Hannah and Powell [8] propose an algorithm for one-stage stochastic combinatorial optimization problems, based on evolutionary policy iteration. Li et al. [13] introduce a method based on random search in the most promising area proposed in [12]. And Sklenar [19] considers an exhaustive local search method which is designed explicitly for noisy loss.

The aim here is to present an alternative method that can fully use the information of the structure of objective functions (e.g. “gradient”) and potentially involve fewer function measurements. The simultaneous perturbation stochastic approximation (SPSA) algorithm [20, 21] was

developed for continuous optimization problems of high dimension and where the loss function is expensive to evaluate. SPSA is a popular algorithm that creates gradient-type information from only two noisy function measurements in each iteration. The increase in efficiency over the finite difference stochastic approximation method, for example, has been shown to be a factor equal to the dimension of the problem [20]. Spall [20] has considered the convergence of SPSA for three times differentiable functions, whereas He et al. [9] have analyzed the convergence for nondifferentiable, but continuous optimization. Also Yousefian et al. [23] have discussed a local randomized smoothing technique for convex nondifferentiable continuous stochastic optimization. We want to use a similar idea of SPSA for the discrete case. Because the usual notion of a gradient does not apply in discrete problems, it is not obvious that the convergence properties demonstrated for SPSA hold for the discrete case. Hill et al. [10] considers a discrete form of SPSA and develops preliminary results associated with convergence for a separable discrete loss function under special conditions. However, this algorithm can be shown to not converge to the optimal solution in simple examples. We introduce a different form of discrete algorithm that applies to a broader range of problems while potentially retaining the essential efficiency advantages of standard SPSA.

In particular, we introduce a middle point discrete simultaneous perturbation stochastic approximation (DSPSA) algorithm that applies in a class of discrete problems. As in conventional SPSA, the method needs only two noisy measurements of the loss function at each iteration. Although a full convergence and convergence rate analysis has not yet been conducted, we show conditions for almost sure convergence of the algorithm to the true parameter value.

The paper is organized as the follows. In Section II, we motivate the general approach by considering the case of one dimensional θ , and describe the basic DSPSA algorithm for general $p \geq 1$. In Section III, we show that the algorithm converges to the optimal solution for some class of function under some conditions. In Section IV, we show how this algorithm compares with the localized random search method in two examples. In Section V, we conclude with a discussion.

Qi Wang is with the Department of Applied Mathematics and Statistics of the Johns Hopkins University, Baltimore, MD 21218 USA (e-mail: qwang29@jhu.edu).

James C. Spall is with The Johns Hopkins University, Applied Physics Laboratory, Laurel, MD 20723-6099 USA and with the Department of Applied Mathematics and Statistics of the Johns Hopkins University, Baltimore, MD 21218 USA (e-mail: james.spall@jhuapl.edu).

II. PROBLEM FORMULATION

A. Motivation: One Dimension Case

Let us first consider one dimensional discrete function $L: \mathbb{Z} \rightarrow \mathbb{R}$, where \mathbb{Z} denotes the set of integers $\{\dots, -2, -1, 0, 1, 2, \dots\}$. We want to find the minimal solution of the loss function L . Let the noisy measurement of the loss function be y , where $y = L + \varepsilon$ and ε indicates the noise. Fig. 1 shows an example of a discrete function in one dimension with a line connecting the neighboring integer points. The line \bar{L} can be regarded as a continuous extension of L , but \bar{L} is a nondifferentiable function at the integer points. For a point $\theta \in \mathbb{R} \setminus \mathbb{Z}$, the gradient is

$$\begin{aligned} g(\theta) &= L(\lceil \theta \rceil) - L(\lfloor \theta \rfloor) \\ &= L\left(\pi(\theta) + \frac{1}{2}\right) - L\left(\pi(\theta) - \frac{1}{2}\right) \\ &= \frac{L\left(\pi(\theta) + \frac{1}{2}\Delta\right) - L\left(\pi(\theta) - \frac{1}{2}\Delta\right)}{\Delta}, \end{aligned}$$

where $\lfloor \cdot \rfloor$ is the floor function, $\lceil \cdot \rceil$ is the ceiling function, and $\pi(\theta) = (2\lfloor \theta \rfloor + 1)/2$ is the middle point between $\lfloor \theta \rfloor$ and $\lceil \theta \rceil$, and Δ is a Bernoulli random variable taking the values ± 1 . Actually $\pi(\theta) = n/2$ and n is an odd number, so $\pi(\theta) \pm \Delta/2$ must be integers. We can see that $g(\theta)$ is also well defined at integer points θ , and it is a subgradient (a vector γ is a subgradient of $L(\cdot)$ at θ if $L(\mu) - L(\theta) \geq \gamma^T(\mu - \theta)$ for all $\mu \in \mathbb{R}^p$) at θ ($\pi(\theta)$ is now the middle point between $\theta = \lfloor \theta \rfloor$ and $\theta + 1$). Then the estimated gradient for noisy function is

$$\hat{g}(\theta) = \frac{y\left(\pi(\theta) + \frac{1}{2}\Delta\right) - y\left(\pi(\theta) - \frac{1}{2}\Delta\right)}{\Delta}.$$

Ref. [9] has shown that SPSA method still converges for nondifferentiable functions when the functions are continuous and convex and the domains are convex and compact sets.

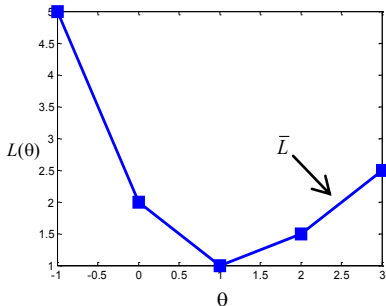


Fig. 1. Example of strictly discrete convex function and \bar{L} is a continuous extension

B. Basic Algorithm of DSPSA

Motivated by the special example shown above, we will consider the case when θ is p -dimensional, $p = 1, 2, 3, \dots$. We have the general basic algorithm as below for function $y = L + \varepsilon$, where $L: \mathbb{Z}^p \rightarrow \mathbb{R}$ and ε is noise.

The basic algorithm is:

Step0: Pick an initial guess $\hat{\theta}_0$.

Step1: Generate $\Delta_k = [\Delta_{k1}, \Delta_{k2}, \dots, \Delta_{kp}]^T$, where the Δ_{ki} are independent Bernoulli random variables taking the values ± 1 with probability $1/2$.

Step2: $\pi(\hat{\theta}_k) = (2\lfloor \hat{\theta}_k \rfloor + \mathbf{1}_p)/2$, where $\mathbf{1}_p$ is a p -dimensional vector with all components equal unity and $\lfloor \hat{\theta}_k \rfloor = [\lfloor \hat{\theta}_{k1} \rfloor, \dots, \lfloor \hat{\theta}_{kp} \rfloor]^T$.

Step3: Evaluate y at $\pi(\hat{\theta}_k) + \Delta_k/2$ and $\pi(\hat{\theta}_k) - \Delta_k/2$, form the estimate of $\hat{g}_k(\hat{\theta}_k)$,

$$\hat{g}_k(\hat{\theta}_k) = \left[y\left(\pi(\hat{\theta}_k) + \frac{1}{2}\Delta_k\right) - y\left(\pi(\hat{\theta}_k) - \frac{1}{2}\Delta_k\right) \right] \Delta_k^{-1},$$

where $\Delta_k^{-1} = [\Delta_{k1}^{-1}, \dots, \Delta_{kp}^{-1}]^T$.

Step4: Update the estimate according to the recursion

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

In the theoretical analysis below, we make use of the following mean gradient-like expression centered at $\pi(\theta)$:

$$\bar{g}(\pi(\theta)) = E \left\{ \left[L\left(\pi(\theta) + \frac{1}{2}\Delta\right) - L\left(\pi(\theta) - \frac{1}{2}\Delta\right) \right] \Delta^{-1} \middle| \theta \right\},$$

where Δ is p -dimensional vector that has the same definition as Δ_k mentioned above and θ may be a random variable in some cases. If each direction is chosen equally, then

$$\bar{g}(\pi(\theta)) = \frac{1}{2^p} \sum_{\Delta} \left[L\left(\pi(\theta) + \frac{1}{2}\Delta\right) - L\left(\pi(\theta) - \frac{1}{2}\Delta\right) \right] \Delta^{-1},$$

where \sum_{Δ} indicates the summation over all possible directions Δ . Note that $\Delta_k^{-1} = \Delta_k$ and $\Delta^{-1} = \Delta$ in the Bernoulli ± 1 case; we use Δ_k^{-1} to accommodate future extension to perturbation distributions other than Bernoulli ± 1 .

III. CONVERGENCE PROPERTIES

We now present an almost sure (a.s.) convergence result for $\hat{\theta}_k$. First we introduce some definitions that are used in

the proof to follow. For any point θ , we denote the set of middle points of all unit hypercubes containing θ as \mathcal{M}_θ . If θ lies strictly inside one unit hypercube, \mathcal{M}_θ contains one point. But if θ lies on the boundary, \mathcal{M}_θ contains multiple points. For any point m_θ in \mathcal{M}_θ , we have $|m_{\theta i} - t_i| \leq 1/2$ for $i=1, \dots, p$, where $\theta = [t_1, t_2, \dots, t_p]^T$ and $m_{\theta i}$ is the i th component of m_θ . Furthermore let $\mathfrak{Z}_k = \{\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k\}$.

Theorem 1. Assume L is a bounded function on \mathbb{Z}^p , and it has unique minimal point θ^* . Assume also (i) $a_k > 0$, $\lim_{k \rightarrow \infty} a_k = 0$, $\sum_{k=0}^{\infty} a_k = \infty$ and $\sum_{k=0}^{\infty} a_k^2 < \infty$; (ii) the components of Δ_k are independently Bernoulli ± 1 distributed; (iii) For all k , $E[(\varepsilon_k^+ - \varepsilon_k^-) | \mathfrak{Z}_k, \Delta_k] = 0$ a.s. and the variance of ε_k^\pm is uniformly bounded; (iv) $\sup_{k \geq 0} \|\hat{\theta}_k\| < \infty$ a.s.; and (v) $-\bar{g}(m_\theta)^T(\theta - \theta^*) < 0$ for all $m_\theta \in \mathcal{M}_\theta$ and all $\theta \in \mathbb{R}^p \setminus \{\theta^*\}$. Then $\hat{\theta}_k \rightarrow \theta^*$ a.s.

Proof. By the algorithm, we have

$$\begin{aligned} \hat{\theta}_{k+1} &= \hat{\theta}_k - a_k \hat{g}_k(\hat{\theta}_k) \\ &= \hat{\theta}_k - a_k \left[y \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - y \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right) \right] \Delta_k^{-1} \\ &= \hat{\theta}_k - a_k \left[L \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - L \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right) + \varepsilon_k^+ - \varepsilon_k^- \right] \Delta_k^{-1}. \end{aligned} \quad (1)$$

By conditions (i), (ii), (iv) and boundedness of L , we have

$$\lim_{k \rightarrow \infty} a_k \left(L(\pi(\hat{\theta}_k) + 1/2 \Delta_k) - L(\pi(\hat{\theta}_k) - 1/2 \Delta_k) \right) \Delta_k^{-1} = \mathbf{0} \quad \text{a.s.} \quad (2)$$

Also suppose the variance of ε_k^\pm are $(\sigma_k^\pm)^2$. Then by Chebyshev's inequality and (iii),

$$\lim_{m \rightarrow \infty} P \left\{ \left| a_k \varepsilon_k^\pm \right| > \eta \text{ for some } k \geq m \right\} \leq \lim_{m \rightarrow \infty} \sum_{k=m}^{\infty} a_k^2 \frac{(\sigma_k^\pm)^2}{\eta^2} \rightarrow 0,$$

implying by [4, Theorem 4.1.] that

$$\lim_{k \rightarrow \infty} a_k \{\varepsilon_k^+ - \varepsilon_k^-\} \Delta_k^{-1} = \mathbf{0} \quad \text{a.s.} \quad (3)$$

Through (1), we have the relationship that $\hat{\theta}_{k+1} - \hat{\theta}_k = -a_k \left[L(\pi(\hat{\theta}_k) + \Delta_k/2) - L(\pi(\hat{\theta}_k) - \Delta_k/2) + \varepsilon_k^+ - \varepsilon_k^- \right] \Delta_k^{-1}$, and by the results of (2) and (3), we get $\hat{\theta}_{k+1} - \hat{\theta}_k \rightarrow \mathbf{0}$ a.s. Hence there exists $\Omega_1 \subseteq \Omega$ such that $\hat{\theta}_{k+1}(\omega) - \hat{\theta}_k(\omega) \rightarrow \mathbf{0}$ and $P(\Omega_1) = 1$. By condition (iv), $\{\hat{\theta}_k(\omega)\}$ is a bounded sequence for any $\omega \in \Omega_1$. Then there exists a subsequence $\{\hat{\theta}_{k_s}(\omega)\}$ and point $\theta'(\omega)$ such that $\{\hat{\theta}_{k_s}(\omega)\} \rightarrow \theta'(\omega)$.

In addition, we can rewrite (1) as

$$\begin{aligned} \hat{\theta}_{k+1} &= \hat{\theta}_k - a_k \bar{g}(\pi(\hat{\theta}_k)) \\ &\quad + a_k \left\{ \bar{g}(\pi(\hat{\theta}_k)) - \left[L \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - L \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right) \right] \Delta_k^{-1} \right\} \\ &\quad - a_k (\varepsilon_k^+ - \varepsilon_k^-) \Delta_k^{-1}. \end{aligned}$$

By the definition of $\bar{g}(\cdot)$, we have

$$\begin{aligned} \bar{g}(\pi(\hat{\theta}_k)) &= E \left\{ \left[L \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - L \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right) \right] \Delta_k^{-1} \middle| \mathfrak{Z}_k \right\} \\ &= \frac{1}{2^p} \sum_{\Delta_k} \left[L \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - L \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right) \right] \Delta_k^{-1}. \end{aligned}$$

Let $b_k = L \left(\pi(\hat{\theta}_k) + \frac{1}{2} \Delta_k \right) - L \left(\pi(\hat{\theta}_k) - \frac{1}{2} \Delta_k \right)$, then for all $i < j$

$$\text{we have} \quad E \left[\left(\bar{g}(\pi(\hat{\theta}_i)) - b_i \Delta_i^{-1} \right)^T \left(\bar{g}(\pi(\hat{\theta}_j)) - b_j \Delta_j^{-1} \right) \right] =$$

$$E \left\{ E \left[\left(\bar{g}(\pi(\hat{\theta}_i)) - b_i \Delta_i^{-1} \right)^T \left(\bar{g}(\pi(\hat{\theta}_j)) - b_j \Delta_j^{-1} \right) \middle| \mathfrak{Z}_j \right] \right\} =$$

$$E \left\{ \left(\bar{g}(\pi(\hat{\theta}_i)) - b_i \Delta_i^{-1} \right)^T E \left[\left(\bar{g}(\pi(\hat{\theta}_j)) - b_j \Delta_j^{-1} \right) \middle| \mathfrak{Z}_j \right] \right\} = 0. \quad \text{Then}$$

due to conditions (i) (ii) and (iv), for any k we have

$$\begin{aligned} E \left[\left\| \sum_{i=k}^{\infty} a_i \left\{ \bar{g}(\pi(\hat{\theta}_i)) - \left[L \left(\pi(\hat{\theta}_i) + \frac{\Delta_i}{2} \right) - L \left(\pi(\hat{\theta}_i) - \frac{\Delta_i}{2} \right) \right] \Delta_i^{-1} \right\} \right\|^2 \right] &= \\ \sum_{i=k}^{\infty} a_i^2 E \left[\left\| \bar{g}(\pi(\hat{\theta}_i)) - \left[L \left(\pi(\hat{\theta}_i) + \frac{\Delta_i}{2} \right) - L \left(\pi(\hat{\theta}_i) - \frac{\Delta_i}{2} \right) \right] \Delta_i^{-1} \right\|^2 \right] &< \infty. \end{aligned} \quad (4)$$

Similarly, due to conditions (i) (ii) and (iii), we have

$$E \left[\left\| \sum_{i=k}^{\infty} a_i (\varepsilon_i^+ - \varepsilon_i^-) \Delta_i^{-1} \right\|^2 \right] = \sum_{i=k}^{\infty} a_i^2 E \left[\left\| (\varepsilon_i^+ - \varepsilon_i^-) \Delta_i^{-1} \right\|^2 \right] < \infty. \quad (5)$$

Since for all k , $\left\{ \sum_{i=k}^n a_i \left(\bar{g}(\pi(\hat{\theta}_i)) - \left[L \left(\pi(\hat{\theta}_i) + \Delta_i/2 \right) - L \left(\pi(\hat{\theta}_i) - \Delta_i/2 \right) \right] \Delta_i^{-1} \right) \right\}_n$ and

$\left\{ \sum_{i=k}^n a_i (\varepsilon_i^+ - \varepsilon_i^-) \Delta_i^{-1} \right\}_n$ are martingales, then by (4), (5) and [3,

Theorem 35.5], we know for all k ,

$\sum_{i=k}^{\infty} a_i \left(\bar{g}(\pi(\hat{\theta}_i)) - \left[L \left(\pi(\hat{\theta}_i) + \Delta_i/2 \right) - L \left(\pi(\hat{\theta}_i) - \Delta_i/2 \right) \right] \Delta_i^{-1} \right)$ exists and

$\sum_{i=k}^{\infty} a_i (\varepsilon_i^+ - \varepsilon_i^-) \Delta_i^{-1}$ exists. Let $M_k = \sum_{i=k}^{\infty} a_i (\varepsilon_i^+ - \varepsilon_i^-) \Delta_i^{-1}$ and

$$N_k = \sum_{i=k}^{\infty} a_i \left\{ \bar{g}(\pi(\hat{\theta}_i)) - \left[L \left(\pi(\hat{\theta}_i) + \Delta_i/2 \right) - L \left(\pi(\hat{\theta}_i) - \Delta_i/2 \right) \right] \Delta_i^{-1} \right\},$$

then $\{M_k\}$ and $\{N_k\}$ are reverse martingales ([2, p.472]),

and by [2, Theorem 35.8], there exist random variables M and N , such that $M_k \rightarrow M$ a.s. and $N_k \rightarrow N$ a.s. Furthermore

due to (4) and (5), we have $\lim_{k \rightarrow \infty} E[\|M_k\|^2] = 0$, and $\lim_{k \rightarrow \infty} E[\|N_k\|^2] = 0$. Also $\lim_{k \rightarrow \infty} E[M_k] = 0$ and $\lim_{k \rightarrow \infty} E[N_k] = 0$. Then $M = 0$ a.s., which indicates $M_k \rightarrow 0$ a.s. and $N = 0$ a.s. which indicates $N_k \rightarrow 0$ a.s. Then there exists $\Omega_2 \subseteq \Omega$ and $\Omega_3 \subseteq \Omega$ such that $P(\Omega_2) = 1$, $P(\Omega_3) = 1$ and such that for any $\omega \in \Omega_2$, $\sum_{i=k}^{\infty} a_i (\varepsilon_i^+(\omega) - \varepsilon_i^-(\omega)) \Delta_i^{-1} \rightarrow 0$ and for any $\omega \in \Omega_3$, $\sum_{i=k}^{\infty} a_i \left\{ \bar{g}(\pi(\hat{\theta}_i(\omega))) - \left[L(\pi(\hat{\theta}_i(\omega)) + \Delta_i/2) - L(\pi(\hat{\theta}_i(\omega)) - \Delta_i/2) \right] \Delta_i^{-1} \right\} \rightarrow 0$. Let $\Omega_4 = \Omega_1 \cap \Omega_2 \cap \Omega_3$ with $P(\Omega_4) = 1$. Then for any $\omega \in \Omega_4$ we have

$$\begin{aligned} \theta'(\omega) &= \hat{\theta}_{k_s}(\omega) - \sum_{i=k_s}^{\infty} a_i \bar{g}(\pi(\hat{\theta}_i(\omega))) \\ &+ \sum_{i=k_s}^{\infty} a_i \left\{ \bar{g}(\pi(\hat{\theta}_i(\omega))) - \left[L\left(\pi(\hat{\theta}_i(\omega)) + \frac{1}{2}\Delta_i\right) - L\left(\pi(\hat{\theta}_i(\omega)) - \frac{1}{2}\Delta_i\right) \right] \Delta_i^{-1} \right\} \\ &- \sum_{i=k_s}^{\infty} a_i (\varepsilon_i^+(\omega) - \varepsilon_i^-(\omega)) \Delta_i^{-1}, \end{aligned}$$

implying $\sum_{i=k_s}^{\infty} a_i (\varepsilon_i^+(\omega) - \varepsilon_i^-(\omega)) \Delta_i^{-1} \rightarrow 0$ and $\sum_{i=k_s}^{\infty} a_i \left\{ \bar{g}(\pi(\hat{\theta}_i(\omega))) - \left[L\left(\pi(\hat{\theta}_i(\omega)) + \Delta_i/2\right) - L\left(\pi(\hat{\theta}_i(\omega)) - \Delta_i/2\right) \right] \Delta_i^{-1} \right\} \rightarrow 0$ as $s \rightarrow \infty$. In addition, we know $\{\hat{\theta}_{k_s}(\omega)\} \rightarrow \theta'(\omega)$, indicating that

$$\sum_{i=k_s}^{\infty} a_i \bar{g}(\pi(\hat{\theta}_i(\omega))) \rightarrow 0 \text{ as } s \rightarrow \infty. \quad (6)$$

Because $\hat{\theta}_{k_s}(\omega) \rightarrow \theta'(\omega)$, then for any $\delta > 0$, there exists $S > 0$ such that when $s > S$, $\|\hat{\theta}_{k_s}(\omega) - \theta'(\omega)\| < \delta$. Thus there exists S' , when $s > S'$, all $\pi(\hat{\theta}_{k_s}(\omega)) \in \mathcal{M}_{\theta'}$. We now show $\theta'(\omega)$ is the optimal point. By way of contradiction, suppose $\theta'(\omega)$ is not the optimal solution. Then by condition (v), we have $-\bar{g}(\mathbf{m}_{\theta'})^T (\theta'(\omega) - \theta^*) < 0$ for all $\mathbf{m}_{\theta'} \in \mathcal{M}_{\theta'}$, which is a contradiction of $\sum_{i=k_s}^{\infty} a_i \bar{g}(\pi(\hat{\theta}_i(\omega))) \rightarrow 0$ when $s > S'$. Then for all $\omega \in \Omega_4$ the limiting point of the sequence $\{\hat{\theta}_k(\omega)\}$ is unique, which is equal to θ^* . Thus $\hat{\theta}_k \rightarrow \theta^*$ a.s.

Comment 1: The inner product condition (v) is a natural extension of the standard inner product condition for continuous problem (e.g. [22, p.106]), which includes convex function as a special case.

Comment 2: Actually some people have considered the discrete convexity. Miller [14] is a forerunner in the early 1970s in the area of discrete convex function. Ref. [14] has introduced the definition of discrete convex function and showed that the local optimal points for discrete convex function are also global optimal solutions. There are other definitions of discrete convex functions [5][15][16][6], but [17] shows that Miller's discrete convexity contains the other classes of discrete convexity. Note that Miller's definition does not include all functions satisfying condition (v), and condition (v) does not include all functions satisfy Miller's definition of discrete convexity. However, for $p = 1$, discrete convex functions satisfying Miller's definition also satisfy (v).

The corollaries below give two common functions satisfying condition (v). Even though we describe the functions in continuous form, for DSPSA we only use their values at multivariate integer points. Strictly convex separable functions mentioned in corollary 1 are discussed in [10].

Corollary 1. Strictly convex separable functions with minimal value at multivariate integer point satisfy the condition (v) in Theorem 1.

Proof. A separable function can be written as $\bar{L}(\theta) = \sum_{i=1}^p \bar{L}_i(t_i)$, where $\theta = [t_1, \dots, t_p]^T$. And L is a discrete function has same values with \bar{L} at multivariate integer points. Suppose the unique minimal point of \bar{L} is θ^* , and θ^* is a multivariate integer point with $\theta^* = [t_1^*, \dots, t_p^*]^T$.

Then θ^* is also the optimal point of L . Because it is strictly convex, then for all $\theta \in \mathbb{R}^p \setminus \{\theta^*\}$ and any subgradient

$\partial \bar{L}_i(t_i)$, we have $\partial \bar{L}_i(t_i)(t_i^* - t_i) < 0$ for $i=1, \dots, p$. Moreover, for any $\mathbf{m}_{\theta} \in \mathcal{M}_{\theta}$, $\bar{g}(\mathbf{m}_{\theta}) =$

$$\frac{1}{2^p} \sum_{\Delta} [L(\mathbf{m}_{\theta} + \Delta/2) - L(\mathbf{m}_{\theta} - \Delta/2)] \Delta^{-1} =$$

$$\frac{1}{2^p} \sum_{\Delta} \left[\sum_{i=1}^p (L_i(m_{\theta i} + \Delta_i/2) - L_i(m_{\theta i} - \Delta_i/2)) \right] \Delta^{-1} =$$

$$\sum_{i=1}^p \frac{1}{2^p} \sum_{\Delta} (L_i(m_{\theta i} + \Delta_i/2) - L_i(m_{\theta i} - \Delta_i/2)) \Delta^{-1} =$$

$$\sum_{i=1}^p (L_i(m_{\theta i} + 1/2) - L_i(m_{\theta i} - 1/2)) \mathbf{e}_i. \text{ Then we have } \bar{g}(\mathbf{m}_{\theta})^T (\theta - \theta^*) = \sum_{i=1}^p (L_i(m_{\theta i} + 1/2) - L_i(m_{\theta i} - 1/2))(t_i - t_i^*).$$

Because the minimal point is a multivariate integer point, then $\bar{L}_i(m_{\theta i} + 1/2) - \bar{L}_i(m_{\theta i} - 1/2)$ has the same sign with one of the subgradient of \bar{L}_i at t_i , indicating that

$$(\bar{L}_i(m_{\theta i} + 1/2) - \bar{L}_i(m_{\theta i} - 1/2))(t_i^* - t_i) < 0 \text{ for all } i=1, \dots, p.$$

Thus $-\bar{\mathbf{g}}(\mathbf{m}_\theta)^T(\boldsymbol{\theta}-\boldsymbol{\theta}^*) < 0$ for all $\mathbf{m}_\theta \in \mathcal{M}_\theta$ and all $\boldsymbol{\theta} \in \mathbb{R}^p \setminus \{\boldsymbol{\theta}^*\}$. Q.E.D.

Corollary 2. \bar{L} is a strictly convex piecewise linear function with minimal value at a multivariate integer point and it is linear in each unit hypercube, then L satisfy the condition (v) in Theorem 1.

Proof. L is a discrete function that has same values with \bar{L} at multivariate integer points. Since \bar{L} is strictly convex function, then for all $\boldsymbol{\theta} \in \mathbb{R}^p \setminus \{\boldsymbol{\theta}^*\}$, and for any subgradient

$\partial\bar{L}(\boldsymbol{\theta})$, we have $\partial\bar{L}(\boldsymbol{\theta})^T(\boldsymbol{\theta}^*-\boldsymbol{\theta}) < 0$. Furthermore for any

$$\mathbf{m}_\theta \in \mathcal{M}_\theta, \quad \bar{\mathbf{g}}(\mathbf{m}_\theta) = \frac{1}{2^p} \sum_{\Delta} [L(\mathbf{m}_\theta + \Delta/2) - L(\mathbf{m}_\theta - \Delta/2)] \Delta^{-1}$$

$$= \frac{1}{2^p} \sum_{\Delta} [\bar{L}(\mathbf{m}_\theta + \Delta/2) - \bar{L}(\mathbf{m}_\theta - \Delta/2)] \Delta^{-1} = \frac{1}{2^p} \sum_{\Delta} (\nabla \bar{L}(\mathbf{m}_\theta)^T \Delta) \Delta^{-1},$$

where ∇ is the notation of gradient. Thus

$$\bar{\mathbf{g}}(\mathbf{m}_\theta)^T(\boldsymbol{\theta}^*-\boldsymbol{\theta}) = \frac{1}{2^p} \sum_{\Delta} (\nabla \bar{L}(\mathbf{m}_\theta)^T \Delta) \Delta^{-T}(\boldsymbol{\theta}^*-\boldsymbol{\theta}) =$$

$$\nabla L(\mathbf{m}_\theta)^T \frac{1}{2^p} \sum_{\Delta} \Delta \Delta^{-T}(\boldsymbol{\theta}^*-\boldsymbol{\theta}) = \nabla \bar{L}(\mathbf{m}_\theta)^T(\boldsymbol{\theta}^*-\boldsymbol{\theta}).$$
 In

addition for any $\mathbf{m}_\theta \in \mathcal{M}_\theta$, there will be one subgradient $\partial\bar{L}(\boldsymbol{\theta})$ at point $\boldsymbol{\theta}$, such that $\nabla \bar{L}(\mathbf{m}_\theta) = \partial\bar{L}(\boldsymbol{\theta})$. Then

$$\bar{\mathbf{g}}(\mathbf{m}_\theta)^T(\boldsymbol{\theta}^*-\boldsymbol{\theta}) = \partial\bar{L}(\boldsymbol{\theta})^T(\boldsymbol{\theta}^*-\boldsymbol{\theta}) < 0, \text{ which indicates}$$

$-\bar{\mathbf{g}}(\mathbf{m}_\theta)^T(\boldsymbol{\theta}-\boldsymbol{\theta}^*) < 0$, for all $\mathbf{m}_\theta \in \mathcal{M}_\theta$ and all $\boldsymbol{\theta} \in \mathbb{R}^p \setminus \{\boldsymbol{\theta}^*\}$. Q.E.D.

IV. COMPARISON WITH LOCALIZED RANDOM SEARCH METHOD

Let us now compare the performance of DSPSA and the localized random search method for two loss functions. The first function considered here is a separable function $\sum_{i=1}^p t_i^2$. The second one is a skewed quartic loss function which is mentioned in [22, Ex 6.6]: $L(\boldsymbol{\theta}) = \boldsymbol{\theta}^T B^T B \boldsymbol{\theta} + 0.1 \sum_{i=1}^p (B\boldsymbol{\theta})_i^3 + 0.01 \sum_{i=1}^p (B\boldsymbol{\theta})_i^4$, where pB is an upper triangular matrix of 1's. Even though the skewed quartic loss function does not satisfy condition (v), we will see that DSPSA still works for this loss function. We consider the high-dimensional case for both functions, where $p = 200$, and the measurement noise ε is i.i.d $N(0,1)$. Since the localized random search method is more efficient in noise-free cases than in noisy cases, then we will consider both the noise-free situation and noisy situation. The localized random search method is described in [22, Sections 2.2–2.3], which consider both noise-free loss functions and noisy loss measurements, where a threshold parameter τ_k is involved. We will restrict the random

search to the closest neighbor points, and all these points are chosen with equal probability. Here for DSPSA, let $a_k = a/(k+1+A)^\alpha$, $a = 0.06$ (for separable); $a = 0.01$ (for skewed quartic), $A = 100$, $\alpha = 0.602$. For the localized random search method, we choose $\tau_k = 2$ for the noisy case after several tuning. The initial guess is set to be $10 \times \mathbf{1}_{200}$ in all runs. Fig. 2 and 3 show the performance of both methods under noise-free and noisy situations for separable function. And Fig. 4 and 5 show the performance of both methods for a skewed quartic function. We can see that DSPSA does better than the random search method for these two examples.

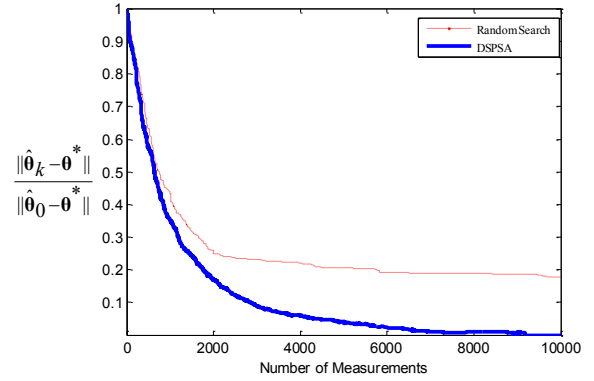


Fig. 2. Performance of localized random search method and DSPSA under noise-free situation for separable function.

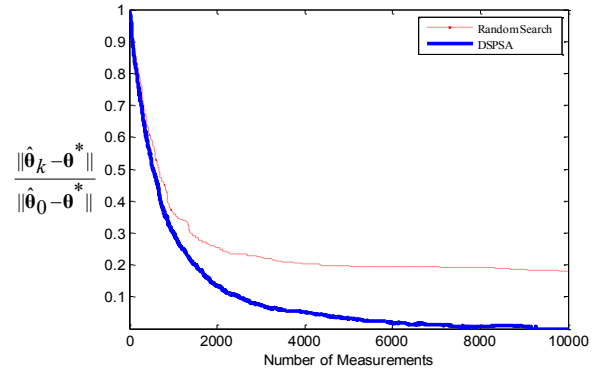


Fig. 3. Performance of localized random search method and DSPSA with noisy measurements for separable function.

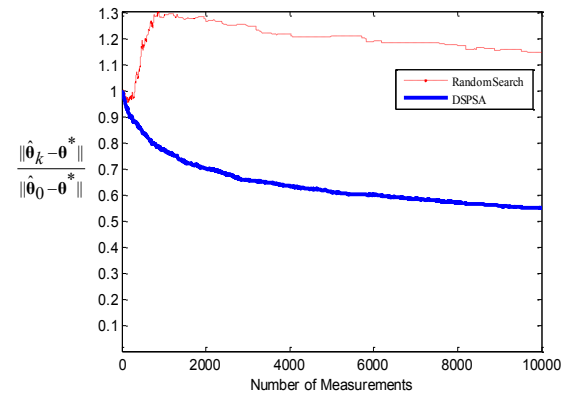


Fig. 4. Performance of localized random search method and DSPSA under noise-free situation for skewed quartic function.

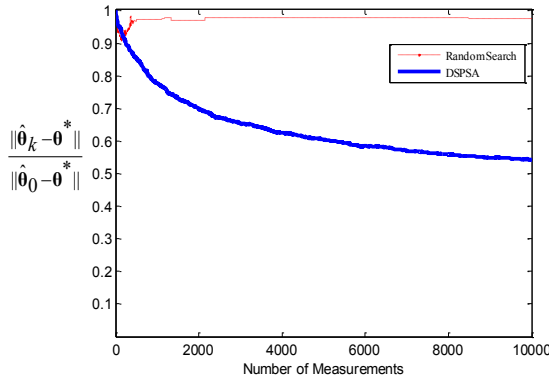


Fig. 5. Performance of localized random search method and DSPSA with noisy measurements for skewed quartic function.

V. CONCLUSION

In this paper, we introduced a discrete SPSA algorithm, and presented some preliminary convergence analysis. A preliminary numerical study shows that DSPSA works well on high-dimensional problems with or without noise in the loss measurements. As part of future work, we plan to formally study the convergence rate of the DSPSA and consider non-Bernoulli random variables for the perturbation vectors.

Also we intend to compare DSPSA with other popular discrete optimization algorithms, including those designed explicitly for handling noisy loss measurements (e.g. [7][13][19]). Two important practical problems of interest that involve stochastic discrete optimization are resource allocation, where a finite amount of a valuable commodity must be optimally allocated, and experimental design, where it is necessary to choose the best subset of input combinations from a large number of possible input combinations in a full-factorial design (e.g., [24]). We intend to explore the application of DSPSA to these or other problems.

ACKNOWLEDGMENT

This work was supported in part by the JHU/APL IRAD Program.

REFERENCES

- [1] M. H. Alrefaie, S. Andradóttir, "A Simulated Annealing Algorithm with Constant Temperature for Discrete Stochastic Optimization," *Management Sci.*, vol. 45, No.5, May 1999, pp. 748–764.
- [2] S. Andradóttir, "A Method for Discrete Stochastic Optimization," *Management Sci.*, vol. 41, No.12, December 1995, pp. 1946–1961.
- [3] P. Billingsley, *Probability and Measure*, Wiley-Interscience, Third Edition, 1995.
- [4] K. L. Chung, *A Course in Probability Theory*, Academic Press, Third Edition, 2001.
- [5] P. Favati, F. Tardella, "Convexity in Nonlinear Integer Programming," *Ricerca Operativa*, vol. 53, 1990, pp. 3–44.
- [6] S. Fujishige, K. Murota, "Notes on L-/M-convex Functions and the Separation Theorems," *Mathematical Programming*, vol. 88, 2000, pp. 129–146.
- [7] W. B. Gong, Y. C. Ho, W. Zhai, "Stochastic Comparison Algorithm for Discrete Optimization with Estimation," *SIAM J. Optim.*, vol. 10, No.2, 2000, pp. 384–404.
- [8] L. A. Hannah and W.B. Powell, "Evolutionary Policy Iteration Under a Sampling Regime for Stochastic Combinatorial Optimization," *IEEE Transactions on Automatic Control*, vol. 55, No.5, May 2010, pp. 1254–1257.
- [9] Y. He, M. C. Fu, and S. I. Marcus, "Convergence of Simultaneous Perturbation Approximation for Nondifferentiable Optimization," *IEEE Transactions on Automatic Control*, vol. 48, No.8, August 2003, pp. 1459–1463.
- [10] S. D. Hill, L. Gerencsér and Z. Vágó, "Stochastic Approximation on Discrete Sets Using Simultaneous Difference Approximations," *Proceeding of the 2004 American Control Conference*, Boston, MA, June 30–July 2, 2004, pp. 2795–2798.
- [11] Y. C. Ho, Q. C. Zhao, and Q. S. Jia, *Ordinal Optimization: Soft Optimization for Hard Problems*. Springer, New York, NY, 2007.
- [12] L. J. Hong and B. L. Nelson, "Discrete Optimization via Simulation Using COMPASS," *Oper. Res.*, vol. 54, No.1, 2006, pp. 115–129.
- [13] J. Li, A. Sava, and X. Xie, "Simulation-Based Discrete Optimization of Stochastic Discrete Event Systems Subject to Non Closed-Form Constraints," *IEEE Transactions on Automatic Control*, vol. 54, No.12, December 2009, pp. 2900–2904.
- [14] B. L. Miller, "On Minimizing Nonseparable Function Defined on the Integer with an Inventory Application," *SIAM Journal on Applied Mathematics*, vol. 21, No.1, July 1971, pp. 166–185.
- [15] K. Murota, "Discrete Convex Analysis," *Mathematical Programming*, vol. 83, 1998, pp. 313–371.
- [16] K. Murota, A. Shioura, "M-convex Function on Generalized Polymatroid," *Mathematics of Operations. Research*, vol. 24, 1999, pp. 95–105.
- [17] K. Murota, A. Shioura, "Relationship of M-/L- Convex Function with Discrete Convex Functions by Miller and Favati-Tardella," *Discrete Applied Mathematics*, vol. 115, 2001, pp. 151–176.
- [18] L. Shi and S. Olafsson, "Nested Partitions Method for Global Optimization," *Oper. Res.*, vol. 48, No.3, 2000, pp. 390–407.
- [19] J. Sklenar, P. Popela "Integer Simulation Based Optimization by Local Search," *Procedia Computer Science*, vol. 1, 2010, pp. 1341–1348.
- [20] J. C. Spall, "Multivariate Stochastic Approximation Using a Simultaneous Perturbation Gradient Approximation," *IEEE Transactions on Automatic Control*, vol. 37, No.3, March 1992, pp. 332–341.
- [21] J. C. Spall, "An Overview of the Simultaneous Perturbation Method for Efficient Optimization," *Johns Hopkins APL Technical Digest*, vol. 19, No.4, 1998, pp. 482–492.
- [22] J. C. Spall, *Introduction to Stochastic Search and Optimization: Estimation, Simulation, and Control*. Wiley, Hoboken, NJ, 2003.
- [23] F. Yousefian, A. Nedić, and U. V. Shanbhag, "Convex Nondifferentiable Stochastic Optimization: A Local Randomized Smoothing Technique," *Proceedings of the American Control Conference*, Baltimore, MD, June 30–July 2, 2010, pp. 4875–4880.
- [24] J. C. Spall, "Factorial Design for Choosing Input Values in Experimentation: Generating Informative Data for System Identification," *IEEE Control Systems Magazine*, vol. 30, no. 5, October 2010, pp. 38–53.