

# Analyzing Convergence and Rates of Convergence of Particle Swarm Optimization Algorithms Using Stochastic Approximation Methods

Quan Yuan and George Yin *Fellow, IEEE*

## Abstract

Recently, much progress has been made on particle swarm optimization (PSO). A number of works have been devoted to analyzing the convergence of the underlying algorithms. Nevertheless, in most cases, rather simplified hypotheses are used. For example, it often assumes that the swarm has only one particle. In addition, more often than not, the variables and the points of attraction are assumed to remain constant throughout the optimization process. In reality, such assumptions are often violated. Moreover, there are no rigorous rates of convergence results available to date for the particle swarm, to the best of our knowledge. In this paper, we consider a general form of PSO algorithms, and analyze asymptotic properties of the algorithms using stochastic approximation methods. We introduce four coefficients and rewrite the PSO procedure as a stochastic approximation type iterative algorithm. Then we analyze its convergence using weak convergence method. It is proved that a suitably scaled sequence of swarms converge to the solution of an ordinary differential equation. We also establish certain stability results. Moreover, convergence rates are ascertained by using weak convergence method. A centered and scaled sequence of the estimation errors is shown to have a diffusion limit.

## Index Terms

Particle swarm optimization, stochastic approximation, weak convergence, rate of convergence.

## I. INTRODUCTION

RECENTLY, optimization using particle swarms have received considerable attention owing to the wide range of applications from networked systems, multi-agent systems, and autonomous systems. Particle swarming refers to a computational method that optimizes a problem by trying recursively to improve a candidate solution with respect to a certain performance measure. Swarm intelligence from bio-cooperation within groups of individuals can often provide efficient solutions for certain optimization problems. When birds are searching food, they exchange and share information. Each member benefits from all other members owing to their discovery and experience based on the information acquired locally. Then each participating member adjusts the next search direction in accordance with the individual's best position currently and the information communicated to this individual by its neighbors. When food sources scattered unpredictably, advantages of such collaboration was decisive. Inspired by this, Kennedy and Eberhart proposed a particle swarm optimization (PSO) algorithm in 1995 [16]. A PSO procedure is a stochastic optimization algorithm that mimics the foraging behavior of birds. The search space of the optimization problem is analogous to the flight space of birds. Using an abstract setup, each bird is modeled as a particle (a point in the space of interest). Finding the optimum is the counterpart of searching for food. A PSO can be carried out effectively by using an iterative scheme. The PSO algorithm simulates social behavior among individuals (particles) "flying" through a multidimensional search space, where each particle represents a point at the intersection of all search dimensions. The particles evaluate their positions according to certain fitness functions at each iteration. The particles share memories of their "best" positions locally, and use the memories to adjust their own velocities and positions. Motivated by this scenario, a model is proposed to represent the traditional dynamics of particles.

To put this in a mathematical form, let  $F : \mathbb{R}^D \rightarrow \mathbb{R}$  be the cost function to be minimized. If we let  $M$  denote the size of the swarm, the current position of particle  $i$  is denoted by  $X^i$  ( $i = 1, 2, \dots, M$ ), and its current velocity is denoted by  $v^i$ . Then, the updating principle can be expressed as

$$\begin{aligned} v_{n+1}^{i,d} &= v_n^{i,d} + c_1 r_{1,n}^{i,d} [\text{Pr}_n^{i,d} - X_n^{i,d}] + c_2 r_{2,n}^{i,d} [\text{Pg}_n^{i,d} - X_n^{i,d}], \\ X_{n+1}^{i,d} &= X_n^{i,d} + v_{n+1}^{i,d}, \end{aligned} \quad (1)$$

where  $d = 1, \dots, D$ ;  $r_1^{i,d} \sim U(0, 1)$  and  $r_2^{i,d} \sim U(0, 1)$  represent two random variables uniformly distributed in  $[0, 1]$ ;  $c_1$  and  $c_2$  represent the acceleration coefficients;  $\text{Pr}_n^i$  represents the best position found by particle  $i$  up to "time"  $n$ , and  $\text{Pg}_n^i$  represents the "global" best position found by particle  $i$ 's neighborhood  $\Pi_i$ , i.e.,

$$\begin{aligned} \text{Pr}_n^i &= \arg \min_{1 \leq k \leq n} F(X_k^i), \\ \text{Pg}_n^i &= \text{Pr}_n^{j^*}, \text{ where } j^* = \arg \min_{j \in \Pi_i} F(\text{Pr}_n^j). \end{aligned} \quad (2)$$

In artificial life and social psychology,  $v_n^i$  in (1) is the velocity of particle  $i$  at time  $n$ , which provides the momentum for particles to pass through the search space. The  $c_1 r_{1,n}^{i,d} [\text{Pr}_n^{i,d} - X_n^{i,d}]$  is known as the ‘‘cognitive’’ component, which represents the personal thinking of each particle. The cognitive component of a particle takes the best position found so far by this particle as the desired input to make the particle move toward its own best positions.  $c_2 r_{2,n}^{i,d} [\text{Pg}_n^{i,d} - X_n^{i,d}]$  is known as the ‘‘social’’ component, which represents the collaborative behavior of the particles to find the global optimal solution. The social component always pulls the particles toward the best position found by its neighbors.

In a nutshell, a PSO algorithm has the following advantages: (1) It has versatility and does not rely on the problem information; (2) it has a memory capacity to retain local and global optimal information; (3) it is easy to implement. Given the versatility and effectiveness of PSO, it is widely used to solve practical problems such as artificial neural networks [13], [29], chemical systems [10], power systems [1], [2], mechanical design [18], communications [55], robotics [22], [47], economy [31], [33], image processing [32], bio-informatics [39], [48], medicine [42], and industrial engineering [26], [44]. Note that swarms have also been used in many engineering applications, for example, in collective robotics where there are teams of robots working together by communicating over a communication network; see [24] for a stability analysis and many related references.

To enable and to enhance further applications, much work has also been devoted to improving the PSO algorithms. Because the original model is similar to a mobile multi-agent system and each parameter describes a special character of natural swarm behavior, one can improve the performance of PSO according to the physical meanings of these parameters [25], [34], [38], [40], [56]. The first significant improvement was proposed by Shi and Eberhart in [43]. They suggested to add a new parameter  $w$  as an ‘‘inertia constant’’, which results in fast convergence. The modified equation of (1) is

$$\begin{aligned} v_{n+1}^{i,d} &= wv_n^{i,d} + c_1 r_{1,n}^{i,d} [\text{Pr}_n^{i,d} - X_n^{i,d}] + c_2 r_{2,n}^{i,d} [\text{Pg}_n^{i,d} - X_n^{i,d}], \\ X_{n+1}^{i,d} &= X_n^{i,d} + v_{n+1}^{i,d}. \end{aligned} \quad (3)$$

Another significant improvement was due to Clerc and Kennedy [9]. They introduced a constriction coefficient  $\chi$  and then proposed to modify (1) as

$$\begin{aligned} v_{n+1}^{i,d} &= \chi(v_n^{i,d} + c_1 r_{1,n}^{i,d} [\text{Pr}_n^{i,d} - X_n^{i,d}] + c_2 r_{2,n}^{i,d} [\text{Pg}_n^{i,d} - X_n^{i,d}]), \\ X_{n+1}^{i,d} &= X_n^{i,d} + v_{n+1}^{i,d}. \end{aligned} \quad (4)$$

This constriction coefficient can control the ‘‘explosion’’ of the PSO and ensure the convergence. Some researchers also considered using ‘‘good’’ topologies of particle connection, in particular adaptive ones (e.g., [7], [28], [30]).

There has been much development on mathematical analysis for the convergence of PSO algorithms as well. Although most researchers prefer to use discrete system [6], [9], [46], [50], there are some works on continuous-time models [11], [27]. Some recent work such as [8], [12], [14], [23], [37], [49] provides guidelines for selecting PSO parameters leading to convergence, divergence, or oscillation of the swarm’s particles. The aforementioned work also gives rise to several PSO variants. Nowadays, it is widely recognized that purely deterministic approach is inadequate in reflecting the exploration and exploitation aspects brought by stochastic variables. However, as criticized by Pedersen [36], the analysis is often oversimplified. For example, the swarm is often assumed to have only one particle; stochastic variables (namely,  $r_{1,n}$ ,  $r_{2,n}$ ) are not used; the points of attraction, i.e., the particle’s best known position  $\text{Pr}$  and the swarm’s best known position  $\text{Pg}$ , are normally assumed to remain constant throughout the optimization process.

In this paper, we study convergence of PSO by using stochastic approximation methods. In the past, some authors have considered using stochastic approximation combined with PSO to enhance the performance or select parameters (e.g., [17]). But to the best of our knowledge, the only paper using stochastic approximation methods to analyze the dynamics of the PSO so far is by Chen and Li [8]. They designed a special PSO procedure and assumed that (i)  $\text{Pr}_n^i$  and  $\text{Pg}_n^i$  are always within a finite domain; (ii) with  $P^*$  representing the global optimal positions in the solution space, and  $\|P^*\| < \infty$ .  $\lim_{n \rightarrow \infty} \text{Pr}_n \rightarrow P^*$  and  $\lim_{n \rightarrow \infty} \text{Pg}_n \rightarrow P^*$ . Using assumption (i), they proved the convergence of the algorithm in the sense of with probability one. With additional assumption (ii), they showed that the swarm will converge to  $P^*$ . Despite the interesting development, their assumptions (i) and (ii) appear to be rather strong. Moreover, they added some specific terms in the PSO procedure. So their algorithm is different from the traditional PSOs (1)-(4). In this paper, we consider a general form of PSO algorithms. We introduce four coefficients  $\varepsilon$ ,  $\kappa_1$ ,  $\kappa_2$ , and  $\chi$  and rewrite the PSOs in a stochastic approximation setup. Then we analyze its convergence using weak convergence method. We prove that a suitably interpolated sequence of swarms converge to the solution of an ordinary differential equation. Moreover, convergence rates are derived by using a sequence of centered and scaled estimation errors.

The rest of the paper is arranged as follows. Section II presents the setup of our algorithm. Section III studies the convergence and Section IV analyzes the rate of convergence. Section V proceeds with several numerical simulation examples to illustrate the convergence of our algorithms. Finally, Section VI provides a few further remarks.

## II. FORMULATION

First, some descriptions on notation are in order. We use  $|\cdot|$  to denote a Euclidean norm. A point  $\theta$  in a Euclidean space is a column vector; the  $i$ th component of  $\theta$  is denoted by  $\theta^i$ ;  $\text{diag}(\theta)$  is a diagonal matrix whose diagonal elements are the

$$\begin{aligned}
\begin{bmatrix} v_{n+1} \\ X_{n+1} \end{bmatrix} &= \begin{bmatrix} v_n \\ X_n \end{bmatrix} + \varepsilon \left\{ \begin{bmatrix} \kappa_1 I & -0.5\chi(c_1 + c_2)I \\ \kappa_2 I & -0.5\chi(c_1 + c_2)I \end{bmatrix} \begin{bmatrix} v_n \\ X_n \end{bmatrix} + \chi \begin{bmatrix} 0.5c_1 I & 0.5c_2 I \\ 0.5c_1 I & 0.5c_2 I \end{bmatrix} \begin{bmatrix} \Pr(\theta_n, \eta_n) \\ \text{Pg}(\theta_n, \eta_n) \end{bmatrix} \right. \\
&+ \chi \begin{bmatrix} 0 & -(c_1 \text{diag}(r_{1,n}) + c_2 \text{diag}(r_{2,n}) - 0.5c_1 I - 0.5c_2 I) \\ 0 & -(c_1 \text{diag}(r_{1,n}) + c_2 \text{diag}(r_{2,n}) - 0.5c_1 I - 0.5c_2 I) \end{bmatrix} \begin{bmatrix} v_n \\ X_n \end{bmatrix} \\
&\left. + \chi \begin{bmatrix} c_1 \text{diag}(r_{1,n}) - 0.5c_1 I & c_2 \text{diag}(r_{2,n}) - 0.5c_2 I \\ c_1 \text{diag}(r_{1,n}) - 0.5c_1 I & c_2 \text{diag}(r_{2,n}) - 0.5c_2 I \end{bmatrix} \begin{bmatrix} \Pr(\theta_n, \eta_n) \\ \text{Pg}(\theta_n, \eta_n) \end{bmatrix} \right\}. \tag{6}
\end{aligned}$$

elements of  $\theta$ ;  $I$  denotes the identity matrix of appropriate dimension;  $z'$  denotes the transposition of  $z$ ; the notation  $O(y)$  denotes a function of  $y$  satisfying  $\sup_y |O(y)|/|y| < \infty$ , and  $o(y)$  denotes a function of  $y$  satisfying  $|o(y)|/|y| \rightarrow 0$ , as  $y \rightarrow 0$ . In particular,  $O(1)$  denotes the boundedness and  $o(1)$  indicates convergence to 0. Throughout the paper, for convenience, we use  $K$  to denote a generic positive constant with the convention that the value of  $K$  may be different for different usage.

In this paper, without loss of generality, we assume that each particle is a one-dimensional scalar. Note that each particle can be a multi-dimensional vector, which does not introduce essential difficulties in the analysis; only the notation is a bit more complex. We introduce four parameters  $\varepsilon$ ,  $\kappa_1$ ,  $\kappa_2$ , and  $\chi$ . Suppose there are  $r$  particles, then the PSO algorithm can be expressed as

$$\begin{aligned}
\begin{bmatrix} v_{n+1} \\ X_{n+1} \end{bmatrix} &= \begin{bmatrix} v_n \\ X_n \end{bmatrix} \\
&+ \varepsilon \left( \begin{bmatrix} \kappa_1 I & -\chi(c_1 \text{diag}(r_{1,n}) + c_2 \text{diag}(r_{2,n})) \\ \kappa_2 I & -\chi(c_1 \text{diag}(r_{1,n}) + c_2 \text{diag}(r_{2,n})) \end{bmatrix} \begin{bmatrix} v_n \\ X_n \end{bmatrix} \right. \\
&\left. + \chi \begin{bmatrix} c_1 \text{diag}(r_{1,n}) & c_2 \text{diag}(r_{2,n}) \\ c_1 \text{diag}(r_{1,n}) & c_2 \text{diag}(r_{2,n}) \end{bmatrix} \begin{bmatrix} \Pr(\theta_n, \eta_n) \\ \text{Pg}(\theta_n, \eta_n) \end{bmatrix} \right), \tag{5}
\end{aligned}$$

where  $c_1$  and  $c_2$  represent the acceleration coefficients,  $X_n = [X_n^1, \dots, X_n^r]' \in \mathbb{R}^r$ ,  $v_n = [v_n^1, \dots, v_n^r]' \in \mathbb{R}^r$ ,  $\theta_n = (X_n, v_n)'$ ,  $r_1, r_2$  are  $r$ -dimensional random vectors in which each component is uniformly distributed in  $(0, 1)$ , and  $\Pr(\theta, \eta)$  and  $\text{Pg}(\theta, \eta)$  are two non-linear functions depending on  $\theta = (X, v)'$  as well as on a ‘‘noise’’  $\eta$ , and  $\varepsilon > 0$  is a small parameter representing the stepsize of the iterations.

*Remark 1:* Note that for a large variety of cases, the structures and the forms of  $\Pr(\theta, \eta)$  and  $\text{Pg}(\theta, \eta)$  are not known. This is similar to the situation in a stochastic optimization problem in which the objective function is not known precisely. Thus, stochastic approximation methods are well suited. As it is well known that stochastic approximation methods are very useful for treating optimization problems in which the form of the objective function is not known precisely, or too complex to compute. The beauty of such stochastic iteratively defined procedures is that one need not know the precise form of the functions. If there is no noise term  $\eta_n$ , let  $\varepsilon = 0.01$ ,  $\chi = 72.9$ ,  $\kappa_1 = -27.1$ , and  $\kappa_2 = 72.9$ , then (5) is equivalent to (3) when  $w = 0.729$  or (4) when  $\chi = 0.729$ . Thus (5) is a generalization of (1)-(4). So a lot of approaches of tuning parameters (e.g., [3], [35], [54]) could also be applied.

*Remark 2:* In the proposed algorithm, we use a constant stepsize. The stepsize  $\varepsilon > 0$  is a small parameter. As is well recognized (see [4], [20]), constant stepsize algorithms have the ability to track slight time variation and is more preferable in many applications. In the convergence and rate of convergence analysis, we let  $\varepsilon \rightarrow 0$ . In the actual computation,  $\varepsilon$  is just a constant. It need not go to 0. This is the same as one carries out any computational problem in which the analysis requires the iteration number going to infinity. However, in the actual computing, one only executes the procedure finitely many steps.

In (5),  $r_1$  and  $r_2$  are used to reflect the exploration of particles. Rearranging terms of (5) and considering that  $E[c_1 \text{diag}(r_{1,n})] = 0.5c_1 I$  and  $E[c_2 \text{diag}(r_{2,n})] = 0.5c_2 I$ , it can be rewritten as (6) (on the top of page 3).

Denote

$$\begin{aligned}
\theta_n &= [v_n, X_n]' \in \mathbb{R}^{2r}, \\
M &= \begin{bmatrix} \kappa_1 I & -0.5\chi(c_1 + c_2)I \\ \kappa_2 I & -0.5\chi(c_1 + c_2)I \end{bmatrix}, \\
P(\theta_n, \eta_n) &= \chi \begin{bmatrix} 0.5c_1 I & 0.5c_2 I \\ 0.5c_1 I & 0.5c_2 I \end{bmatrix} \begin{bmatrix} \Pr(\theta_n, \eta_n) \\ \text{Pg}(\theta_n, \eta_n) \end{bmatrix}, \tag{7}
\end{aligned}$$

and  $W(\theta_n, r_{1,n}, r_{2,n}, \eta_n)$  to be the sum of the last two terms in the curly braces of (6). Then (6) can be expressed as a stochastic approximation algorithm

$$\theta_{n+1} = \theta_n + \varepsilon [M\theta_n + P(\theta_n, \eta_n) + W(\theta_n, r_{1,n}, r_{2,n}, \eta_n)]. \tag{8}$$

One of the challenges in analyzing the convergence of PSO is that the concrete forms of  $\Pr(\theta_n, \eta_n)$  and  $\text{Pg}(\theta_n, \eta_n)$  are unknown. However, this will not concern us. As mentioned before, stochastic approximation methods are known to have advantages in treating such situations. We shall use the following assumptions.

(A1) The  $\Pr(\cdot, \eta)$  and  $\text{Pg}(\cdot, \eta)$  are continuous for each  $\eta$ . For each bounded  $\theta$ ,  $E|P(\theta, \eta_m)|^2 < \infty$  and  $E|W(\theta, r_{1,n}, r_{2,n}, \eta_n)|^2 < \infty$ . There exist continuous functions  $\overline{\Pr}(\theta)$  and  $\overline{\text{Pg}}(\theta)$  such that

$$\begin{aligned} \frac{1}{n} \sum_{j=m}^{n+m-1} E_m \Pr(\theta, \eta_j) &\rightarrow \overline{\Pr}(\theta) \quad \text{in probability,} \\ \frac{1}{n} \sum_{j=m}^{n+m-1} E_m \text{Pg}(\theta, \eta_j) &\rightarrow \overline{\text{Pg}}(\theta) \quad \text{in probability,} \end{aligned} \quad (9)$$

where  $E_m$  denotes the conditional expectation on the  $\sigma$ -algebra  $\mathcal{F}_m = \{\theta_0, r_{i,j}, i = 1, 2, \eta_j : j < m\}$ . Moreover, for each  $\theta$  in a bounded set,

$$\begin{aligned} \sum_{j=n}^{\infty} |E_n \Pr(\theta, \eta_j) - \overline{\Pr}(\theta)| &< \infty, \\ \sum_{j=n}^{\infty} |E_n \text{Pg}(\theta, \eta_j) - \overline{\text{Pg}}(\theta)| &< \infty. \end{aligned} \quad (10)$$

(A2) Define

$$\overline{P}(\theta) = \chi \begin{bmatrix} 0.5c_1 I & -0.5c_2 I \\ 0.5c_1 I & -0.5c_2 I \end{bmatrix} \begin{bmatrix} \overline{\Pr}(\theta) \\ \overline{\text{Pg}}(\theta) \end{bmatrix}.$$

The ordinary differential equation

$$\frac{d\theta(t)}{dt} = M\theta(t) + \overline{P}(\theta(t)) \quad (11)$$

has a unique solution for each initial condition  $\theta(0) = (\theta_0^1, \dots, \theta_0^{2r})'$ .

(A3) For  $i = 1, 2$ ,  $\{r_{i,n}\}$  and  $\{\eta_n\}$  are mutually independent;  $\{r_{i,n}\}$  are i.i.d. sequences of random variables with each component being uniformly distributed in  $(0, 1)$ .

*Remark 3:* Condition (A1) is satisfied by a large class of functions and random variables. The continuity is assumed for convenience. In fact, only weak continuity is needed so we can in fact deal with indicator type of functions whose expectations are continuous.

In fact, (9) mainly requires that  $\{\Pr(\theta, \eta_n)\}$  is a sequence that satisfies a law of large number type of condition, although it is weaker than the usual weak law of large numbers. Condition (10) is modeled by the mixing type condition. For instance, we may assume that for each bounded random vector  $\theta$  and each  $T < \infty$ , either

$$\begin{aligned} \lim_{j \rightarrow \infty, \Delta \rightarrow 0} E \sup_{|Y| \leq \Delta} |\Pr(\theta + Y, \eta_j) - \Pr(\theta, \eta_j)| &= 0, \quad \text{or} \\ \lim_{n \rightarrow \infty, \Delta \rightarrow 0} \frac{1}{n} \sum_{j=m}^{m+n-1} E \sup_{|Y| \leq \Delta} |\Pr(\theta + Y, \eta_j) - \Pr(\theta, \eta_j)| &= 0. \end{aligned}$$

Apparently, the second alternative is even weaker. With either of this assumption, all of the subsequent development follows, but the argument is more complex. Under the above condition, one can treat discontinuity involving sign function or indicator function among others. For the corresponding stochastic approximation algorithms, see [19, p. 100]; the setup in [20] is even more general, which allows in addition to the discontinuity, the functions involved to be time dependent. Inserting the conditional expectation is much weaker than without. For example, if  $\{\eta_n\}$  is a sequence of i.i.d. random variables with distribution function  $F_\eta$ , then for each  $\theta$ ,  $\overline{\Pr}(\theta) = E\Pr(\theta, \eta_1) = \int \Pr(\theta, \zeta) F_\eta(d\zeta)$ , so (9) is easily verified. Likewise, if  $\{\eta_n\}$  is a martingale difference sequence, the condition is also satisfied. Next, if  $\{\eta_n\}$  is a moving average sequence driven by a martingale difference sequence, (9) is also satisfied. In addition, if  $\{\eta_n\}$  is a mixing sequence [5, p.166] with the mixing rate decreasing to 0, the condition is also satisfied. Note that in a mixing sequence, there can be infinite correlations and the remote past and distant future are only asymptotically uncorrelated.

In the simplest additive noise case, i.e.,  $\Pr(\theta, \eta) = \Pr(\theta) + \eta$ , then the condition is mainly on the noise sequence  $\{\eta_n\}$ . Condition (10) is modeled after the so-called mixing inequality; see [19, p.82] and references therein. Suppose that  $\{\Pr(\theta, \eta_n)\}$  is a stationary mixing sequence with mean  $\overline{\Pr}(\theta)$  and mixing rate  $\phi_n$  such that  $\sum_n \phi_n^{1/2} < \infty$ , then (10) is satisfied.

With these assumptions, we proceed to analyze the convergence and rates of convergence of PSO algorithms with general form (8). The scheme is a constant-step-size stochastic approximation algorithm with step size  $\varepsilon$ . Our interest lies in obtaining convergence and rates of convergence as  $\varepsilon \rightarrow 0$ . We emphasize that in the actual computation, it is not necessary to modify it as the generalized PSO form (8). This generalized PSO form is simply a convenient form that allows us to analyze the algorithm by using methods of stochastic approximation.

### III. CONVERGENCE

This section is devoted to obtaining asymptotic properties of algorithm (8). In relation to PSO the word ‘‘convergence’’ typically means one of two things, although it is often not clarified which definition is meant and sometimes they are mistakenly thought to be identical. (i) Convergence may refer to the swarm’s best known position  $P_g$  approaching (converging to) the optimum of the problem, regardless of how the swarm behaves. (ii) Convergence may refer to a swarm collapse in which all particles have converged to a point in the search space, which may or may not be the optimum. Since the convergence may rely on structure of the cost function if we use the first definition of convergence, we use the second one as the definition of convergence in this study. Our first result concerns the property of the algorithm as  $\varepsilon \rightarrow 0$  through an appropriate continuous-time interpolation. We define

$$\theta^\varepsilon(t) = \theta_n \text{ for } t \in [\varepsilon n, \varepsilon n + \varepsilon).$$

Then  $\theta^\varepsilon(\cdot) \in D([0, T] : \mathbb{R}^{2r})$ , which is the space of functions that are defined on  $[0, T]$  taking values in  $\mathbb{R}^{2r}$ , and that are right continuous and have left limits endowed with the Skorohod topology [20, Chapter 7].

*Theorem 4:* Under (A1)-(A3),  $\theta^\varepsilon(\cdot)$  is tight in  $D([0, T] : \mathbb{R}^{2r})$ . Moreover, as  $\varepsilon \rightarrow 0$ ,  $\theta^\varepsilon(\cdot)$  converges weakly to  $\theta(\cdot)$ , which is a solution of (11).

*Remark 5:* An equivalent way of stating the ODE limit (11) is to consider its associated martingale problem [19, pp. 15-16]. Consider the differential operator associated with  $\theta(\cdot)$

$$\mathcal{L}f(\theta) = (\nabla f(\theta))'(M\theta + \overline{P}(\theta)),$$

and define

$$\widetilde{M}_f(t) = f(\theta(t)) - f(\theta(0)) - \int_0^t \mathcal{L}f(\theta(s))ds.$$

If  $\widetilde{M}_f(\cdot)$  is a martingale for each  $f(\cdot) \in C_0^1$  ( $C^1$  function with compact support), then  $\theta(\cdot)$  is said to solve a martingale problem with operator  $\mathcal{L}$ . Thus, an equivalent way to state the theorem is to prove that  $\theta^\varepsilon(\cdot)$  converges weakly to  $\theta(\cdot)$ , which is a solution of the martingale problem with operator  $\mathcal{L}$ .

**Proof of Theorem 4.** To prove the tightness in  $D([0, T] : \mathbb{R}^{2r})$ , we first need to show

$$\lim_{K \rightarrow \infty} \limsup_{\varepsilon \rightarrow 0} P\{\sup_{t \leq T} |\theta^\varepsilon(t)| \geq K\} = 0 \quad (12)$$

To avoid verifying (12), we define a process  $\theta^{\varepsilon, N}(\cdot)$  satisfies  $\theta^{\varepsilon, N}(t) = \theta^\varepsilon(t)$  up until the first exit from  $S_N = \{x \in \mathbb{R}^{2r} : |x| \leq N\}$  and satisfies (12), the  $\theta^{\varepsilon, N}(\cdot)$  is said to be an  $N$ -truncation of  $\theta^\varepsilon(\cdot)$ . Introduce a truncation function  $q^N(\cdot)$  that is smooth and that satisfies  $q^N(\theta) = 1$  for  $|\theta| \leq N$ ,  $q^N(\theta) = 0$  for  $|\theta| \geq N + 1$ . Then the discrete system (8) is defined as

$$\theta_{n+1}^N = \theta_n^N + \varepsilon[M\theta_n^N + P(\theta_n^N, \eta_n) + W(\theta_n^N, r_{1,n}, r_{2,n}, \eta_n)]q^N(\theta_n^N), \quad (13)$$

using the  $N$ -truncation. Moreover, the  $N$ -truncated ODE and the operator  $\mathcal{L}^N$  of the associated martingale problem can be defined as

$$\frac{d\theta^N(t)}{dt} = [M\theta^N(t) + \overline{P}(\theta^N(t))]q^N(\theta(t)), \quad (14)$$

and

$$\mathcal{L}^N f(\theta) = (\nabla f(\theta))'[M\theta + \overline{P}(\theta)]q^N(\theta), \quad (15)$$

respectively.

To prove the theorem, we proceed to verify the following claims: (a) for each  $N$ ,  $\{\theta^{\varepsilon, N}(\cdot)\}$  is tight. By virtue of the Prokhorov theorem [20, p.229], we can extract a weakly convergent subsequence. For notational simplicity, we still denote the subsequence by  $\{\theta^{\varepsilon, N}(\cdot)\}$  with limit denoted by  $\theta^N(\cdot)$ .

(b)  $\theta^N(\cdot)$  is a solution of the martingale problem with operator  $\mathcal{L}^N$ . Using the uniqueness of the limit, passing to the limit as  $N \rightarrow \infty$ , and by the corollary in [19, p.44],  $\{\theta^\varepsilon(\cdot)\}$  converges weakly to  $\theta(\cdot)$ .

Now we start to prove claims (a) and (b).

(a) **Tightness.** For any  $\delta > 0$ , let  $t > 0$  and  $s > 0$  such that  $s \leq \delta$ , and  $t, t + \delta \in [0, T]$ . Note that

$$\begin{aligned} & \theta^{\varepsilon, N}(t+s) - \theta^{\varepsilon, N}(t) \\ & \stackrel{(t+s)/\varepsilon - 1}{=} \varepsilon \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon - 1} (M\theta_k^N + P(\theta_k^N, \eta_k) \\ & \quad + W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k))q^N(\theta_k^N). \end{aligned}$$

In the above and hereafter, we use the conventions that  $t/\varepsilon$  and  $(t+s)/\varepsilon$  denote the corresponding integer parts  $\lfloor t/\varepsilon \rfloor$  and  $\lfloor (t+s)/\varepsilon \rfloor$ , respectively. For notational simplicity, in what follows, we will not use the floor function notation unless it is necessary.

Using the Cauchy-Schwarz inequality,

$$\begin{aligned} \varepsilon^2 E_t^\varepsilon \left| \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon-1} M\theta_k^N q^N(\theta_k^N) \right|^2 \\ \leq \varepsilon K s \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon-1} E_t^\varepsilon |\theta_k^N q^N(\theta_k^N)|^2. \end{aligned} \quad (16)$$

where  $E_t^\varepsilon$  denotes the expectation conditioned on the  $\sigma$ -algebra  $\mathcal{F}_t^\varepsilon$ . Likewise,

$$\varepsilon^2 E_t^\varepsilon \left| \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon-1} W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k) q^N(\theta_k^N) \right|^2 \leq K s^2, \quad (17)$$

and

$$\varepsilon^2 E_t^\varepsilon \left| \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon-1} P(\theta_k^N, \eta_k) q^N(\theta_k^N) \right|^2 \leq K s^2. \quad (18)$$

So we have

$$\begin{aligned} E_t^\varepsilon \left| \theta^{\varepsilon, N}(t+s) - \theta^{\varepsilon, N}(t) \right|^2 \\ \leq K \varepsilon s \sum_{k=t/\varepsilon}^{(t+s)/\varepsilon-1} \sup_{t/\varepsilon \leq k \leq (t+s)/\varepsilon-1} E_t^\varepsilon |\theta_k^N q^N(\theta_k^N)|^2 + K s^2. \end{aligned} \quad (19)$$

As a result, there is a  $\zeta^\varepsilon(\delta)$  such that

$$E_t^\varepsilon |\theta^{\varepsilon, N}(t+s) - \theta^{\varepsilon, N}(t)|^2 \leq E_t^\varepsilon \zeta^\varepsilon(\delta) \quad \text{for all } 0 \leq s \leq \delta,$$

and that  $\lim_{\delta \rightarrow 0} \limsup_{\varepsilon \rightarrow 0} E \zeta^\varepsilon(\delta) = 0$ . The tightness of  $\{\theta^{\varepsilon, N}(\cdot)\}$  then follows from [19, p.47].

(b) Characterization of the limit. To characterize the limit process, we need to work with a continuously differentiable function with compact support  $f(\cdot)$ . Choose  $m_\varepsilon$  so that  $m_\varepsilon \rightarrow \infty$  as  $\varepsilon \rightarrow 0$  but  $\delta_\varepsilon = \varepsilon m_\varepsilon \rightarrow 0$ . Using the recursion (13),

$$\begin{aligned} f(\theta^{\varepsilon, N}(t+s)) - f(\theta^{\varepsilon, N}(t)) \\ = \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} [f(\theta_{lm_\varepsilon+m_\varepsilon}^N) - f(\theta_{lm_\varepsilon}^N)] \\ = \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} [M\theta_k^N + \bar{P}(\theta_k^N)] q^N(\theta_k^N) \\ + \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} [P(\theta_k^N, \eta_k) - \bar{P}(\theta_k^N)] q^N(\theta_k^N) \\ + \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k) q^N(\theta_k^N) \\ + \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} \left\{ (\nabla f(\theta_{lm_\varepsilon}^{N+}) - \nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ \left. \times \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} [M\theta_k^N + P(\theta_k^N, \eta_k) \right. \\ \left. + W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k)] q^N(\theta_k^N) \right\}, \end{aligned} \quad (20)$$

where  $\theta_{lm_\varepsilon}^{N+}$  is a point on the line segment joining  $\theta_{lm_\varepsilon}^N$  and  $\theta_{lm_\varepsilon+m_\varepsilon}^N$ .

Our focus here is to characterize the limit. By the Skorohod representation [20, p.230], with a slight abuse of notation, we may assume that  $\theta^{\varepsilon, N}(\cdot)$  converges to  $\theta^N(\cdot)$  with probability one and the convergence is uniform on any bounded time interval. To show that  $\{\theta^{\varepsilon, N}(\cdot)\}$  is a solution of the martingale problem with operator  $\mathcal{L}^N$ , it suffices to show that for any  $f(\cdot) \in C_0^1$ , the class of functions that are continuously differentiable with compact support,

$$\widetilde{M}_f^N(t) = f(\theta^N(t)) - f(\theta^N(0)) - \int_0^t \mathcal{L}^N f(\theta^N(u)) du$$

is a martingale. To verify the martingale property, we need only show that for any bounded and continuous function  $h(\cdot)$ , any positive integer  $\kappa$ , any  $t, s > 0$ , and  $t_i \leq t$  with  $i \leq \kappa$ ,

$$\begin{aligned} & Eh(\theta^N(t_i) : i \leq \kappa)[\widetilde{M}_f^N(t+s) - \widetilde{M}_f^N(t)] \\ &= Eh(\theta^N(t_i) : i \leq \kappa) \\ &\quad \times [f(\theta^N(t+s)) - f(\theta^N(t)) - \int_t^{t+s} \mathcal{L}^N f(\theta^N(u)) du] \\ &= 0. \end{aligned} \tag{21}$$

To verify (21), we begin with the process indexed by  $\varepsilon$ . For notational simplicity, denote

$$\widetilde{h} = h(\theta^N(t_i) : i \leq \kappa), \quad \widetilde{h}^\varepsilon = h(\theta^{\varepsilon, N}(t_i) : i \leq \kappa). \tag{22}$$

Then the weak convergence and the Skorohod representation together with the boundedness and the continuity of  $f(\cdot)$  and  $h(\cdot)$  yield that as  $\varepsilon \rightarrow 0$ ,

$$\begin{aligned} & E\widetilde{h}^\varepsilon[f(\theta^{\varepsilon, N}(t+s)) - f(\theta^{\varepsilon, N}(t))] \\ & \rightarrow E\widetilde{h}[f(\theta^N(t+s)) - f(\theta^N(t))]. \end{aligned}$$

For the last term of (20), as  $\varepsilon \rightarrow 0$ , since  $f(\cdot) \in C^1_0$ ,

$$\begin{aligned} & E\widetilde{h}^\varepsilon \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} \left\{ (\nabla f(\theta_{lm_\varepsilon}^{N+}) - \nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ & \quad \times \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} [M\theta_k^N + P(\theta_k^N, \eta_k) \\ & \quad \left. + W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k)] q^N(\theta_k^N) \right\} \\ &= O(\varepsilon) \rightarrow 0. \end{aligned} \tag{23}$$

For the next to the last term of (20),

$$\begin{aligned} & \lim_{\varepsilon \rightarrow 0} E\widetilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ & \quad \times \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k) q^N(\theta_k^N) \left. \right] \\ &= \lim_{\varepsilon \rightarrow 0} E\widetilde{h}^\varepsilon \left[ \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ & \quad \times \frac{\delta_\varepsilon}{m_\varepsilon} \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} E_{lm_\varepsilon} W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k) q^N(\theta_k^N) \left. \right]. \end{aligned} \tag{24}$$

Using (A1) and (A3),

$$\frac{1}{m_\varepsilon} \sum_{j=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} E_{lm_\varepsilon} W(\theta_{lm_\varepsilon}^N, r_{1,j}, r_{2,j}, \eta_j) q^N(\theta_{lm_\varepsilon}^N) \rightarrow 0$$

in probability, we obtain that

$$\begin{aligned} & E\widetilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ & \quad \times \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} W(\theta_k^N, r_{1,k}, r_{2,k}, \eta_k) q^N(\theta_k^N) \left. \right] \rightarrow 0. \end{aligned} \tag{25}$$

Using (A1), we obtain

$$\begin{aligned} & E\widetilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\ & \quad \times \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} (P(\theta_k^N, \eta_k) - \overline{P}(\theta_k^N)) q^N(\theta_k^N) \left. \right] \rightarrow 0. \end{aligned} \tag{26}$$

Next, we consider the first term. We have

$$\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\
& \quad \times \left. \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} (M\theta_k^N + \overline{P}(\theta_k^N))q^N(\theta_k^N) \right] \\
& = \lim_{\varepsilon \rightarrow 0} E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\
& \quad \times \left. \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} (M\theta_{lm_\varepsilon}^N + \overline{P}(\theta_{lm_\varepsilon}^N))q^N(\theta_{lm_\varepsilon}^N) \right].
\end{aligned} \tag{27}$$

Thus, to get the desired limit, we need only examine the last two lines above. Let  $\varepsilon lm_\varepsilon \rightarrow u$  as  $\varepsilon \rightarrow 0$ . Then for all  $k$  satisfying  $lm_\varepsilon \leq k \leq lm_\varepsilon + m_\varepsilon - 1$ ,  $\varepsilon k \rightarrow u$  since  $\delta_\varepsilon \rightarrow 0$ . As a result,

$$\begin{aligned}
& \lim_{\varepsilon \rightarrow 0} E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} (\nabla f(\theta_{lm_\varepsilon}^N))' \right. \\
& \quad \times \left. \sum_{k=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} (M\theta_{lm_\varepsilon}^N + \overline{P}(\theta_{lm_\varepsilon}^N))q^N(\theta_{lm_\varepsilon}^N) \right] \\
& = E\tilde{h} \left[ \int_t^{t+s} (\nabla f(\theta^N(u)))' (M(\theta^N(u)) \right. \\
& \quad \left. + \overline{P}(\theta(u)))q^N(\theta(u))du \right].
\end{aligned} \tag{28}$$

The desired result then follows.  $\square$

To proceed, consider (11). For simplicity, suppose that there is a unique stationary point  $\theta^*$ . Denote  $\overline{\text{Pr}}(\theta^*) = \text{Pr}^*$  and  $\overline{\text{Pg}}(\theta^*) = \text{Pg}^*$ . By the inversion formula of partitioned matrix [45], solving  $M\theta^* + \overline{P}(\theta^*) = 0$  yields that the equilibrium point of the ODE satisfies

$$\begin{aligned}
\theta^* & = \begin{bmatrix} \kappa_1 I & -0.5\chi(c_1 + c_2)I \\ \kappa_2 I & -0.5\chi(c_1 + c_2)I \end{bmatrix}^{-1} \\
& \quad \times \begin{bmatrix} -0.5\chi(c_1 \text{Pr}^* + c_2 \text{Pg}^*) \\ -0.5\chi(c_1 \text{Pr}^* + c_2 \text{Pg}^*) \end{bmatrix} \\
& = \begin{bmatrix} 0 \\ \frac{c_1 \text{Pr}^* + c_2 \text{Pg}^*}{c_1 + c_2} \end{bmatrix}.
\end{aligned} \tag{29}$$

*Corollary 6:* Suppose that the stationary point  $\theta^*$  is asymptotically stable in the sense of Lyapunov and that  $\{\theta_n\}$  is tight. Then for any  $t_\varepsilon \rightarrow \infty$  as  $\varepsilon \rightarrow 0$ ,  $\theta^\varepsilon(\cdot + t_\varepsilon)$  converges weakly to  $\theta^*$ .

**Proof.** Define  $\tilde{\theta}^\varepsilon(\cdot) = \theta^\varepsilon(\cdot + t_\varepsilon)$ . Let  $T > 0$  and consider the pair  $\{\tilde{\theta}^\varepsilon(\cdot), \tilde{\theta}^\varepsilon(\cdot - T)\}$ . Using the same argument as in the proof of Theorem 4,  $\{\tilde{\theta}^\varepsilon(\cdot), \tilde{\theta}^\varepsilon(\cdot - T)\}$  is tight. Select a convergent subsequence with limit denoted by  $(\theta(\cdot), \theta_T(\cdot))$ . Then  $\theta(0) = \theta_T(T)$ . The value of  $\theta_T(0)$  is not known, but all such  $\theta_T(0)$ , over all  $T$  and convergent subsequences, belong to a tight set. This together with the stability and Theorem 4 implies that for any  $\Delta > 0$ , there is a  $T_\Delta$  such that for  $T > T_\Delta$ ,  $P(\theta_T(T) \in U_\Delta(\theta^*)) > 1 - \Delta$ , where  $U_\Delta(\theta^*)$  is a  $\Delta$ -neighborhood of  $\theta^*$ . The desired result then follows.  $\square$

In Corollary 6, we used the tightness of the set  $\{\theta_n\}$ , which can be proved using the argument of Lemma 9. The result indicates that as the stepsize  $\varepsilon \rightarrow 0$  and  $n \rightarrow \infty$  with  $n\varepsilon \rightarrow \infty$ ,  $\theta_n$  converges to  $\theta^*$  in the sense in probability. Note that if  $\theta^*$  turns out to be the optimum of the search space, then  $\theta_n$  converges to the optimum.

*Remark 7:* Note that for notational simplicity, we have assumed that there is a unique stationary point of (11). As far as the convergence is concerned, one need not assume that there is only one  $\theta^*$ . See how multimodal cases can be handled in the related stochastic approximation problems in [20, Chpater 5, 6, 8]. In fact, for the multimodal cases, we can show that  $\theta^\varepsilon(\cdot + t_\varepsilon)$  converges in an appropriate sense to the set of the stationary points. Thus Corollary 6 can be modified. In the rate of convergence study, [15] suggested an approach using conditional distribution, which is a modification of a single stationary point. If multiple stationary points are involved, we can simply use the approach of [15] combined with our weak convergence analysis. The notation will be a bit more complex, but main idea still rest upon the basic analysis method to be presented in the next section. It seems to be more instructive to present the main ideas, so we choose the current setting.

#### IV. RATE OF CONVERGENCE

Once the convergence of a stochastic approximation algorithm is established, the next task is to ascertain the convergence rate. To study the convergence rate, we take a suitably scaled sequence  $z_n = (\theta_n - \theta^*)/\varepsilon^\alpha$ , for some  $\alpha > 0$ . The idea is

to choose  $\alpha$  such that  $z_n$  converges (in distribution) to a nontrivial limit. The scaling factor  $\alpha$  together with the asymptotic covariance of the scaled sequence gives us the rate of convergence. That is, the scaling tells us the dependence of the estimation error  $\theta_n - \theta^*$  on the step size, and the asymptotic covariance is a mean of assessing “goodness” of the approximation. Here the factor  $\alpha = 1/2$  is used. To some extent, this is dictated by the well-known central limit theorem. For related work on convergence rate of various stochastic approximation algorithms, see [21], [51].

As mentioned above, by using the definition of the rate of convergence, we are effectively dealing with convergence in the distributional sense. In lieu of examining the discrete iteration directly, we are again taking continuous-time interpolations. Three assumptions are provided in what follows.

(A4) The following conditions hold:

- (i) in a neighborhood of  $\theta^*$ ,  $\Pr(\cdot, \eta)$  and  $\text{Pg}(\cdot, \eta)$  are continuously differentiable for each  $\eta$ , and the second derivatives (w.r.t.  $\theta$ ) of  $W(\cdot, r_1, r_2, \eta)$  and  $P(\cdot, \eta)$  exist and are continuous.
- (ii) denoting by  $E_m$  the conditional expectation on the  $\sigma$ -algebra  $\mathcal{F}_m = \{\theta_0, r_{1j}, r_{2j}, \eta_j : j < m\}$ , and by  $\zeta_\theta$  the first partial derivative w.r.t.  $\theta$  of  $\zeta = W$  or  $P$ , resp., for each positive integer  $m$ , as  $n \rightarrow \infty$ ,

$$\begin{aligned} \frac{1}{n} \sum_{j=m}^{n+m-1} E_m \Pr_\theta(\theta, \eta_j) &\rightarrow \overline{\Pr}_\theta(\theta) \text{ in probability,} \\ \frac{1}{n} \sum_{j=m}^{n+m-1} E_m \text{Pg}_\theta(\theta, \eta_j) &\rightarrow \overline{\text{Pg}}_\theta(\theta) \text{ in probability,} \\ \sum_{j=m}^{\infty} |E_m W_\theta(\theta^*, r_{1,j}, r_{2,j}, \eta_j)| &< \infty, \\ \sum_{j=m}^{\infty} |E_m P_\theta(\theta^*, \eta_j) - \overline{P}_\theta(\theta^*)| &< \infty. \end{aligned} \tag{30}$$

(iii) The matrix  $M + \overline{P}_\theta(\theta^*)$  is stable in that all of its eigenvalues are on the left half of the complex plane.

(iv) There is a twice continuously differentiable Lyapunov function  $V(\cdot) : \mathbb{R}^{2r} \rightarrow \mathbb{R}$  such that

- $V(\theta) \rightarrow \infty$  as  $|\theta| \rightarrow \infty$ , and  $V_{\theta\theta}(\cdot)$  is uniformly bounded.
- $|V_\theta(\theta)| \leq K(1 + V^{1/2}(\theta))$ .
- $|M\theta + \overline{P}(\theta)|^2 \leq K(1 + V(\theta))$  for each  $\theta$ .
- $V'_\theta(\theta)(M\theta + \overline{P}(\theta)) \leq -\lambda V(\theta)$  for some  $\lambda > 0$  and each  $\theta \neq \theta^*$ .

(A5)  $\sum_{j=m}^{\infty} |E\widetilde{W}'(\theta_m, r_{1,m}, r_{2,m}, \eta_m)\widetilde{W}(\theta_j, r_{1,j}, r_{2,j}, \eta_j)| < \infty$ , where  $\widetilde{W}(\theta, r_1, r_2, \eta) = P(\theta, \eta) - \overline{P}(\theta) + W(\theta, r_1, r_2, \eta)$ .

(A6) The sequence  $B^\varepsilon(t) = \sqrt{\varepsilon} \sum_{j=0}^{t/\varepsilon-1} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j)$  converges weakly to  $B(\cdot)$ , a Brownian motion whose covariance  $\Sigma t$  with  $\Sigma \in \mathbb{R}^{2r \times 2r}$  given by

$$\begin{aligned} \Sigma &= E\widetilde{W}(\theta^*, r_{1,0}, r_{2,0}, \eta_0)\widetilde{W}'(\theta^*, r_{1,0}, r_{2,0}, \eta_0) \\ &+ \sum_{k=1}^{\infty} E\widetilde{W}(\theta^*, r_{1,0}, r_{2,0}, \eta_0)\widetilde{W}'(\theta^*, r_{1,k}, r_{2,k}, \eta_k) \\ &+ \sum_{k=1}^{\infty} E\widetilde{W}(\theta^*, r_{1,k}, r_{2,k}, \eta_k)\widetilde{W}'(\theta^*, r_{1,0}, r_{2,0}, \eta_0). \end{aligned} \tag{31}$$

*Remark 8:* Note that (A4)(ii) is another noise condition. The motivation is similar to Remark 3. The main difference of (9) and (10) and (30) is that (30) is on the derivative of the functions evaluated at the point  $\theta^*$ . In fact, we only need the derivative exists in a neighborhood of this point only. This is because that we are analyzing the asymptotic normality locally. In view of this condition and condition of  $\{r_{i,n}\}$ ,

$$\begin{aligned} \frac{1}{n} \sum_{j=m}^{m+n-1} E_m W_\theta(\theta^*, r_{1,j}, r_{2,j}, \eta_j) &\rightarrow 0 \text{ in probability,} \\ \frac{1}{n} \sum_{j=m}^{m+n-1} E_m P_\theta(\theta^*, \eta_j) &\rightarrow \overline{P}_\theta(\theta^*) \text{ in probability.} \end{aligned}$$

The traditional PSO algorithms do not allow non-additive noise, here we are treating a more general problem. Nonadditive noise can be allowed.

(A4)(iv) assumes the existence of a Lyapunov function. Only the existence is needed; its precise form need not be known. For simplicity, we have assumed the convergence of the scaled sequence to a Brownian motion in (A6); sufficient conditions are well known; see for example, [20, Section 7.4]. Before proceeding further, we first obtain a moment bound of  $\theta_n$ .

*Lemma 9:* Assume that (A1)-(A6) hold. Then there is an  $N_\varepsilon$  such that for all  $n > N_\varepsilon$ ,  $EV(\theta_n) = O(\varepsilon)$ .

**Proof.** To begin, it can be seen that

$$\begin{aligned} E_n V(\theta_{n+1}) - V(\theta_n) \\ \leq -\varepsilon\lambda V(\theta_n) + \varepsilon E_n V'_\theta(\theta_n) \widetilde{W}(\theta_n, r_{1,n}, r_{2,n}, \eta_n) \\ + O(\varepsilon^2)(1 + V(\theta_n)), \end{aligned} \quad (32)$$

where  $\theta_n^+$  is on the line segment joining  $\theta_n$  and  $\theta_{n+1}$ . The bound in (32) follows from the growth condition in (A4)(iv), the last inequality follows from (A1). To proceed, we use the methods of perturbed Lyapunov functions, which entitles to introduce small perturbations to a Lyapunov function in order to make desired cancelation. Define a perturbation

$$V_1^\varepsilon(\theta, n) = \varepsilon \sum_{j=n}^{\infty} E_n V'_\theta(\theta) \widetilde{W}(\theta, r_{1,j}, r_{2,j}, \eta_j).$$

Note that

$$|V_1^\varepsilon(\theta, n)| = K\varepsilon(1 + V(\theta)). \quad (33)$$

Moreover,

$$\begin{aligned} E_n V_1^\varepsilon(\theta_{n+1}, n+1) - V_1^\varepsilon(\theta_n, n) \\ = O(\varepsilon^2)(V(\theta_n) + 1) - \varepsilon E_n V'_\theta(\theta_n) \widetilde{W}(\theta_n, r_{1,n}, r_{2,n}, \eta_n). \end{aligned} \quad (34)$$

Define  $V^\varepsilon(\theta, n) = V(\theta) + V_1^\varepsilon(\theta, n)$ . Using (32) and (34), we obtain

$$E_n V^\varepsilon(\theta_{n+1}, n+1) \leq (1 - \varepsilon\lambda)V^\varepsilon(\theta_n, n) + O(\varepsilon^2)(1 + V^\varepsilon(\theta_n, n)). \quad (35)$$

Choosing  $N_\varepsilon$  to be a positive integer such that  $(1 - (\lambda\varepsilon/2))^{N_\varepsilon} \leq K\varepsilon$ . Iterating on the recursion (35), taking expectation, and using the order of magnitude estimate (33), we can then obtain

$$\begin{aligned} EV^\varepsilon(\theta_{n+1}, n+1) \\ \leq (1 - \varepsilon\lambda)EV^\varepsilon(\theta_n, n) + O(\varepsilon^2)(1 + V^\varepsilon(\theta_n, n)) \\ \leq (1 - \frac{\varepsilon\lambda}{2})^n EV^\varepsilon(\theta_0, 0) + O(\varepsilon) = O(\varepsilon). \end{aligned} \quad (36)$$

when  $n > N_\varepsilon$ . The second line of (36) follows from  $1 - \lambda\varepsilon + O(\varepsilon^2) \leq 1 - \frac{\lambda\varepsilon}{2}$  for sufficiently small  $\varepsilon$ . Now using (33) again, we also have  $EV(\theta_{n+1}) = O(\varepsilon)$ . Thus the desired estimate follows.  $\square$

Define  $z_n = (\theta_n - \theta^*)/\sqrt{\varepsilon}$ . Then it is readily verified that

$$\begin{aligned} z_{n+1} = z_n + \varepsilon(M + \overline{P}_\theta(\theta^*))z_n \\ + \sqrt{\varepsilon}(P(\theta^*, \eta_n) - \overline{P}(\theta^*) + W(\theta^*, r_{1,n}, r_{2,n}, \eta_n)) \\ + \varepsilon(P_\theta(\theta^*, \eta_n) - \overline{P}_\theta(\theta^*) \\ + W_\theta(\theta^*, r_{1,n}, r_{2,n}, \eta_n))z_n + o(|z_n|^2). \end{aligned} \quad (37)$$

*Corollary 10:* Assume that (A1)-(A6) hold. If the Lyapunov function is locally quadratic, i.e.,

$$V(\theta) = (\theta - \theta^*)'Q(\theta - \theta^*) + o(|\theta - \theta^*|^2).$$

Then  $EV(z_n) = O(1)$  for all  $n > N_\varepsilon$ .

Now we are in a position to study the asymptotic properties through weak convergence of appropriately interpolated sequence of  $z_n$ . Define  $z^\varepsilon(t) = z_n$  for  $t \in [(n - N_\varepsilon)\varepsilon, (n - N_\varepsilon)\varepsilon + \varepsilon]$ . We can introduce a truncation sequence. That is, in lieu of  $z^\varepsilon(\cdot)$ , we let  $N$  be a fixed but otherwise arbitrary large positive integer and define  $z^{\varepsilon, N}(\cdot)$  as an  $N$ -truncation of  $z^\varepsilon(\cdot)$ . That is, it is equal to  $z^\varepsilon(\cdot)$  up until the first exit of the process from the sphere  $S_N = \{|z| \leq N\}$  with radius  $N$ . Also define a truncation function  $q^N(z) = 1$  if  $z \in S_N$ ,  $= 0$  if  $z \in \mathbb{R}^r - S_{N+1}$ , and is smooth. Corresponding to such a truncation, we also have a modified operator with truncation (i.e., the functions used in the operator are all modified by use of  $q^N(z)$ ). Then we proceed to establish the convergence of  $z^{\varepsilon, N}(\cdot)$  as a solution of a martingale problem with the truncated operator. Then finally, letting  $N \rightarrow \infty$ , we use the uniqueness of the martingale problem to conclude the proof. The argument is similar to that of Section III. For further technical details, we refer the reader to [20, pp. 284-285]. Such a truncation device is also widely used in the analysis of partial differential equations. For notational simplicity, we choose to simply assume the boundedness rather than go with the truncation route. Thus merely for notational simplicity, we suppose  $z^\varepsilon(t)$  is bounded. For the rate of convergence, our focus is on the convergence of the sequence  $z^\varepsilon(\cdot)$ . We shall show that it converges to a diffusion process whose covariance matrix together with the scaling factor will provide us with the desired convergence rates. Although more complex than Theorem 4, we still use the martingale problem setup. To keep the presentation relatively brief, we shall only outline the main steps needed.

For any  $t, s > 0$ ,

$$\begin{aligned}
z^\varepsilon(t+s) - z^\varepsilon(t) &= \varepsilon \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon-1} (M + \bar{P}_\theta(\theta^*)) z_j \\
&\quad + \sqrt{\varepsilon} \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) \\
&\quad + \varepsilon \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon} \widetilde{W}_\theta(\theta^*, r_{1,j}, r_{2,j}, \eta_j) z_j.
\end{aligned} \tag{38}$$

Note that for any  $\delta > 0$ ,  $t, s > 0$  with  $s < \delta$ ,

$$\begin{aligned}
&E_t^\varepsilon \left| \sqrt{\varepsilon} \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon-1} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) \right|^2 \\
&\leq K\varepsilon \left( \frac{t+s}{\varepsilon} - \frac{t}{\varepsilon} \right) = Ks \leq K\delta.
\end{aligned}$$

and similarly,

$$E_t^\varepsilon \left| \varepsilon \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon-1} \widetilde{W}_\theta(\theta^*, r_{1,j}, r_{2,j}, \eta_j) z_j \right|^2 \leq Ks \leq K\delta.$$

Using Corollary 10 and similar argument as that of Theorem 4, we have the following result.

*Lemma 11:* Assume conditions of Corollary 10,  $\{z^\varepsilon(\cdot)\}$  is tight on  $D([0, T] : \mathbb{R}^{2r})$ .

Next we can extract a convergent subsequence of  $\{z^\varepsilon(\cdot)\}$ . Without loss of generality, still denote the subsequence by  $z^\varepsilon(\cdot)$  with limit  $z(\cdot)$ . For any  $t, s > 0$ , (38) holds. The way to derive the limit is similar to that of Theorem 4 using martingale problem formulation although the analysis is more involved. We proceed to show that the limit is the unique solution for the martingale problem with operator

$$Lf(z) = \frac{1}{2} \text{tr}(\Sigma f_{zz}(z)) + (\nabla f(z))'(M + \bar{P}(\theta_*)), \tag{39}$$

for  $f \in C_0^2$ ,  $C^2$  functions with compact support.

Using similar notation as that of Section III, redefine

$$\tilde{h} = h(z(t_i) : i \leq \kappa), \quad \tilde{h}^\varepsilon = h(z^\varepsilon(t_i) : i \leq \kappa). \tag{40}$$

By (A4) (ii), as  $\varepsilon \rightarrow 0$

$$\begin{aligned}
&E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) z_j \right] \\
&= E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} \sum_{j=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) z_j \right] \rightarrow 0.
\end{aligned}$$

Using the notation as in Section III,

$$\begin{aligned}
&E\tilde{h}^\varepsilon \left[ \sum_{l=t/\delta_\varepsilon}^{(t+s)/\delta_\varepsilon} \frac{\delta_\varepsilon}{m_\varepsilon} \sum_{j=lm_\varepsilon}^{lm_\varepsilon+m_\varepsilon-1} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) [z_j - z_{lm_\varepsilon}] \right] \\
&\rightarrow 0 \quad \text{as } \varepsilon \rightarrow 0.
\end{aligned}$$

Moreover, by (A6) we have

$$\sqrt{\varepsilon} \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon} \widetilde{W}(\theta^*, r_{1,j}, r_{2,j}, \eta_j) \rightarrow \int_t^{t+s} dB(u)$$

as  $\varepsilon \rightarrow 0$ . For the first term of (38), we have

$$\begin{aligned}
&E\tilde{h}^\varepsilon \left[ \varepsilon \sum_{j=t/\varepsilon}^{(t+s)/\varepsilon-1} (M + \bar{P}_\theta(\theta^*)) z_j \right] \\
&\rightarrow E\tilde{h} \left[ \int_t^{t+s} (M + \bar{P}_\theta(\theta^*)) z(u) du \right]
\end{aligned}$$

as  $\varepsilon \rightarrow 0$ . Putting the aforementioned arguments together, we have the following theorem.

*Theorem 12:* Under conditions (A1)-(A7),  $\{z^\varepsilon(\cdot)\}$  converges to  $z(\cdot)$  such that  $z(\cdot)$  is a solution of the following stochastic differential equation

$$dz = [M + \bar{P}_\theta(\theta^*)]zdt + \Sigma^{1/2}d\widehat{B}(t), \quad (41)$$

where  $\widehat{B}(\cdot)$  is a standard Brownian motion.

*Remark 13:* To see what kind of functions and the associated ODE and SDE we are working with, we look at two simple examples. In the first example we use  $F(x) = x^2$ , take 2 particles,  $\chi = 1$ ,  $\kappa_1 = -0.271$ ,  $\kappa_2 = 1$ ,  $c_1 = c_2 = 1.5$ , and assume  $\{\eta_k\}$  is an i.i.d. sequence with mean  $[0, 0, 0, 0]'$  and variance  $I$ . Then

$$M = \begin{bmatrix} -0.271 & 0 & -1.5 & 0 \\ 0 & -0.271 & 0 & -1.5 \\ 1 & 0 & -1.5 & 0 \\ 0 & 1 & 0 & -1.5 \end{bmatrix}, \quad (42)$$

and the limit ODE is given by

$$\dot{\theta}(t) = M\theta(t).$$

Thus  $\theta^* = [0, 0, 0, 0]'$  is the minimizer of the swarm, and  $P_\theta(\theta^*) = 0 \in \mathbb{R}^{4 \times 4}$  (a  $4 \times 4$  matrix with all entries being 0). In the standard optimization algorithm, one processor is running to approximate the optimum. Here, we have two particles running simultaneously. Note that  $\theta$  has four components. Two of them represent the particles' positions, and the other two are the particles' speeds. At the end, both of the particles reach the minimum, representing something that might be called "overlapping." In addition, eventually the speeds of both particles reach 0 (or at resting point). As far as the rate of convergence is concerned, we conclude that  $\theta_n - \theta^*$  decays in the order of  $\sqrt{\varepsilon}$  (in the sense of convergence in distribution). Not only is the mean squares error of  $(\theta_n - \theta^*)$  of the order  $\varepsilon$ , but also the interpolation of the scaled sequence  $z_n$  has a limit represented by a stochastic differential equation

$$dz = Mzdt + d\widehat{B}(t).$$

That is, (41) is satisfied with  $P_\theta(\theta^*) = 0$  and  $\Sigma = I$ . As illustrated in [20], the scaling factor  $\sqrt{\varepsilon}$  together with stationary covariance of the SDE gives us the rate of convergence. In terms of the swarm, loosely, we have  $\theta_n - \theta^* \sim N(0, \varepsilon \Xi_0)$  [that is,  $(\theta_n - \theta^*)$  is asymptotically normal with mean  $0 \in \mathbb{R}^4$  and covariance matrix  $\varepsilon \Xi_0$ ], where  $\Xi_0$  is the asymptotic covariance matrix that is the solution of the Lyapunov equation  $M\Xi_0 + \Xi_0M' = -I$ .

Likewise, in the second example,  $F(x) = \sin x$  with  $x \in [0, 1]$ . We still take 2 particles, same parameters setting, and assume  $\{\eta_k\}$  is the same i.i.d. sequence as before. Then  $M$  is as in (42), and

$$P_\theta(\theta^*) = \begin{bmatrix} 0.75 & 0 & 0.75 & 0 \\ 0 & 0.75 & 0 & 0.75 \\ 0.75 & 0 & 0.75 & 0 \\ 0 & 0.75 & 0 & 0.75 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}.$$

It follows that (41) holds with

$$M + \bar{P}_\theta(\theta^*) = \begin{bmatrix} 1.229 & 0 & -1.5 & 0 \\ 0 & 1.229 & 0 & -1.5 \\ 2.5 & 0 & -1.5 & 0 \\ 0 & 2.5 & 0 & -1.5 \end{bmatrix}$$

and  $\Sigma = I$ . Similar to the previous example, we have that  $\theta_n - \theta^*$  is asymptotically normal with mean 0 and covariance  $\varepsilon \tilde{\Xi}$ , where  $\tilde{\Xi}$  is the asymptotic covariance satisfying the Lyapunov equation  $(M + \bar{P}_\theta(\theta^*))\tilde{\Xi} + \tilde{\Xi}(M + \bar{P}_\theta(\theta^*))' = -I$ .

## V. NUMERICAL EXAMPLES

We use two simulation examples to illustrate the convergence properties. Using (5), we take  $\varepsilon = 0.01$ ,  $\chi = 1$ ,  $\kappa_1 = -0.271$ ,  $\kappa_2 = 1$ ,  $c_1 = c_2 = 1.5$ . For simplicity, we take the additive noise  $\Pr(\theta_n, \eta_n) = \Pr(\theta_n) + \eta_n$  and  $\text{Pg}(\theta_n, \eta_n) = \text{Pg}(\theta_n) + \eta_n$ , where  $\eta_n$  is a sequence of i.i.d. random variables with a standard normal distribution  $\mathcal{N}(0, 1)$ . In addition, we set the number of swarms to be 5.

*Example 14:* Consider the sphere function:

$$F_1(x) = \sum_{i=1}^D x_i^2, \quad (43)$$

where  $D$  is the dimension of the variable  $x$ . Its global optimum is  $(0, 0, \dots, 0)'$ . First, the dimension of  $X$  is set to be 1. Figures 1 and 2 show the state trajectories (top) and the centered and scaled errors of the first component  $\theta_n^1$  (bottom). The graphs of  $\Pr$  (top) and  $\text{Pg}$  (bottom) are also provided.

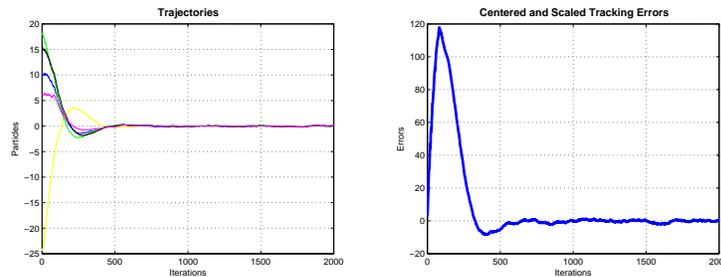


Fig. 1. Particle swarm of one-dimensional  $X$  using  $F_1$  defined in (43).

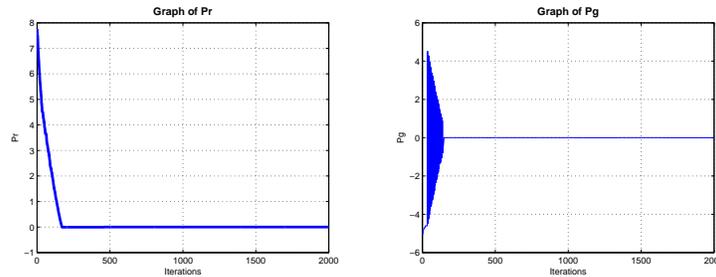


Fig. 2. Graphs of  $Pr$  and  $Pg$  using  $F_1$  defined in (43).

Next, we consider the 2-dimension case of  $X$ . Figures 3 and 4 illustrate the state trajectories (top) and the centered and scaled errors of the first component  $\theta_n^1$  (bottom), and the graph of  $Pr$  (top) and  $Pg$  (bottom), respectively.

*Example 15:* Consider the Rastrigin function [41]

$$F_2(x) = 10D + \sum_{i=1}^D [x_i^2 - 10 \cos(2\pi x_i)], \quad (44)$$

where  $D$  is the dimension of the variable  $x$ .

This function has many local minima. Its global optimum is given by  $(0, 0, \dots, 0)'$ . Same as Example 14, we set the dimension of  $X$  to be 1 and 2, respectively. The particle swarm trajectories, the centered and scaled errors of the first component, and graphs of  $Pr$  and  $Pg$  are given in Figures 5 to 8, respectively.

From these figures, we can conclude that all the swarms converge to a point in the searching space. These results were obtained without assuming that  $r_1$ ,  $r_2$ ,  $Pr$ , and  $Pg$  are fixed. Our numerical results confirm our theoretical findings in Sections III and IV.

*Remark 16:* We use the definition of convergence here that a swarm collapse in which all particles have converged to a point in the search space. Sometimes we observe (e.g., in the second example) that the convergence point is not the global or even local optimum. This problem, referred to as premature in literatures, occurs commonly in evolutionary algorithms such as PSOs, genetic algorithms, evolutionary strategies, etc. Based on our numerical experiments, we found that if the cost function is unimodal and with low dimensions, the equilibrium coincides with the proper parameter choice. The problem of under what conditions the equilibrium coincides with the optimum deserves to be carefully studied in the future.

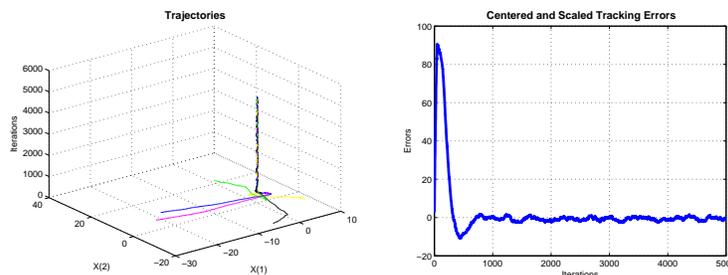


Fig. 3. Particle swarm of two-dimensional  $X$  using  $F_1$  defined in (43).

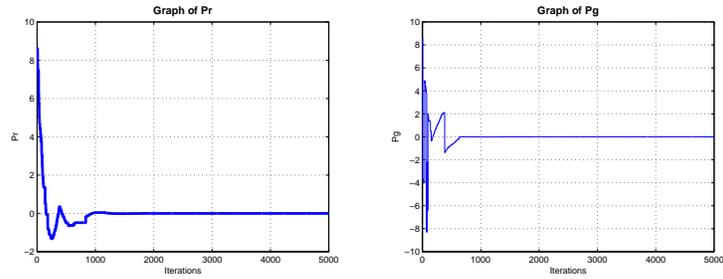


Fig. 4. Graphs of Pr and Pg using  $F_1$  defined in (43).

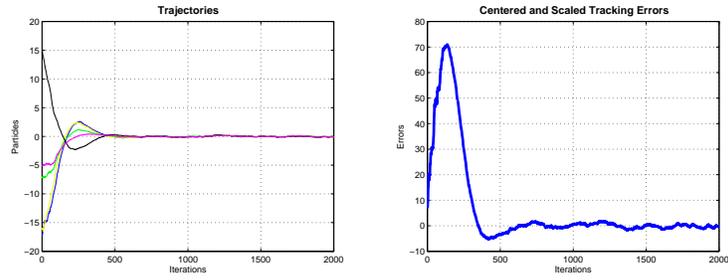


Fig. 5. Particle swarm of one-dimensional  $X$  using  $F_2$  defined in (44).

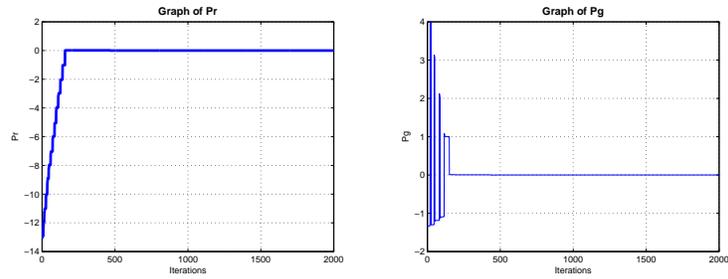


Fig. 6. Graphs of Pr and Pg using  $F_2$  defined in (44).

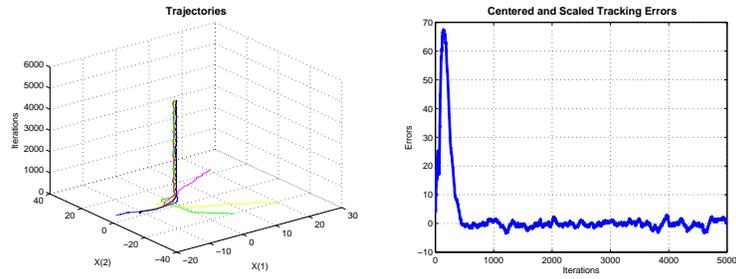


Fig. 7. Particle swarm of two-dimensional  $X$  using  $F_2$  defined in (44).

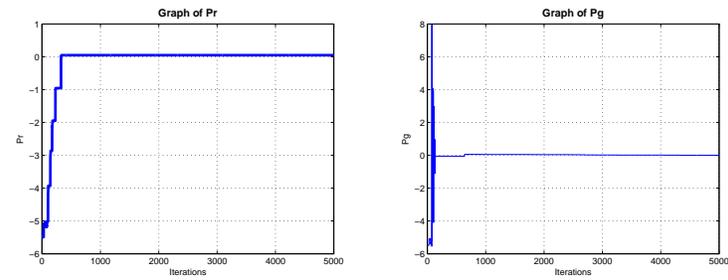


Fig. 8. Graphs of Pr and Pg using  $F_2$  defined in (44).

## VI. FURTHER REMARKS

In this paper, we considered a general form of PSO algorithms using a stochastic approximation scheme. Different from the existing results in the literature, we have used weaker assumptions and obtained more general results without depending on empirical work. In addition, we obtained rates of convergence for the PSO algorithms for the first time.

Several research directions may be pursued in the future. We can use stochastic approximation methods to analyze other schemes of PSO, for example, the SPSO2011 considered in [53]. We can set up a stochastic approximation similar to (8) and analyze its convergence and convergence rate. Finding ways to systematically choose the parameter values  $\kappa_1$ ,  $\kappa_2$ ,  $c_1$ , and  $c_2$  is a practically challenging problem. One thought is to construct a level two (stochastic) optimization algorithm to select best parameter value in a suitable sense. To proceed in this direction requires careful thoughts and consideration. In addition, we can consider that some parameters such as  $\chi$ ,  $\kappa_1$ , etc. are not fixed but change randomly during iterations or change owing to some random environment change (for example, see [52]). The problem to study is to analyze the convergence and convergence rates in such a case. Furthermore, using another definition of convergence, i.e., the swarm's best known position  $P_g$  approaching (converging to) the optimum of the problem, is another possible study direction.

To conclude, this paper demonstrated convergence properties of a class of general PSO algorithms and derived the rates of convergence by using a centered and scaled sequence of the iterates. This study opens new arenas for subsequent studies on determining convergence capabilities of different PSO algorithms and parameters.

## REFERENCES

- [1] M.A. Abido, "Particle swarm optimization for multimachine power system stabilizer design", in *Proc. Power Eng. Soc. Summer Meeting*, 2001, pp. 1346-1351.
- [2] M. R. AlRashidi, M. E. El-Hawary, "A Survey of Particle Swarm Optimization Applications in Electric Power Systems", *IEEE Trans. Evolutionary Comp.*, vol. 13, no. 4, pp. 913-918, 2009.
- [3] T. Beielstein, K.E. Parsopoulos, and M.N. Vrahatis, "Tuning pso parameters through sensitivity analysis", in: *Technical Report, Reihe Computational Intelligence CI 124/02*, Dept. Computer Sci., Univ. of Dortmund, 2002.
- [4] A. Benveniste, M. Metivier, and P. Priouret, *Adaptive Algorithms and Stochastic Approximations*, Springer-Verlag, Berlin, 1990.
- [5] P. Billingsley, *Convergence of Probability Measures*, J. Wiley, New York, NY, 1968.
- [6] B. Brandstaier and U. Baumgartner, "Particle swarm optimization Mass-spring system analogon", *IEEE Trans. Magn.*, vol. 38, no. 2, pp. 997-1000, 2002.
- [7] D. Bratton, and J. Kennedy, "Defining a standard for particle swarm optimization", in *Proc IEEE Swarm Intell. Symp*, 2007, pp. 120-127.
- [8] X. Chen and Y. Li, "A modified PSO structure resulting in high exploration ability with convergence guaranteed". *IEEE Trans. Syst. Man Cybern. B Cybern.*, vol. 37, no. 5, pp. 1271-1289, 2007.
- [9] M. Clerc and J. Kennedy, "The particle swarm: explosion, stability, and convergence in a multidimensional complex space". *IEEE Trans. Evolut. Comput.*, vol. 6, no. 1, pp. 58-73, 2002.
- [10] Jr. E.F. Costa, P.L.C Lage, and Jr. E.C. Biscaia. "On the numerical solution and optimization of styrene polymerization in tubular reactors". *Comput. Chem. Eng.*, vol. 27, no. 11, pp. 1591-1604, 2003.
- [11] H.M. Emarra and H.A. Fattah, "Continuous swarm optimization technique with stability analysis", in *Proc. Amer. Control Conf.*, 2004, vol. 3, pp. 2811-2817.
- [12] M. Jiang, Y.P. Luo, and S.Y. Yang. "Stochastic Convergence Analysis and Parameter Selection of the Standard Particle Swarm Optimization Algorithm". *Inf. Process Lett.*, vol. 102, no. 1, pp. 8-16, 2007.
- [13] C.F. Juang, "A hybrid of genetic algorithm and particle swarm optimization for recurrent network design", *IEEE Trans. Syst., Man, Cybern. B, Cybern.*, vol. 34, no. 2, pp. 997-1006, 2004.
- [14] S. Kamisetty, J. Garg, J.N. Tripathi, and J. Mukherjee. "Optimization of Analog RF Circuit Parameters Using Randomness in Particle Swarm Optimization". in: *WICT*, pp. 274-278, 2011.
- [15] Yu.M. Kaniovskii, "Limit distribution of processes of stochastic approximation type when the regression function has several roots", (in Russian) *Dokl. Akad. Nauk SSSR* **301** vol. 6, pp. 1308-1309, 1988.
- [16] J. Kennedy, R.C. Eberhart, "Particle swarm optimization", in: *Proc. IEEE Conf. on Neural Networks, IV*, Piscataway, NJ, 1995, pp. 1942-1948.
- [17] S. Kiranyaz, T. Ince, and M. Gabbouj. "Stochastic Approximation Driven Particle Swarm Optimization". in: *IIT'09*, 2009, pp. 40-44.
- [18] G. Kovacs, A. Groenwold, and K. Jarmai, et al., "Analysis and optimum design of fibereinforced composite structures", *Struct. Multidiscip. Opti.*, vol. 28, no. 2-3, pp. 170-179, 2004.
- [19] H.J. Kushner, *Approximation and Weak Convergence Methods for Random Processes, with Applications to Stochastic Systems Theory*, MIT Press, Cambridge, MA, 1984.
- [20] H.J. Kushner and G. Yin, *Stochastic Approximation and Recursive Algorithms and Applications*, 2nd Ed., Springer-Verlag, New York, NY, 2003.
- [21] P. L'Ecuyer and G. Yin, "Budget-dependent convergence rate of stochastic approximation", *SIAM J. Optim.* vol. 8, no. 1, pp. 217-247, 1998.
- [22] Y. Li and X. Chen, "Mobile robot navigation using particle swarm optimization and adaptive NN", in: *Proc. 1st Int. Conf. Nat. Comput.*, Changsha, China, Lecture Notes in Computer Science, vol. 3612. Berlin, Germany: Springer-Verlag, 2005, pp. 554-559.
- [23] H. Liu, A. Abraham, and V. Snasel. "Convergence Analysis of Swarm Algorithm". in *NaBIC2009*, pp. 1714-1719, 2009.
- [24] Y. Liu and K.M. Passino, "Stable social foraging swarms in a noisy environment", *IEEE Trans. Automat. Contr.*, vol. 49, no. 1, pp. 30-44, 2004.
- [25] Y. Liu, Z. Qin, and Z. Shi, "Hybrid particle swarm optimizer with line search", in *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2004, vol. 4, pp. 3751-3755.
- [26] B. Liu, L. Wang, Y. Jin, "An effective PSO-based memetic algorithm for flowshop scheduling". *IEEE Trans. Syst. Man. Cy. B.-Cybernetics*, vol. 37, no. 1, pp. 18-27, 2007.
- [27] J.L. Fernandez-Martinez, E. Garcıa-Gonzalo, "Stochastic Stability Analysis of the Linear Continuous and Discrete PSO Models". *IEEE Trans. Evol. Comput.*, vol. 15, no. 3, pp. 405-423, 2011.
- [28] R. Mendes, "Population Topologies and Their Influence in Particle Swarm Performance", Ph.D. Thesis, Universidade do Minho, 2004.
- [29] L. Messerschmidt and A.P. Engelbrecht, "Learning to play games using a PSO-based competitive learning approach", *IEEE Trans. Evol. Comput.*, vol. 8, no. 3, pp. 280-288, 2004.
- [30] V. Miranda, H. Keko, and A.J. Duque, "Stochastic Star Communication Topology in Evolutionary Particle Swarms (EPSO)", *Inter. J. Comput. Intelligent Research*, vol. 4, no. 2, pp. 105-116, 2008.
- [31] J. Nenoraite, R. Simutis, "Stocks' trading system based on the particle swarm optimization algorithm", in: *ICCS 2004*, 2004, pp. 843-850.
- [32] K.E. Parsopoulos, E.I. Papageorgiou, and P.P. Groumpos, et al. "Evolutionary computation techniques for optimizing fuzzy cognitive maps in radiation therapy systems". In: *Proc. GECCO*. 2004, pp. 402-413.

- [33] N.G. Pavlidis, K.E. Parsopoulos, and M.N. Vrahatis, "Computing Nash equilibria through computational intelligence methods". *J. Comput. Appl. Math.*, vol. 175, no. 1, pp. 113-136, 2005.
- [34] K.E. Parsopoulos and M.N. Vrahatis, "Recent approaches to global optimization problems through particle swarm optimization", *Nat. Comput.*, vol. 1, no. 2-3, pp. 235-306, 2002.
- [35] M. Pedersen, "Good parameters for particle swarm optimization". *Hvass Laboratories Technical Report no HL1001*, 2010.
- [36] M.E.H. Pedersen and A.J. Chipperfield, "Simplifying Particle Swarm Optimization", *Appl. Soft Computing*, vol. 10, no. 2, pp. 618-628, 2010.
- [37] R. Poli, "Dynamics and Stability of the Sampling Distribution of Particle Swarm Optimisers via Moment Analysis", *J. Artif. Evol. Appl.*, 2008. doi:10.1155/2008/761459.
- [38] R. Poli, C.D. Chio, and W.B. Langdon, "Exploring extended particle swarms: A genetic programming approach", in: *Proc. Conf. Genet. & Evol. Comput.*, Washington DC, 2005, pp. 169-176.
- [39] T.K. Rasmussen, T. Krink, "Improved hidden Markov model training for multiple sequence alignment by a particle swarm optimization-evolutionary algorithm hybrid". *Biosystems*, vol. 72, no. 1-2, pp. 5-17, 2003.
- [40] A. Ratnaweera, S.K. Halgamuge, and H. C. Watson, "Self-organizing hierarchical particle swarm optimizer with time-varying acceleration coefficients", *IEEE Trans. Evol. Comput.*, vol. 8, no. 3, pp. 240-255, 2004.
- [41] R.G. Reynolds and C.-J. Chung, "Knowledge-based self-adaption in evolutionary programming using cultural algorithms", in: *Proc. IEEE Int. Conf. Evolutionary Computation*, Indianapolis, IN, 1997, pp. 71-76.
- [42] Q. Shen, J. Jiang, and C. Jiao, et al., "Modified particle swarm optimization algorithm for variable selection in MLR and PLS modeling: QSAR studies of antagonism of angiotensin II: antagonists". *Eur. J. Pharm. Sci.*, vol. 22, no. 2-3, pp. 145-152, 2004.
- [43] Y. Shi and R. Eberhart, "A modified particle swarm optimizer", in: *Proc. IEEE World Congr. Comput. Intell.*, 1998, pp. 69-73.
- [44] M.F. Tasgetiren, M. Sevkli, and Y. Liang, et al., "Partical swarm optimization algorithm for permutation flowshop sequencing problem". *Lecture Notes in Comput. Sci.*, vol. 3172, pp. 382-389, 2004.
- [45] Y. Tian and Y. Takane, "The inverse of any two-by-two nonsingular partitioned matrix and three matrix inverse completion problems", *Comp. & Math. Appl.*, vol. 57, no. 8, pp. 1294-1304, 2009.
- [46] I.C. Trelea, "The Particle Swarm Optimization Algorithm: convergence analysis and parameter selection". *Inf. Process Lett.* vol. 85, no. 6, pp. 317-325, 2003.
- [47] H. Wu, F. Sun, and Z. Sun, et al. "Optimal trajectory planning of a flexible dual-arm space robot with vibration reduction". *J. Intell. Robot. Syst.*, vol. 40, no. 2, pp. 147-163, 2004.
- [48] X. Xiao, E.R. Dow, R. Eberhart, et al., "Hybrid self-organizing maps and particle swarm optimization approach". *Concurr. Comp-Pract. E.*, vol. 16, no. 9, pp. 895-915, 2004.
- [49] R. Xiao, B. Li, and X. He, "The Particle Swarm: Parameter Selection and Convergence", *CCIS*, vol. 2, pp. 396-402, 2007.
- [50] K. Yasuda, A. Ide, and N. Iwasaki, "Adaptive particle swarm optimization", in: *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2003, pp. 1554-1559.
- [51] G. Yin, "Rates of convergence for a class of global stochastic optimization algorithms", *SIAM J. Optim.*, vol. 10, no. 1, pp. 99-120, 1999.
- [52] G. Yin and C. Zhu, *Hybrid Switching Diffusions: Properties and Applications*, Springer, New York, 2010.
- [53] M. Zambrano-Bigiarini, M. Clerc, and R. Rojas, "Standard particle swarm optimization 2011 at CEC-2013: A baseline for future PSO improvements", *IEEE Congr. Evol. Comp. (CEC)*, Cancun, 2013, pp. 2337-2344.
- [54] L. Zhang, H. Yu, and S. Hu, "Optimal choice of parameters for particle swarm optimization". *J. Zhejiang Univ. Sci.* vol. 6A, no. 6, pp. 528-534, 2005.
- [55] X. Zhang, L. Yu, and Y. Zheng, et al. "Two-stage adaptive PMD compensation in a 10 Gbit/s optical communication system using particle swarm optimization". *Opt. Commun.*, vol. 231, no. 1-6, pp. 233-242, 2004.
- [56] W.J. Zhang and X.F. Xie, "DEPSO: Hybrid particle swarm with differential evolution operator", in: *Proc. IEEE Int. Conf. Syst., Man, Cybern.*, 2003, pp. 3816-3821.