# Working Paper

## Strong Convergence of Stochastic Approximation Without Lyapunov Functions

*Yuri M. Kaniovski*

# Strong Convergence of Stochastic Approximation Without Lyapunov Functions

*Yuri M. Kaniovski*

WP-95-19
February 1995

# Summary

We prove convergence with probability one of a multivariate Markov stochastic approximation procedure of the Robbins – Monro type with several roots. The argument exploits convergence of the corresponding system of ordinary differential equations to its stationary points. If the points are either linearly stable or linearly unstable, we prove convergence with probability 1 of the procedure to a random vector whose distribution concentrates on the set of stable stationary points. This generalizes for procedures with several roots the approach suggested by L. Ljung for processes with a single root.

Along with stochastic approximation processes as such, the result can be applied to generalized urn schemes and stochastic models of technological and economic dynamics based on them, in particular, evolutionary games with incomplete information.

**Key words:** stochastic approximation procedure with several roots, strong convergence, Bendixon theorem, cycle, phase polygon, generalized urn scheme, evolutionary game.

# Contents

# Strong Convergence of Stochastic Approximation Without Lyapunov Functions

*Yuri M. Kaniovski*

## 1 Motivation and formulation of the problem

Suppose we are interested in finding a root of a (Lipschitz) continuous $m$-dimensional regression function $\vec{g}(\cdot)$ given on $\mathbf{R}^m$. We cannot observe $\vec{g}(\cdot)$ itself, but only

$$\vec{\Xi}(t, \vec{x}) = \vec{g}(\vec{x}) + \vec{W}^t(\vec{x}) + \vec{Z}^t(\vec{x}), \ t \geq 1,$$

independent in $t$ observations of $\vec{g}(\cdot)$ together with deterministic $\vec{W}^t(\cdot)$ and random $\vec{Z}^t(\cdot)$ noise. It is assumed that the random noise has zero mean, i.e. $\mathbf{E}\vec{Z}^t(\vec{x}) = \vec{0}$, finite variance, i.e. $\mathbf{E}||\vec{Z}^t(\vec{x})||^2 \leq c$, and there is a measurable with respect to $\mathcal{B}_{2m}$ mapping $\vec{n}(t, \cdot, \cdot)$ such that $\vec{Z}^t(\vec{x}) = \vec{n}(t, \vec{x}, \vec{\xi}^t)$. Here $|| \cdot ||$ stands for the Euclidean norm in $\mathbf{R}^m$, $\mathcal{B}_l$ designates the $\sigma$-algebra of Borel sets in $\mathbf{R}^l$. Also $\vec{\xi}^t$, $t \geq 1$, is a sequence of independent $m$-dimensional random vectors on a probability space $\{\Omega, \mathcal{F}, P\}$. The deterministic noise with respect to $\mathcal{B}_m$ is a sequence of measurable two-dimensional vector-functions. The Robbins – Monro procedure gives successive approximations $\vec{X}^t$, $t \geq 1$, for finding a root of $\vec{g}(\cdot)$ in the following form

$$\vec{X}^{t+1} = \vec{X}^t + \gamma_t \vec{\Xi}(t, \vec{X}^t), \tag{1}$$

where $\vec{X}^1$ is a deterministic vector, $\gamma_t$ stands for the step-sizes, i.e. positive numbers such that

$$\sum_{t \geq 1} \gamma_t = \infty, \quad \sum_{t \geq 1} \gamma_t^2 < \infty, \quad \sum_{t \geq 1} \gamma_t \sup_{\vec{x}} ||\vec{W}^t(\vec{x})|| < \infty. \tag{2}$$

Since we assumed that the random noise is formed by independent random vectors, the random process $\vec{X}^t$ is *Markov*.

Traditionally the regression function $\vec{g}(\cdot)$ is assumed to have a single root. But it is a hypothesis to be checked. In many cases, like generalized urn schemes [1] or evolutionary games with incomplete information [3,6], it does not hold. We shall study here the case when $\vec{g}(\cdot)$ has *several roots*.

Usually, proving almost sure convergence for such processes, one needs a Lyapunov function [9]. But it is rather an exception than a rule, when the function is known. In particular, one hardly can expect existence of a Lyapunov function for evolutionary games [3,6]. On the

other hand, in [6] convergence of a two-dimensional process like (1) was proved for a class of evolutionary games without using a Lyapunov function. The argument in [6] exploits the fact that asymptotically (1) behaves like the following system of ordinary differential equations

$$\frac{d}{dt}\vec{x} = \vec{g}(\vec{x}). \tag{3}$$

The main problem is to show that (3) does not have cycles and phase polygons. This, due to the theorem of Bendixon [4, p. 66], implies convergence of trajectories of (3) to its stationary points. The attractors in this case are either linearly stable or linearly unstable. An argument suggested in [6] allows to prove almost sure convergence of (1) to its (linearly) stable attractors. To prove that there are no cycles and phase polygons Dulac's criterion [4, p. 66] is used in [6].

Here we *prove convergence with probability 1 of a general stochastic approximation procedure without using a Lyapunov function*. We assume that the *system (3) converges to one of its stationary points*. In a sense the approach proposed here can be thought as an extension of the one suggested by L. Ljung [7]. Studying procedures like (1) with a single root, he assumed that the process (1) belongs to the domain of attraction of (3) to this root.

Now we shall introduce further hypotheses about the procedure (1). We assume that:

A1. There is a compact set $K \subset \mathbf{R}^m$ such that every solution of (3) originating at time 0 at a point from $K$ remains in $K$ for $t > 0$ (existence and uniqueness of the solution follow from the Lipschitz continuity of $\vec{g}(\cdot)$ that we required earlier);

A2. The equation $\vec{g}(\vec{x}) = \vec{0}$ has a finite number of roots $\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l, \vec{\eta}^1, \vec{\eta}^2, \ldots, \vec{\eta}^s$ in $K$; at each of these points $\vec{g}(\cdot)$ is differentiable and the Jakobian $J(\cdot)$ has eigenvalues with negative real parts at $\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l$ (i.e. these points are *linearly stable* for (3)) and at least one of its eigenvalues has a positive real part at each of $\vec{\eta}^1, \vec{\eta}^2, \ldots, \vec{\eta}^s$ (thus these points are *linearly unstable* for (3));

A3. There is a non-empty set of initial approximations $D$ such that for every $\vec{X}^1 \in D$ and for every $\tau > 0$ one can find a time instant $t(\vec{X}^1, \tau)$ for which

$$P\{\vec{X}^t \in K, \ t \geq t(\vec{X}^1, \tau)\} \geq 1 - \tau;$$

A4. $P\{\vec{X}^t \to \vec{\eta}^i\} = 0$ for every $i = 1, 2, \ldots, s$;

A5. Every trajectory of (3) belonging to $K$ converges to one of its stationary points (i.e. $\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l$ or $\vec{\eta}^1, \vec{\eta}^2, \ldots, \vec{\eta}^s$) as $t \to \infty$.

We shall prove that under assumptions A1 – A5, for every $\vec{X}^1 \in D$ the successive approximation $\vec{X}^t$ converges with probability 1 as $t \to \infty$ to a random vector $\vec{X}^0$ whose distribution concentrates on the set $\{\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l\}$[1].

---

[1]Conditions when $P\{\vec{X}^0 = \vec{\theta}^i\} > 0$ for $i = 1, 2, \ldots, l$ are known both for general stochastic approximation procedures [5] and for generalized urn processes [1].

Let us discuss conditions A1 – A5. A1 holds when $K$ is *invariant* for the system (3). For $2 \times 2$ evolutionary games considered in [3,6] it holds and $K = [0,1] \times [0,1]$. The same is true with $K = T_m$ for those of generalized urn process [1] (with balls of $m + 1$ colors) where the dynamic is given by Lipschitz functions. Here $T_m = \{\vec{x} \in \mathbf{R}^m : x_i \geq 0, \; x_1 + x_2 + \ldots + x_m \leq 1\}$. A2 is quite natural from the point of view of stability theory (so-called "the stability in the first approximation"). The assumption concerning unstable roots is coupled with A4, since only for linearly unstable roots one can apply the results on non-attainability [1,8,9,10]. What we really need is that stable roots are isolated and that all others (which might not necessarily be singleton) are attained with zero probability. Assumption A3 holds for the evolutionary games mentioned above and for generalized urn processes (since all trajectories belong correspondingly to $[0,1] \times [0,1]$ or $T_m$ with certainty). For a general stochastic approximation procedure one can use a projection (truncation) mechanism [7] or some global (i.e. for $\mathbf{R}^m \setminus K$) criterion of strong convergence to $K$ based on a Lyapunov function [9]. The most fundamental question concerns convergence if (3) to the stationary points, i.e. A5. In the case of $\mathbf{R}^2$, the Bendixon theorem [4, p. 66] implies that, if one can exclude cycles and phase polygons, trajectories of (3) converge to stationary points. The simplest result for excluding cycles and phase polygons is Dulac's criterion [4, p. 66]:

$$\frac{\partial}{\partial x_1}[h(\vec{x})g_1(\vec{x})] + \frac{\partial}{\partial x_2}[h(\vec{x})g_2(\vec{x})]$$

preserves its sign in $K$, where $h(\cdot)$ stands for a continuously differentiable in $K$ scalar function. In requires continuous differentiability of $\vec{g}(\cdot)$ in $K$. For higher dimensions we do not know anything as universal as the Bendixon theorem and Dulac's criterion. But for a particular system one can possibly suggest a specific criterion for ensuring convergence of trajectories to stationary points.

Now we shall prove convergence with probability 1 of stochastic approximation procedure (1).

## 2   Strong convergence of stochastic approximation procedures

The main result of this paper is given by the following theorem.

**Theorem 1** *Assume that conditions A1 – A5 hold true for a stochastic approximation process (1) with a Lipschitz regression function $\vec{g}(\cdot)$. Then $\vec{X}^t$ converges with probability 1 for every $\vec{X}^1 \in D$ to a random vector $\vec{X}^0$ whose distribution is concentrated on the set $\{\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l\}$ of stable roots of $\vec{g}(\cdot)$ in $K$.*

*Proof.* Since the random noise has uniformly bounded variance, the martingale convergence theorem implies that

$$\sum_{t \geq 1} \gamma_t \vec{Z}^t(\vec{X}^t)$$

exists with probability 1. Designate by $\Omega_0$ the joint event that it exists and $\vec{X}^t$ does not converge to $\vec{\eta}^i$, $i = 1, 2, \ldots, s$. Owing to A4, $P\{\Omega_0\} = 1$. Fix $\tau > 0$ and $\vec{X}^1 \in D$. Set $\Omega_\tau = \{\vec{X}^t \in K, \, t \geq t(\vec{X}^1, \tau)\} \bigcap \Omega_0$. Due to A3, one has

$$P\{\Omega_\tau\} \geq 1 - \tau. \tag{4}$$

Fix an elementary outcome $\omega$ from $\Omega_\tau$. Then the stochastic sequence (1) converts to the following deterministic sequence

$$\vec{x}^{t+1} = \vec{x}^t + \gamma_t \vec{g}(\vec{x}^t) + \vec{\sigma}^t, \, t \geq 1, \, \vec{x}^1 = \vec{X}^1, \tag{5}$$

where $\vec{x}^t$ and $\vec{\sigma}^t$ stand for the realizations of $\vec{X}^t$ and $\gamma_t[\vec{W}^t(\vec{X}^t) + \vec{Z}^t(\vec{X}^t)]$. Since $\omega \in \Omega_0$, the above infinite series containing $\vec{Z}^t(\vec{X}^t)$ converges. Hence, due to (2), one has

$$\lim_{t \to \infty} \|\sum_{i=t}^{\infty} \vec{\sigma}^i\| = 0. \tag{6}$$

What has to be shown that $\{\vec{x}^t\}$ converges to one of $\vec{\theta}^i$, $i = 1, 2, \ldots, l$. Since $\tau$ in (4) can be arbitrary small, this implies that with probability 1 the limit of $\vec{X}^t$ exists and its support coincides with $\{\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l\}$.

Assume to the contrary that there exists a subsequence of $\{\vec{x}^t\}$ converging to a point different from $\vec{\theta}^i$, $i = 1, 2, \ldots, l$. We shall show that this assumption leads to a contradiction.

Suppose there is a subsequence $\{n_p\}$ such that $\vec{x}^{n_p} \to \vec{z}$ as $p \to \infty$ and $\vec{z} \notin \{\vec{\theta}^1, \vec{\theta}^2, \ldots, \vec{\theta}^l\}$.

For all positive integers $n$ and all real $t \geq 0$ set

$$\vec{x}^n(t) = \vec{x}^i \quad \text{where} \quad \sum_{j=n}^{i} \gamma_t \leq t < \sum_{j=n}^{i+1} \gamma_t.$$

Let $\vec{x}(\vec{z}, \cdot)$ stand for the solution of the system of ordinary differential equations (3) satisfying the initial condition $\vec{x}(\vec{z}, 0) = \vec{z}$. Using that $\vec{g}(\cdot)$ is a Lipschitz function, one can show (see, for example, [2], p.p. 230-231) that for every finite $T > 0$

$$\lim_{p \to \infty} \sup_{t \in [0,T]} \|\vec{x}^{n_p}(t) - \vec{x}(\vec{z}, t)\| = 0. \tag{7}$$

We assumed that trajectories of (3) converge to its stationary points. Due to A4 and construction of the event $\Omega_\tau$, we can exclude the unstable points from this set. Without loss of generality let us consider the case when $\lim_{t \to \infty} \vec{x}(\vec{z}, t) = \vec{\theta}^1$.

We conclude that there is a subsequence $\{m_p\}$ such that $n_p < m_p < n_{p+1}$ and $\vec{x}^{m_p} \to \vec{\theta}^1$ as $p \to \infty$.

4

Due to condition A2 $\max(Re\lambda_1(\vec{\theta^1}), Re\lambda_2(\vec{\theta^1})) = \lambda(\vec{\theta^1}) < 0$, where $\lambda_i(\vec{\theta^1})$ stand for eigenvalues of $J(\vec{\theta^1})$. For every $\lambda \in (0, -\lambda(\vec{\theta^1}))$ a lemma of Lyapunov [9, p.173] guarantees existence of a symmetric positive definite matrix $C_\lambda$ such that

$$\langle C_\lambda J(\vec{\theta^1})\vec{x}, \vec{x} \rangle \leq -\lambda \langle C_\lambda \vec{x}, \vec{x} \rangle, \tag{8}$$

where $\langle \cdot, \cdot \rangle$ stands for the Euclidean scalar product in $\mathbf{R^m}$ and $\vec{x}$ is an arbitrary vector. Introduce a new scalar product $\langle \cdot, \cdot \rangle_C = \langle C_\lambda \cdot, \cdot \rangle$. Designate by $||\cdot||_C$ the corresponding norm. (This norm is equivalent to $||\cdot||$). We shall use it from now on. Owing to inequality (8), there exists $\epsilon^0$ such that $||\vec{x} - \vec{\theta^1}||_C \leq \epsilon^0$ implies

$$\langle \vec{g}(\vec{x}), \vec{x} - \vec{\theta^1} \rangle_C \leq -\frac{1}{2}\lambda ||\vec{x} - \vec{\theta^1}||_C^2. \tag{9}$$

Fix $\epsilon > 0$ such that $\epsilon < \min(\epsilon^0, ||\vec{z} - \vec{\theta^1}||_C)$. There is a subsequence $\{l_p\}$ such that $l_p = \max n > m_p : n < n_{p+1}$ and $||\vec{x}^n - \vec{\theta^1}||_C \leq \epsilon$. Selecting a subsequence of $\{l_p\}$ if necessary, we can assume that $\vec{x}^{l_p} \to \vec{x}(\epsilon)$ as $p \to \infty$, where $||\vec{x}(\epsilon) - \vec{\theta^1}||_C = \epsilon$. Fix small enough $\epsilon' > 0$ and define a subsequence $\{j_p\}$ such that $j_p = \max n > l_p : n < n_{p+1}$ and $||\vec{x}^n - \vec{x}^{l_p}||_C \leq \epsilon'$. Then from (5)

$$\vec{x}^{j_p} - \vec{x}^{l_p} = \sum_{i=l_p}^{j_p-1} \gamma_i \vec{g}(\vec{x}^i) + \vec{\sum}^p, \tag{10}$$

where

$$\vec{\sum}^p = \sum_{i=l_p}^{j_p-1} \vec{\sigma}^i.$$

Using the Lipschitz property of $\vec{g}(\cdot)$, one obtains from (10)

$$||\vec{x}^{j_p} - \vec{x}^{l_p}||_C \leq [||\vec{g}(\vec{x}^{l_p})||_C + L\epsilon']T^p + ||\vec{\sum}^p||_C$$

and

$$||\vec{x}^{j_p} - \vec{x}^{l_p}||_C \geq [||\vec{g}(\vec{x}^{l_p})||_C - L\epsilon']T^p + ||\vec{\sum}^p||_C,$$

where $L$ stands for the Lipschitz constant and

$$T^p = \sum_{i=l_p}^{j_p-1} \gamma_i.$$

Due to equivalence of the norms, from (6) it follows that, for small enough $\epsilon'$ and all sufficiently large $p$, there are positive constants $c_1$ and $c_2$ such that

$$c_1\epsilon' \leq T^p \leq c_2\epsilon'. \tag{11}$$

Using (6), (9) - (11), we obtain, for large enough $p$,

$$||\vec{x}^{j_p} - \vec{\theta^1}||_C^2 = ||\vec{x}^{l_p} - \vec{\theta^1}||_C^2 + 2T^p \langle \vec{x}^{l_p} - \vec{\theta^1}, \vec{g}(\vec{x}^{l_p}) \rangle_C +$$

5

$$2\langle \vec{x}^{l_p} - \vec{\theta^1}, \sum_{i=l_p}^{j_p-1} \gamma_i[\vec{g}(\vec{x}^i) - \vec{g}(\vec{x}^{l_p})]\rangle_C + 2\langle \vec{x}^{l_p} - \vec{\theta^1}, \vec{\sum}^p\rangle_C +$$

$$2\langle \sum_{i=l_p}^{j_p-1} \gamma_i\vec{g}(\vec{x}^i), \vec{\sum}^p\rangle_C + \|\sum_{i=l_p}^{j_p-1} \gamma_i\vec{g}(\vec{x}^i)\|_C^2 + \|\vec{\sum}^p\|_C^2 \le$$

$$\|\vec{x}^{l_p} - \vec{\theta^1}\|_C^2[1 - c_1\lambda\epsilon' + L^2c_2^2(\epsilon')^2] + 2\|\vec{x}^{l_p} - \vec{\theta^1}\|_C Lc_2(\epsilon')^2 + o_p(1),$$

where $o_p(1) \to 0$ as $p \to \infty$. Passing to the limit as $p \to \infty$ we conclude that

$$\limsup_{p\to\infty} \|\vec{x}^{j_p} - \vec{\theta}\|_C^2 \le \epsilon^2[1 - c_1\lambda\epsilon' + L^2c_2^2(\epsilon')^2] + 2\epsilon Lc_2(\epsilon')^2.$$

Consequently, if $\epsilon'$ is so small that $\epsilon^2[1 - c_1\lambda\epsilon' + L^2c_2^2(\epsilon')^2] + 2\epsilon Lc_2(\epsilon')^2 < \epsilon^2$, then

$$\limsup_{p\to\infty} \|\vec{x}^{j_p} - \vec{\theta^1}\|_C < \epsilon.$$

However, $j_p > l_p$ and $l_p$ is the last time instant before $n_{p+1}$ when the sequence is inside the $\epsilon$-neighborhood of $\vec{\theta^1}$. Hence $\vec{x}^{j_p}$ must lie outside the $\epsilon$-neighborhood of $\vec{\theta^1}$, that is,

$$\liminf_{p\to\infty} \|\vec{x}^{j_p} - \vec{\theta^1}\|_C \ge \epsilon.$$

This contradiction shows that there is no subsequence of $\{\vec{x}^t\}$ converging to a limit different from $\vec{\theta^i}$, $i = 1, 2, \ldots, l$.

The theorem is proved.

# 3 Conclusions

Possible development of the approach given here lies in two directions. One consists in developing particular conditions ensuring convergence of trajectories of ordinary differential equations to its stationary points. As we mentioned before, for two-dimensional systems this reduces to finding conditions which cancel cycles and phase polygons. Another consists in considering non-Markov procedures. We did not exploit in our argument explicitly that the noise is formed by independent random vectors. The only place where we used this assumption implicitly is A4. This is since we do not know any unattainability conditions for non-Markov procedures.

# References

[1] Arthur, W. B., Y. M. Ermoliev, and Y. M. Kaniovski (1987). Adaptive Growth Process Modeled by Urn Schemes, Kibernetika, No. 6, 49–57 (in Russian). (Translated into English in *Cybernetics*, **23**, 779–789.)

[2] Benveniste, A., M. Métivier, and P. Priouret (1990): *Adaptive Algorithms and Stochastic Approximations.* Springer–Verlag, New York.

[3] Dosi, G., and Y. Kaniovski (1994). On "Badly Behaved" Dynamics *Some Applications of Generalized Urn Schemes to Technological and Economic Change.* Journal of Evolutionary Economics, 4, pp. 93–123.

[4] Hahn, W. (1967). *Stability of Motion.* Springer–Verlag, New York.

[5] Kaniovski, Yu. M. (1988). Limit Theorems for Processes of Stochastic Approximation when the Regression Function Has Several Roots. Kibernetika, No. 2, pp. 136–138 (in Russian).

[6] Kaniovski, Yu. M., and H. P. Young (1994). *Learning Dynamics in Games with Stochastic Perturbations.* Working paper WP-94-30, International Institute for Applied Systems Analysis.

[7] Ljung, L. and T. Söderström (1983). *Theory and Practice of Recursive Identification.* M.I.T. Press, Cambridge, MA.

[8] Ljung, L. (1978). Strong Convergence of a Stochastic Approximation Algorithm. Ann. Statist., **6**, pp. 680–696.

[9] Nevel'son, M. B. and R. Z. Has'minski (1972). *Stochastic Approximation and Recurrent Estimation*, Nauka, Moscow (in Russian). (Translated into English: Nevelson, M. B., and R. Z. Hasminskii (1976). *Stochastic Approximation and Recursive Estimation.* 47 American Math. Soc., Providence, RI.)

[10] Pemantle, R. (1990). Nonconvergence to Unstable Points in Urn Models and Stochastic Approximations. Ann. Prob. **18**, pp. 698–712.