# ASYNCHRONOUS STOCHASTIC APPROXIMATIONS[*]

## VIVEK S. BORKAR[†]

**Abstract.** The asymptotic behavior of a distributed, asynchronous stochastic approximation scheme is analyzed in terms of a limiting nonautonomous differential equation. The relation between the latter and the relative values of suitably rescaled relative frequencies of updates of different components is underscored.

**Key words.** distributed algorithms, asynchronous algorithms, communication delays, stochastic approximation, ODE limit

**AMS subject classifications.** 62L20, 93E25

**PII.** S0363012995282784

**1. Introduction.** There has been a resurgence of interest in stochastic approximation algorithms, particularly as mechanisms for learning systems. They can, for example, be a learning algorithm for neural networks [13] or a model of learning by boundedly rational agents in a macroeconomic system [20], in addition to their traditional applications in adaptive engineering systems [2]. These applications call for a distributed, asynchronous implementation of stochastic approximation schemes. In engineering applications, this is a natural consequence of dealing with large interconnected systems. In macroeconomics, it is simply the reality of life. It is not, however, apparent that the traditional analysis of these schemes, extensively dealt with in [2], automatically holds ground in the new scenario. Prompted by these and similar concerns, there have been studies of distributed implementations of these algorithms [17, 18, 21, 22]. (See [3] for an extensive account of parallel distributed algorithms in general). The present work is in the same spirit, but with some crucial differences.

1. Our model of asynchronism postulates a set-valued random process that marks the indices to be updated at each iteration. This clumping of indices into sets can be an artifice as long as causal relationships are not violated; thus the set-up is very general indeed. We impose on this process a natural condition that requires all indices to be updated comparably often in a precise sense.

2. In addition, we allow random, possibly nonstationary and unbounded delays that are required to satisfy a mild conditional moment condition.

3. The analysis depends on proving that the algorithm asymptotically tracks a nonautonomous ODE, in contrast to the traditional autonomous "ODE limit." In particular, it gives a handle on situations when the latter may not be feasible.

4. The ODE in question differs from the traditional one in that each component of the driving vector field is now weighted by a time-varying nonnegative scalar. These scalars add to 1 and may be interpreted as relative frequencies of updates of different components after suitable time-scaling. This clearly brings out the desired relationships between update schemes for different components and paves the way for analyzing situations where they are not desirable (see remark 4 of the conclusion).

---

[†]Department of Computer Science and Automation, Indian Institute of Science, Bangalore 560012, India (borkar@csa.iisc.ernet.in).

The paper is organized as follows. The remainder of this section describes the problem framework. The next section states the key assumptions and their immediate consequences. The third section provides the convergence analysis. The final section highlights some further possibilities.

Let TS (for tapering stepsize) denote the set of sequences $\{a(n)\}$ in (0,1) satisfying

$$(1.1) \qquad \sum_n a(n) = \infty, \qquad \sum_n a(n)^2 < \infty.$$

The standard stochastic approximation algorithm is the recursive scheme in $R^d$, $d \geq 1$, described by

$$(1.2) \qquad X(n+1) = X(n) + a(n)F(X(n), \xi(n)),$$

where $\{a(n)\} \in$ TS, $X(n) = [X_1(n), \ldots, X_d(n)]^T \in R^d$ with a prescribed $X(0)$, $F(\cdot, \cdot) : R^d \times R^k \to R^d$, and $\{\xi(n)\}$ is a stationary random process in $R^k$. For simplicity, we take $\{\xi(n)\}$ to be independently and identically distributed (i.i.d.) with a common law $\psi$ (say). The $i$th row of this vector iteration reads

$$(1.3) \qquad X_i(n+1) = X_i(n) + a(n)F_i(X_1(n), \ldots, X_d(n), \xi(n)).$$

A distributed but synchronous version of (1.2) could be as follows. Let $I = \{1, 2, \ldots, d\}$ and $S$ be a collection of nonempty subsets of $I$ that cover $I$. Let $\{Y_n\}$ be an $S$-valued random process that selects the coordinates to be updated at time $n$, and for each $n$, let $\tau_{ij}(n)$, $i \neq j \in I$, be random variables taking values in $\{0, 1, \ldots, n\}$ that represent communication delays. We set $\tau_{ii}(n) = 0 \ \forall i, n$. The synchronous distributed version of (1.3) is then

$$(1.4) \quad X_i(n+1) = X_i(n) + a(n)F_i(X_1(n - \tau_{i1}(n)), \ldots, X_d(n - \tau_{id}(n)), \xi(n))I\{i \in Y_n\}$$

for $i \in I$, $n \geq 0$. Special instances of (1.4) were studied in [5, 6]. The reason this is a synchronous version is that the decision to use stepsize $a(n)$ at time $n$ by the processor updating the $i$th component (say) presupposes the availability of a global clock to all processors. This is not reasonable in an asynchronous environment. The asynchronous version we propose is as follows. Let $\{a(n, i)\} \in$ TS, $i \in I$, and define

$$\nu(n, i) = \sum_{m=0}^{n} I\{i \in Y_m\}, \quad i \in I, \quad n \geq 0,$$

$$\bar{a}(n, i) = a(\nu(n, i), i), \quad i \in I, \quad n \geq 0.$$

The first of these is the total number of times the $i$th component was updated up until time $n$. Assuming that each component of the iteration is assigned to one and only one processor once and for all, $\bar{a}(n, i)$ is a random variable known to the $i$th processor at time $n$. The proposed algorithm is

$$(1.5) \quad X_i(n+1) = X_i(n) + \bar{a}(n, i)F_i(X_1(n - \tau_{i1}(n)), \ldots, X_d(n - \tau_{id}(n)), \xi(n))I\{i \in Y_n\}$$

for $i \in I$, $n \geq 0$. This is the algorithm analyzed in this paper, under the assumptions stipulated in the next section. We conclude this section with the remark that even the implicit presence of an unobserved global clock in the background in (1.5) is not really needed. The clumping of updated coordinates into $Y_m$'s could be a complete artifice as long as causal relationships are not violated and the additional assumptions of the next section (notably (A3)) remain valid.

**2. Preliminaries.** The additional assumptions and their consequences that we present in this section concern, respectively, the stepsize routines $\{a(n, i)\}$, the sampling process $\{Y_n\}$, the communication delays $\{\tau_{ij}(n)\}$, and the function $F$. We proceed in that order. These assumptions, (A1)–(A5), are enforced throughout the paper without further mention.

Let ITS (for "ideal tapering stepsize") denote the subset of TS consisting of $\{a(n)\}$ satisfying:

(i) $a(n + 1) \leq a(n)$ from some $n$ onwards;

(ii) there exists $r \in (0, 1)$ such that

$$(2.1) \qquad \sum_n a(n)^{1+q} < \infty, \qquad q \geq r;$$

(iii) for $x \in (0, 1)$,

$$(2.2) \qquad \sup_n a([xn])/a(n) < \infty,$$

where $[\cdots]$ stands for the integer part of "$\cdots$";

(iv) for $x \in (0, 1)$ and $A(n) \triangleq \sum_{i=0}^{n} a(i)$,

$$(2.3) \qquad A([yn])/A(n) \to 1$$

uniformly in $y \in [x, 1]$.

By (i), (2.2) may be strengthened to

$$(2.4) \qquad \sup_n \sup_{y \in [x,1]} a([yn])/a(n) < \infty.$$

It is easy to construct examples of $\{a(n)\}$ in TS which violate (2.2). Condition (iv) can be given an alternative formulation. Let $h : R^+ \to R^+$ be an eventually nonincreasing function satisfying $h(n) = a(n)$, $n \geq 0$. Then (2.3) is equivalent to

$$(2.5) \qquad \lim_{t \to \infty} \frac{\int_0^{yt} h(s)ds}{\int_0^t h(s)ds} = 1,$$

which, by l'Hôpital's rule, reduces to

$$\lim_{t \to \infty} yh(yt)/h(t) = 1.$$

One needs this to hold uniformly in $y \in [x, 1]$. One sufficient condition for this would be that the derivative of the left-hand side of (2.5) in $y$, which is $th(yt)/\int_0^t h(s)ds$, be bounded uniformly in $y$, $t$, ensuring the equicontinuity in $y$ for the ratio in (2.5). It is not clear whether (iv) is implied by (i)–(iii). Examples of $\{a(n)\}$ satisfying (i)–(iv) are $\{1/n\}$, $\{1/n \log n\}$, and $\{\log n/n\}$, with suitable modification for $n = 0, 1$ where needed.

One property of $\{a(n)\} \in$ TS that we shall need later is the following.

LEMMA 2.1. *For $s \in (0, 1)$, $a(n)^{-s}/n \to 0$.*

*Proof.* It suffices to prove that $(a(n)n^x)^{-1} \to 0$ for $x = 1/s > 1$, or equivalently, that $(a(n) + n^{-x})/a(n) \to 1$. Let $h_1, h_2 : R^+ \to R^+$ be continuous functions linearly interpolated from $h_1(n) = a(n) + n^{-x}$, $h_2(n) = a(n)$, $n \geq 0$. Since $\int_0^t h_1(y)dy \to \infty$ as $t \to \infty$ and $\int_1^\infty t^{-x}dt < \infty$, we have

$$\lim_{t \to \infty} \frac{\int_0^t h_2(y)dy}{\int_0^t h_1(y)dy} = 1.$$

The claim now follows from l'Hôpital's rule.    □

Our first assumption then is:

(A1) $\{a(n,i)\} \in$ ITS for $i \in I$.

Next, introduce for $n \geq 0$ the $\sigma$-fields $\mathcal{F}_n = \sigma(X(m), Y(m), m \leq n, \tau_{ij}(m), \xi(m), m < n, i, j \in I)$ and $\mathcal{G}_n = \sigma(X(m), Y(m), \tau_{ij}(m), \xi(m), m \leq n, i, j \in I)$. Our assumption concerning $\{Y_n\}$ is as follows.

(A2) There exists a $\delta > 0$ such that for any $A, B \in S$, the quantity

$$(2.6) \qquad\qquad P(Y_{n+1} = B / Y_n = A, \mathcal{G}_n)$$

is either always zero almost surely (a.s.) or always exceeds $\delta$ a.s. That is, having picked $A$ at time $n$, picking $B$ at time $n+1$ is either improbable or probable with a conditional probability of at least $\delta$, regardless of $n$ and the "history" $\mathcal{G}_n$. Furthermore, if we draw a directed graph with node set $S$ and an edge from $A$ to $B$ whenever (2.6) exceeds $\delta$ a.s., the graph is irreducible; i.e., there is a directed path from any $A \in S$ to any $B \in S$. (As will become apparent later, this may be replaced by the weaker requirement that every communicating class of the directed graph comprises sets that together cover $I$.)

This has the following important consequence. Let $\mathcal{P}(\cdots)$ denote the space of probability vectors on "$\cdots$."

LEMMA 2.2. *There exists a deterministic constant $\Delta > 0$ such that for any $A \in S$,*

$$(2.7) \qquad\qquad \liminf_{n\to\infty} \frac{1}{n} \sum_{m=0}^{n-1} I\{Y_m = A\} \geq \Delta \quad a.s.$$

*Proof.* For $A \in S$, let $D_A = \{B \in S | (2.6) \text{ exceeds } \delta \text{ a.s.}\}$ and $V_A = \{u \in \mathcal{P}(D_A) | u(B) \geq \delta \ \forall B \in D_A\}$, $V = \prod_A V_A$. Define $p : S \times S \times V \to [0,1]$ by $p(A, B, u) = u_A(B)$, where $u_A$ is the $A$th component of $u$. Define $V$-valued random variables $\{Z^n\}$ by

$$Z_A^n(B) = P(Y_{n+1} = B / \mathcal{G}_n)I\{Y_n = A\} + \psi_A I\{Y_n \neq A\},$$

where $\psi_A$ is a fixed element of $V_A$ for $A \in S$. Then (2.6) equals $p(A, B, Z^n)$ and $\{Y_n\}$ may be viewed as an $S$-valued controlled Markov chain with action space $V$ and transition probability function $p$. (This is a pure artifice for the sake of the proof. It is in no way implied that $\{Z^n\}$ is an actual control process.) In particular, this allows us to conceive of a stationary policy $\pi$ associated with a map $\pi \colon S \to V$. (A1) implies, in particular, that $\{Y_n\}$ will be an ergodic Markov chain under a stationary policy $\pi$ with a corresponding stationary distribution $\nu_\pi \in \mathcal{P}(S)$. Then the left-hand side of (2.7) a.s. exceeds $\min_\pi \nu_\pi(A) > 0$ by Lemmas 1.2 and 2.1 of [4, pp. 56, 60]. $\quad\square$

For the communication delays, we assume the following. (Recall the $r$ in (2.1).)

(A3) $\tau_{ij}(n) \in \{0, 1, \ldots, n\}, \tau_{ii}(n) = 0 \ \forall i, a$. There exist $b > r/(1-r), C > 0$ such that

$$(2.8) \qquad\qquad E[(\tau_{ij}(n))^b / \mathcal{F}_n] \leq C \quad \text{a.s.} \quad \forall i, j, n.$$

(In particular, we do not require the delays to be either bounded or stationary.) Also, $\{\xi(n)\}$ is i.i.d. and independent of $\{X_0, \xi(m), \tau_{ij}(m), m < n\}$ for all $n$.

Next come the conditions on $F$.

(A4) $F$ is assumed to be measurable and uniformly Lipschitz in the first argument; i.e., for some $K > 0$,

$$\|F(x, z) - F(y, z)\| \leq K\|x - y\| \quad \forall \ x, y, z.$$

Other conditions on $F$ will be given in terms of the function $f : R^d \to R^d$ defined by

$$(2.9) \qquad\qquad f(x) = \int F(x,y)\psi(dy).$$

Under our conditions on $F$, $f$ is Lipschitz with Lipschitz constant $K$. The traditional analysis of (1.2) [2] proceeds by showing that it asymptotically tracks the ODE

$$(2.10) \qquad\qquad \dot{x}(t) = f(x(t)),$$

which in turn has trajectories converging to $J = \{x | f(x) = 0\}$.

(A5) $J$ is assumed to be compact and nonempty.

We shall also have reason to consider a related nonautonomous ODE. Let $\mathcal{D}$ denote the set of diagonal $d \times d$ matrices with nonnegative diagonal entries that add to 1. For $a > 0$, say that $M = \text{diag}(m_1, \ldots, m_d)$ is $a$-thick if $m_i \geq a \; \forall i$. The ODE in question is

$$(2.11) \qquad\qquad \dot{x}(t) = M(t)f(x(t)),$$

where $t \to M(t)$ is a $\mathcal{D}$-valued measurable process.

We consider two scenarios.

*Case* 1: *Strict Liapunov systems.* A continuously differentiable function $V : R^d \to R^+$ is said to be a strict Liapunov function for (2.10) if $\nabla V . f < 0$ outside $J$. Call (2.10) a strict Liapunov system if it has bounded trajectories and a strict Liapunov function $V$ exists. The latter implies the former if $V(x) \to \infty$ as $\|x\| \to \infty$, which we assume to hold. (Call this assumption (A6).) Examples of such systems can be found among gradient systems and their variants, certain systems arising in neural networks [14], and analog fixed point algorithms wherein $f(x) = g(x) - x$ and $g$ is either a contraction under a $\|.\|_p$-norm for $p \in [1, \infty]$ or nonexpansive under a $\|.\|_p$- norm for $p \in (1, \infty)$. (Here $V(.) = \|. - x^*\|_p$, where $x^* \in J$, will do. For $p = \infty$, this is not continuously differentiable, but this does not pose any problems for contractions [7].)

Finally, a strict Liapunov system as above will be said to be $a$-robust for some $a > 0$ if $\nabla V . Mf < 0$ outside $J$ for any $a$-thick $M \in \mathcal{D}$.

Given $T, \delta > 0$, a $(T, \delta)$-perturbation of (2.10) (resp., (2.11)) is a function $y : R^+ \to R^d$ such that there exist $0 = T_0 < T_1 < \cdots < T_n \uparrow \infty$ and solutions $x^j(t)$, $t \in [T_j, T_{j+1}]$, $j \geq 0$, of (2.10) (resp., (2.11)) such that $T_{j+1} - T_j \geq T$ for $j \geq 0$ and

$$\|y(t) - x^j(t)\| < \delta, \quad T_j \leq t \leq T_{j+1}, j \geq 0.$$

For $\epsilon > 0$, let $J^\epsilon = \{x \in R^d | \|x - y\| < \epsilon \text{ for some } y \in J\}$.

LEMMA 2.3. *Under* (A6), *we have:* (a) *For any* $T, \epsilon > 0$, *there exists a* $\delta_0 = \delta_0(T, \epsilon) > 0$ *such that for* $\delta \in (0, \delta_0)$, *any* $(T, \delta)$-*perturbation of* (2.10) *converges to* $J^\epsilon$. *(In particular, solutions of* (2.10) *converge to* $J$.)

(b) *Suppose that* (2.10) *is* $a$-*robust for some* $a > 0$ *and* $M(t)$ *in* (2.11) *is* $a$-*thick for almost every* $t$. *Then, for any* $T, \epsilon > 0$, *there exists a* $\delta_0 = \delta_0(T, \epsilon, a) > 0$ *such that for* $\delta \in (0, \delta_0)$, *any* $(T, \delta)$-*perturbation of* (2.11) *converges to* $J^\epsilon$. *(In particular, the solutions of* (2.11) *converge to* $J$.)

These are straightforward adaptations of Theorem 1 of [14, p. 339].

*Case* 2: $\infty$-*nonexpansive maps.* In this case $f(x) = g(x) - x$, where $g$ is $\infty$-nonexpansive, i.e., $\|g(x) - g(y)\|_\infty \leq \|x - y\|_\infty, x, y \in R^d$. Thus $J$ is the set of fixed points of $g$. This case is important in dynamic programming applications [3, 7].

For this case, we have the following analog of Lemma 2.3.

LEMMA 2.4. *The conclusions of Lemma* 2.3(a) *continue to hold. Those of Lemma* 2.3(b) *hold if* $M(t)$ *is* a-*thick for almost every* t, *for some* $a > 0$.

This is proved in Theorem 2.1 and Corollary 2.2 of [5].

**3. Convergence.** We start by establishing a link between (1.4) and (2.11). Our first observation is that we may equivalently consider the recursion

$$(3.1) \quad X_i(n+1) = X_i(n) + \bar{a}(n,i)F_i(X_i(n-\tau_{ij}(n)), \ldots, X_d(n-\tau_{id}(n)), \tilde{\xi}(n))I\{\varphi_n = i\},$$

where $\{\varphi_n\}$ is an $I$-valued random process satisfying the following statement. There exists a deterministic constant $\eta > 0$ such that

$$(3.2) \qquad\qquad \liminf_{n\to\infty} \frac{1}{n} \sum_{m=0}^{n-1} I\{\varphi_n = i\} \geq \eta \quad \text{a.s.} \quad \forall i \in I$$

and $\tilde{\xi}(n) = \xi(k(n))$ for a nondecreasing map $n \to k(n)$, satisfying $k(n+1) - k(n) \in \{0,1\}$.

This is achieved simply by unfolding each iteration as follows.

Let $Y_n = \{i_1, \ldots, i_{c(n)}\}$ (say) with the elements arranged in ascending order. Replace the iteration (1.4) by $c(n)$ distinct iterations such that the $j$th iteration among them updates only the $i_j$th component in accordance with (1.4). Next, relabel the iteration index and the delays to obtain a correspondence with (3.1). Then (3.2) is an immediate consequence of Lemma 2.2. Note that this blows up the delays at most $d$ fold, thus still retaining (A3). Note also that for $m > n$, $\varphi_m = \varphi_n$ implies $k(m) > k(n)$. With these considerations, we proceed to analyze (3.1). We start with some preliminaries.

Let $U$ be the space of $\mathcal{P}(I)$-valued trajectories $\bar{\mu} = \{\mu_t, t \geq 0\}$ with the coarsest topology that renders continuous the maps $\bar{\mu} \to \int_0^T h(t)\mu_t(i)dt$ for $T \geq 0$, $i \in I$, $h \in L_2[0,T]$. $U$ is compact metrizable. Say that $\mu \in \mathcal{P}(I)$ is $\alpha$-thick for some $\alpha > 0$ if $\mu(i) \geq \alpha \ \forall i$. Say that $\bar{\mu} \in U$ is $\alpha$-thick, $\alpha > 0$, if $\mu_t$ is so for almost every $t$. Say that $\mu$ (resp., $\bar{\mu}$) is thick if it is $\alpha$-thick for some $\alpha > 0$.

LEMMA 3.1. (a) *For* $\alpha > 0$, $\{\bar{\mu}|\bar{\mu}$ *is* $\alpha$-*thick*$\}$ *is compact in* $U$. (b) *The map* $(\bar{\mu}, x) \to x(.) : U \times R^d \to C([0,\infty); R^d)$ *defined via*

$$(3.3) \qquad\qquad \dot{x}(t) = M^{\bar{\mu}}(t)f(x(t)), \qquad x(0) = x,$$

*with* $M^{\bar{\mu}}(t) = \text{diag}(\mu_t(1), \ldots, \mu_t(d))$, *is continuous.*

*Proof.* (i) For $i \in I, t > s, n \geq 1$, and any $\alpha$-thick $\bar{\mu}$, $\alpha > 0$,

$$\int_s^t \mu_y(i)dy \geq \alpha(t-s).$$

The relation is preserved under limits in $U$, implying the claim.

(ii) Let $(\bar{\mu}^n, x_n) \to (\bar{\mu}^\infty, x_\infty)$. For $n \geq 1$, let $x^n(.)$ satisfy

$$(3.4) \qquad\qquad \dot{x}^n(t) = M^{\bar{\mu}^n}(t)f(x^n(t)), \qquad x^n(0) = x_n.$$

Using the Gronwall lemma and the Arzela–Ascoli theorem, one verifies that $\{x^n(.)\}$ is relatively compact in $C([0,\infty); R^d)$, and a straightforward limiting argument (keeping in mind our topology on $U$) shows that any limit $x^\infty(.)$ thereof must satisfy (3.4) with $n = \infty$. The claim follows.   □

Let $\tilde{a}(n) = \bar{a}(n, \varphi_n)$ and rewrite (3.1) as

$$X(n+1) = X(n) + \tilde{a}(n)W(n)$$

for appropriately defined $W(n) = [W_1(n), \dots, W_d(n)]^T$. Redefine $\mathcal{F}_n$ by $\mathcal{F}_n = \sigma(X(m), m \leq k^{-1}(n), \xi(m), m < k^{-1}(n), \tau_{ij}(m), m < n, \varphi_m, m \leq n)$, where $k^{-1}(n) = \min\{j | k(j) = n\}$. Set $\hat{W}(n) = E[W(n)/\mathcal{F}_n], n \geq 0$, the conditioning bring componentwise. Write $\hat{W}(n) = [\hat{W}_1(n), \dots, \hat{W}_d(n)]^T$. Define $f^i : R^d \to R^d$ by $f^i_j(x) = f_i(x)\delta_{ij}$, $i, j \in I$, $\delta_{ij}$ being the Kronecker delta. Let $b(n) = \max_i \bar{a}(n, i), n \geq 0$. Let $Q$ denote the set of sample points for which $\hat{X} \triangleq \sup_n \|X(n)\| < \infty$.

LEMMA 3.2. $\{b(n)\}$ *satisfies* $\sum_n b(n)^{1+r} < \infty$ *a.s., and for* $a \in (0, 1]$,

$$\sup_n \ \sup_{\alpha \in [a,1]} b([\alpha n])/b(n) < \infty \quad a.s.$$

*Proof.* By (2.4) and (3.2), $\sup_n \bar{a}(n, i)/a(n, i) < \infty$ a.s., $i \in I$. Combining this with property (2.1) for $\{a(n, i)\}$, we have $\sum \bar{a}(n, i)^{1+r} < \infty$ a.s. The first claim follows. The second follows easily from (2.4) applied to $\{a(n, i)\}$. $\quad\square$

LEMMA 3.3. *Almost surely on* $Q$, *there exist* $K_1 > 0$, $N \geq 1$ *(random) such that for* $n \geq N$,

$$\|f^{\varphi_n}(X(n)) - \hat{W}(n)\| < K_1 b(n)^r.$$

*Proof.* Consider $\omega \in Q$. Let $K_2$ be an upper bound on $\{\|f(x)\|_\infty \mid \|x\| \leq \hat{X}\}$. Let $\tilde{W}_i(n) = f^{\varphi_n}_i(X_1(n - \tau_{i1}(n)), \dots, X_d(n - \tau_{id}(n)))$. Let $c = 1 - r$. For $i \in I$, we have

(3.5)

$$|f^{\varphi_n}_i(X(n)) - \hat{W}_i(n)| \leq E[|f^{\varphi_n}_i(X(n)) - \tilde{W}_i(n)|I\{\tau_{ij}(n) \leq b(n)^{-c} \ \forall i, j\}/\mathcal{F}_n]$$
$$+ E[|f^{\varphi_n}_i(X(n)) - \tilde{W}_i(n)|I\{\tau_{ij}(n) > b(n)^{-c} \text{ for some } i, j\}/\mathcal{F}_n] \quad a.s.$$

By (A3) and the conditional Chebyshev inequality, the second term is a.s. bounded by $2K_2 C d^2 b(n)^{bc}$. Let $\bar{n} = [b(n)^{-c}]$. By Lemma 2.1, $\bar{n}$ is $o(1)$ a.s. as $n \to \infty$, and outside a zero probability set, we may pick $n$ large enough so that $n > \bar{n}$. Then for $m \leq \bar{n}$,

$$\|X(n) - X(n-m)\| \leq 2K_2 d \sum_{j=n-\bar{n}}^{n} b(j) \leq K_3 b(n)^{1-c}$$

for a suitable (random) $K_3 > 0$, by the above lemma. Thus the first term in (3.5) is bounded by $K_4 b(n)^r$ for a suitable (random) $K_4 > 0$. Since $b > r/(1 - r)$, the claim follows. $\quad\square$

Let $T > 0$. Define $t_0 = T_0 = 0, t_n = \sum_{m=0}^{n-1} \tilde{a}(m)$, $n \geq 1$, and $T_n = \min\{t_m | t_m \geq T_{n-1} + T\}$, $n \geq 1$. Then $T_n = t_{m(n)}$ for a strictly increasing sequence $\{m(n)\}$. Let $I_n = [T_n, T_{n+1}]$, $n \geq 0$. Define $\bar{x}^n(t)$, $t \in I_n$, by $\bar{x}^n(T_n) = X(m(n))$ and

$$\bar{x}^n(t_{m(n)+k+1}) = \bar{x}^n(t_{m(n)+k}) + \tilde{a}(m(n) + k)f^{\varphi_{m(n)+k}}(\bar{x}^n(t_{m(n)+k})),$$

with linear interpolation on each interval $[t_{m(n)+k}, t_{m(n)+k+1}]$. Define $x(t)$, $t \geq 0$, by $x(t_n) = X(n)$ with linear interpolation on each interval $[t_n, t_{n+1}]$.

LEMMA 3.4. $\lim_{n\to\infty} \sup_{t\in I_n} \|x(t) - \overline{x}^n(t)\| = 0$ a.s. on $Q$.

*Proof.* Let $n \geq 1$. For $i \geq m(n)$, we have

$$x(t_{i+1}) = x(t_i) + \tilde{a}(i)f^{\varphi_i}(x(t_i)) + \tilde{a}(i)(\hat{W}(i) - f^{\varphi_i}(x(t_i))) + \tilde{a}(i)(W(i) - \hat{W}(i)).$$

Let $\overline{M}_i = \sum_{j=0}^{i} \tilde{a}(j)(W(j) - \hat{W}(j))$ and $\lambda_i = \overline{M}_i - \overline{M}_{m(n)}, i \geq m(n)$. Also, let $M_i^k = \sum_{j=0}^{k} \tilde{a}(j)(W(j) - \hat{W}(j))I\{\varphi_j = k\}, 1 \leq k \leq d$. Recall that $m > n$ and $\varphi_m = \varphi_n$ implies $k(m) > k(n)$. Then, for each $k$, $\{M_i^k, \mathcal{F}_i\}$ is a zero mean–bounded increment vector martingale, and the quadratic variation process of each of its component martingales is a.s. convergent on $Q$. By Proposition VII-3-(c) of [19, pp. 149–150], each $\{M_i^k\}$ and hence $\{\overline{M}_i\}$ converges a.s. on $Q$. Fix a sample point for which this convergence holds and let $\epsilon > 0$. Then $\sup_{i\geq m(n)} \|\lambda_i\| < \epsilon/2$ for sufficiently large $n$. Let $\hat{x}_{i+1} = x(t_{i+1}) - \lambda_i, i \geq m(n)$, with $\hat{x}_{m(n)} = X(m(n))$. Then, for $i \geq m(n)$, we have

$$\hat{x}_{i+1} = \hat{x}_i + \tilde{a}(i)f^{\varphi_i}(\hat{x}_i) + \tilde{a}(i)(f^{\varphi_i}(\hat{x}_i + \lambda_{i-1}) - f^{\varphi_i}(\hat{x}_i)) + \tilde{a}(i)(\hat{W}(i) - f^{\varphi_i}(x(t_i))).$$

Also,

$$\overline{x}^n(t_{i+1}) = \overline{x}^n(t_i) + \tilde{a}(i)f^{\varphi_i}(\overline{x}^n(t_i)).$$

Fix $\omega \in Q$, where the foregoing and Lemma 3.3 hold. Subtracting and using Lemma 3.3, we have, for $n$ sufficiently large,

$$\|\hat{x}_{i+1} - \overline{x}^n(t_{i+1})\| \leq (1 + K\tilde{a}(i))\|\hat{x}_i - \overline{x}^n(t_i)\| + \tilde{a}(i)\|\lambda_{i-1}\|K + K_1\tilde{a}(i)b(i)^{1+r}.$$

By increasing $n$ if necessary, we may suppose that

$$\sum_{i\geq n} b(i)^{1+r} < \epsilon/2.$$

Then using the inequality $1 + x \leq \exp(x)$ and iterating, we have

$$\sup_{m(n)\leq i\leq m(n+1)} \|\hat{x}_i - \overline{x}^n(t_i)\| \leq e^{K(T+1)}(K_1 + K(T+1))\epsilon$$

for sufficiently large $n$. Since $\|\hat{x}_i - x(t_i)\| < \epsilon/2$, $i \geq m(n)$ for sufficiently large $n$, $\sup_{m(n)\leq i\leq m(n+1)} \|x(t_i) - \overline{x}^n(t_i)\| \leq \tilde{K}\epsilon$ for a suitable $\tilde{K} > 0$. Since $\epsilon > 0$ was arbitrary, the claim follows on noting that both $x(.)$ and $\overline{x}^n(.)$ are linearly interpolated from their values at $\{t_i\}$. $\square$

Next, define $\overline{\mu} \in U$ by $\mu_t = $ the Dirac measure at $\varphi_n$ for $t \in [t_n, t_{n+1}), n \geq 0$. Define $\tilde{x}^n(t), t \in I_n$, by $\tilde{x}^n(t_{m(n)}) = x(t_{m(n)})$ and

$$(3.6) \qquad \dot{\tilde{x}}^n(t) = M^{\overline{\mu}}(t)f(\tilde{x}^n(t)), \qquad t \in I_n.$$

LEMMA 3.5. $\lim_{n\to\infty} \sup_{t\in I_n} \|\tilde{x}^n(t) - \overline{x}^n(t)\| = 0$ a.s.

*Proof.* This follows easily from the Gronwall inequality. $\square$

For $\overline{\mu}$ as above, define $\overline{\mu}^t = \{\mu_{t+s}, i \geq 0\} \in U$ for $t \geq 0$. Combining the foregoing with Lemmas 2.3 and 2.4, we have the following theorem.

THEOREM 3.1. (a) *Suppose there exists an $a > 0$ such that (2.10) is an $a$-robust strict Liapunov system, (A6) applies, and all limit points of $\overline{\mu}^t$ in $U$ as $t \to \infty$ are $a$-thick a.s. Then the algorithm converges to $J$ a.s. on $Q$.*

(b) *For the $\infty$-nonexpansive case (Case 2), suppose all limit points of $\overline{\mu}^t$ in $U$ as $t \to \infty$ are $a$-thick for some $a > 0$, a.s. Then the algorithm converges to $J$ a.s. on $Q$.*

*Remark.* For Case 1 without the $a$-robustness hypothesis, the above analysis still gives some clue about the convergence of the algorithm: if all limit points of $\overline{\mu}^t$ are $a$-thick a.s., the algorithm will converge to the smallest closed set outside which $\nabla V. Mf < 0$ for $a$-thick $M$.

Clearly, one would like $P(Q) = 1 = P(\hat{X} < \infty)$. One observes that the boundedness of $\hat{X}$ is used twice: in Lemma 3.3 and to prove almost sure convergence on $Q$ of $\{\overline{M}_n\}$. In either case, it is unnecessary if $f$ (or, in Case 2, $g$) is bounded. If not, the problem of establishing $P(Q) = 1$ remains. This is so even for the traditional "centralized" algorithm, and it is not unusual to find results that state convergence if the iterates remain bounded or visit a neighborhood of the desired attractor infinitely often, a.s. There is no general scheme for showing $P(Q) = 1$. There are, however, problem-specific techniques for special problem classes. We list a few recent ones below without details, referring the reader to the original works for those.

(i) *Martingale methods.* These usually take the form of establishing the "almost supermartingale" property [19, p. 33] for $\{V(X(n))\}$, where $V : R^d \to R^+$ is a continuously differentiable "stochastic Liapunov function" satisfying $V(x) \to \infty$ as $\|x\| \to \infty$. This leads to the almost sure boundedness of $\{V(X(n))\}$, hence of $\{X(n)\}$. For strict Liapunov systems, the Liapunov function therein will itself suffice in most cases. The adaptation of this approach to the asynchronous case, however, is rendered difficult by the presence of delays. Specific instances of it have been worked out, a good example being the stochastic gradient schemes discussed in [3, section 7.8]. For the "centralized" case without delays, see [2, p. 239].

(ii) *Projection and related schemes.* One way to escape the boundedness issue is to alter the algorithm by projecting the iterates back onto a prescribed, large bounded set whenever they exit from the same. The trade-off is that the limiting ODE becomes more complicated. It is now confined to the said set and thus involves a "reflection at the boundary" of the same in an appropriate sense. The analysis of such schemes for the centralized case is by now standard, and an excellent exposition appears in [16, pp. 191–194]. It seems possible to extend it to the present case. (See [1] for a specific instance.)

In an ingenious boundedness proof for the case when $F(x, y)$ is homogeneous of degree 1 in its first argument (important in certain "learning" algorithms), Jaakola, Jordan, and Singh [15] use the almost sure convergence of the algorithm with rescaling to deduce the almost sure boundedness of the one without. See [1] for some extensions of this idea and application to a specific asynchronous situation.

In a somewhat similar spirit, but using different techniques, Chen [8] discusses stabilization of the (centralized) algorithm by truncating the iterates while slowly increasing the truncation bounds.

(iii) *Tsitsiklis conditions.* For Case 2 (nonexpansive maps) studied above, Tsitsiklis [21] gives a remarkable set of conditions for almost sure boundedness when additional structure is available, such as an appropriate monotonicity property of the map or contraction property under a suitable norm. These are very useful for applications arising from dynamic programming.

In some special cases (e.g., when $F(\cdot, y)$ has a common fixed point), one may adapt the conditions of [3, p. 433], for deterministic algorithms to prove almost sure boundedness. See [5] for an instance of this.

In [5], almost sure $a$-thickness of the limit points of $\{\overline{\mu}^t\}$ for a suitable $a > 0$ is established for the synchronous case. That argument does not follow for the asynchronous case. In fact, it will soon become clear that such a result need not hold in general, and whether it does depends crucially on the relationships between the sequences $\{a(n,i)\} \in$ ITS, $i \in I$. We now consider an important special case where things work out.

Say that the family $\{a(n,i)\}, i \in I$, is balanced if there exist $a_{ij} > 0, i, j \in I$, such that

$$\lim_{n\to\infty} \frac{\sum_{m=0}^n a(m,j)}{\sum_{m=0}^n a(m,i)} = a_{ij}.$$

Equivalently, if $h_i, h_j$ are continuous, eventually nonincreasing functions $R^+ \to R^d$ that restrict to $\{a(n,i)\}, \{a(n,j)\}$, respectively, at integer values of their arguments, then

(3.7)
$$\lim_{t\to\infty} \frac{\int_0^t h_j(s)ds}{\int_0^t h_i(s)ds} = a_{ij}.$$

Certain relations between $a_{ij}$'s are obvious: $a_{ii} = 1$, $a_{ik} = a_{ij}a_{jk}$, $a_{ji} = 1/a_{ij}$. An important special case is $a_{ij} = 1\ \forall i, j$, which would be true, e.g., when all $\{a(n,i)\}, i \in I$, are identical. Let $\beta(i) = a_{1i}/a_{11}$ and $\overline{\beta}(i) = \beta(i)/\sum_j \beta(j)$. Then $\overline{\beta}(i) \in (0,1)\ \forall i$ and $\sum_i \overline{\beta}(i) = 1$. Also $a_{ij} = \overline{\beta}(j)/\overline{\beta}(i)\ \forall i, j$. Set $\overline{a} = \min \overline{\beta}(i)$.

THEOREM 3.2. *If $\{a(n,i)\}, i \in I$, are balanced, the conclusions of Theorem* 3.1(b) *hold. Those of Theorem* 3.1(a) *hold if, in addition, $a \le \overline{a}$.*

*Proof.* For $i \in I$, let $q(i,n) = \sum_{m=0}^n I\{\varphi_m = i\}$. By (3.2),

(3.8)
$$\liminf_{n\to\infty} q(i,n)/n \ge \eta \quad \text{a.s.}, \quad i \in I.$$

Fix $i, j \in I$. Then, for $z > 0$,

$$\lim_{t\to\infty} \frac{\int_z^t \mu_s(j)ds}{\int_z^t \mu_s(i)ds} = \lim_{n\to\infty} \frac{\sum_{m=0}^{q(j,n)} a(m,j)}{\sum_{m=0}^{q(i,n)} a(m,i)}$$

$$= \lim_{n\to\infty} \frac{\int_0^{q(j,n)} h_j(s)ds}{\int_0^{q(i,n)} h_i(s)ds}$$

$$= \lim_{n\to\infty} \frac{\int_0^{q(j,n)} h_j(s)ds}{\int_0^n h_j(s)ds} \cdot \frac{\int_0^n h_j(s)ds}{\int_0^n h_i(s)ds} \cdot \frac{\int_0^n h_i(s)ds}{\int_0^{q(i,n)} h_i(s)ds}$$

$$= a_{ij} \quad \text{a.s.}$$

uniformly in $z$ in a compact interval, by (2.5) and (3.8).

Thus, for $x > 0$,

$$\lim_{t\to\infty} \frac{\int_0^x \int_0^t \mu_{s+y}(j)dsdy}{\int_0^x \int_0^t \mu_{s+y}(i)dsdy} = \lim_{t\to\infty} \frac{\int_0^t \int_0^x \mu_{s+y}(j)dsdy}{\int_0^t \int_0^x \mu_{s+y}(i)dsdy} = a_{ij} \quad \text{a.s.}$$

By l'Hôpital's rule,

$$\lim_{t\to\infty} \frac{\int_0^x \mu_{t+y}(j)dy}{\int_0^x \mu_{t+y}(i)dy} = a_{ij} \quad \text{a.s.}$$

It follows that, a.s., any limit point $\overline{\mu}^*$ of $\{\overline{\mu}^t\}$ in $U$ as $t \to \infty$ must satisfy $\int_0^x \mu_t^*(j)dt / \int_0^x \mu_t^*(i)dt = a_{ij}$. Then so will $\overline{\mu}^{*t}$, $t \geq 0$. Since $x > 0$ was arbitrary, we have $\mu_t^*(j)/\mu_t^*(i) = a_{ij}$ for almost every $t$, where we may drop the "almost every $t$" by taking a suitable modification. Then we must have $\mu_t^*(i) = \overline{\beta}(i) \; \forall i, \; t$, and the matrix $M^{\overline{\mu}^*}(t)$ is the constant diagonal matrix $M^* = \mathrm{diag}(\overline{\beta}(1), \ldots, \overline{\beta}(d))$. The rest is easy.    □

*Remark.* In the latter case, one may in fact replace the $a$-robustness condition and the condition $a \leq \overline{a}$ by the simpler condition $\nabla V. M^* f < 0$ outside $J$.

In particular, if $\{a(n, i)\}$ are identical, $M^{\overline{\mu}^*}$ is $1/d$ times the identity matrix, implying that the rescaled time axis is apportioned equally to all components. One may dub this the "asymptotic equipartition of time."

**4. Conclusions.** The foregoing analysis raises several interesting issues, which are listed below.

1. We have not presented any results on the convergence rate. For the ODE, the rate of convergence to $J^\epsilon$ for $\epsilon > 0$ could be gleaned from the Liapunov function and would be eventually mimicked by the interpolated algorithm $x(\cdot)$. There are two catches here. One is that "eventually" could be a long way into the future. Second, the passage from $\{X(n)\}$ to $x(\cdot)$ involves a time-scaling $n \to t(n)$, which has to be inverted to obtain the actual convergence rate of $\{X(n)\}$. These aspects need further study.

2. It seems plausible that one could retain the above results if (3.7) were replaced by the weaker requirement that the corresponding liminf be bounded away from zero. One cannot then expect $\{\overline{\mu}^t\}$ to a.s. converge to a fixed element, but it is conjectured that one will still retain the property that all limit points of $\{\overline{\mu}^t\}$ in $U$ are $a$-thick for some $a > 0$.

3. In engineering applications, $\{a(n, i)\}$ are design parameters and can be chosen to be balanced. This may not, however, be so in "emergent" computations or when (1.4) is merely a computational paradigm for a natural process such as a macroeconomic learning system. An interesting problem, then, is to let each agent (processor) "learn" its stepsize scheme in real time based on observations of stepsizes used by, say, "neighboring" agents. One may then try to show that under reasonable conditions, this leads to balanced schemes.

4. If we had allowed some of the $a_{ij}$'s to be zero, it is clear that the corresponding diagonal elements of $M^*$ will be zero and $M^*$ is no longer thick. This reflects different time scales in the speed of adjustment of different learners. It would be interesting to analyze this situation using the theory of singularly perturbed differential equations.

5. If (1.2) had no extraneous randomness, i.e., $F(X(n), \xi(n)) = H(X(n)) \; \forall n$ for a suitable $H$, the foregoing shows that a stepsize scheme from ITS suppresses the effects of communication delays in deterministic recursions under a mild conditional moment condition (A3). This is in contrast to the usual role of ITS as a pure noise-suppressing mechanism. Compare this with the fact that even linear recursions with constant stepsize show very complex behavior in the presence of communication delays [12].

6. Yet another possibility to explore is the use of the Wentzell–Freidlin theory of small noise asymptotics [10] to get a dynamic picture of the behavior of the algorithm in the vicinity of $J$, in particular, to see if it favors certain points in $J$. This is in the spirit of some recent work on annealing algorithms [11] and equilibrium selection in evolutionary games [9].

## REFERENCES

[1] J. ABOUNADI, D. BERTSEKAS, AND V.S. BORKAR, *O.D.E. Analysis for Q-Learning Algorithms*, preprint, Laboratory for Information and Decision Systems, M.I.T., Cambridge, MA, 1996.

[2] A. BENVENISTE, M. METIVIER, AND P. PRIOURET, *Adaptive Algorithms and Stochastic Approximations*, Springer-Verlag, Berlin, Heidelberg, 1990.

[3] D.P. BERTSEKAS AND J.N. TSITSIKLIS, *Parallel and Distributed Computation*, Prentice-Hall, Englewood Cliffs, NJ, 1989.

[4] V.S. BORKAR, *Topics in Controlled Markov Chains*, Pitman Res. Notes in Math. Ser. 240, Longman Scientific and Technical, Harlow, UK, 1991.

[5] V.S. BORKAR, *Distributed computation of fixed points of $\infty$-nonexpansive maps*, Proc. Indian Acad. Sci. Math. Sci., 106 (1996), pp. 289–300.

[6] V.S. BORKAR AND V.V. PHANSALKAR, *Managing interprocessor delays in distributed recursive algorithms*, Sādhanā, 19 (1994), pp. 995–1003.

[7] V.S. BORKAR AND K. SOUMYANATH, *A new analog parallel scheme for fixed point computation*, Part I: Theory, IEEE Trans. Circuits Systems I Fund. Theory Appl., 44 (1997), pp. 509–522.

[8] H.F. CHEN, *Stochastic approximation and its new applications*, in Proc. 1994 Hong Kong Internat. Workshop on New Directions in Control and Manufacturing, 1994, pp. 2–12.

[9] D. FOSTER AND P. YOUNG, *Stochastic evolutionary game dynamics*, Theoret. Population Biol., 38 (1990), pp. 229–232.

[10] M. FREIDLIN AND A. WENTZELL, *Random Perturbations of Dynamical Systems*, Springer-Verlag, Berlin, New York, 1984.

[11] S.B. GELFAND AND S.K. MITTER, *Recursive stochastic algorithms for global optimization in $R^d$*, SIAM J. Control Optim., 29 (1992), pp. 999–1018.

[12] R. GHARAVI AND V. ANANTHARAM, *A Structure Theorem for Partially Asynchronous Relaxations with Random Delays*, ERL Memo. No. M92/143, Electronics Research Laboratory, University of California, Berkeley, 1993.

[13] S. HAYKIN, *Neural Networks: A Comprehensive Foundation*, McMillan, New York, 1994.

[14] M. HIRSCH, *Convergent activation dynamics in continuous time networks*, Neural Networks, 2 (1987), pp. 331–349.

[15] T. JAAKOLA, M. JORDAN, AND S.P. SINGH, *On the convergence of stochastic iterative dynamic programming algorithms*, Neural Computation, 6 (1994), pp. 1185–1201.

[16] H. KUSHNER AND D. CLARK, *Stochastic Approximation Methods for Constrained and Unconstrained Systems*, Springer-Verlag, New York, 1978.

[17] H. KUSHNER AND G. YIN, *Stochastic approximation algorithms for parallel and distributed processing*, Stochastics, 22 (1987), pp. 219–250.

[18] H. KUSHNER AND G. YIN, *Asymptotic properties of distributed and communicating stochastic approximation algorithms*, SIAM J. Control Optim., 25 (1987), pp. 1266–1290.

[19] J. NEVEU, *Discrete Parameter Martingales*, North-Holland, Amsterdam, 1975.

[20] T. SARGENT, *Bounded Rationality in Macroeconomics*, Clarendon Press, Oxford, UK, 1993.

[21] J. TSITSIKLIS, *Asynchronous stochastic approximation and Q-learning*, Machine Learning, 16 (1994), pp. 185–202.

[22] J. TSITSIKLIS, D. BERTSEKAS, AND M. ATHANS, *Distributed asynchronous deterministic and stochastic gradient optimization algorithms*, IEEE Trans. Automatic. Control, AC-31 (1986), pp. 803–812.