

ANALYSIS OF A STOCHASTIC APPROXIMATION ALGORITHM FOR COMPUTING QUASI-STATIONARY DISTRIBUTIONS

J. BLANCHET,^{**} *Columbia University*

P. GLYNN,^{***} *Stanford University*

S. ZHENG,^{****} *Columbia University*

Abstract

We study the convergence properties of a Monte Carlo estimator proposed in the physics literature to compute the quasi-stationary distribution on a transient set of a Markov chain (see De Oliveira and Dickman (2005), (2006), and Dickman and Vidigal (2002)). Using the theory of stochastic approximations we verify the consistency of the estimator and obtain an associated central limit theorem. We provide an example showing that convergence might occur very slowly if a certain eigenvalue condition is violated. We alleviate this problem using an easy-to-implement projection step combined with averaging.

Keywords: Quasi-stationary distribution; stochastic approximation; Markov chain; central limit theorem

2010 Mathematics Subject Classification: Primary 60J22
Secondary 60J10

1. Introduction

A quasi-stationary distribution can be characterized as the principal (left) eigenvector of a substochastic matrix. Consequently, it is natural to use numerical linear algebra methods [15] for computing such a principal eigenvector. Nevertheless, the application of these methods, and related Monte Carlo variations, such as [13], [18], [19], and [22], becomes difficult when the underlying matrix is large.

In [9], [10], and [12] the authors proposed an iterative Monte Carlo procedure to estimate the quasi-stationary distribution of very large Markov chains. A key feature of the procedure is that in each iteration only a small portion of the matrix is used.

The Fleming–Viot (FV) method, [6], [14], [16], [21], provides a powerful alternative. It consists of N particles evolving according to suitable dynamics in continuous time. As both time, t , and the number of particles, N , tend to ∞ , the empirical measure of the positions of the particles at time t converges almost surely to the underlying quasi-stationary distribution. A significant advantage of the FV method is that it can be run in parallel. A disadvantage is that, for a fixed value t , if only N is sent to ∞ , the method will be biased.

Received 21 October 2014; revision received 10 August 2015.

* Postal address: Department of Industrial Engineering & Operations Research, Columbia University, 500 West 120th Street, New York, NY 10027, USA.

** Email address: jose.blanchet@columbia.edu

*** Postal address: Department of Management Science & Engineering, Stanford University, Huang Engineering Center, 475 Via Ortega, Stanford, CA 94035-4121, USA. Email address: glynn@stanford.edu

**** Email address: johnz622@gmail.com

In contrast to the FV method, the method in [9], [10], and [12], which we analyze here, is asymptotically unbiased as the number of iterations, N , tends to ∞ . Moreover, as we shall show, a small modification can be made to the method to ensure convergence at rate $N^{-1/2}$. Note that such a rate of convergence is impossible to achieve in the FV method because of the presence of the bias appearing by truncating the evolution at time t .

The method suggested in [9], [10], and [12], it turns out, is equivalent to a class of algorithms studied in the context of urn processes, [1], [2], [23]. So, the convergence of the sampling procedure has been rigorously established in the urn process literature. Moreover, in [2], some results on rates of convergence have been obtained. These results involve a central limit theorem (CLT) for the inner product of the following two quantities: any nonprincipal eigenvector (or linear combinations thereof) and the estimated quasi-stationary vector (i.e. the estimated principal eigenvector). One of our contributions in this paper is the development of a multidimensional CLT for the estimated quasi-stationary vector. Therefore, we can obtain a CLT for the inner product between the estimated quasi-stationary vector and any other vector. More generally, our main contributions are as follows.

- Our paper recognizes the algorithm in [9], [10], and [12] as a stochastic approximation algorithm, [25], (Section 3.2).
- Using the stochastic approximation connection we prove the convergence of the underlying estimator and provide sufficient conditions for a multidimensional CLT (Theorem 3.1).
- We illustrate the very slow convergence rate that might occur if the CLT conditions fail to apply (Section 4).
- More importantly, using Polyak–Ruppert averaging [24], we suggest an improved algorithm (Section 4.2) which exhibits a valid CLT with optimal rate of convergence assuming only irreducibility of the underlying substochastic matrix.
- We provide an estimator which allows us to compute the variance in the CLT, see Section 4.2.1.

The vanilla version of the algorithm analyzed here (without averaging) has independently been studied in [3]. In contrast to [3], our focus is more algorithmic. In particular, our emphasis on exploring the close connection of the algorithm to stochastic approximation leads naturally to a Polyak–Ruppert averaging variant that generally exhibits an optimal square root convergence rate, in contrast to the original algorithm which may display sub-square root convergence rates. Numerical experiments showing the dramatic improvement obtained by introducing averaging are given in [4], and [26]; see also Section 5.

We study discrete-time Markov chains. The adaptation to continuous-time Markov chains is relatively straightforward and was given in [26]. The convergence of the estimator for uniformly ergodic Markov chains taking values on a compact space was also studied in [26]. The work of [16] provides a modern survey of stochastic approximations for countable state-spaces.

The rest of the paper is organized as follows. In Section 2 we introduce notation and the definition of a quasi-stationary distribution. In Section 3 we review stochastic approximation methods and sketch the proof of Theorem 3.1 (the full proof is given in Sections 6.1 and 6.2). In Section 4.2 we discuss the averaging improvement. All the technical results are elaborated in Section 6.

2. Quasi-stationary distribution: basic notions

Let $\{X_n : n \geq 0\}$ be a discrete-time, finite-state-space, Markov chain $\{X_n : n \geq 0\}$ with transition matrix \mathcal{P} . We assume that 0 is an absorbing state and that $1, \dots, d$ are nonabsorbing. We shall write $\mathcal{P}(i, j) = Q(i, j)$ for $i, j \in \{1, \dots, d\}$ so that Q defines a substochastic transition matrix of size d by d . We impose the following assumption on Q .

Assumption 2.1. *Throughout the rest of the paper we shall assume that Q is irreducible and that $Q^n \rightarrow 0$ as $n \rightarrow \infty$.*

Under this assumption the Perron–Frobenius theorem [17] guarantees the existence of a unique positive probability vector μ_* such that $\mu_*^\top Q = \lambda_* \mu_*^\top$, where $\lambda_* \in (0, 1)$. (Throughout the rest of the paper we use ‘ \top ’ to denote transposition.) Precisely, the vector μ_* is the quasi-stationary distribution of \mathcal{P} . Early reference for quasi-stationary distributions are [7] and [8]; see also [16] and [21] for a discussion on quasi-stationary distributions on infinite spaces.

3. Stochastic approximations analysis of the algorithm

3.1. Basic notions of stochastic approximations

Given θ_0 , the standard stochastic approximations recursion takes the form

$$\theta_{n+1} = \theta_n + \varepsilon_n W_n \quad \text{for } n \geq 0, \tag{3.1}$$

where $\{\varepsilon_n\}_{n \geq 0}$ is a (deterministic) step-size sequence of nonnegative numbers satisfying $\sum \varepsilon_n = \infty$ but $\sum \varepsilon_n^2 < \infty$. And the n th noise observation, W_n , is measurable with respect to $\mathcal{F}_n = \sigma\{(\theta_k, W_{k-1}) : 1 \leq k \leq n\}$. Moreover, it is assumed that $\mathbb{E}(W_n \mid \mathcal{F}_n) = g(\theta_n)$ for some function $g(\cdot)$. Under mild regularity conditions, to be reviewed momentarily in our setting, we have the fact that θ_n converges almost surely to the stable attractors of the ordinary differential equation (ODE) $\dot{\theta}(t) = g(\theta(t))$; see, for example, [20, Theorem 5.2.1].

3.2. The precise algorithm in stochastic approximation form

Suppose that we have d bins (one for each element in the underlying transient set). At the beginning of the n th iteration we have a certain distribution of balls across the bins and we select an initial position according to such distribution. For example, if $d = 2$ and there are three balls in the first bin and five balls in the second bin, then state 1 is selected with probability $\frac{3}{8}$ and state 2 is selected with probability $\frac{5}{8}$. The n th iteration then proceeds by running a Markov chain starting from the selected state $i^* \in \{1, \dots, d\}$ according to the underlying dynamics, until absorption (i.e. until hitting state 0) and we call such a trajectory a tour. We count the number of times state j (for $j \in \{1, \dots, d\}$) is visited during such a tour; note, for example, that the state i^* is visited at least once. We then update the distribution of balls across bins by adding these counts. So, for instance, back to the earlier example with $d = 2$, if during the n th tour state 1 was visited twice, while state 2 was visited four times, then the distribution of balls at the beginning of the $(n + 1)$ th iteration will be $(5 = 3 + 2, 9 = 5 + 4)$. The output of the algorithm is the normalized distribution of balls (in order to obtain a probability vector) after many iterations.

We now explain how this procedure can be described in terms of a stochastic approximation recursion.

Notation. We set μ_n as the sequence of probability vectors of the transient set $\{1, \dots, d\}$ obtained at the n th iteration of the algorithm. This vector will store the cumulative empirical measure up to, and including, the n th iteration of the algorithm. We use $\mu_n(x)$ to denote the particular value at the transient state x .

We set $\{X_k^{(n)}\}_{k \geq 0}$ as the Markov chain that is realized during the n th iteration of the algorithm. These Markov chains are conditionally independent (given the values X_0^n). The n th Markov chain has an initial position drawn from the vector μ_n .

We define $\tau^{(n)} = \inf\{k \geq 0 : X_k^{(n)} = 0\}$. (Recall that 0 is the underlying absorbing state.)

We are interested in analyzing the recursion

$$\mu_{n+1}(x) = \frac{(\sum_{k=0}^n \tau^{(k)})\mu_n(x) + (\sum_{k=0}^{\tau^{(n+1)}-1} \mathbf{1}\{X_k^{(n+1)} = x \mid X_0^{(n)} \sim \mu_n\})}{(\sum_{k=0}^{n+1} \tau^{(k)})} \tag{3.2}$$

for all $x \in \{1, \dots, d\}$, where the notation $\mathbf{1}\{X_k^{(n+1)} = x \mid X_0^{(n)} \sim \mu_n\}$ describes the indicator of the event $\{X_k^{(n+1)} = x\}$ and we emphasize that $X_0^{(n)}$ is sampled using the distribution μ_n . We may select the initial probability distribution μ_0 supported on $\{1, \dots, d\}$ in an arbitrary way and set $\tau^{(0)} = 0$ by convention.

We transform μ_n into a more familiar stochastic approximation form by writing

$$\mu_{n+1}(x) = \mu_n(x) + \frac{1}{n+1} \left(\frac{\sum_{k=0}^{\tau^{(n+1)}-1} (\mathbf{1}\{X_k^{(n+1)} = x \mid X_0^{(n)} \sim \mu_n\} - \mu_n(x))}{(\sum_{j=0}^{n+1} \tau^{(j)})/(n+1)} \right).$$

Compared to the standard form in (3.1) we recognize that $\varepsilon_n = 1/(n+1)$; however, if we attempt to make a direct translation into (3.1) we see that the denominator is somewhat problematic because its conditional expectation (given the whole history of the algorithm up to the end of the n th iteration) is not only a function of μ_n . To address this issue, we add another variable, T_n , leading to the recursions (assuming $T_0 = 0$),

$$T_{n+1} = T_n + \frac{1}{n+2} (\tau^{(n+1)} - T_n) = \frac{1}{n+1} \sum_{j=0}^n \tau^{(j)}, \tag{3.3}$$

$$\mu_{n+1}(x) = \mu_n(x) + \frac{1}{n+1} \left(\frac{\sum_{k=0}^{\tau^{(n+1)}-1} (\mathbf{1}\{X_k^{(n+1)} = x \mid X_0^{(n)} \sim \mu_n\} - \mu_n(x))}{T_n + \tau^{(n+1)}/(n+1)} \right).$$

In order to provide a more succinct notation let us define

$$Y_n(\mu^\top, T)(x) := \frac{\sum_{k=0}^{\tau-1} (\mathbf{1}\{X_k = x \mid X_0 \sim \mu\} - \mu(x))}{T + \tau/(n+1)}, \quad Z(\mu^\top, T) := (\tau - T), \tag{3.4}$$

where $\{X_l : l \geq 0\}$ denotes a generic Markov chain with transition matrix P , X_0 is distributed according to μ (supported on $\{1, \dots, n\}$), and τ corresponds to the first hitting time to 0 of the chain $\{X_l : l \geq 0\}$. We also write

$$Y(\mu^\top, T)(x) := \sum_{k=0}^{\tau-1} \frac{\mathbf{1}\{X_k = x \mid X_0 \sim \mu\} - \mu(x)}{T}.$$

Note that Y is time homogeneous whereas Y_n is not. Then, we can write the stochastic approximation recursion in distribution via

$$\mu_{n+1}(x) = \mu_n(x) + \frac{1}{n+1} Y_n(\mu_n, T_n)(x), \quad T_{n+1} = T_n + \frac{1}{n+2} Z(\mu_n, T_n).$$

If we let $\theta_n = (\mu_n^\top, T_n)$ we now have a setting very close to that described in (3.1), except for the fact that $g(\cdot)$ is time homogeneous (i.e. of the form $g_n(\cdot)$).

We note the following.

- As mentioned earlier, Y_n is not time homogeneous because of the presence of the term $\tau/(n + 1)$. It is not difficult to argue (as we shall do in Lemma 6.1) that such a term is asymptotically negligible because $\tau^{(n)} = O(\log(n))$ almost surely.
- During the course of the algorithm, each μ_n is a probability vector; that is, $\mu_n \in H := \{x \in \mathbb{R}_+^d : \mathbf{e}^\top x = 1\}$, where \mathbf{e} is the vector of 1s. So, the boundedness requirement in [20, Theorem 5.2.1] holds automatically for μ_n . We only need to argue boundedness for the coordinate T_n .

3.2.1. *Convergence result: consistency and CLT.* We now state the main result of this section.

Theorem 3.1. *Suppose that μ_0 is any probability vector supported on $\{1, \dots, n\}$ and pick $T_0 \geq 1$. Then $(\mu_n^\top, T_n) \rightarrow \theta_* := (\mu_*^\top, 1/(1 - \lambda_*))$ almost surely (a.s.), where the left principal eigenvector μ_*^\top of Q is normalized so that $\mu_*^\top \mathbf{e} = 1$. Finally, if $\bar{\lambda}$ is any nonprincipal eigenvalue of Q (i.e. $\bar{\lambda} \neq \lambda_*$) and*

$$\operatorname{Re}((1 - \bar{\lambda})^{-1}) < \frac{1}{2}(1 - \lambda_*)^{-1}, \tag{3.5}$$

then $n^{1/2}(\mu_n - \mu_*) \xrightarrow{D} N(0, V_0)$ for some V_0 , explicitly characterized by (6.11), where \xrightarrow{D} denotes convergence in distribution.

Proof of Theorem 3.1 (sketch). The full proof is in Sections 6.1 and 6.2, here we outline the main ideas. We use the ODE method [20, Theorem 5.2.1], which involves studying the behavior, as $t \rightarrow \infty$, of the pair $(\mu(t), T(t))$ satisfying

$$\begin{aligned} \dot{T}(t) &= \mathbb{E}(\tau \mid X_0 \sim \mu(t)) - T(t) = \mu(t)^\top R \mathbf{e} - T(t), \\ \dot{\mu}(t)T(t) &= \mathbb{E}\left(\sum_{k=0}^{\tau-1} [\mathbf{1}\{X_k = \cdot \mid X_0 \sim \mu(t)\} - \mu(t)]\right) = (\mu(t)^\top R - (\mu(t)^\top R \mathbf{e})\mu(t)^\top)^\top, \end{aligned}$$

where $R = (I - Q)^{-1}$. In Section 6.1 we are able to show, using Duhamel’s principle, that for a given initial position in the probability simplex H , the solution to a suitably reduced dynamical system (obtained by assuming that $T(t) = 1$) exists and converges as $t \rightarrow \infty$ to its stationary point. This stationary point is the unique solution to the eigenvalue problem $\mu_*^\top R = \rho_* \mu_*^\top$, $\mu_*^\top \mathbf{e} = 1$, and $\mu_* \geq 0$, where $\rho_* = 1/(1 - \lambda_*)$. Uniqueness follows from Perron–Frobenius’ theorem. The complete dynamical system, for $\mu(t)$ and $T(t)$, is a time change of the reduced one, so we can connect them via a simple transformation.

Thus, applying [20, Theorem 5.2.1] we can conclude that μ_n converges to the quasi-stationary distribution for all initial configurations $(\mu_0, T_0) \in H \times [1, \infty)$.

For the CLT we invoke [20, Theorem 10.2.1]. Because the recursion in (3.3) uses step size $\varepsilon_n = 1/(n + 1)$, we need to verify that the Jacobian matrix of the ODE vector field, evaluated at the stability point, has spectral radius less than $-\frac{1}{2}$. As we show in Section 6.2 this is equivalent to requiring (3.5). The expression for V_0 is extracted from the variance of an associated Ornstein–Uhlenbeck process as in [20, p. 332]. □

4. Variations on the algorithm with improved rate of convergence

We study the (dramatic) deterioration that can occur in the rate of convergence if (3.5) is not satisfied. We focus on a simple example consisting of two states, the phenomenon is not unique to this example but to the implicit selection of $\varepsilon_n = O(1/n)$ in the stochastic approximations

form, and the importance of the constant multiplying $1/n$ in these cases. Moreover, we note there are natural algorithmic ‘tricks’ that one might attempt to use, and which are likely to induce a violation of (3.5). For example, observe that the eigenvectors of Q and ϕQ coincide for any $\phi \in (0, 1)$. So we can choose $\phi > 0$ small in order to shorten the expected size of each iterate, since absorption now occurs with probability at least $1 - \phi$ in each step. However, because of the nonlinearity of the $1/(1 - \bar{\lambda})$ as a function of $\bar{\lambda}$ and its presence in (3.5), we can see that choosing $\phi > 0$ too small might result in a significant deterioration in the rate of convergence (despite the gain in speed at each iteration).

4.1. Counterexample to square root convergence

Consider the Markov chain with states $\{0, 1, 2\}$, the state 0 is absorbing and the matrix Q satisfies

$$Q = \begin{pmatrix} \frac{1 - \varepsilon}{2} & \frac{1 - \varepsilon}{2} \\ \frac{1 - \varepsilon}{2} & \frac{1 - \varepsilon}{2} \end{pmatrix}$$

By symmetry the recursion that we analyze, namely (3.2), can be tracked by a simple process, $\{\bar{X}_m : m \geq 0\}$, which we describe now. Assume that the distribution of \bar{X}_0 is given. At step m , the value of \bar{X}_m is decided as follows. First we sample a Bernoulli trial which we call the *type*. The type has success with probability equal to $1 - \varepsilon$.

If the type is a success, we sample a second Bernoulli trial with probability $\frac{1}{2}$ of the success, if the second trial is successful we let $\bar{X}_m = 1$, if it is a failure we let $\bar{X}_m = 2$.

If the type is a failure (which occurs with probability ε), we sample state 1 or 2 according to the empirical measure of $\{\bar{X}_k : 0 \leq k \leq m - 1\}$.

Let T_n be the time at which the n th failure type occurs. Then (3.2) is equivalent to studying $\mu_n(x) = \sum_{k=0}^{T_n} \mathbf{1}\{\bar{X}_k = x\} / T_n$. The process $\{\bar{X}_m : m \geq 0\}$ is known as a self-interacting Markov chain, see [11]. Equation (3.5) in Theorem 3.1 applied to this case corresponds to requiring that $\varepsilon < \frac{1}{2}$. Of course, $\mu = (\frac{1}{2}, \frac{1}{2})$. Del Moral and Miclo [11] is applicable to this example and shows that if $f \neq 0$ there exists $\delta > 0$ such that, for $n \geq 1$, $\delta n^{-2(1-\varepsilon)} \leq \mathbb{E}((\mu_n^\top f - \mu^\top f)^2) \leq \delta^{-1} n^{-2(1-\varepsilon)}$. So the rate of convergence is not $O(n^{-1/2})$ but rather $O(n^{-(1-\varepsilon)})$.

4.2. Projection and averaging

Now that the method is under the stochastic approximation umbrella, we can modify the algorithm to enforce an optimal CLT rate regardless of the eigenvalues of Q . So, we consider the recursion

$$\bar{\mu}_{n+1} = \Pi_H \left(\bar{\mu}_n + \varepsilon_n \left(\sum_{k=0}^{\tau^{(n+1)}-1} (\mathbf{1}\{X_k^{(n+1)} = \cdot \mid X_0^{(n)} \sim \mu_n\} - \bar{\mu}_n(\cdot)) \right) \right),$$

where Π_H denotes the L_2 -projection into the probability simplex H . We still require $\sum \varepsilon_n = \infty$ and $\sum \varepsilon_n^2 < \infty$. As we now explain, we only need to perform a small number of projections. The vector inside the projection operator is equal to

$$\bar{\mu}_n(1 - \varepsilon_n \tau^{(n+1)}) + \varepsilon_n \sum_{k=0}^{\tau^{(n+1)}-1} \mathbf{1}\{X_k^{(n+1)} = \cdot \mid X_0^{(n)} \sim \bar{\mu}_n\},$$

so it always has components which add up to 1. Moreover, for a component to become negative, it is necessary that $\tau^{(n+1)} > \varepsilon_n^{-1}$. We shall argue in Lemma 6.1(iii), that there exists $\delta > 0$ such that $\tau^{(n+1)} > \delta \log(n)$ for only finitely many values of $n \geq 1$ with probability 1; thus, $\tau^{(n+1)} > \varepsilon_n^{-1}$ occurs only finitely many times if $\varepsilon_n = O(n^{-\alpha})$ for $\alpha > 0$. Finally, it is quite easy to perform the L_2 projection into a probability simplex. In particular,

$$\bar{\mu}_{n+1} = \left(\bar{\mu}_n(1 - \varepsilon_n \tau^{(n+1)}) + \varepsilon_n \sum_{k=0}^{\tau^{(n+1)}-1} \mathbf{1}\{X_k^{(n+1)} = \cdot \mid X_0^{(n)} \sim \bar{\mu}_n\} - u_{n+1} \mathbf{e} \right)_+$$

where $u_{n+1} > 0$ is the unique constant such that $\bar{\mu}_{n+1}^\top \mathbf{e} = 1$ (see [5]).

The advantage of the projection version is that we are free to choose slower step sizes so that we can weaken the condition for the required CLT to hold. In particular, when $\varepsilon_n = n^{-\alpha}$ and $\alpha \in (\frac{1}{2}, 1)$ we always obtain a $\varepsilon_n^{-1/2}$ -CLT. We summarize this observation in the following result proved in Section 7.

Proposition 4.1. *If $\varepsilon_n = n^{-\alpha}$ for $\alpha \in (\frac{1}{2}, 1)$, it follows that $\varepsilon_n^{-1/2}(\bar{\mu}_n - \mu_*) \xrightarrow{D} N(0, V_1)$, where V_1 can be characterized via (7.1).*

The Polyak–Ruppert averaging technique [24] can be applied jointly with the projection algorithm to ensure ‘square root convergence’, regardless of whether (3.5) holds or not, as the next theorem shows. Its proof, given in Section 7, is based on [24] and it uses the analyses behind Proposition 4.1 and Theorem 3.1.

Theorem 4.1. *Suppose that μ_0 is any probability vector supported on $\{1, \dots, n\}$ and pick $T_0 \geq 1$. Selecting $\varepsilon_n = n^{-\alpha}$ for $\alpha \in (\frac{1}{2}, 1)$, let $v_n = \sum_{k=1}^n \bar{\mu}_k/n$. Then, $n^{1/2}(v_n - \mu_*) \xrightarrow{D} N(0, \bar{V}_1)$, where \bar{V}_1 is given in (8.1).*

We can apply Theorem 4.1 in the estimation of quasi-stationary expectations of the form $\mathbb{E}(s(X) \mid X \sim \mu_*) = \mu_*^\top s$, using the estimator $\mu_n^\top s$ (note that we are encoding the function $s(\cdot)$ as a column vector). As a consequence of our CLT, we have the following corollary.

Corollary 4.1. *Under the notation defined in Theorem 3.1, we have $n^{1/2}(v_n^\top s - \mu_*^\top s) \xrightarrow{D} N(0, \sigma_s^2)$, where $\sigma_s^2 = s^\top \bar{V}_1 s$.*

4.2.1. *Estimating the asymptotic variance.* In the next result (proved in Section 8) we indicate how to estimate V_1 using the outcomes of the improved algorithm, which we ultimately advocate using.

Proposition 4.2. *Set v_n as in Theorem 4.1. For $\varepsilon_n = n^{-\alpha}$ with $\alpha \in (\frac{1}{2}, 1)$, let $n_k = \lceil k^{\beta/\alpha} \rceil$ and $N_n = \lceil n^{\alpha/\beta} \rceil$, with $\beta \in (\alpha, 1)$. Then*

$$\frac{1}{N_n} \sum_{k=0}^{N_n} \varepsilon_{n_k}^{-1} (\bar{\mu}_{n_k} - v_{n_k})(\bar{\mu}_{n_k} - v_{n_k})^\top \xrightarrow{\mathbb{P}} \bar{V}_1 \quad \text{as } n \rightarrow \infty,$$

where ‘ $\xrightarrow{\mathbb{P}}$ ’ denotes convergence in probability.

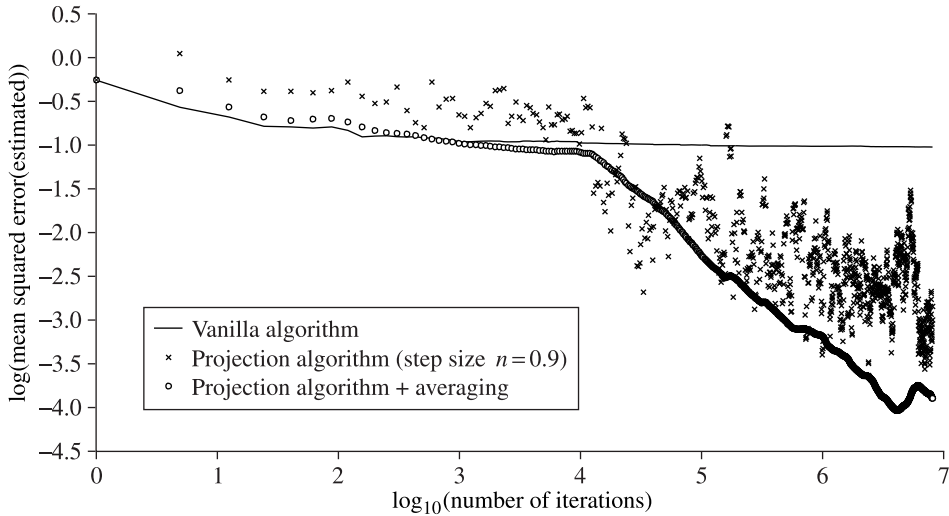


FIGURE 1: The M/M/1/100 queue with $\rho = 1.25$. Averaging plus projection (circles) significantly outperforms the original algorithm (solid line), and the projection-only algorithm (crosses).

In the context of Corollary 4.1, we can take advantage of Proposition 4.2 in order to estimate the asymptotic variance $\sigma_f^2 = s^\top \bar{V}_1 s$ using the relationship

$$\frac{1}{N_n} \sum_{k=0}^{N_n} \varepsilon_{n_k}^{-1} s^\top (\bar{\mu}_{n_k} - v_{n_k})(\bar{\mu}_{n_k} - v_{n_k})^\top s = \frac{1}{N_n} \sum_{k=0}^{N_n} \varepsilon_{n_k}^{-1} ((\bar{\mu}_{n_k} - v_{n_k})^\top s)^2.$$

5. Numerical example

We simulated the embedded discrete-time Markov chain induced by an absorbing M/M/1/100 queue. That is, a birth–death chain with a reflecting boundary at 100 and an absorbing state at 0. The traffic intensity of the system $\rho = 1.25$ (i.e. up rate is 1.25, down rate is 1). The expected time to absorption $\mathbb{E}(\tau)$ is large, so we introduce the trick of replacing Q by $0.95 \times Q$; so the transformed chain does not satisfy (3.5). In Figure 1 we show that the projection-with-averaging algorithm significantly outperforms both the original algorithm and the projection-only algorithm.

6. Proofs of main results

6.1. Proof of Theorem 3.1: convergence

We first restate a series of assumptions and notations that are used in [20, Theorem 5.2.1]. We adopt the abstract form of the recursion $\theta_{n+1} = \theta_n + \varepsilon_n W_n$. In our setting $\theta_n = (\mu_n^\top, T_n)$ and $W_n = (Y_n(\theta_n), Z(\theta_n))$ as defined in (3.4). Recall that \mathcal{F}_n is the σ -field generated by the iterates of the algorithm, namely, $\mathcal{F}_n = \sigma(\theta_0, \theta_i, W_{i-1} : 1 \leq i \leq n)$.

To show that $\mu_n \xrightarrow{\mathbb{P}} \mu_*$ a.s. we must verify that

(C.1) $\sum \varepsilon_n = \infty, \sum \varepsilon_n^2 < \infty$. This is immediately satisfied with the choice $\varepsilon_n = 1/(n + 1)$, as in our case. Moreover, define $t_n = \sum_{j=1}^n \varepsilon_j$ (with $t_0 = 0$), and let $m(s) = \max\{n : t_n \leq s\}$.

(C.2) Uniformly bounded variance: $\sup_n \mathbb{E} \|W_n\|_\infty^2 < \infty$. This is shown in Lemma 6.4 below. (We can use any norm, but we choose the norm $\|x\|_\infty = \max_i |x_i|$.)

(C.3) Local averaging condition: define $g_n(\theta_n) := \mathbb{E}(W_n \mid \mathcal{F}_n)$. The family of functions $\{g_n(\cdot)\}_{n \geq 0}$ must be uniformly equicontinuous, and there must exist a continuous function $g(\cdot)$ such that, for each θ , and each $t > 0$, $\sum_{k=n}^{m(t_n+t)} \varepsilon_k(g_k(\theta) - g(\theta)) \xrightarrow{\mathbb{P}} 0$, a.s. This local averaging condition is proved in Lemma 6.3 below.

In our setting, we write $\theta = (\mu^\top, T)$ and define $g(\theta) = (f^\top(\theta), h(\theta))$, where f and h are given via

$$f(\mu^\top, T) = \frac{1}{T} \mathbb{E} \left(\sum_{k=0}^{\tau-1} (\mathbf{1}\{X_k = \cdot\} - \mu) \mid X_0 \sim \mu \right) = \frac{1}{T} (\mu^\top R - (\mu^\top R e) \mu^\top)^\top, \tag{6.1}$$

$$h(\mu^\top, T) = \mathbb{E}(\tau \mid X_0 \sim \mu) - T = \mu^\top R e - T, \tag{6.2}$$

with $R = (I - Q)^{-1}$. We also define

$$f_n(\mu^\top, T) = \mathbb{E} \left(\sum_{k=0}^{\tau-1} \frac{(\mathbf{1}\{X_k = \cdot\} - \mu)}{T + \tau/(n+1)} \mid X_0 \sim \mu \right)$$

and set $g_n(\theta) = (f_n^\top(\theta), h(\theta))$.

Under (C.1)–(C.3), [20, Theorem 5.2.1] indicates that if the ODE $\dot{\theta}(t) = g(\theta(t))$ has an attractor (asymptotically stable point) in some domain \mathcal{D} and the sequence $\{\theta_n\}$ visits a compact subset within the domain \mathcal{D} infinitely often with probability 1, then θ_n converges to the attractor with probability 1.

In our situation it turns out that the entire probability simplex H is the domain of attraction for an attractor which is precisely the quasi-stationary vector. So we just need to study $\{T_n\}$. We will compute the functions $\{g_n(\cdot)\}$, verify (C.3), then (C.2), and finally the asymptotic stability behavior of the ODE. We will show that $\{T_n\}$ stays within a compact set throughout the course of the algorithm, and the uniform continuity of the functions $\{g_n(\cdot)\}$ holds for every compact set in $H \times [1, \infty)$. First, however, we obtain auxiliary results for $\{\tau^{(n)}\}$.

6.1.1. *Auxiliary results.* Define $\bar{\tau}(x)$ to be a random variable with the distribution of the first passage time to the absorbing state 0 given that the initial condition of the chain is the transient state $x \in \{1, \dots, d\}$. Suppose that the random variables $\{\bar{\tau}(j) : 1 \leq j \leq d\}$ are all independent. Then, let $\bar{\tau} = \max\{\bar{\tau}(j) : 1 \leq j \leq d\}$. We have the following simple but useful result.

Lemma 6.1. *The following claims hold:*

- (i) $\tau^{(n+1)}$ is stochastically bounded by $\bar{\tau}$;
- (ii) there exists $\delta > 0$ such that $\mathbb{E} \exp(\delta \bar{\tau}) < \infty$;
- (iii) $\mathbb{P}(\tau^{(n+1)} > \log(n) \text{ infinitely often}) = 0$;
- (iv) almost surely we have $1 \leq \overline{\lim} \sum_{k=1}^n \tau^{(k)} / n \leq \mathbb{E}(\bar{\tau}) < \infty$.

Proof. It is clear that (i) follows regardless of any assumption. For (ii)–(iv), we need $Q^n \rightarrow 0$ as $n \rightarrow \infty$ because this ensures that $\bar{\tau}(x)$ has a finite moment generating function in a neighborhood of the origin. □

We also have a useful expression for $f_n(\mu^\top, T)$. Define $v(x, s) := \mathbb{E}(\exp(-s\tau) \mid X_0 = x)$.

Lemma 6.2. For each $(\mu^\top, T) \in H \times (0, \infty)$, the x th component of $f_n(\mu^\top, T)$, namely, $f_n(\mu^\top, T)(x)$, is equal to

$$\int_0^\infty e^{-Tu} \mu^\top \left[v\left(x, \frac{u}{n+1}\right) (I - e^{-u/(n+1)} Q)^{-1} e_x - (I - e^{-u/(n+1)} Q)^{-1} v\left(\cdot, \frac{u}{n+1}\right) \mu(x) \right] du,$$

where e_x denotes the vector which has 1 in the x th coordinate and zeroes elsewhere.

Proof. First note that $(T + \tau/(n+1))^{-1} = \int_0^\infty e^{-(T+\tau/(n+1))u} du$. So, applying Fubini's theorem, since $\int_0^\infty \mathbb{E}(e^{-T\tau} \tau \mid X_0 \sim \mu) du < \infty$,

$$f_n(\mu^\top, T)(x) = \int_0^\infty \mathbb{E}\left(e^{-(T+\tau/(n+1))u} \sum_{k=0}^{\tau-1} (\mathbf{1}\{X_k = x\} - \mu(x)) \mid X_0 \sim \mu \right) du.$$

Again, another application of Fubini's theorem (also valid because $\mathbb{E}(\tau \mid X_0 \sim \mu) < \infty$) yields that the previous expression equals

$$\begin{aligned} & \int_0^\infty \mathbb{E}\left(e^{-(T+\tau/(n+1))u} \sum_{k=0}^\infty \mathbf{1}\{\tau > k\} (\mathbf{1}\{X_k = x\} - \mu(x)) \mid X_0 \sim \mu \right) du \\ &= \int_0^\infty e^{-Tu} \sum_{k=0}^\infty e^{-ku/(n+1)} \mathbb{E}(e^{-u((\tau-k)/(n+1))} \mathbf{1}\{\tau > k\} \\ & \quad \times (\mathbf{1}\{X_k = x\} - \mu(x)) \mid X_0 \sim \mu) du \\ &= \int_0^\infty e^{-Tu} \sum_{k=0}^\infty e^{-ku/(n+1)} \mathbb{E}\left(v\left(X_k, \frac{u}{n+1}\right) (\mathbf{1}\{X_k = x\} - \mu(x)) \mid X_0 \sim \mu \right) du \\ &= \int_0^\infty e^{-Tu} \sum_{k=0}^\infty e^{-ku/(n+1)} \left[v\left(x, \frac{u}{n+1}\right) \mu^\top Q^k e_x - \left(\mu^\top Q^k v\left(\cdot, \frac{u}{n+1}\right) \right) \mu(x) \right] du. \end{aligned}$$

The infinite series inside the integral of the last term can be simplified following the same logic behind the first equality in the display; hence, we obtain the conclusion of the lemma. \square

6.1.2. *Local averaging and uniformly bounded variance.* We first verify the uniformly bounded variance condition.

Lemma 6.3. It holds that $\{g_n(\cdot)\}$ is a sequence of uniformly equicontinuous functions on $H \times [1, \infty)$ and we have, for each $t > 0$,

$$\lim_{n \rightarrow \infty} \sum_{k=n}^{m(t_n+t)} \varepsilon_k (g_k(\mu^\top, T) - g(\mu^\top, T)) = 0.$$

Proof. Clearly, it holds that $|f_k(\mu^\top, T)(x) - f(\mu^\top, T)(x)|$ is bounded by

$$\mathbb{E}\left(\sum_{k=0}^{\tau-1} \frac{|\mathbf{1}\{X_k = x\} - \mu(x)| \tau}{T(T(k+1) + \tau)} \mid X_0 \sim \mu \right) \leq \mathbb{E}\left(\frac{\bar{\tau}^2}{k+1} \right).$$

Therefore,

$$\sum_{k=n}^{m(t_n+t)} \varepsilon_k |(f_k(\mu^\top, T) - f(\mu^\top, T))| \leq \mathbb{E} \left(\frac{\bar{\tau}^2}{n+1} \right)^{m(t_n+t)} \sum_{k=n}^{m(t_n+t)} \varepsilon_k \leq t \mathbb{E} \left(\frac{\bar{\tau}^2}{n+1} \right) \rightarrow 0$$

as $n \rightarrow \infty$. Finally, we need to argue that $g_n(\cdot)$ is uniformly equicontinuous on compact sets in $H \times [1, \infty)$. This follows easily by noting, from the expression obtained in Lemma 6.2, that the Jacobian $(Df_n)(\mu^\top, T)$ is uniformly bounded for a neighborhood around any point $(\mu^\top, T) \in H \times [1, \infty)$.

The coordinate of g_n corresponding to h does not depend on n and, thus, the result follows immediately in this case. □

Now we turn our attention to (C.2), namely, uniformly bounded variance.

Lemma 6.4. *We have that $\sup_n \mathbb{E} \|W_n\|_\infty^2 < \infty$.*

Proof. It holds that $\|Y_n(\mu_n, T_n)\|_\infty \leq \tau^{(n+1)}$. By Lemma 6.1(i), we have $\mathbb{E}|Z(\mu_n, T_n)| \leq 2\mathbb{E}\bar{\tau}^2$. Thus, Lemmas 6.1(i) and 6.1(ii) yield $\mathbb{E}\|W_n\|_\infty^2 \leq 3\mathbb{E}\bar{\tau}^2 < \infty$, thus concluding the proof. □

6.1.3. *Stability of the dynamical system and final convergence argument.* The dynamical system of interest, namely, $\dot{\theta}(t) = g(\theta(t))$ takes the form

$$\begin{aligned} \dot{\mu}(t)^\top &= f(\mu(t), T(t)) = \frac{1}{T(t)} (\mu(t)^\top R - (\mu(t)^\top R e) \mu^\top(t)), & (6.3) \\ \dot{T}(t) &= h(\mu(t), T(t)) = \mu(t)^\top R e - T(t). \end{aligned}$$

In the proof of [20, Theorem 5.2.1] it was shown that any converging subsequence of the suitably normalized iterates converges to a function $\theta(\cdot, \omega)$ which is a solution to the ODE $\dot{\theta}(t) = g(\theta(t))$. We only need to show that these solutions converge as $t \rightarrow \infty$ to $(\mu_*, 1/(1 - \lambda_*))$.

Define an associated reduced ODE as

$$\dot{v}(t) = (v(t)^\top R - (v(t)^\top R e) v^\top(t))^\top, \quad v(0) = \mu_0 \in H. \tag{6.4}$$

Note that the gradient of the vector field in (6.4) is continuously differentiable in H , therefore $v(\cdot)$ has a unique solution given $v(0) \in H$.

Suppose that $T_0 \geq 1$ and let $(\mu(\cdot), T(\cdot))$ be a solution to (6.3) obtained by the subsequence procedure in [20, Theorem 5.2.1], then define $\Gamma(t) = \int_0^t (1/T(s)) ds$. It follows by formal differentiation that $v(t) = \mu(\Gamma^{-1}(t))$ solves (6.4). The following result ensures regularity properties of $\Gamma(\cdot)$.

Lemma 6.5. *It holds that $\Gamma(t) > 0$, $\Gamma(\cdot)$ is strictly increasing, and $\Gamma(t) \rightarrow \infty$ as $t \rightarrow \infty$.*

Proof. Clearly, we have $T(s) \geq 1$ so $\Gamma(\cdot)$ is strictly increasing and nonnegative. Now, suppose that $\Gamma(\infty) < \infty$, this implies that

$$\begin{aligned} T(t) &= T_0 \exp \left(\int_0^t \frac{\mathbb{E}(\tau \mid X_0 \sim \mu(s))}{T(s)} ds - t \right) \\ &\leq T_0 \exp \left(\mathbb{E}(\bar{\tau}) \int_0^t \frac{ds}{T(s)} - t \right) \\ &\leq T_0 \exp(\mathbb{E}(\bar{\tau})\Gamma(\infty)) \in (0, \infty). \end{aligned}$$

This bound would imply that there exists $\delta > 0$ so that $1/T(t) > \delta$, obtaining a contradiction to the assumption that $\Gamma(\infty) < \infty$. Thus, we must have $\Gamma(\infty) = \infty$. \square

Lemma 6.6. *Any solution to the reduced ODE in (6.4) (regardless of $v(0) \in H$) converges to the quasi-stationary distribution μ as $t \rightarrow \infty$.*

Proof. By applying the inner product with the vector e we can see that the entire trajectory stays in H . We let $R = (I - Q)^{-1} = \sum_{k=0}^{\infty} Q^k$ which has only nonnegative entries (actually, the entries are strictly positive if Q is irreducible). By Duhamel’s principle all the solutions $v(\cdot)$ can be represented by

$$v(t)^\top = v(0)^\top \exp\left(Rt - \int_0^t (v(s)^\top R e) ds\right). \tag{6.5}$$

Because R has nonnegative entries it follows that $\exp(Rt)$ has nonnegative entries and, thus, $v(t) \geq 0$ (in fact the entries are strictly positive if Q is irreducible). Rearranging (6.5), we arrive at

$$v(t)^\top \exp\left(\int_0^t (v(s)^\top R e) ds - \lambda_R t\right) = v(0)^\top \exp(Rt - \lambda_R t), \tag{6.6}$$

where $\lambda_R > 0$ is the principal eigenvalue of R . The matrix $\exp(R)/\exp(\lambda_R)$ is a matrix with strictly positive entries and it has spectral radius equal to 1. By the Perron–Frobenius theorem (see [17]) it follows that there exists a strictly positive vector w such that $v(0)^\top \exp(Rn - \lambda_R n) \rightarrow w^\top$ as $n \rightarrow \infty$, where w is a principal eigenvector of $R - \lambda_R I$. The convergence holds also along real numbers t (not only natural numbers n) by virtue of a continuity argument noting that $(R - \lambda_R I)/m$ has the same eigenvectors regardless of the value of $m > 0$.

Now, take the inner product with e in both sides of (6.6) to obtain (because $v(t) \in H$)

$$\exp\left(\int_0^t (v(s)^\top R e) ds - \lambda_R t\right) \rightarrow \gamma := w^\top e \in (0, \infty).$$

Finally, write (6.5) as

$$v(t)^\top = v(0)^\top \exp(Rt - \lambda_R t) \exp\left(-\int_0^t (v(s)^\top R e) ds + \lambda_R t\right) \rightarrow \frac{w^\top}{\gamma}.$$

The lemma follows by noting that the Perron–Frobenius eigenvectors of R and Q are identical, so we see that $v(t)$ converges to the quasi-stationary distribution as $t \rightarrow \infty$. \square

Now we are ready to conclude the consistency portion of Theorem 3.1 by invoking [20, Theorem 5.2.1] together with the following proposition.

Proposition 6.1. *Any subsequence solution (obtained as in [20, Theorem 5.2.1]) of the system (6.3) satisfies $\mu_0 \in H$ and $T_0 \geq 1$ and it is such that $\mu(t) \rightarrow \mu_*$ and $T(t) \rightarrow (1 - \lambda_*)^{-1} = \mathbb{E}(\tau \mid X_0 \sim \mu_*)$ as $t \rightarrow \infty$. The sequence $\{(\mu_n^\top, T_n)\}$ stays in a compact set of the attractor domain $H \times [1, \infty)$ eventually. Therefore, $\mu_n \rightarrow \mu_*$ and $T_n \rightarrow (1 - \lambda_*)^{-1}$ with probability 1.*

Proof. We have $\mu_n \in H$ and $T_n \geq 1$ because T_n is the average of the $\tau^{(n)}$ which are greater than 1, so the subsequence procedure in [20, Theorem 5.2.1] produces trajectories that lie in $H \times [1, \infty)$ and which are solutions to (6.3). Now, we have noted that $\Gamma(\cdot)$ is nonnegative and strictly increasing, according to Lemma 6.5 and, thus, Lemma 6.6 implies that $v(t) = \mu(\Gamma^{-1}(t)) \rightarrow \mu_*$. Moreover, Lemma 6.5 indicates that $\Gamma^{-1}(t) \rightarrow \infty$; therefore, we have $\mu(t) \rightarrow \mu_*$ as $t \rightarrow \infty$.

Note that

$$T(t) = e^{-t} \int_0^t \mathbb{E}(\tau \mid X_0 \sim \mu(s)) e^s ds + e^{-t} T_0, \quad \mathbb{E}(\tau \mid X_0 \sim \mu(t)) \rightarrow \mathbb{E}(\tau \mid X_0 \sim \mu_*).$$

So, we can use l'Hôpital's rule to conclude that

$$\lim_{t \rightarrow \infty} T(t) = \mathbb{E}(\tau \mid X_0 \sim \mu(t)) = \mathbb{E}(\tau \mid X_0 \sim \mu_*) = (1 - \lambda_*)^{-1}.$$

The fact that $\{(\mu_n^\top, T_n)\}$ stays in a compact set of $H \times [1, \infty)$ follows from Lemma 6.1(iv). □

6.2. Proof of Theorem 1: CLT

In order to prove the CLT portion of Theorem 3.1, we shall invoke [20, Theorem 10.2.1]; this requires verifying the following conditions.

- (C.4) The sequence $\{W_n \mathbf{1}\{\|\theta_n - \theta_*\| \leq \delta\}\}$ is uniformly integrable. (This follows immediately because due to Lemma 6.4 we have the fact that W_n is L_2 -bounded.)
- (C.5) That θ_* is an isolated stable point of the ODE. (This follows from the Perron–Frobenius theorem and the analysis in Proposition 6.1.)
- (C.6) The expansion $g_n(\theta) = g_n(\theta_*) + (Dg_n)(\theta_*)(\theta - \theta_*) + o(\|\theta - \theta_*\|)$ holds uniformly in n . (This estimate will be elaborated in Lemma 6.7.)
- (C.7) We must have $\lim_{n,m \rightarrow \infty} m^{-1/2} \sum_{k=n}^{n+mt-1} (Dg_k)(\theta_*) = 0$, uniformly over t in compact sets. (See Lemma 6.9.)
- (C.8) There exists a matrix A such that $\lim_{n,m} \sum_{k=n}^{n+m-1} ((Dg_k)(\theta_*) - A) = 0$. (Let $A = (Dg)(\theta_*)$, then this condition will hold by Lemma 6.10 which shows that $(Dg_n)(\theta_*) \rightarrow (Dg)(\theta_*)$.)
- (C.9) The matrix A must also be such that $A + I/2$ is Hurwitz (i.e. all its eigenvalues have a negative real part). (This corresponds precisely to (3.5) and it will be established in Proposition 6.2.)
- (C.10) The sequence $\{(\theta_n - \theta_*)/\varepsilon_n^{1/2}\}$ is tight. (See Lemma 6.12.)
- (C.11) Define $\delta M_n = W_n - g_n(\theta_n)$, then there exists $p > 0$ such that $\sup_n \mathbb{E} \|\delta M_n\|_\infty^{2+p} < \infty$, and a nonnegative matrix Σ such that $\mathbb{E}(\delta M_n (\delta M_n)^\top \mid \mathcal{F}_n) \xrightarrow{p} \Sigma$ as $n \rightarrow \infty$. (This is established in Lemma 6.13.)

6.2.1. Analysis of the Jacobian: (C.6)–(C.8).

Lemma 6.7. *We have*

$$\lim_{\|\theta - \theta_*\|_\infty \rightarrow 0} \sup_{n \geq 1} \left| \frac{g_n(\theta) - g_n(\theta_*) - (Dg_n)(\theta_*)(\theta - \theta_*)}{\|\theta - \theta_*\|_\infty} \right| = 0.$$

Proof. We consider the analysis only for $f_n(\mu^\top, T)$ because h is a simpler quantity and does not depend on n . The analysis follows as an application of the representation derived in Lemma 6.2. It is easy to justify the interchange of differentiation and integration in the representation given in Lemma 6.2 because the integrand consists of products of a second degree polynomial in μ , the exponential factor $\exp(-uT)$ on the region of interest which is $T \geq 1$, and the term including $v(x, s) \in (0, 1]$. Thus, the second derivatives of f_n will be bounded uniformly in n around a neighborhood of the stationary point θ_* . □

Next we turn to (C.7), but first we have an auxiliary result.

Lemma 6.8. *It follows that $g_n(\theta_*) = O(1/n)$.*

Proof. Note that $h(\theta_*) = 0$, so we only focus on $f_n(\theta_*)$. On the other hand, we observed in the proof of Lemma 6.3 that $|f_n(\mu^\top, T)(x) - f(\mu^\top, T)(x)| \leq \mathbb{E}(\tau^2)/(n + 1)$, but (with $R = (I - Q)^{-1}$)

$$f(\mu_*^\top, T_*) = \frac{1}{T_*} (\mu_*^\top R - (\mu_*^\top R e) \mu_*^\top)^\top = \frac{1}{T_*} \left(\frac{1}{1 - \lambda} \mu_*^\top - \frac{1}{1 - \lambda} \mu_*^\top \right)^\top = 0.$$

Hence, Lemma 6.8 follows. □

Lemma 6.9. *We have $\lim_{n,m \rightarrow \infty} m^{-1/2} \sum_{k=n}^{n+mt-1} (Dg_k)(\theta_*) = 0$, uniformly over compact sets in t .*

Proof. Using Lemma 6.8, it follows that there exists $c \in (0, \infty)$ (independent of n) such that

$$m^{-1/2} \sum_{k=n}^{n+mt-1} \|(Dg_k)(\theta_*)\|_\infty \leq cm^{-1/2} \log \left(1 + \frac{mt}{n} \right).$$

Let $u = mt/n$, then

$$cm^{-1/2} \log \left(1 + \frac{mt}{n} \right) \leq ct^{1/2} n^{-1/2} \sup_{u \geq 0} u^{-1/2} \log(1 + u).$$

Since $\sup_{u \geq 0} u^{-1/2} \log(1 + u)/u^{1/2} < \infty$, we can send $n \rightarrow \infty$ to conclude the statement of Lemma 6.9. □

Lemma 6.10. *We have $\lim_{n \rightarrow \infty} (Dg_n)(\theta_*) = (Dg)(\theta_*)$.*

Proof. Once again, it suffices to concentrate on f_n . Because of Lemma 6.7, we know that $f_n(\theta) = f_n(\theta_*) + (Df_n)(\theta_*)(\theta - \theta_*) + o(\theta - \theta_*)$. Taking the limit as $n \rightarrow \infty$, we arrive at $f(\theta) = \lim_{n \rightarrow \infty} (Df_n)(\theta_*)(\theta - \theta_*) + o(\theta - \theta_*)$. Expanding the left-hand side, we have $f(\theta) = (Df)(\theta_*)(\theta - \theta_*) + o(\theta - \theta_*)$. Matching these terms and noting that $\theta - \theta_*$ can have any direction, we conclude the result of the lemma. □

6.2.2. *The Hurwitz property: (C.9).*

Proposition 6.2. *Let $A = (Dg)(\theta_*)$, then $A + I/2$ is Hurwitz assuming that the eigenvalues of Q satisfy (3.5).*

Proof. Recall that $g(\mu^\top, T) = (f^\top(\mu^\top, T), h(\mu^\top, T))$ and (6.1) and (6.2). Letting $B = R^\top = (I - Q^\top)^{-1}$, it follows that the Jacobians are given by (using D_μ and D_T to denote the derivatives with respect to μ and T , respectively),

$$\begin{aligned} D_\mu f(\mu^\top, T) &= \frac{1}{T} [B - (\mu^\top B e)I - \mu e^\top B], \\ D_T f(\mu^\top, T) &= -\frac{1}{T^2} [B\mu - (\mu^\top B e)\mu], \\ D_\mu h(\mu^\top, T) &= e^\top B, \quad D_T h(\mu^\top, T) = -1. \end{aligned}$$

We consider the stationary point and note that $\mu_*^\top B e = T_*$. Then, define

$$J := D_\mu f(\mu_*^\top, T_*) = \frac{1}{T_*} [B - T_* I - \mu_* e^\top B]. \tag{6.7}$$

Also, note that $D_T f(\mu_*^\top, T_*) = 0$. We now establish a one-to-one correspondence between the eigenvectors of J and the eigenvectors of B . The overall Jacobian in block form takes the form

$$A = \begin{bmatrix} J & 0 \\ e^\top B & -1 \end{bmatrix}.$$

This matrix has the same eigenvalues as J , with the addition of the eigenvalue -1 which has no effect on the Hurwitz property. Hence, we only need to ensure that $J + I/2$ is Hurwitz.

Let y be any vector such that $Jy = \lambda_J y$ and such that y is linearly independent of μ_* . Note that $Jy = \lambda_J y$ is equivalent to

$$By = T_* \lambda_J y + T_* y + (e^\top B y) \mu_*,$$

and, therefore, if we let $x = y + r \mu_*$, for some r to be characterized momentarily, we have

$$Bx = By + r T_* \mu_* = T_* (\lambda_J + 1) y + (e^\top B y + r T_*) \mu_*.$$

So, the value of r that would make x an eigenvector of B is such that $r T_* \lambda_J = e^\top B y$. Since $T_* > 0$ the existence of r is guaranteed if $\lambda_J \neq 0$ and the corresponding eigenvalue for B would be $\lambda_B = T_* (1 + \lambda_J)$. On the other hand, if y is a multiple of μ_* , its eigenvalue for the matrix J equals -1 (the eigenvalue for B is, of course, T_*). In Lemma 6.11 below we will argue that λ_J cannot be 0, so the argument just given shows that every eigenvector of J is an eigenvector of B .

Conversely, given any vector z , such that $Bz = \lambda_B z$, we can define $u = z + r \mu_*$. If we choose $r = (T_* + \lambda_B e^\top z) / (T_* - \lambda_B)$ then we conclude that $Ju = (\lambda_B / T_* - 1)u$. This selection of r is valid if z is not the principal right eigenvector of B (i.e. in case z is not μ_*). In case we select $z = \mu_*$, then trivially $Jz = -z$.

Consequently, we conclude that there is a one-to-one correspondence between the eigenvectors (and eigenvalues) of J and B , and the relationship between the eigenvalues is

$$\lambda_J = \frac{\lambda_B}{T_*} - 1 \quad \text{for } \lambda_B \neq T_* \text{ or } \lambda_J \neq 0, \quad \lambda_J = -1 \quad \text{if } \lambda_B = T_*.$$

Therefore, in order to ensure that $J + I/2$ is Hurwitz, we must have $\text{Re}(\lambda_B) < T_*/2$ for all $\lambda_B \neq T_*$, which is precisely (3.5). □

We finish the analysis of the Hurwitz condition, with the following result invoked in the previous proof.

Lemma 6.11. *We have $\lambda_J \neq 0$.*

Proof. Assume that y is such that $Jy = 0$. This implies that

$$By = T_* y + \mu_* (e^\top B y) = T_* y + (\mu_* e^\top) B y.$$

Therefore,

$$(I - (\mu_* e^\top)) B y = T_* y.$$

We recognize that $\bar{P} = (I - (\mu_* e^\top))$ is a (nonorthogonal) projection in the sense that $\bar{P}^2 = \bar{P}$. Also we have $\bar{P} \mu_* = 0$ and $e^\top \bar{P} = 0$. This means that T_* is an eigenvalue of $\bar{P} B$, which in turn implies that there exists a left eigenvector x such that $x^\top \bar{P} B = T_* x^\top$, or

$$x^\top \bar{P} = T_* x^\top B^{-1}. \tag{6.8}$$

Now let $x^\top \bar{P} = z^\top$ and consider all possible solutions w such that $w^\top \bar{P} = z^\top$ which must be written as the sum of an element of the null space and a particular solution. Observe, because $\bar{P}^2 = \bar{P}$, that $z^\top \bar{P} = z^\top$ is a particular solution and, therefore, any solution x (i.e. any eigenvector corresponding to T_* for $\bar{P}B$) must take the form $x = c\mathbf{e} + z$ for some constant c . Observe from (6.8), multiplying by μ_* from the right, that

$$0 = T_* x^\top B^{-1} \mu_* = T_* x^\top (I - Q^\top) \mu_* = x^\top \mu_*$$

Therefore, we have $c = 0$ because $x^\top \mu_* = 0$ and $0 = x^\top \mu_* = c + z^\top \mu_* = c + 0$. Consequently, $x^\top = z^\top = T_* x B^{-1}$, which implies that $x^\top B = T_* x^\top$, therefore concluding that x is the principal left eigenvector of B . Consequently, x must have strictly positive entries and, in turn, we must have $x^\top \mu_* > 0$, thus, arriving at a contradiction. So, there is no eigenvalue T_* for the matrix $\bar{P}B$ and, thus, $\lambda_J = 0$ is not possible. \square

6.2.3. *Tightness: (C.10).*

Lemma 6.12. *The sequence $\{(\theta_n - \theta_*)/\varepsilon_n^{1/2}\}$ is tight.*

Proof. We use the local techniques discussed in [20, Section 10.5.2] and apply them as in the proof of [20, Theorem 10.4.1], albeit with some modifications. We shall use the Lyapunov function $V(\theta) = \|\theta - \theta_*\|_2^2$. However, we now have to deal with the gradient of g_n as opposed to the gradient of g as in the proof of [20, Theorem 10.4.1]. We shall control the changes in g_n by expanding around the stationary point. We have

$$\begin{aligned} \mathbb{E}(V(\theta_{n+1}) \mid \mathcal{F}_n) - V(\theta_n) &= 2\varepsilon_n(\theta_n - \theta_*)^\top g_n(\theta_n) + O(\varepsilon_n^2) \\ &= 2\varepsilon_n(\theta_n - \theta_*)^\top g_n(\theta_*) + 2\varepsilon_n(\theta_n - \theta_*)^\top A_n(\theta_n - \theta_*)^\top + \varepsilon_n o(\|\theta_n - \theta_*\|_2^2) + O(\varepsilon_n^2), \end{aligned}$$

where the first equality uses an idea similar to that of the proof of Lemma 6.4 to arrive at the error term $O(\varepsilon_n^2)$, and the second inequality is just an expansion of g_n around θ_* followed by an application of Lemma 6.7. Since $2A_n$ has eigenvalues with negative real part less than -1 (i.e. $A_n + I/2$ is Hurwitz) for large enough n , we conclude that there exists $\delta > 0$ such that, for all sufficiently large n ,

$$(\theta_n - \theta_*)^\top A_n(\theta_n - \theta_*)^\top < -(1 + 2\delta)\|\theta_n - \theta_*\|_2^2.$$

Moreover, because $g_n(\theta_*) = O(1/n)$ due to Lemma 6.8, we have

$$2\varepsilon_n(\theta_n - \theta_*)^\top g_n(\theta_*) \leq O(\varepsilon_n^2(\theta_n - \theta_*))$$

and $o(\|\theta_n - \theta_*\|_2^2) \leq \delta V(\theta_n)$. We then conclude that

$$\mathbb{E}(V(\theta_{n+1}) \mid \mathcal{F}_n) - V(\theta_n) \leq -\varepsilon_n(1 + \delta)V(\theta_n) + O(\varepsilon_n^2).$$

The rest of the proof can now be concluded as in [20, Theorem 10.4.1]. \square

6.2.4. *Quadratic variation of the martingales: (C.11).*

Lemma 6.13. *Let $\delta M_n = W_n - g_n(\theta_n)$, then*

$$\sup_{n \geq 0} \mathbb{E} \|\delta M_n\|_2^4 < \infty. \tag{6.9}$$

Moreover,

$$\mathbb{E}(\delta M_n \delta M_n^\top \mid \mathcal{F}_n) \xrightarrow{\mathbb{P}} \Sigma \text{ for some matrix } \Sigma. \tag{6.10}$$

Proof. We use the notation $\mathbb{E}_n(\cdot)$ for $\mathbb{E}(\cdot \mid \mathcal{F}_n)$,

$$\begin{aligned} \|\delta M_n\|_2^4 &\leq 2(\|Y_n(\theta_n) - \mathbb{E}_n Y_n(\theta_n)\|_2^4 + |Z_n(\theta_n) - \mathbb{E}_n Z_n(\theta_n)|^4) \\ &\leq 16(\|Y_n(\theta_n)\|_2^4 + \|\mathbb{E}_n Y_n(\theta_n)\|_2^4 + |Z_n(\theta_n)|^4 + |\mathbb{E}_n Z_n(\theta_n)|^4). \end{aligned}$$

An argument similar to Lemma 6.3 yields that $\|Y_n(\theta_n)\|_2^4 \leq \bar{\tau}^4$ and $|Z_n(\theta_n)|^4 \leq \bar{\tau}^4$ (stochastically) and, therefore, by Lemma 6.1 we conclude bound (6.9).

To establish (6.10) note that $\delta M_n \delta M_n^\top$ is equal to

$$\begin{bmatrix} (Y_n(\theta_n) - f_n(\theta_n))(Y_n(\theta_n) - f_n(\theta_n))^\top & (Y_n(\theta_n) - f_n(\theta_n))(Z(\theta_n) - h(\theta_n)) \\ (Z(\theta_n) - h(\theta_n))(Y_n(\theta_n) - f_n(\theta_n))^\top & (Z(\theta_n) - h(\theta_n))^2 \end{bmatrix}.$$

By Lemma 6.8, we have $f_n(\theta_n) \rightarrow 0$, and $h(\theta_n) \rightarrow 0$. Note that the distribution of $Z(\theta)$ and $Y(\theta)$ can be written in a way that is continuous in θ (as a mixture of the initial distribution); therefore, $Z(\theta_n) \xrightarrow{D} Z(\theta_*)$ and $Y(\theta_n) \xrightarrow{D} Y(\theta_*)$; consequently, we also have $Y_n(\theta_n) \xrightarrow{D} Y(\theta_*)$. We observe that each entry of the matrix is stochastically dominated by $2\bar{\tau}$ and, thus, we can apply Lemma 6.1 to conclude uniform integrability (since $\mathbb{E}(\bar{\tau}^2) < \infty$), thereby obtaining

$$\mathbb{E}_n(\delta M_n \delta M_n^\top) \rightarrow \Sigma := \mathbb{E}([Y(\theta_*)^\top \quad Z(\theta_*)]^\top [Y(\theta_*)^\top \quad Z(\theta_*)]).$$

In turn, due to the equality in [20, p. 332], the asymptotic covariance equals

$$V_0 = \int_0^\infty \exp\left(\left(J + \frac{I}{2}\right)t\right) \Sigma \exp\left(\left(J^\top + \frac{I}{2}\right)t\right) dt, \tag{6.11}$$

concluding the proof. □

7. Proof of Proposition 4.1

Proof of Proposition 4.1. The proof is almost identical to that of Theorem 3.1. In fact, the analysis is somewhat simpler because there is no denominator and so we just need to analyze the reduced ODE in (6.4). The proof of tightness also follows similar steps as the argument given in [20, Theorem 10.4.1], which distinguishes the cases $\varepsilon_n = 1/n$ and the case $\varepsilon_n = n^{-\alpha}$ for $\alpha \in (\frac{1}{2}, 1)$ as we do here.

Now, recall that J is the Jacobian of the vector field obtained in (6.7), evaluated at the unique stability point μ_* . We need to ensure that J is Hurwitz (i.e. all the eigenvalues have a strictly negative real part). This is in contrast to requiring that $J + I/2$ is Hurwitz — which is a stronger condition because then one needs that all the eigenvalues have a real part less than $-\frac{1}{2}$, which leads to (3.5). Instead, requiring that J be Hurwitz is equivalent to the condition that $\text{Re}(\bar{\lambda}) < \lambda_*$ for all nonprincipal eigenvalue $\bar{\lambda}$ of the matrix Q , which is automatic by the Perron–Frobenius theorem [17]. Hence, we conclude the result by invoking [20, Theorem 10.2.1]. The asymptotic covariance matrix in this case takes the form

$$V_1 = \int_0^\infty \exp(Jt) \Sigma_0 \exp(J^\top t) dt, \tag{7.1}$$

where $\Sigma_0 = \mathbb{E}(\bar{Y}(\mu_*^\top) \bar{Y}(\mu_*^\top)^\top)$, with $\bar{Y}(\mu^\top) = \sum_{k=0}^{\tau-1} (\mathbf{1}\{X_k = \cdot \mid X_0 \sim \mu\} - \mu(\cdot))$. □

We are now ready to discuss the proof of Theorem 4.1

Proof of Theorem 4.1. We shall verify the conditions in [24, Theorem 2]. First, define $\bar{f}(\mu) = f(\mu^\top, 1)$ and recall that $f(\cdot)$, introduced in (6.1), coincides with $f(\mu^\top, 1) = \mathbb{E}\bar{Y}(\mu^\top)$. We must first verify that

(C.12) there exists a function L a (globally) Lipschitz continuous function $V(\cdot)$ such that $L(\mu_*) = 0$, and for some positive definite matrix G ,

$$DL(\mu)G\bar{f}(\mu^\top) < 0 \quad \text{for } \mu \neq \mu_*, \tag{7.2}$$

there exists $\varepsilon, \delta > 0$ such that

$$DL(\mu)G\bar{f}(\mu^\top) \leq -\delta V(\mu) \tag{7.3}$$

if $\|\mu - \mu_*\| \leq \varepsilon$, and $L(\mu - \mu_*) \geq \delta\|\mu - \mu_*\|_2^2$ for some $\delta > 0$.

This condition is satisfied if we construct $L(\cdot)$ by noting that μ_* is the unique root of $\bar{f}(\mu_*) = 0$ and we have established in Lemma 6.2 that $\bar{J} := (D\bar{f})(\mu_*) = JT_*$ is Hurwitz (in fact $J + I/2 = \bar{J}/T_* + I/2$ is Hurwitz); therefore,

$$\bar{f}(\mu) = \bar{J}(\mu - \mu_*) + O(\|\mu - \mu_*\|_2^2). \tag{7.4}$$

Now, it is standard in stability of dynamical systems that given a Hurwitz matrix \bar{J} and a given symmetric positive definite matrix G (which we might take as $G = I$, the identity, here), there is a unique symmetric positive definite matrix K such that $\bar{J}^\top K + K\bar{J} = -G = -I$, and there is a positive constant $\gamma > 0$ such that $I - \delta K$ is symmetric and positive definite. We first set $\bar{L}(\mu) = (\mu - \mu_*)^\top K(\mu - \mu_*)$ so that

$$\begin{aligned} D\bar{L}(\mu)G\bar{f}(\mu^\top) &= 2(\mu - \mu_*)^\top K\bar{J}(\mu - \mu_*) + O(\|\mu - \mu_*\|_2^3) \\ &= (\mu - \mu_*)^\top (\bar{J}^\top K + K\bar{J})(\mu - \mu_*) + O(\|\mu - \mu_*\|_2^3) \\ &= -\|\mu - \mu_*\|_2^2 + O(\|\mu - \mu_*\|_2^3). \end{aligned}$$

From this bound we obtain the existence of L satisfying (7.2) and (7.3) by modifying \bar{L} outside a neighborhood of μ_* inside the compact set H .

The second condition [24, Theorem 2] follows directly from (7.4) and the fact that \bar{J} is Hurwitz. There are two more conditions to check for the application of [24, Theorem 2], the fourth condition is trivially satisfied for $\varepsilon_n = n^{-\alpha}$ with $\alpha \in (\frac{1}{2}, 1)$ and the third condition involves the martingale difference process and its quadratic variation, the verification is completely parallel to that of Proposition 6.2. □

8. Proof of Proposition 4.2

In this section we provide the proof of consistency for the estimator of

$$\bar{V}_1 = J^{-1}\Sigma_0J^{-1}. \tag{8.1}$$

Equation (8.1) follows from [20, Theorem 11.1.1].

Proof of Proposition 4.2. We write

$$\begin{aligned} \varepsilon_n^{-1}(\bar{\mu}_n - v_n)(\bar{\mu}_n - v_n)^\top &= \varepsilon_n^{-1}(\bar{\mu}_n - \mu_*)(\bar{\mu}_n - \mu_*)^\top + \varepsilon_n^{-1}(v_n - \mu_*)(v_n - \mu_*)^\top \\ &\quad - \varepsilon_n^{-1}(\bar{\mu}_n - \mu_*)(v_n - \mu_*)^\top - \varepsilon_n^{-1}(\bar{\mu}_n - \mu_*)(v_n - \mu_*)^\top. \end{aligned} \tag{8.2}$$

Note that $\varepsilon_n^{-1/2}(\bar{\mu}_n - \mu_*)(\varepsilon_n n)^{-1/2}n^{1/2}(v_n - \mu_*)^\top \xrightarrow{\mathbb{P}} 0$ as $n \rightarrow \infty$ because of Proposition 4.1 and Theorem 4.1, since $\varepsilon_n n \rightarrow \infty$ as $n \rightarrow \infty$. A similar argument applies to all the terms in (8.2) involving $(v_n - \mu_*)$; so it suffices to study the limit of

$$N_n^{-1} \sum_{n_k=1}^{N_n} \varepsilon_{n_k}^{-1} (\bar{\mu}_{n_k} - \mu_*)(\bar{\mu}_{n_k} - \mu_*)^\top \quad \text{as } n \rightarrow \infty.$$

The rest of the calculation is similar to the analysis of the asymptotic covariance in [20, Theorem 11.3.1]. The idea is that the sequence $\{(\bar{\mu}_{n_k} - \mu_*)\varepsilon_{n_k}^{-1/2}\}$ is weakly dependent and each term is asymptotically normal (as $k \rightarrow \infty$) with variance V_1 . \square

Acknowledgements

We acknowledge support from the National Science Foundation (grant numbers CMMI-1069064 and DMS-1320550). We also thank the anonymous referee for the careful review of our paper.

References

- [1] ALDOUS, D., FLANNERY, B. AND PALACIOS, J. L. (1988). Two applications of urn processes to the fringe analysis of search trees and the simulation of quasi-stationary distributions of Markov chains. *Prob. Eng. Inf. Sci.* **2**, 293–307.
- [2] ATHREYA, K. B. AND KARLIN, S. (1968). Embedding of urn schemes into continuous time Markov branching processes and related limit theorems. *Ann. Math. Statist.* **39**, 1801–1817.
- [3] BENAÏM, M. AND CLOEZ, B. (2015). A stochastic approximation approach to quasi-stationary distributions on finite spaces. *Electron. Commun. Probab.* **20**, 14 pp.
- [4] BLANCHET, J., GLYNN, P. AND ZHENG, S. (2013). Empirical analysis of a stochastic approximation approach for computing quasi-stationary distributions. In *EVOLVE*, Springer, Berlin, 19–37.
- [5] BOYD, S. AND VANDENBERGHE, L. (2004). *Convex Optimization*. Cambridge University Press.
- [6] BURZDY, K., HOEYST, R. AND MARCH, P. (2000). A Fleming–Viot particle representation of the Dirichlet Laplacian. *Commun. Math. Phys.* **214**, 679–703.
- [7] DARROCH, J. N. AND SENETA, E. (1965). On quasi-stationary distributions in absorbing discrete-time finite Markov chains. *J. Appl. Prob.* **2**, 88–100.
- [8] DARROCH, J. N. AND SENETA, E. (1967). On quasi-stationary distributions in absorbing continuous-time finite Markov chains. *J. Appl. Prob.* **4**, 192–196.
- [9] DE OLIVEIRA, M. M. AND DICKMAN, R. (2005). How to simulate the quasistationary state. *Phys. Rev. E* **71**, 016129.
- [10] DE OLIVEIRA, M. M. AND DICKMAN, R. (2006). Quasi-stationary simulation: the subcritical contact process. *Brazilian J. Phys.* **36**, 685–689.
- [11] DEL MORAL, P. AND MICLO, L. (2006). Self-interacting Markov chains. *Stoch. Anal. Appl.* **24**, 615–660.
- [12] DICKMAN, R. AND VIDIGAL, R. (2002). Quasi-stationary distributions for stochastic processes with an absorbing state. *J. Phys. A* **35**, 1147–1166.
- [13] DIMOV, I. T., KARAIANOVA, A. N. AND YORDANOVA, P. I. (1998). Monte Carlo algorithms for calculating eigenvalues. In *Monte Carlo and Quasi-Monte Carlo Methods 1996* (Salzburg; Lecture Notes Statist. **127**), Springer, New York, pp. 205–220.
- [14] FERRARI, P. A. AND MARIĆ, N. (2007). Quasistationary distributions and Fleming–Viot processes in countable spaces. *Electron. J. Probab.* **12**, 684–702.
- [15] GOLUB, G. H. AND VAN LOAN, C. F. (1996). *Matrix Computations*, 3rd edn. Johns Hopkins University Press, Baltimore, MD.
- [16] GROISMAN, P. AND JONCKHEERE, M. (2013). Simulation of quasi-stationary distributions on countable spaces. *Markov. Process. Relat. Fields.* **19**, 521–542.
- [17] KARLIN, S. AND TAYLOR, H. M. (1975). *A First Course in Stochastic Processes*, 2nd edn. Academic Press, New York.
- [18] KRASULINA, T. P. (2013). The method of stochastic approximation for the determination of the least eigenvalue of a symmetrical matrix. *USSR Comput. Math. Math. Phys.* **9**, 189–195.

- [19] KRASULINA, T. P. (1970). Method of stochastic approximation in the determination of the largest eigenvalue of the mathematical expectation of random matrices. *Automat. Rem. Contr.* **2**, 215–221.
- [20] KUSHNER, H. J. AND YIN, G. (2003). *Stochastic Approximation and Recursive Algorithms and Applications*, 2nd edn. Springer, New York.
- [21] MÉLÉARD, S. AND VILLEMONAIS, D. (2012). Quasi-stationary distributions and population processes. *Prob. Surveys* **9**, 340–410.
- [22] OJA, E. AND KARHUNEN, J. (1985). On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. *J. Math. Anal. Appl.* **106**, 69–84.
- [23] PERMANTLE, R. (2007). A survey of random processes with reinforcement. *Prob. Surveys* **4**, 1–79.
- [24] POLYAK, B. T. AND JUDITSKY, A. B. (1992). Acceleration of stochastic approximation by averaging. *SIAM J. Control Optimization* **30**, 838–855.
- [25] ROBBINS, H. AND MONRO, S. (1951). A stochastic approximation method. *Ann. Math. Statist.* **22**, 400–407.
- [26] ZHENG, S. (2014). Stochastic approximation algorithms in the estimation of quasi-stationary distribution of finite and general state space Markov chains. Doctoral Thesis, Columbia University.