

# Oja's Algorithm for Graph Clustering and Markov Spectral Decomposition

V. Borkar<sup>\*</sup>

School of Technology and Computer Science  
Tata Institute of Fundamental Research, Homi  
Bhabha Road, Mumbai 400005, India.  
borkar@tifr.res.in

S.P. Meyn<sup>†</sup>

Department of Electrical and Computer Engg.  
and the Coordinated Sciences Laboratory  
University of Illinois at Urbana-Champaign,  
Urbana, IL 61801, USA.  
meyn@uiuc.edu

## ABSTRACT

Given a positive definite matrix  $M$  and an integer  $N_m \geq 1$ , Oja's subspace algorithm will provide convergent estimates of the first  $N_m$  eigenvalues of  $M$  along with the corresponding eigenvectors. It is a common approach to principal component analysis. This paper introduces a normalized stochastic-approximation implementation of Oja's subspace algorithm, as well as new applications to the spectral decomposition of a reversible Markov chain. Stability and convergence are established under conditions far milder than assumed in previous work. Applications to graph clustering and Markov spectral decomposition are surveyed, along with numerical results.

## Keywords

Graph algorithms, Oja's algorithm, stochastic approximation, Markov chains, spectral theory of Markov chains

**2000 AMS Subject Classification:** 05C85, 94C15, 68W20, 62L20, 60J22, 60J10, 37A30, 92B20

## 1. INTRODUCTION

Spectral decomposition is a classical approach to model reduction for systems that are complex due to dimension or randomness. This technique is known as principal component analysis or the Karhunen-Loève decomposition, depending on the context [17, 15, 14]. The same technique has been developed more recently as an alternative to the

<sup>\*</sup>V.B. was supported in part by a J. C. Bose Fellowship of Dept. of Science and Technology, Govt. of India, and a grant from General Motors India Science Lab.

<sup>†</sup>S.P.M. was supported in part by National Science Foundation grants ECS-0523620 and CCF 07-29031. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the National Science Foundation.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

ValueTools 2008, October 21 – 23, 2008, Athens, GREECE.  
Copyright © 2008 ICST ISBN # 978-963-9799-31-8.

min-cut max-flow theorem for network decomposition [23, 24, 19].

Given a symmetric  $N \times N$  matrix  $w$ , its spectral decomposition amounts to the computation of its  $N$  real eigenvalues and corresponding eigenvectors. In the Karhunen-Loève decomposition the matrix  $w$  is a covariance matrix, and the decomposition leads to a representation of a stationary process as a moving-average of white noise. In the graph clustering problem the elements of this matrix represent positive edge weights:  $w_{ij} = w_{ji}$  is the weight of the link connecting nodes  $i$  and  $j$ . The first decomposition of a connected graph is obtained by computation of the eigenvector corresponding to the second eigenvalue. It can be shown that the eigenvector possesses positive and negative entries, and this sign structure is used to define a generalized network cut in [23, 24, 19].

Oja's subspace algorithm is one approach to computation of the leading eigenvalues and eigenvectors of the matrix  $w$  [20, 7]. Fix an integer  $N_m \leq N$ , and let  $m(t)$  denote an  $N \times N_m$  matrix whose columns are intended to approximate an  $N_m$ -dimensional eigenspace corresponding to the  $N_m$  largest of the  $N$  eigenvalues of  $w$ . A deterministic version of Oja's algorithm is expressed as the polynomial differential equation,

$$\frac{d}{dt}m(t) = [I - m(t)m^T(t)]wm(t) \quad (1)$$

where  $m(0)$  is given as initial condition. If the matrix  $w$  is positive definite then the analysis of [7] establishes convergence of  $m$  for almost every initial condition.

This paper introduces a normalized implementation of Oja's algorithm, that is also a multi-dimensional generalization of the one-dimensional algorithm of Krasulina [16]. Stability and convergence of the normalized algorithm are established under conditions far milder than assumed in previous work. Applications to graph clustering are surveyed, as well as new applications to the spectral decomposition of a reversible Markov chain.

In the following section we introduce the stochastic approximation algorithm, and present the main result establishing convergence of the algorithm. Applications to spectral graph theory are surveyed in Section 3, and Section 4 contains extensions of the algorithm to compute the spectrum of a reversible Markov chain. Examples are contained in Section 5, and conclusions may be found in Section 6

## 2. STOCHASTIC APPROXIMATION AND OJA'S ALGORITHM

Oja's 1985 paper [21] introduces a stochastic approximation algorithm based on the o.d.e. (1). Suppose that  $\mathbf{X}$  is an  $\mathbb{R}^n$ -valued stationary process with covariance matrix  $w = \mathbb{E}[X(t)X(t)^T]$ . We can express Oja's stochastic approximation algorithm as the matrix recursion,

$$M(n+1) - M(n) = a(n)[I - M(n)M^T(n)]\widehat{W}(n)M(n) \quad (2)$$

where  $\widehat{W}(n) = X(n)X^T(n)$ . Convergence is established by applying current stochastic approximation techniques. However, these techniques require Lipschitz continuity of the right hand side of the recursion in the variable  $M(n)$ , which is violated in this case. This issue is addressed by imposing additional conditions on  $\mathbf{X}$ .

The lack of Lipschitz continuity presents problems even in deterministic approximations of (1) in discrete time. One such algorithm is introduced in [25] through sampling the o.d.e. to obtain,

$$m(n+1) - m(n) = a(n)[I - m(n)m^T(n)]wm(n) \quad (3)$$

While convergence is established for the deterministic algorithm, the proof is complex. Complexity is due in large part to the cubic nonlinearity seen here just as in the stochastic approximation algorithm.

To enforce the Lipschitz continuity assumption and thereby place the algorithm within the framework of [3, 2, 1] we introduce a normalization. The normalized o.d.e. is given by

$$\begin{aligned} \frac{d}{dt}m(t) &= a(t)[I - m(t)m^T(t)]wm(t) \\ a(t) &= (1 + \text{trace}(m(t)m(t)^T))^{-1}. \end{aligned} \quad (4)$$

The right hand side of the differential equation is Lipschitz in the variable  $m(t)$ . Solutions to this differential equation are simply time-scaled versions of the solutions to (1). In particular, from each initial condition the set of limit points are identical.

The stochastic approximation algorithm considered in this paper is again of the form (2) in which the gain sequence is modified through the choice of a non-negative gain sequence  $\{b(n) : n \geq 0\}$ :

$$a(n) = b(n)(1 + \text{trace}(m(n)m(n)^T))^{-1}. \quad (5)$$

It is assumed throughout that the following assumptions hold for the gain sequence  $\mathbf{b}$ : It is non-negative, with

$$\sum_{n=0}^{\infty} b(n) = \infty, \quad \sum_{n=0}^{\infty} b(n)^2 < \infty, \quad \sup_{n \geq 0} \left( \frac{\sum_{k \geq n} b(k)^2}{b(n)} \right) < \infty \quad (6)$$

An example is  $b(n) = (1+n)^{-1}$ ,  $n \geq 0$ .

Under these conditions the algorithm is stable. To guarantee consistency we modify the algorithm slightly through the introduction of white noise,

$$M(n+1) - M(n) = a(n)([I - M(n)M^T(n)]\widehat{W}(n)M(n) + \xi(n+1)) \quad (7)$$

where  $\xi$  is an i.i.d.  $N(0, I)$  sequence. Proposition 2.1 states that this recursion shares the best possible convergence properties observed in the o.d.e. (1). While the deterministic algorithm can become trapped in an arbitrary eigenspace of  $w$ , the stochastic algorithm (7) is strongly consistent from each initial condition.

PROPOSITION 2.1. Consider the algorithms (2) or (7), where  $\mathbf{a}$  is given in (5), and with  $\mathbf{b}$  satisfying the conditions in (6). Suppose that the process  $\mathbf{X}$  is i.i.d., with covariance  $w > 0$ , and that it is independent of the i.i.d.  $N(0, I)$  sequence  $\xi$ .

Then, the following conclusions hold for each initial  $M(0)$ :

- (i) Stability: For either of the algorithms (2) or (7),

$$\limsup_{n \rightarrow \infty} \|M(n)\| < \infty \quad a.s.$$

- (ii) Convergence: For the algorithm (7), with probability one, any limit point  $M(\infty)$  of the sequence of matrices  $\{M(n)\}$  has columns that lie in the eigenspace spanned by the first  $m$  eigenvalues of  $w$ .

PROOF. First we establish that the solutions to either stochastic approximation recursion are bounded a.s. by applying Theorem 7 of [2, Ch. 3] (see also [3]). This result constructs an "o.d.e. at infinity" that approximates the behavior of the recursion for large initial conditions. Based on the recursion (2) or (7) we obtain the o.d.e.,

$$\frac{d}{dt}m^\infty(t) = - \left[ \frac{m^\infty(t)m^{\infty T}(t)}{\text{trace}(m^\infty(t)m^{\infty T}(t))} \right]wm^\infty(t) \quad (8)$$

where  $m^\infty(0) \in \mathbb{R}^{N \times N_m}$  is given as initial condition. Define the real valued function  $V : \mathbb{R}^{N \times N_m} \rightarrow \mathbb{R}_+$  as the quadratic,

$$V(m) := \text{trace}(m^T w m), \quad m \in \mathbb{R}^{N \times N_m}.$$

Under the positivity assumption on  $w$  this function vanishes only when  $m$  is identically zero. This property combined with the following drift condition implies that  $V$  serves as a Lyapunov function,

$$\frac{d}{dt}V(m^\infty(t)) = -2 \left[ \frac{\text{trace}([m^{\infty T}(t)wm^\infty(t)]^2)}{\text{trace}(m^\infty(t)m^{\infty T}(t))} \right] < 0, \quad m^\infty(t) \neq 0.$$

It follows that the origin is the unique asymptotically stable equilibrium for (8). Theorem 7 of [2, Ch. 3] completes the proof of (i).

We now restrict to the algorithm (7). From the analysis of [7] it follows that the eigenspace spanned by the first  $m$  eigenvectors of  $w$  is a locally stable invariant set for (4), whereas the remaining eigenvectors are unstable invariant sets. The introduction of the i.i.d. process  $\xi$  combined with the assumptions on the gain sequence ensure that the results of section 4.3 of [2] apply, and the iterates avoid these unstable invariant sets with probability one. In turn, Theorem 19 of [2, Ch. 4] then ensures the desired convergence with probability one.  $\square$

## 3. SPECTRAL GRAPH CLUSTERING

We now show how these methods can be adapted to spectral graph clustering, following [23, 24, 19]. The algorithms described here are variants of stochastic approximation based on the construction of a Markov chain evolving on the nodes of the graph.

Suppose that  $w$  is a symmetric matrix with non-negative entries that defines weights in a graph with adjacency matrix  $A_{ij} = A_{ji} = \mathbf{1}\{w_{ij} > 0\}$ . Throughout this section we impose the following assumptions on the matrix  $w$ :

- (i) Symmetry:  $w = w^T$
- (ii) Probabilistic normalization:  $\sum_{i,j} w_{ij} = 1$ .
- (iii) Irreducibility:  $\sum_{k=1}^{\infty} w_{ij}^k > 0$  for each  $i, j$ , where  $w^k$  denotes the  $k$ -fold matrix product.

The normalization can be assumed without loss of generality by scaling, and irreducibility is equivalent to connectedness of the graph.

Oja's technique is not directly applicable because  $w$  is not necessarily positive definite. One approach to enforce positivity is to add a scaled identity matrix to obtain  $w^{(r)} := w + rI$ . This matrix is positive definite for  $r \geq 0$  sufficiently large. The relationship between the spectrum of  $w$  and  $w^{(r)}$  is obvious, and the eigenvectors coincide. We henceforth assume that this scaling has been performed so that the matrix  $w$  is positive definite.

A stochastic approximation algorithm is obtained by constructing a Markov chain on the state space  $\mathbf{X} := \{1, \dots, N\}$ . Under the normalization assumption, the matrix  $w$  can be interpreted as a probability measure on the product space  $\mathbf{X} \times \mathbf{X}$ . Its common marginal distribution is denoted  $\pi(i) = \sum_j w_{ij}$ , and a transition matrix is defined as the ratio,

$$P(i, j) = \frac{w_{ij}}{\pi(i)}.$$

The detailed balance equations hold,  $\pi(i)P(i, j) = \pi(j)P(j, i)$ ,  $1 \leq i, j \leq N$ , so that  $\pi$  is invariant for  $P$ . The transition matrix is irreducible since the graph is connected, which implies that the invariant measure  $\pi$  is unique. Denote the Markov chain with this transition matrix by  $\mathbf{X} = \{X(n) : n \geq 0\}$ .

In the applications considered in this section we redefine the matrix  $\widehat{W}$  by,

$$\widehat{W}_{ij}(n) = r\mathbf{1}\{i = j\} + \mathbf{1}\{X(n) = i, X(n+1) = j\}, \quad 1 \leq i, j \leq N, \quad (9)$$

so that we obtain  $\mathbb{E}[\widehat{W}(n)] = w^{(r)}$  for each  $n$ . A stochastic approximation algorithm is obtained by applying (2) using this matrix sequence.

If the second eigenvalue of  $P$  is close to unity then the mixing rate of the Markov chain  $\mathbf{X}$  will be slow, and this may adversely affect the convergence rate of (2) (see [9], [18, Ch. 20], and the discussion in Section 4.) In this case the following variant can be used, known as *split sampling* [1]. Let  $\mathbf{X}^1$  denote an i.i.d. sequence with marginal  $\pi$ . Construct a second stochastic process as follows: For each  $n = 1, 2, \dots$  the random variable  $X^2(n)$  is chosen in two stages. First, the value  $j = X^1(n-1)$  is observed. Next, the value of  $X^2(n)$  is chosen according to the distribution  $P(j, \cdot)$ , independent of  $\{X^1(r), X^2(k) : r \in \mathbb{Z}_+, k \leq n-1\}$ . Based on this pair of stochastic processes, the algorithm is then defined by (2) using

$$\widehat{W}_{ij}(n) = r\mathbf{1}\{i = j\} + \mathbf{1}\{X^1(n) = i, X^2(n+1) = j\}, \quad 1 \leq i, j \leq N. \quad (10)$$

Analogous of Proposition 2.1 can be formulated for each of these algorithms. Once again we can establish global consistency only for a perturbed algorithm, as in (7).

## 4. SPECTRAL DECOMPOSITION OF A MARKOV CHAIN

It is known that the rate of convergence to equilibrium for a finite state-space Markov chain is determined by the

second largest eigenvalue of its transition matrix. Based on this observation, there is a large and growing literature on rates of convergence of Markov chains based on spectral theory and related methods.

For a reversible chain with finite state-space each of the eigenvalues is real. Diaconis and Stroock in [9] obtain bounds on the second largest eigenvalue in this setting. A striking conclusion is the following explicit bound on the rate of convergence, as defined by the total-variation norm distance:

$$\|P^n(x, \cdot) - \pi\| \leq \sqrt{\frac{1-\pi(x)}{\pi(x)}} \lambda_*^n \quad (11)$$

where  $\lambda_*$  is the magnitude of the second largest eigenvalue,  $\lambda_*^2 = \max\{\lambda^2 : \lambda \neq 1\}$ , and  $\|\cdot\|$  denotes the total-variation norm. Bounds on the rate of convergence for chains that are not necessarily reversible are obtained in [11], again in the finite state-space case. The bounds are based on spectral theory, but the spectrum of the symmetrized kernel  $P\tilde{P}$  is considered, where  $\tilde{P}$  is the transition kernel for the time-reversed chain.

Just as eigenvectors are used for clustering in graph models, the use of eigenvectors or eigenfunctions (in general state space models) can be used to decompose a Markov model. This is a component of the classical Wentzell-Freidlin theory for model reduction. Much of this work concerns Markov processes that are reversible [22, 8, 4, 12, 5, 6]. Extensions to non-reversible processes appeared for the first time in [13]. The foundation of this paper is the theory of quasi-stationarity, building on the work of [10].

In this section we restrict to the simpler reversible setting. Our goal is to obtain a variant of the stochastic approximation algorithm that will provide estimates of the spectrum of  $P$  rather than a symmetric matrix  $w$ .

Our starting point is a finite state space Markov chain  $\mathbf{X}$  on the state space  $\{1, \dots, N\}$  with transition matrix  $P$ , and invariant measure  $\pi$ . It is assumed that  $\mathbf{X}$  is irreducible and reversible. We write  $\Pi = \text{diag}(\pi)$  and  $w := \Pi P$ . Recall that reversibility implies that  $w$  is symmetric:

$$w = \Pi P = P^T \Pi = w^T \quad (12)$$

Consider the matrix defined by the transformation,

$$w^\circ = \Pi^{-\frac{1}{2}} w \Pi^{-\frac{1}{2}} = \Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}} \quad (13)$$

This matrix remains symmetric. Suppose that  $v^\circ$  is an eigenvector,

$$w^\circ v^\circ = \lambda v^\circ.$$

Then by definition the vector  $v = \Pi^{-\frac{1}{2}} v^\circ$  is an eigenvector of  $P$ .

The Oja o.d.e. to compute the spectrum of  $w^\circ$  is given by,

$$\frac{d}{dt} m(t) = \Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}} m(t) - m(t) m^T(t) \Pi^{\frac{1}{2}} P \Pi^{-\frac{1}{2}} m(t),$$

This is not attractive from the point of view of stochastic approximation. Letting  $g = \Pi^{-\frac{1}{2}} m$  we eliminate the square-root in the o.d.e.,

$$\frac{d}{dt} g(t) = [I - g(t)g^T(t)\Pi]Pg(t) \quad (14)$$

This is very similar to the original Oja o.d.e. using the matrix  $P$ . The point of all this is that this construction implies that  $g(t)$  converges to the maximal eigenspace of  $P$  even though  $P$  is not symmetric.

To ensure convergence of the o.d.e. (1) or its stochastic approximation counterparts we must also assume that  $w^\circ$  is

positive definite. In this setting we are interested in calculating the eigenvalues that are maximal in modulus, so that adding the matrix  $rI$  will not solve the problem of interest. Instead, we work with the two-step transition matrix  $P^2$ . Its eigenvalues are the square of those of  $P$ , so that they are non-negative. We can then replace  $P$  with the transition matrix  $P_\varepsilon := \varepsilon I + (1-\varepsilon)P^2$  where  $\varepsilon \in (0, 1)$  is arbitrary. This matrix has strictly positive eigenvalues, which implies that  $w^\circ$  is positive definite. To simplify notation we assume that  $P$  has been transformed in this way so that it is positive.

A discrete-time implementation of (14) is given by,

$$G(n+1) = G(n) + a(n)[I - G(n)G(n)^T \Pi] P G(n) \quad (15)$$

where  $a$  is redefined by,

$$a(n) = b(n)(1 + \text{trace}(G(n)G(n)^T))^{-1}. \quad (16)$$

A stochastic approximation algorithm is obtained once more by mimicking the deterministic recursion. One algorithm is expressed in matrix form by,

$$G(n+1) - G(n) = a(n)[I - G(n)G(n)^T \hat{\Pi}(n)] \hat{P}(n) G(n) \quad (17)$$

based on the following definitions:  $\hat{\pi}(n)$  is the empirical distribution of  $\mathbf{X}$  based on the first  $n$  samples,  $\hat{\Pi}(n) = \text{diag}(\hat{\pi}(n))$ , and  $[\hat{P}(n)]_{ij} = [\hat{W}(n)]_{ij} / [\hat{\pi}(n)]_i$ , with  $[\hat{W}(n)]_{ij} = \mathbf{1}(X(n) = i, X(n+1) = j)$ .

To obtain a version of the split sampling algorithm we recall the notation introduced in Section 2:  $\mathbf{X}^1$  is i.i.d. with marginal  $\pi$ , and  $\mathbf{X}^2$  is constructed based on the transition matrix  $P$ . We then apply the recursion (17) in which the random quantities are redefined by  $[\hat{W}(n)]_{ij} = \mathbf{1}(X^1(n) = i, X^2(n+1) = j)$ , and  $\hat{\pi}(n)$  is the true marginal  $\pi$ . There is no need to estimate the marginal since it is required in the construction of  $\mathbf{X}^1$ .

In experiments it is found that the multidimensional algorithm in which  $N_m \geq 2$  is slow. To compute the second eigenvector an alternative algorithm is given as follows: The  $N$ -dimensional vector sequence  $\mathbf{G}$  is constructed recursively,

$$G(n+1) - G(n) = a(n)[I - G(n)G(n)^T \hat{\Pi}(n)] (\hat{P}(n) - \varrho \mathbf{1} \otimes \hat{\pi}(n)) G(n) \quad (18)$$

where  $\varrho \in (0, 1)$  is chosen near unity, and we adopt the same conventions as above in the i.i.d. or Markovian versions. The associated o.d.e. is given by

$$\frac{d}{dt} g(t) = [I - g(t)g(t)^T \Pi](P - \varrho \mathbf{1} \otimes \pi) g(t) \quad (19)$$

which is a transformation of the o.d.e. (1) using the positive-definite matrix

$$w^\circ = \Pi^{-\frac{1}{2}}(w - \varrho \pi \otimes \pi) \Pi^{-\frac{1}{2}} = \Pi^{\frac{1}{2}}(P - \varrho \mathbf{1} \otimes \pi) \Pi^{-\frac{1}{2}}$$

Once again, analogs of Proposition 2.1 can be formulated for each of these algorithms, subject to the same caveats stated at the end of section 3.

## 5. EXAMPLES

In most of the applications envisioned we are primarily interested in the sign structure of eigenvectors rather than their values. In such cases we judge the value of an estimate  $\hat{v}$  of an eigenvector  $v$  by the error criterion,

$$\mathcal{E}(\hat{v}) = \min \|\text{sign}(v) - \text{sign}(r\hat{v})\|_1 \quad (20)$$

where the sign is computed pointwise,  $\|\cdot\|_1$  denotes the  $\ell_1$  norm, and the minimum is over  $r = \pm 1$ .

We did not introduce the noise term  $\xi$  in any of our experiments. We found that the algorithms were globally convergent without this modification.

In our first set of examples we apply Oja's subspace algorithm for network decomposition.

### 5.1 Spectral graph clustering

Figure 1 shows the two graphs used in experiments using the deterministic algorithm, and its stochastic counterpart based on i.i.d. observations. The weighting matrix was chosen to coincide with the adjacency matrix, so that each weight was either one or zero.

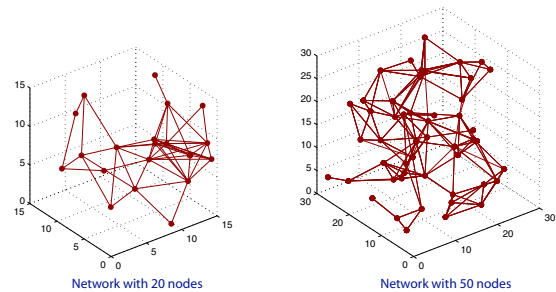


Figure 1: The two networks considered in experiments.

Figure 2 shows results from several experiments using the normalized deterministic algorithm (3) and its stochastic approximation counterpart (2). These plots illustrate the transient behavior of the algorithm for each of the two graphs. In each plot, the vertical axis shows the error  $\mathcal{E}(\hat{v}(n))$  for  $n = 0, 2, \dots$ , where  $\hat{v}$  is the estimate of the second eigenvector of  $w$  obtained from  $m(n)$  defined by (3) (red dashed line), and  $M(n)$  defined in (2) (blue solid). In these experiments the algorithm was run using  $N_m = 2$ . In each case the algorithm was run for 100,000 iterations. For  $N = 20$  the initial 10,000 samples are shown together with the eigenvector approximation obtained after this many samples. Two sets of plots are shown for  $N = 50$ . These results are based on the initial 10,000 iterations, and also the final results after 100,000 iterations.

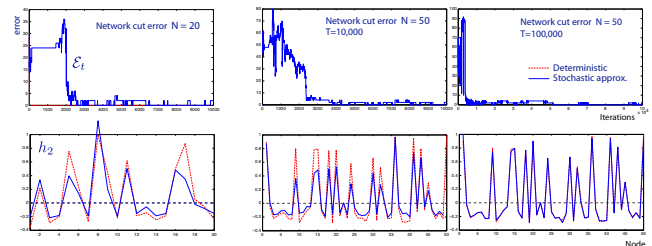
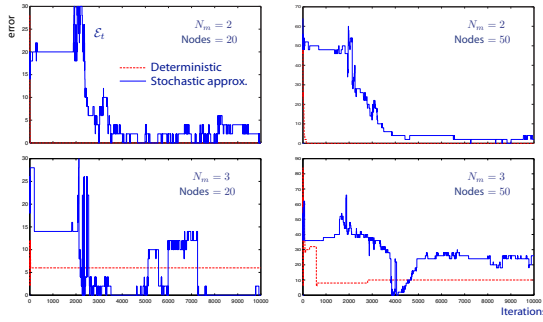


Figure 2: Computation of the first spectral cut for the two networks shown in Figure 1.

For either graph, the sign structure of the eigenvector is identified after approximately 3,000 iterations in the stochastic approximation algorithm. Convergence of the sign structure for the deterministic algorithm was nearly instantaneous.

Note that the slower rate of convergence in the stochastic



**Figure 3: Computation of the first and spectral cut for the 20 and the 50-node network after 10,000 samples. The rate of convergence is slowed significantly when  $N_m$  is increased from 2 to 3.**

algorithm is misleading since the required computation in each iteration is much smaller in the stochastic algorithm.

Convergence slowed considerably when  $N_m$  was increased. Figure 3 shows a comparison of the algorithm using  $N_m = 2$  and  $N_m = 3$ . The convergence rate might be improved by first applying the algorithm with  $N_m = 2$  to find the second eigenvector  $v^2$ , normalized so that its  $L_2$ -norm is unity. Replacing  $w$  by  $w' = w - \lambda_2 v^2 v^{2T}$ , the algorithm can be re-run with  $N_p = 2$  to compute  $v^3$ .

The next examples illustrate computation of the spectrum of a Markov transition matrix.

## 5.2 Markovian spectral clustering

To compute eigenvectors of the transition matrix we applied three approaches, each based on the recursion (18):

- (i) The deterministic algorithm in which  $\hat{\Pi}(n) \equiv \Pi$ ,  $\hat{P}(n) \equiv P$ , and  $\hat{\pi}(n) \equiv \pi$
- (ii) The Markovian algorithm.
- (iii) The i.i.d. algorithm.

In each case the gain sequence was taken of the form (5) in which  $b(n) = (1 + n)^{-1}$  for  $n \geq 0$ .

In experiments we found that the split sampling approach converges much more quickly than the Markovian approach. Note however that the Markovian algorithm can be run using observations of the queue process  $\mathbf{X}$ , without knowledge of the model.

### 5.2.1 Queueing model

Following uniformization, the M/M/1/b model is a doubly reflected random walk,

$$Q(t+1) = [Q(t) + \Delta(t+1)]_{0,b} \quad (21)$$

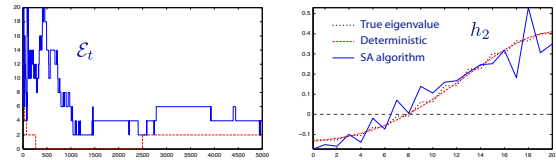
where  $[x]_{0,b} = \min(\max(x, 0), b)$  is a projection onto the interval  $[0, b]$ , and  $\Delta$  is an i.i.d. process. Letting  $\alpha$  denote the arrival rate, and  $\mu$  the service rate, scaled so that  $\alpha + \mu = 1$ , the marginal distribution of  $\Delta$  is given by,

$$P\{\Delta(t) = 1\} = \alpha, \quad P\{\Delta(t) = -1\} = \mu$$

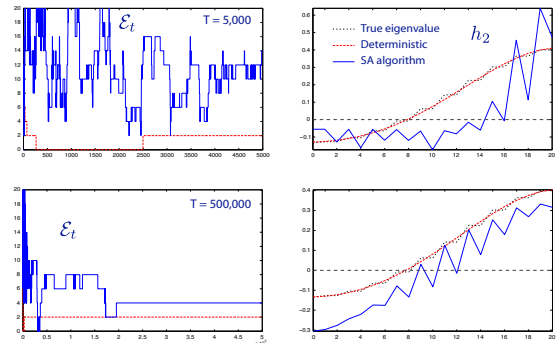
Hence its Markov transition matrix is given by,

$$P(x, \min(x+1, b)) = \alpha, \quad P(x, \max(x-1, 0)) = \mu, \quad x \in \mathbf{X} \quad (22)$$

To ensure that the matrix (13) is positive semi-definite we chose the Markov chain to be sampled at even integer values  $X(k) = Q(2k)$ ,  $k \geq 0$ , in the stochastic approximation algorithm (18) based on Markovian observations. In the split sampling algorithm, the i.i.d. process  $\mathbf{X}^1$  was constructed with geometric marginal distribution on  $\mathbf{X}$ . For  $n = 1, 2, \dots$  the random variable  $X^2(n)$  was chosen based on  $X^1(n-1)$  using  $P^2$ , with  $P$  defined in (22).



**Figure 4: Error trajectories defined by (20) and the final second eigenvalue estimates using the i.i.d and deterministic algorithms for the M/M/1/20 queue.**



**Figure 5: Error trajectories defined by (20) and the final second eigenvalue estimates using the Markovian and deterministic algorithms. After 500,000 steps the Markovian algorithm provides a good approximation of the sign structure of the second eigenvector, but the absolute error remains high.**

### 5.2.2 Statistical mechanics model

A running example in [13] is the Smoluchowski equation, defined by the Itô equation,

$$dX(t) = \nabla U(X(t)) dt + \sigma dN(t)$$

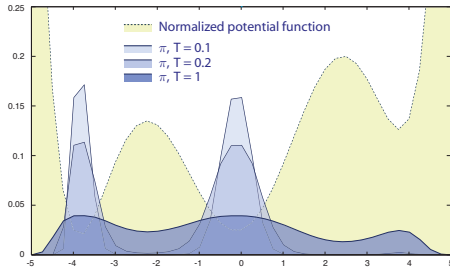
where  $\mathbf{N}$  is standard Brownian motion on  $\mathbb{R}$ , and the function  $U: \mathbb{R} \rightarrow \mathbb{R}$  is the polynomial,

$$U(x) = \frac{1}{200} \left( \frac{1}{2}x^6 - 15x^4 + 119x^2 + 28x + 50 \right)$$

Eigenfunctions of this diffusion were used to construct metastable subsets of  $\mathbb{R}$ .

Here we consider a related discrete-time Markov chain, and compute the spectrum of the transition matrix using the algorithms introduced in Section 4.

The Markov chain is constructed by restricting to a finite subset of  $\mathbb{R}$ :  $x$  restricted to  $N$  equally spaced values between  $-5$  and  $5$ , denoted  $\mathbf{X} = \{-5, -5 + \delta, -5 + 2\delta, \dots, 5 - \delta, 5\}$  where  $\delta = 10/(N-1)$ . We fix a scalar  $T_e > 0$  called the



**Figure 6:** Plot of  $\pi$  for  $N = 41$  and  $T_e = 0.1, 0.2$  and  $1$ . Also shown is a plot of the normalized potential function  $U/10$ .

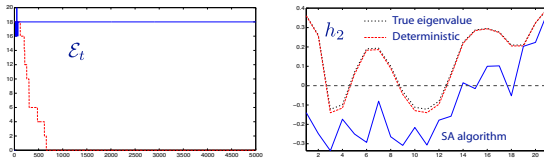
temperature, and define a bivariate distribution  $w$  on  $\mathbf{X} \times \mathbf{X}$  as follows:

$$w(x, y) = \frac{1}{\zeta} \exp(-(\max(U(x), U(y))/T_e))$$

where  $\zeta$  is the normalizing factor,  $\zeta := \sum_{x', y'} \exp(-(\max(U(x'), U(y'))/T_e))$ . As in the general construction described in Section 2, we define  $\pi(x) = \sum_y w(x, y)$ ,  $x \in \mathbf{X}$ , and a transition matrix is defined by  $P(x, y) = w(x, y)/\pi(x)$ ,  $x, y \in \mathbf{X}$ . The discretized Smoluchowski equation is then defined as the Markov chain with transition matrix,

$$P(x, y) = \frac{1}{\alpha(x)} \exp(-(\max(U(y) - U(x), 0)/T_e))$$

with  $\alpha(x) := \sum_{y'} \exp(-(\max(U(y') - U(x), 0)/T_e))$ .



**Figure 7:** Results for the discretized Smoluchowski equation based on the potential (6) using  $T_e = 1$ . Shown on the right is a plot of the second eigenvector for  $P$  and the approximation obtained from 5,000 iterations of the deterministic and split sampling algorithms. Shown on the left is the error process (20).

Shown on the right in Figure 7 is the resulting eigenfunction approximation after 5,000 iterations using the deterministic and split sampling algorithms. Shown on the left is the error process (20) using this algorithm. Convergence of the SA algorithm is slow, but note that a time horizon of 5,000 steps is very short. In this model, convergence to within 1% occurred after approximately 100,000 iterations.

## 6. CONCLUSIONS

We have introduced several stochastic-approximation variants of Oja's subspace algorithm for principal component analysis, Markov spectral theory, and spectral graph clustering. Convergence of these algorithms has been established through recent stochastic approximation techniques combined with stability theory from [7] that establishes convergence of the associated o.d.e..

Questions in current research include extensions to Markov chains that are not reversible, and variance reduction techniques for these algorithms.

## 7. REFERENCES

- [1] V. S. Borkar. Reinforcement learning - a bridge between numerical methods and Markov Chain Monte Carlo. In N. S. N. Sastry, B. Rajeev, M. Delampady, and T. S. S. R. K. Rao, editors, *Springer Lecture Notes in Control and Info. Sci.*, Perspectives in Mathematical Sciences. World Scientific, 2008.
- [2] V. S. Borkar. *Stochastic Approximation: A Dynamical Systems Viewpoint*. Hindustan Publishing Agency and Cambridge University Press (jointly), Dehli, India and Cambridge, UK, 2008.
- [3] V. S. Borkar and S. P. Meyn. The O.D.E. method for convergence of stochastic approximation and reinforcement learning. *SIAM J. Control Optim.*, 38(2):447–469, 2000. (also presented at the *IEEE CDC*, December, 1998).
- [4] A. Bovier, M. Eckhoff, V. Gayrard, and M. Klein. Metastability in stochastic dynamics of disordered mean-field models. *Probab. Theory Related Fields*, 119(1):99–161, 2001.
- [5] A. Bovier, M. Eckhoff, V. Gayrard, and M. Klein. Metastability in reversible diffusion processes. I. Sharp asymptotics for capacities and exit times. *J. Eur. Math. Soc. (JEMS)*, 6(4):399–424, 2004.
- [6] A. Bovier, V. Gayrard, and M. Klein. Metastability in reversible diffusion processes. II. Precise asymptotics for small eigenvalues. *J. Eur. Math. Soc. (JEMS)*, 7(1):69–99, 2005.
- [7] T. Chen, Y. Hua, and W.-Y. Yan. Global convergence of Oja's subspace algorithm for principal component extraction. *IEEE Trans. Neural Networks*, 9(1):58–67, Jan. 1998.
- [8] P. Deuffhard, W. Huisinga, A. Fischer, and C. Schütte. Identification of almost invariant aggregates in reversible nearly uncoupled Markov chains. *Linear Algebra Appl.*, 315(1-3):39–59, 2000.
- [9] P. Diaconis and D. Stroock. Geometric bounds for eigenvalues of Markov chains. *Ann. Appl. Probab.*, 1:36–61, 1991.
- [10] P. A. Ferrari, H. Kesten, and S. Martínez.  $R$ -positivity, quasi-stationary distributions and ratio limit theorems for a class of probabilistic automata. *Ann. Appl. Probab.*, 6:577–616, 1996.
- [11] J. A. Fill. Eigenvalue bounds on convergence to stationarity for nonreversible Markov chains, with an application to the exclusion process. *Ann. Appl. Probab.*, 1(1):62–87, 1991.
- [12] W. Huisinga. *Metastability of Markovian Systems: A transfer operator approach to molecular dynamics*. PhD thesis, Free University Berlin, 2001.
- [13] W. Huisinga, S. Meyn, and C. Schütte. Phase transitions and metastability in Markovian and molecular systems. *Ann. Appl. Probab.*, 14(1):419–458, 2004. Presented at the 11TH INFORMS Applied Probability Society Conference, NYC, JULY 25 - July 27, 2001.
- [14] A. Hyvärinen. Survey on independent component analysis. *Neural Computing Surveys*, 2:94–128, 1999.

- [15] I. T. Jolliffe. *Principal component analysis*. Springer Series in Statistics. Springer-Verlag, New York, second edition, 2002.
- [16] T. P. Krasulina. Method of stochastic approximation in the determination of the largest eigenvalue of the mathematical expectation of random matrices. *Automat. Remote Control*, 2:215–221, 1970.
- [17] M. Loève. *Probability Theory II*. Springer-Verlag, New York, Heidelberg, Berlin, 4th edition, 1978.
- [18] S. P. Meyn and R. L. Tweedie. *Markov Chains and Stochastic Stability*. Springer-Verlag, London, second edition, 1993. 2008 Edition to appear, Cambridge University Press, Cambridge Mathematical Library. 1993 edition online: <http://black.csl.uiuc.edu/~meyn/pages/book.html>.
- [19] B. Nadler, S. Lafon, R. Coifman, and I. Kevrekidis. Diffusion maps, spectral clustering and eigenfunctions of fokker-planck operators. In Y. Weiss, B. Schölkopf, and J. Platt, editors, *Advances in Neural Information Processing Systems 18*, pages 955–962. MIT Press, Cambridge, MA, 2006.
- [20] E. Oja. A simplified neuron model as a principal component analyzer. *J. Math. Biol.*, 15(3):267–273, 1982.
- [21] E. Oja and J. Karhunen. On stochastic approximation of the eigenvectors and eigenvalues of the expectation of a random matrix. *J. Math. Anal. Appl.*, 106(1):69–84, 1985.
- [22] L. Rey-Bellet and L. E. Thomas. Asymptotic behavior of thermal nonequilibrium steady states for a driven chain of anharmonic oscillators. *Comm. Math. Phys.*, 215(1):1–24, 2000.
- [23] B. Schölkopf, A. J. Smola, and K.-R. Müller. Nonlinear component analysis as a kernel eigenvalue problem. *Neural Computation*, 10(5):1299–1319, 1998.
- [24] Y. Weiss. Segmentation using eigenvectors: a unifying view. *Computer Vision, 1999. The Proceedings of the Seventh IEEE International Conference on*, 2:975–982 vol.2, 1999.
- [25] Z. Yi, M. Ye, J. C. Lv, and K. K. Tan. Convergence analysis of a deterministic discrete time system of Oja’s PCA learning algorithm. *IEEE Trans. Neural Networks*, 16(6):1318–1328, Nov. 2005.