

## OPTIMAL ESTIMATION OF FREE ENERGIES AND STATIONARY DENSITIES FROM MULTIPLE BIASED SIMULATIONS\*

HAO WU<sup>†</sup> AND FRANK NOÉ<sup>†</sup>

**Abstract.** When studying high-dimensional dynamical systems such as macromolecules, quantum systems, and polymers, a prime concern is the identification of the most probable states and their stationary probabilities or free energies. Often, these systems have metastable regions or phases, prohibiting the estimation of the stationary probabilities by direct simulation. Efficient sampling methods such as umbrella sampling, metadynamics, and conformational flooding have been developed that perform a number of simulations where the system’s potential is biased so as to accelerate the rare barrier crossing events. A joint free energy profile or stationary density can then be obtained from these biased simulations with the weighted histogram analysis method. This approach (a) requires a few essential order parameters to be defined in which the histogram is set up, and (b) assumes that each simulation is in global equilibrium. Both assumptions make the investigation of high-dimensional systems with previously unknown energy landscape difficult. Here, we introduce the transition matrix based unbiasing method (TMU), a simple and efficient estimation method which dismisses both assumptions. The configuration space is discretized into sets, but these sets are not only restricted to a preselected slow coordinate but can be clusters that form a partition of high-dimensional state space. The assumption of global equilibrium is replaced by requiring only local equilibrium within the discrete sets, and the stationary density or free energy is extracted from the transitions between clusters. We prove the asymptotic convergence and normality of TMU, give an efficient approximate version of it, and demonstrate its usefulness in numerical examples.

**Key words.** Markov chain, maximum likelihood estimation, error analysis, free energy, stationary density, simulation

**AMS subject classifications.** 60J20, 62M09, 65C20

**DOI.** 10.1137/120895883

**NOTATION.** For convenience, we summarize some notation used throughout the paper as follows:

$\mathbf{1}_\omega$	Indicator function of event $\omega$ , taking value 1 if $\omega$ holds and 0 otherwise
$\bar{x}, \hat{x}, \tilde{x}$	“True” value, estimate obtained from simply “counting,” and maximum likelihood estimate of an unknown variable $x$ (except if stated otherwise)
$\mathbf{0}, \mathbf{1}$	Vectors of zeros and ones in appropriate dimensions
$\mathbf{I}$	Identity matrix
$\mathcal{N}(u, \Sigma)$	Multivariate normal distribution with mean $u$ and covariance matrix $\Sigma$
$\mathcal{V}(G)$	Vector $(G_{11}, G_{12}, \dots, G_{mn})^T$ which consists of elements of $G = [G_{ij}] \in \mathbb{R}^{m \times n}$
$\mathcal{V}(G_1, \dots, G_m)$	$(\mathcal{V}(G_1)^T, \dots, \mathcal{V}(G_m)^T)^T$
$\text{tr}(G)$	Trace of square matrix $G$
$\mu(v)$	Arithmetic mean of elements of vector $v$ , i.e., $\mu(v) = \frac{1}{m} \sum_i v_i$ for $v = [v_i] \in \mathbb{R}^m$
$\nabla_x y$	Jacobian matrix $[\partial y_i / \partial x_j]$ of $y = [y_i]$ with respect to $x = [x_i]$
$\nabla_y x z$	$\nabla_y (\nabla_x z)^T$ for $z \in \mathbb{R}$
$G > 0 (\geq 0)$	Each element of matrix $G$ is positive (nonnegative)
$G < 0 (\leq 0)$	Matrix $G$ is negative-definite (negative-semidefinite)
$G \succ 0 (\succeq 0)$	Matrix $G$ is positive-definite (positive-semidefinite)
$G^+$	Moore–Penrose pseudoinverse of $G$
$x_n \xrightarrow{d} x$	$x_n$ converges in distribution to $x$ w.r.t. $n$ ; i.e., $\lim_{n \rightarrow \infty} \Pr(x_n \in \mathcal{B}) = \Pr(x \in \mathcal{B})$ for all continuity sets $\mathcal{B}$
$x_n \xrightarrow{p} x$	$x_n$ converges in probability to $x$ w.r.t. $n$ ; i.e., $\lim_{n \rightarrow \infty} \Pr(\ x_n - x\  \geq \epsilon) = 0$ for all $\epsilon > 0$

\*Received by the editors October 22, 2012; accepted for publication (in revised form) October 3, 2013; published electronically January 16, 2014.

<http://www.siam.org/journals/mms/12-1/89588.html>

<sup>†</sup>Mathematics Institute, Department of Mathematics and Computer Science, Free University of Berlin, 14195 Berlin, Germany (hwu@zedat.fu-berlin.de, frank.noe@fu-berlin.de). The first author was supported by DFG grants NO 825/2-1 and WU 744/1-1. The second author was supported by DFG research center Matheon and ERC grant 307494 “pcCell.”

$\ G\ $	Frobenius norm of $G$
$\ G\ _{\max}$	$\max_{i,j}  G_{ij} $ for $G = [G_{ij}]$

**1. Introduction.** Stochastic simulations of chemical, physical, or biological processes often involve rare events that render the exploration of relevant states, or the calculation of expectation values by direct numerical simulation, difficult or impossible. Examples include phase transitions in spin systems [24, 3], transitions between different chemical states in quantum dynamics simulations [14], and conformational transitions in biomolecules [8]. For this reason, many enhanced sampling techniques have been developed to modify the dynamics of the original simulation system such that the relevant rare events become more frequent and can be accessed by direct numerical simulation of the modified simulation system. Such an approach is, of course, reasonable only if there exists a way to reliably compute at least some quantities of interest of the original simulation system from the realizations of the modified simulation system.

In this paper we focus on processes that are asymptotically stationary and ergodic, and on enhanced sampling approaches that use bias potentials (or, equivalently, conservative bias force fields) that attempt to modify the original system’s dynamics so as to avoid rare events. Well-known examples of such approaches are umbrella sampling [29], conformational flooding [10], and metadynamics and its variants [13, 1]. These approaches assume that one has some prior knowledge of coordinates or order parameters which are “slow”; i.e., the rare event dynamics of the system is resolved by state transitions in these selected coordinates.

Umbrella sampling defines a series of biased simulations, each of which uses the forces from the original dynamics and the forces arising from a specific harmonic potential. These potentials restrain the biased simulations to stay close to positions in the selected coordinates which are the centers of the umbrella potentials. The force constant(s) of these potentials must be chosen such that the corresponding biased stationary densities overlap significantly and the unification of all biased stationary densities covers the part of state space in which the original stationary density is significantly greater than zero. In this case, all the biased simulations, together with the knowledge of the umbrella potentials, can be used in order to estimate the original stationary density in the selected coordinates (or the corresponding free energy landscape).

Metadynamics is based on an opposite philosophy. Rather than constraining the simulation to a set of points, it adds bias potentials to drive the simulation away from regions that it has sampled sufficiently well. In practice this is often done by adding Gaussian hat functions to the biased potential every constant number of simulation steps, centered at the current simulation state. We consider that this happens a set number of times, leading to the same number of simulation snippets, each with a different biasing potential. Due to limitations of filling high-dimensional space volumes, these bias potentials also usually live in a few predefined coordinates. Since the sequence of added bias potentials depends on the simulation history, metadynamics is usually used to first “fill up” the free energy wells until the states that cause the rare event waiting times have been destabilized and the corresponding free energy landscape is approximately “flat.” It can be shown that at this point continuing the metadynamics simulation will sample bias potentials that are the negative free energy landscapes of the original system, up to an arbitrary additive constant. Since metadynamics does not require the modeler to know the relevant states along the slow coordinates, it not only is an approach to quantifying the stationary distribution/free

energy landscape of the original system but has been very successful in terms of exploring the state space in complex systems [19]. Unfortunately, this approach of using metadynamics also appears to suggest that all simulation effort that has been spent until the free energy surface is approximately flat cannot be used for quantitative estimations.

Here we concentrate on the step of unbiasing the modified dynamics so as to obtain the stationary distribution of the original dynamical system. For both umbrella sampling and metadynamics, the step of “unbiasing” is usually done with the weighted histogram analysis method (WHAM) [7]. WHAM uses a discretization of the selected coordinates in which the biased simulation was done and collects a set of histograms, one for each of the biased simulations. These biased histograms are then combined into a single unbiased histogram by solving a set of self-consistent equations to minimize the statistical error. The assumption used by WHAM is that each of the biased simulations done at different conditions is sufficiently long such that they generate unbiased samples of the corresponding biased stationary density. In other words, each subsimulation is assumed to be in global equilibrium at its conditions. We will see that this assumption is unnecessary, and a method that does not rely on this assumption can provide estimates that are substantially more precise (or, equivalently, require substantially less simulation effort for a given level of precision), even in trivial double-well examples.

This paper develops the transition matrix based unbiasing method (TMU), which replaces the assumption that the biased simulations are in global equilibria by the much weaker assumption that each simulation is only in local equilibrium in the discrete states on which the stationary distribution is estimated. TMU has been motivated by the recent progress in Markov modeling [28, 6, 18, 21] and constructs the joint unbiased stationary distribution from a series of transition count matrices estimated from the biased simulations. However, it is important to note that TMU does *not* need the discrete dynamics to be Markovian.

Subsequently, we describe the basic mathematical assumptions underlying our method, then describe TMU in its most general form and show that the method always has a solution that is asymptotically normal and convergent. We then provide an approximate TMU that is efficient for very large state spaces and a large number of subsimulations. The method is demonstrated in conjunction with umbrella sampling and metadynamics on double-well potentials, and its performance is compared with that of the standard WHAM and a recently introduced method, the multiple Markov transition matrix method (MMMM), that had a similar motivation [23].

**2. Background.** In this section, we briefly review the mathematical background of biased simulation techniques. Let us consider a reference system on the finite state space  $\mathcal{S} = \{1, \dots, n\}$  with free energy  $V = [V_i]$ , where  $V_i$  is the energy of state  $i$ . If we denote the system state at time  $t$  by  $x_t$ , the state sequence  $\{x_t\}$  is then a stochastic process. In this paper, we focus on processes  $\{x_t\}$  with properties of asymptotic stationarity, wide-sense ergodicity, and detailed balance, which are relevant for many physical simulation processes. The detailed descriptions of these properties are listed as follows:

1. *Asymptotic stationarity* means that the sequence  $\{x_t\}_{t \geq \tau}$  is approximately stationary if  $\tau$  is large enough. More formally,  $\{x_t\}$  is said to be asymptotically stationary if the limits  $\lim_{t \rightarrow \infty} \Pr(x_{t+1} = s_1, \dots, x_{t+m} = s_m)$  exist for all  $m \geq 1$  and  $s_1, \dots, s_m \in \mathcal{S}$ . Specifically, the marginal distribution of  $x_t$

converges to the Boltzmann distribution, i.e.,

$$(2.1) \quad \pi_i = \frac{\exp(-\beta V_i)}{\sum_j \exp(-\beta V_j)},$$

where  $\pi = [\pi_i]$  denotes the system stationary distribution defined by  $\pi_i = \lim_{t \rightarrow \infty} \Pr(x_t = i)$  and  $\beta$  is a constant and generally proportional to the inverse temperature in physical systems. Furthermore, we define

$$(2.2) \quad T = [T_{ij}] = \left[ \lim_{t \rightarrow \infty} \Pr(x_{t+1} = j | x_t = i) \right].$$

It is easy to see that the matrix  $T \in \mathbb{R}^{n \times n}$  represents the stationary state transition probabilities and satisfies the condition that each row sums to 1, so here, for simplicity, we call  $T$  the *transition matrix* of  $\{x_t\}$  even if  $\{x_t\}$  is not a Markov chain.

2. *Wide-sense ergodicity* states that

$$\frac{1}{\tau + 1} \sum_{t=0}^{\tau} \mathbf{1}_{x_t=i} \xrightarrow{P} \pi_i$$

and

$$\frac{1}{\tau} \sum_{t=0}^{\tau-1} \mathbf{1}_{x_t=i} \cdot \mathbf{1}_{x_{t+1}=j} \xrightarrow{P} \pi_i T_{ij} \quad \text{as } \tau \rightarrow \infty,$$

and implies that  $\pi$  and  $T$  can be consistently estimated by time averages of  $\mathbf{1}_{x_t=i}$  and  $\mathbf{1}_{x_t=i} \cdot \mathbf{1}_{x_{t+1}=j}$ .

3. The *detailed balance* condition can be written as

$$\lim_{t \rightarrow \infty} \Pr(x_t = i, x_{t+1} = j) = \lim_{t \rightarrow \infty} \Pr(x_t = j, x_{t+1} = i)$$

or, equivalently,  $\pi_i T_{ij} = \pi_j T_{ji}$  for all  $i, j$ , which means that each state transition has the same unconditional probability as its reverse. Note that this property clearly holds for systems that are time-reversible at equilibrium.

*Remark 2.1.* In the present study we work at constant thermodynamic temperature. For convenience, we measure energies  $V_i$  in units of thermal energy  $\beta^{-1} = k_B \mathcal{T}$  with  $k_B$  the Boltzmann constant and  $\mathcal{T}$  the thermodynamic temperature, yielding  $\beta = 1$  in (2.1). Furthermore, we assume without loss of generality that all involved free energies in this paper have zero mean. Then  $\sum_i V_i = 0$ , and we can construct a bijection between stationary distributions and free energies with  $\pi(V)$  defined by (2.1) and

$$(2.3) \quad V(\pi) = -\log \pi + \mu(\log \pi),$$

where  $\log \pi = [\log \pi_i] \in \mathbb{R}^n$  and  $\mu(\cdot)$  denotes the mean operator defined in the notation list.

*Remark 2.2.* In general, the discrete state space  $\mathcal{S}$  may correspond to a set partition of an underlying dynamic system with a large-scale or continuous state space. In this case,  $V_i$  corresponds to the configurational free energy difference of the  $i$ th partition with respect to an arbitrary reference state.

For now, our goal is to estimate  $V$ , or equivalently  $\pi$ , from simulations in the case that it is unknown. Due to the wide-sense ergodicity, when sufficient simulation data can be generated, we can simply carry out one or multiple simulations of the reference system and get the estimate of  $V$  through computing the histogram of simulation data. This approach, however, is very inefficient when the reference system has multiple metastable states, because the simulation process is very likely to get stuck in some local minima of the energy landscape for a long time. To alleviate this drawback, biased simulation techniques, such as umbrella sampling [26] and metadynamics [13], were developed to solve this problem; they perform simulations for a set of biased potentials so that the energy landscape can be explored more efficiently.

Although many practical algorithms use a different approach, we can roughly summarize the estimation of stationary distributions through biased simulations in terms of the following pseudocode:

**Step 1.** Design a set of biasing potentials  $\{U^{(k)}\}_{k=1}^K$ , where  $U^{(k)} = [U_i^{(k)}] \in \mathbb{R}^n$ .

**Step 2.** Repeat Steps 2.1 and 2.2 for  $k = 1, \dots, K$ :

**Step 2.1.** Change the system potential as

$$(2.4) \quad V^{(k)} = V + U^{(k)} - \mu \left( V + U^{(k)} \right),$$

where  $V^{(k)}$  is called the biased potential and the last term is used to shift the mean of  $V^{(k)}$  to zero.

**Step 2.2.** Perform a biased simulation with length  $M$  using the same simulation model as the reference system except that the potential energy is changed from  $V$  to  $V^{(k)}$ , and record the simulation trajectory  $\{x_t^{(k)}\}_{t=0}^M$ .

**Step 3.** Estimate the reference (unbiased) free energy  $V$  or stationary distribution  $\pi$  from  $K$  biased simulation trajectories.

In this paper, we will focus on the estimation problem in Step 3. We start with the assumption that each simulation  $\{x_t^{(k)}\}_{t=0}^M$  is a Markov chain, and the developed estimation method will then be proved to be applicable to more general simulation models.

### 3. Maximum likelihood estimation from multiple simulations.

**3.1. Maximum likelihood estimation.** In this section, we investigate a maximum likelihood approach to the estimation problem described in section 2 under the assumption of the Markovity of biased simulations. For the description of the estimation method, it is convenient to denote the biased stationary distribution, transition matrix, and count matrix of the  $k$ th simulation by  $\pi^{(k)} = [\pi_i^{(k)}]$ ,  $T^{(k)} = [T_{ij}^{(k)}]$ , and  $C^{(k)} = [C_{ij}^{(k)}]$ , where

$$(3.1) \quad C_{ij}^{(k)} := \sum_{t=1}^M \mathbf{1}_{x_{t-1}=i}^{(k)} \cdot \mathbf{1}_{x_t=j}^{(k)}$$

i.e.,  $C_{ij}^{(k)}$  is the number of transitions from state  $i$  to state  $j$  in the  $k$ th simulation, and the relationship between  $\pi^{(k)}$  and  $\pi$  can be written as

$$(3.2) \quad \pi_i^{(k)} = \frac{\exp\left(-U_i^{(k)}\right) \pi_i}{\sum_{j=1}^n \exp\left(-U_j^{(k)}\right) \pi_j}.$$

Suppose that each simulation  $\{x_t^{(k)}\}$  is a time-homogeneous and reversible Markov chain. The maximum likelihood estimation (MLE) of the unbiased stationary distribution  $\pi$  can then be obtained by solving the following optimization problem:

$$(3.3) \quad \begin{aligned} & \max_{\pi, T^{(1)}, \dots, T^{(K)}} && L = \sum_k L^{(k)}(T^{(k)} | C^{(k)}) \\ & \text{subject to (s.t.)} && \forall i, j = 1, \dots, n, \quad k = 1, \dots, K, \\ & && \pi \text{ is a probability vector,} \\ & && T^{(k)} \text{ is a transition matrix,} \\ & && \pi_i^{(k)} T_{ij}^{(k)} = \pi_j^{(k)} T_{ji}^{(k)}, \end{aligned}$$

where

$$(3.4) \quad L^{(k)}(T^{(k)} | C^{(k)}) = \log \Pr(x_1^{(k)}, \dots, x_M^{(k)} | x_0^{(k)}, T^{(k)}) = \sum_{i,j} C_{ij}^{(k)} \log T_{ij}^{(k)}$$

denotes the log-likelihood function of the  $k$ th simulation, and the last constraint is the detailed balance constraint. (Here we set  $0 \log 0 = 0$  and  $a \log 0 = -\infty$  if  $a > 0$ .) After performing the MLE of  $\pi$ , the optimal estimate of  $V$  can also be obtained by using (2.3).

*Remark 3.1.* It is a nontrivial task to solve (3.3) even if the problem is of small size because the detailed balance constraint is a highly nonlinear equality constraint. To overcome this difficulty, we can replace the detailed balance constraint by the following equivalent expression:

$$(3.5) \quad \exp(-U_i^{(k)}) \pi_i T_{ij}^{(k)} = \exp(-U_j^{(k)}) \pi_j T_{ji}^{(k)}.$$

(The equivalence can be simply proved by (3.2).) Note that both sides of (3.5) can be represented as a difference of two convex functions by using  $\pi_i T_{ij}^{(k)} = \frac{1}{4}(\pi_i + T_{ij}^{(k)})^2 - \frac{1}{4}(\pi_i - T_{ij}^{(k)})^2$ . Therefore the constrained concave-convex procedure [25] can be used to perform the MLE.

For the MLE problem (3.3), we have the following theorem.

**THEOREM 3.2.** *The optimization problem (3.3) has at least one optimal solution satisfying*

1.  $T_{ij}^{(k)} = 0$  for  $(i, j, k) \in \{(i, j, k) | C_{ij}^{(k)} + C_{ji}^{(k)} = 0 \text{ and } i \neq j\}$ ;
2.  $1_{T_{ij}^{(k)} > 0} \equiv 1_{C_{ij}^{(k)} > 0}$  if  $C_{ii}^{(k)} > 0$  and  $1_{C_{ij}^{(k)} > 0} = 1_{C_{ji}^{(k)} > 0}$  for all  $i, j, k$ .

*Proof.* See Appendix A.  $\square$

According to Theorem 3.2, the dimension of the optimization variable of (3.3) can be significantly reduced by setting  $T_{ij}^{(k)} = 0$  for  $(i, j, k)$  belonging to  $\{(i, j, k) | C_{ij}^{(k)} + C_{ji}^{(k)} = 0, i \neq j\}$  when count matrices are sparse. However, even if each  $C^{(k)}$  is sparse with  $O(n)$  nonzero elements, the reduced problem involves  $O(nK)$  decision variables and nonlinear equality constraints. (Note that  $\pi$  and  $T^{(k)}$  are both unknown in the last constraint.) It is still inefficient to search the optimal solution by direct methods. In section 4, we will adopt an approximate MLE method to improve the efficiency.

**3.2. Convergence analysis.** The MLE method of stationary distribution in section 3.1 is motivated by the assumption that  $\{x_t^{(k)}\}$  is a Markov chain. Interestingly, it turns out that the Markov property is not necessary for the convergence of MLE. In this section we will prove the convergence of MLE under more general conditions.

First, we discuss the relationship between the MLE and the counting based estimation by using the Kullback–Leibler (KL) divergence rate proposed in [22], and we provide an intuitive explanation for why the MLE can work for non-Markovian stochastic processes.

**DEFINITION 3.3.** *Let  $T' = [T'_{ij}]$  and  $T'' = [T''_{ij}]$  be two transition matrices of the same dimension; then the KL divergence rate between  $T'$  and  $T''$  w.r.t. the probability vector  $\pi' = [\pi'_i]$  is*

$$(3.6) \quad \text{KLR}_{\pi'}(T'|T'') = \sum_{i,j} \pi'_i T'_{ij} (\log T'_{ij} - \log T''_{ij}).$$

*Remark 3.4.* It is easy to see that  $\text{KLR}_{\pi'}(T'|T'')$  can be employed as a “pseudo-metric” to measure the distance between  $T'$  and  $T''$  with  $\text{KLR}_{\pi'}(T'|T'') = 0 \Leftrightarrow T' = T''$  in the case of  $\pi' > 0$ . (The KL divergence rate is not a true metric because neither the symmetry nor the triangle inequality is satisfied.) Furthermore, for two Markov chains  $\{x'_t\}$  and  $\{x''_t\}$  with transition matrices  $T', T''$ , it can be proved that  $\text{KLR}_{\pi'}(T'|T'') = \frac{1}{t+1} \lim_{t \rightarrow \infty} \text{KL}(x'_0, \dots, x'_t | x''_0, \dots, x''_t)$  if  $\pi'$  is the stationary distribution of  $T'$  [22], where  $\text{KL}(\cdot|\cdot)$  denotes the KL divergence.

Generally speaking, if there is no other knowledge available,  $T^{(k)}$  can be estimated as  $\hat{T}^{(k)} = [\hat{T}_{ij}^{(k)}]$  with  $\hat{T}_{ij}^{(k)}$  being the fraction of observed transitions from the  $i$ th state to the  $j$ th state:

$$(3.7) \quad \hat{T}_{ij}^{(k)} = C_{ij}^{(k)} / \left( \sum_l C_{il}^{(k)} \right).$$

But the transition matrix estimates obtained from (3.7) generally do not satisfy the detailed balance condition and do not share the same unbiased stationary distribution for finite-time simulations. Therefore we search for the feasible transition matrices which are the closest to  $\hat{T}^{(1)}, \dots, \hat{T}^{(K)}$  based on the KL divergence rate:

$$(3.8) \quad \begin{aligned} & \min_{\pi, T^{(1)}, \dots, T^{(K)}} \sum_k \text{KLR}_{\hat{\pi}^{(k)}} \left( \hat{T}^{(k)} || T^{(k)} \right) \\ & \text{s.t.} \quad \forall i, j = 1, \dots, n, \quad k = 1, \dots, K, \\ & \quad \pi \text{ is a probability vector,} \\ & \quad T^{(k)} \text{ is a transition matrix,} \\ & \quad \pi_i^{(k)} T_{ij}^{(k)} = \pi_j^{(k)} T_{ji}^{(k)}, \end{aligned}$$

where  $\hat{\pi}^{(k)} = [\hat{\pi}_i^{(k)}]$  is the counting estimate of  $\pi^{(k)}$  given by

$$(3.9) \quad \hat{\pi}_i^{(k)} = \left( \sum_l C_{il}^{(k)} \right) / \left( \sum_{j,l} C_{jl}^{(k)} \right).$$

Note that the KL divergence rate  $\text{KLR}_{\hat{\pi}^{(k)}}(\hat{T}^{(k)} || T^{(k)})$  can be decomposed as

$$(3.10) \quad \text{KLR}_{\hat{\pi}^{(k)}}(\hat{T}^{(k)} || T^{(k)}) = -\frac{1}{M} L^{(k)} \left( T^{(k)} | C^{(k)} \right) + \frac{1}{M} L^{(k)} \left( \hat{T}^{(k)} | C^{(k)} \right)$$

and the optimal solution of (3.8) is  $T^{(k)} = \hat{T}^{(k)}$  if  $\hat{T}^{(1)}, \dots, \hat{T}^{(K)}$  satisfy all the constraints. Therefore (3.8) is equivalent to (3.3), and the MLE can be considered to be a projector which projects the counting estimates (3.7) onto the feasible space.

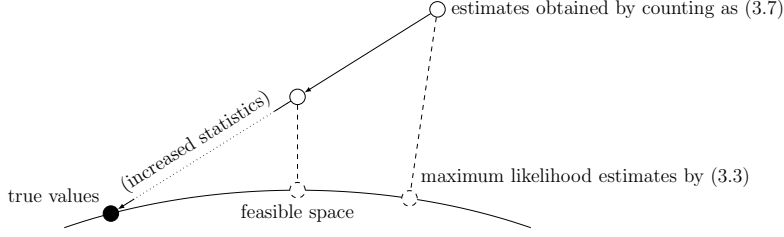


FIG. 3.1. Relationship between the estimates obtained by counting, the maximum likelihood estimates, and the true values of transition matrices.

Figure 3.1 shows the relationship between different estimates and the true values of  $(T^{(1)}, \dots, T^{(K)})$ , where the estimates obtained by counting converge to their true values due to the wide-sense ergodicity of simulations. Obviously, if it can be shown that the error of the MLE converges to zero when  $\hat{T}^{(1)}, \dots, \hat{T}^{(K)}$  converge to their true values, the consistency of the MLE can hold without assuming the Markov property.

We now present a formal analysis of the convergence of the MLE. Before proving the main theorems, we make some assumptions and introduce some notation.

Unless stated otherwise, in this paper all convergence statements are made w.r.t.  $M \rightarrow \infty$ .

ASSUMPTION 3.5.  $U^{(1)}, \dots, U^{(K)}$  and  $\bar{V}$  are finite, where  $\bar{V}$  denotes the true value of  $V$ .

ASSUMPTION 3.6.  $\{x_t^{(1)}\}, \dots, \{x_t^{(K)}\}$  are all asymptotically stationary and wide-sense ergodic processes with detailed balance.

ASSUMPTION 3.7. For any  $i, j, k$ , if there exists some  $t$  such that  $\Pr(x_t^{(k)} = i, x_{t+1}^{(k)} = j) > 0$ , then  $\lim_{\tau \rightarrow \infty} \Pr(x_\tau^{(k)} = i, x_{\tau+1}^{(k)} = j) > 0$ .

ASSUMPTION 3.8.  $\pi = \bar{\pi}$  is the unique solution of the following set of equations and inequalities:

$$(3.11) \quad \begin{cases} \pi_i^{(k)}(\pi) \cdot \bar{T}_{ij}^{(k)} = \pi_j^{(k)}(\pi) \cdot \bar{T}_{ji}^{(k)} & \text{for } i, j = 1, \dots, n \text{ and } k = 1, \dots, K, \\ \mathbf{1}^\top \pi = 1 & \text{and } \pi \geq 0, \end{cases}$$

where  $\bar{\pi}$  and  $\bar{T}^{(k)} = [\bar{T}_{ij}^{(k)}]$  denote true values of  $\pi$  and  $T^{(k)}$ .

The above assumption means that the unbiased stationary distribution can be uniquely determined if all the transition matrices are given.

Furthermore, here we let  $\theta = \mathcal{V}(\pi, T^{(1)}, \dots, T^{(K)})$  be the vector consisting of elements of the unbiased stationary distribution  $\pi$  and transition matrices,  $\tilde{\theta} = \mathcal{V}(\bar{\pi}, \bar{T}^{(1)}, \dots, \bar{T}^{(K)})$  be the solution of (3.3),  $X^{(k)} = [X_{ij}^{(k)}] = [\pi_i^{(k)} T_{ij}^{(k)}]$  be the matrix of unconditional transition probabilities,  $\bar{X}^{(k)} = [\bar{X}_{ij}^{(k)}] = [\bar{\pi}_i^{(k)} \bar{T}_{ij}^{(k)}]$  denote the true value of  $X^{(k)}$ , and  $\hat{X}^{(k)} = [\hat{X}_{ij}^{(k)}] = [\hat{\pi}_i^{(k)} \hat{T}_{ij}^{(k)}]$  be the estimate of  $X^{(k)}$  obtained by counting. (The definition of  $\mathcal{V}(\cdot)$  is given in the list of notation.)

Based on the above assumptions and notation, we have the following theorems on the convergence of  $\hat{\theta}$ .

THEOREM 3.9. If Assumptions 3.5–3.8 hold, then  $\hat{\theta} \xrightarrow{P} \bar{\theta}$ .

Proof. See Appendix B.  $\square$

THEOREM 3.10. If Assumptions 3.5–3.8 hold and the conditions



1. for each  $k$ , there exists a  $\Sigma_X^{(k)}$  such that

$$(3.12) \quad \sqrt{M} \left( \mathcal{V} \left( \hat{X}^{(k)} \right) - \mathcal{V} \left( \bar{X}^{(k)} \right) \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \Sigma_X^{(k)} \right);$$

2.  $\tilde{T}_{ij}^{(k)} = 0$  if  $C_{ij}^{(k)} + C_{ji}^{(k)} = 0$ ;  
 3. diagonal elements of  $\bar{X}^{(1)}, \dots, \bar{X}^{(K)}$  are positive;  
 4. all  $K$  simulations are statistically independent;  
 5.  $H = \sum_k \nabla_{\theta_r} L^{(k)}(T^{(k)}(\bar{\theta}_r) | \bar{X}^{(k)})$  is nonsingular with  $\theta_r$  the vector consisting of  $\{T_{ij}^{(k)} | \bar{X}_{ij}^{(k)} > 0, i < j\}$  and  $\{\pi_1, \dots, \pi_{n-1}\}$  and  $\bar{\theta}_r$  the corresponding true value of  $\theta_r$

are satisfied, then

$$(3.13) \quad \sqrt{M} \left( \tilde{\theta} - \bar{\theta} \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \nabla_{\theta_r} \theta(\theta_r) \cdot H^{-1} \Sigma H^{-1} \cdot \left( \nabla_{\theta_r} \theta(\theta_r) \right)^T \right),$$

where  $\Sigma = \left( \nabla_{\theta_r} \Phi(\theta(\bar{\theta}_r)) \right)^T \Sigma_X \nabla_{\theta_r} \Phi(\theta(\bar{\theta}_r))$ ,  $\Sigma_X = \text{diag}(\Sigma_X^{(1)}, \dots, \Sigma_X^{(K)})$ , and  $\Phi(\theta) = \mathcal{V}(\log T^{(1)}, \dots, \log T^{(K)})$ .

*Proof.* See Appendix C.  $\square$

*Remark 3.11.* Note that  $\hat{X}^{(k)}$  can be expressed as  $\hat{X}^{(k)} = \frac{1}{M} \sum_{t=1}^M \Delta C_t^{(k)}$  with  $\Delta C_t^{(k)} = [\Delta C_{t,ij}^{(k)}] = [1_{x_{t-1}^{(k)}=i} \cdot 1_{x_t^{(k)}=j}]$ , so condition 1 stated in Theorem 3.10 means that the central limit theorem holds for  $\{\Delta C_t^{(k)}\}$  and can be justified by using Markov chain central limit theorems [12] in many practical situations. For example, if for each simulation  $k$ ,  $\{x_t^{(k)}\}$  is obtained by coarse-graining a stochastic process  $\{y_t^{(k)}\}$  on a large-scale or continuous state space as mentioned in Remark 2.2 and  $\{y_t^{(k)}\}$  is a geometrically ergodic Markov chain, then (3.12) can be simply proved by Theorem 1.2 in [11].

*Remark 3.12.* In this section we only characterize the convergence of  $\tilde{\pi}$ . For the MLE of free energy  $V$ , the consistency and asymptotic normality are immediate consequences of Theorems 3.9 and 3.10 by considering that  $V$  is a smooth function of  $\pi$  (see (2.3)). Here we omit the detailed description and proof as they are trivial.

**4. Approximate MLE.** In this section, we develop an approximate MLE method based on a decomposition strategy in order to improve the efficiency of MLE, and the convergence of the method is also shown.

For convenience of analysis and computation, here we introduce two new variables,  $\underline{C}^{(k)}$  and  $Z^{(k)}$ .  $\underline{C}^{(k)} = [\underline{C}_{ij}^{(k)}]$  is a modified count matrix used to replace  $C^{(k)}$  to avoid singularity in the approximate MLE and is assumed to satisfy the following assumption.

**ASSUMPTION 4.1.**  $\underline{C}^{(1)}, \dots, \underline{C}^{(K)}$  are irreducible matrices with positive diagonal elements and satisfy  $1_{\underline{C}_{ij}^{(k)} > 0} = 1_{\underline{C}_{ji}^{(k)} > 0}$  and  $1_{\underline{C}^{(k)} = C^{(k)}} \xrightarrow{P} 1$  for all  $i, j, k$ .

One way to perform the count matrix modification is as follows:

$$(4.1) \quad \underline{C}_{ij}^{(k)} = \begin{cases} \max \left\{ C_{ij}^{(k)}, \delta \right\}, & C_{ji}^{(k)} > 0 \text{ or } i = j, \\ C_{ij}^{(k)} & \text{otherwise,} \end{cases}$$

where  $\delta \in (0, 1)$  is a small number. (This approach is similar to the so-called neighbor prior used in [20, 2].)

**THEOREM 4.2.** *If Assumptions 3.5–3.7 hold, and if  $\bar{X}_{ii}^{(k)} > 0$  and  $\sum_{t=0}^M 1_{x_t^{(k)}=i} > 0$  for all  $i, k$ , then the modified count matrices defined in (4.1) satisfy Assumption 4.1.*

*Proof.* The proof is omitted because it is trivial by contradiction.  $\square$

The variable  $Z^{(k)} = [Z_{ij}^{(k)}]$  is defined by

$$(4.2) \quad \exp(-Z_{ij}^{(k)}) \propto X_{ij}^{(k)},$$

which can be interpreted as the “free energy matrix” of state transitions in the  $k$ th simulation because  $\exp(-Z_{ij}^{(k)}) \propto \lim_{t \rightarrow \infty} \Pr(x_t^{(k)} = i, x_{t+1}^{(k)} = j)$ , and the relationship between the free energy matrix and the free energy can be expressed as

$$(4.3) \quad V_Z(Z^{(k)}) = V^{(k)}$$

with

$$(4.4) \quad V_Z(Z^{(k)}) = (Z_1^{(k)}, \dots, Z_n^{(k)})^T - \frac{1}{n} \sum_i Z_i^{(k)},$$

where

$$(4.5) \quad Z_i^{(k)} = -\log \sum_j \exp(-Z_{ij}^{(k)})$$

denotes the potential of state  $i$  derived from  $Z^{(k)}$ . Like the free energy  $V$ , we also assume that  $\sum_{(i,j) \in \{(i,j) | X_{ij}^{(k)} > 0\}} Z_{ij}^{(k)} = 0$  such that all the finite elements of  $Z^{(k)}$  have zero mean and there is a one-to-one correspondence between  $Z^{(k)}$  and  $X^{(k)}$ . (Note that  $Z_{ij}^{(k)} = \infty$  if  $X_{ij}^{(k)} = 0$ .)

Under the above assumption and variable definitions and replacing  $C^k$  by  $\underline{C}^k$ , (3.3) can be written as

$$(4.6) \quad \begin{aligned} & \max_{V, \{Z_{ij}^{(k)} | \underline{C}_{ij}^{(k)} > 0\}} L = \sum_k L_Z^{(k)}(Z^{(k)} | \underline{C}^{(k)}) \\ & \text{s.t.} \quad \forall k = 1, \dots, K, \\ & \quad Z^{(k)} = Z^{(k)T}, \\ & \quad \sum_{(i,j) \in \{(i,j) | \underline{C}_{ij}^{(k)} > 0\}} Z_{ij}^{(k)} = 0, \\ & \quad \mathbf{1}^T V = 0, \\ & \quad V_Z(Z^{(k)}) = V^{(k)}(V), \end{aligned}$$

where

$$(4.7) \quad L_Z^{(k)}(Z^{(k)} | \underline{C}^{(k)}) = -\sum_{i,j} C_{ij}^{(k)} Z_{ij}^{(k)} + \sum_i \underline{C}_i^{(k)} Z_i^{(k)},$$

$0 \cdot \infty$  is set to be 0, and  $\underline{C}_i^{(k)} = \sum_j \underline{C}_{ij}^{(k)}$ . A brief description of the objective function and constraints of (4.6) follows.

1. Each term of the objective function is the log-likelihood of the free energy matrix  $Z^{(k)}$ , given  $\underline{C}^{(k)}$  with  $L_Z^{(k)}(Z^{(k)} | \underline{C}^{(k)}) = L^{(k)}(T^{(k)} | \underline{C}^{(k)})$ , and the objective function is a concave function of since  $Z_i^{(k)}$  is a “log-sum-exp” function [5] of  $Z^{(k)}$ .
2. According to the second conclusion of Theorem 3.2, we set  $Z_{ij}^{(k)} = \infty$  when  $\underline{C}_{ij}^{(k)} = 0$ .

3. The first constraint is the detailed balance constraint which is equivalent to the third constraint in (3.3).
4. The last constraint means that the state potential obtained from  $Z^{(k)}$  is consistent with  $V^{(k)} = V^{(k)}(V)$ , where  $V^{(k)}(V)$  is given by (2.4).

**4.1. Approximate MLE problem and its solution.** It is easy to see that (4.6) is a convex optimization problem if we drop the free energy constraint  $V_Z(Z^{(k)}) = V^{(k)}(V)$ . This motivates an approximation method for solving (4.6) which is based on Taylor expansions and consists of two steps: First, we solve (4.6) without consideration of the free energy constraint and get the optimal solutions of free energy matrices which are denoted by  $\check{Z}^{(k)} = [\check{Z}_{ij}^{(k)}]$  with  $k = 1, \dots, K$ . Then, we replace the free energy constraint in (4.6) with its Taylor expansion around  $\check{Z}^{(k)}$  and search for the corresponding approximate optimal solution of the reference free energy. In what follows, we describe the two steps in more detail.

**4.1.1. Optimization without the free energy constraint.** Note that both  $L$  and  $Z^{(k)}$  are independent of  $V$  in (4.6) if the free energy constraint is omitted. Thus we can eliminate the variable  $V$  from the entire problem and decompose (4.6) into  $K$  subproblems:

$$(4.8) \quad \begin{aligned} & \max_{\{Z_{ij}^{(k)} | \underline{\mathcal{C}}_{ij}^{(k)} > 0\}} L_Z^{(k)} \left( Z^{(k)} | \underline{\mathcal{C}}^{(k)} \right) \\ & \text{s.t.} \quad Z^{(k)} = Z^{(k)\top}, \\ & \quad \quad \sum_{(i,j) \in \{(i,j) | \underline{\mathcal{C}}_{ij}^{(k)} > 0\}} Z_{ij}^{(k)} = 0 \end{aligned}$$

for  $k = 1, \dots, K$ . It is clear that (4.8) is a convex optimization problem with linear constraints, and the gradient and Hessian matrix of  $L_Z^{(k)}(Z^{(k)} | \underline{\mathcal{C}}^{(k)})$  can be simply obtained by the following equations:

$$(4.9) \quad \frac{\partial Z_i^{(k)}}{\partial Z_{lj}^{(k)}} = \begin{cases} T_{ij}^{(k)}, & l = i, \\ 0, & l \neq i, \end{cases}$$

$$(4.10) \quad \frac{\partial^2 Z_i^{(k)}}{\partial Z_{jm}^{(k)} \partial Z_{lm'}^{(k)}} = \begin{cases} T_{im}^{(k)} T_{im'}^{(k)} - 1_{m=m'} T_{im'}^{(k)}, & j = l = i, \\ 0 & \text{otherwise.} \end{cases}$$

Therefore, (4.8) can be efficiently solved by standard convex optimization numerical methods, and we will employ a conjugate gradient algorithm [30] to find the solution  $Z^{(k)} = \check{Z}^{(k)}$  of (4.8) in our numerical experiments.

*Remark 4.3.* It can be seen that  $\check{Z}^{(k)}$  is, in fact, the maximum likelihood estimate of  $Z^{(k)}$  given the data of the  $k$ th simulation and (4.8) can be solved independently of the other simulations, which is also the reason why we select  $\check{Z}^{(k)}$  of (4.8) as the central point of the Taylor expansions in the next step.

*Remark 4.4.* It is worth pointing out that the transition matrix estimation under the reversibility constraint is a very important problem in the Markov state modeling of stochastic simulations, but the existing MLE methods (see, e.g., [4] and [21]) often suffer from undesirable local optima and slow convergence. Here we show that the MLE of the reversible transition matrix can be formulated as a convex optimization problem, where the global optimum can be found in polynomial time.

**4.1.2. Optimization with the approximate free energy constraint.** The first-order Taylor expansion of the free energy constraint around  $\check{Z}^{(k)}$  can be expressed as

$$(4.11) \quad \check{V}^{(k)} + \nabla_{\mathcal{V}(Z^{(k)})} V_Z \left( \check{Z}^{(k)} \right) \cdot \mathcal{V} \left( Z^{(k)} - \check{Z}^{(k)} \right) = V^{(k)}(V)$$

with  $\check{V}^{(k)} = V_Z(\check{Z}^{(k)})$ , which is a linear equality. Through approximating the free energy constraint by (4.11), (4.6) can be simplified to a convex problem. However, the direct computation of the simplified problem is still too time consuming in many practical cases because it involves a large number of decision variables and cannot be decomposed into a set of small-sized subproblems as (4.8). Hence we use again the Taylor expansion to approximate the objective function, and further simplify (4.6) to a quadratic optimization problem with equality constraints:

$$(4.12) \quad \begin{aligned} & \max_{V, \{Z_{ij}^{(k)} | \underline{C}_{ij}^{(k)} > 0\}} && \sum_k \check{L}_Z^{(k)} \left( Z^{(k)} | \underline{C}^{(k)} \right) \\ & \text{s.t.} && \forall k = 1, \dots, K, \\ & && Z^{(k)} = Z^{(k)\top}, \\ & && \sum_{(i,j) \in \{(i,j) | \underline{C}_{ij}^{(k)} > 0\}} Z_{ij}^{(k)} = 0, \\ & && \mathbf{1}^\top V = 0, \\ & && \check{V}^{(k)} + \nabla_{\mathcal{V}(Z^{(k)})} V_Z \left( \check{Z}^{(k)} \right) \cdot \mathcal{V} \left( Z^{(k)} - \check{Z}^{(k)} \right) = V^{(k)}(V) \end{aligned}$$

with

$$(4.13) \quad \begin{aligned} \check{L}_Z^{(k)} \left( Z^{(k)} | \underline{C}^{(k)} \right) &= \check{L}_Z^{(k)} \left( \check{Z}^{(k)} | \underline{C}^{(k)} \right) + \nabla_{\mathcal{V}(Z^{(k)})} L_Z^{(k)} \left( \check{Z}^{(k)} | \underline{C}^{(k)} \right) \cdot \mathcal{V} \left( Z^{(k)} - \check{Z}^{(k)} \right) \\ &+ \frac{1}{2} \mathcal{V} \left( Z^{(k)} - \check{Z}^{(k)} \right)^\top \nabla_{\mathcal{V}(Z^{(k)})\mathcal{V}(Z^{(k)})} L_Z^{(k)} \left( \check{Z}^{(k)} | \underline{C}^{(k)} \right) \mathcal{V} \left( Z^{(k)} - \check{Z}^{(k)} \right). \end{aligned}$$

Applying a bilevel optimization procedure to (4.12), we can obtain the closed-form solution of  $V$ :

$$(4.14) \quad \check{V} = \Xi^{(k)} \left( \underline{C}^{(k)}, \check{\rho}^{(k)} \right) \check{V}^{(k)} + b^{(k)} \left( \underline{C}^{(k)}, \check{\rho}^{(k)} \right),$$

where  $\check{\rho}^{(k)}$  is a vector which consists of elements of  $\{\check{Z}_{ij}^{(k)} | \underline{C}_{ij}^{(k)} > 0, i \leq j, (i, j) \neq (n, n)\}$ , and  $\Xi^{(k)}(\underline{C}^{(k)}, \check{\rho}^{(k)})$  and  $b^{(k)}(\underline{C}^{(k)}, \check{\rho}^{(k)})$  are defined in (D.7) and (D.8). (The detailed optimization procedure is given in Appendix D.)

*Remark 4.5.* We can now explain why variables  $Z^{(k)}, V$  are used instead of  $T^{(k)}, \pi$  in the approximate MLE. First, the MLE problem w.r.t.  $T^{(k)}, \pi$  is a nonconvex optimization problem even if we drop the free energy constraint. Second, the quadratic approximation of the MLE problem w.r.t.  $T^{(k)}, \pi$  involves inequality constraints and therefore has no analytic solution.

**4.2. Convergence analysis.** In this section, we will analyze consistency and asymptotic normality of the approximate MLE as the exact MLE under the assumptions stated in section 3.2 and Assumption 4.1.

Before introducing the main theorem, some definitions and a lemma are needed. Let  $\rho^{(k)}$  be a vector consisting of  $\{Z_{ij}^{(k)} | \bar{X}_{ij}^{(k)} > 0, i \leq j, (i, j) \neq (n, n)\}$  with  $\bar{\rho}^{(k)}$  consisting of  $\{\bar{Z}_{ij}^{(k)} | \bar{X}_{ij}^{(k)} > 0, i \leq j, (i, j) \neq (n, n)\}$ , where  $\bar{Z}^{(k)} = [\bar{Z}_{ij}^{(k)}]$  denotes the

true value of  $Z^{(k)}$ . In this section both  $Z^{(k)}$  and  $Z_i^{(k)}$  are viewed as functions of  $\rho^{(k)}$  by considering the symmetry of  $Z^{(k)}$  and the zero-mean property of finite elements of  $Z^{(k)}$ .

LEMMA 4.6. *Provided that Assumptions 3.5–3.8 and 4.1 hold, we have the following:*

1.  $\check{\rho}^{(k)} \xrightarrow{P} \bar{\rho}^{(k)}$  and  $\check{V}^{(k)} \xrightarrow{P} \bar{V}^{(k)}$ , where  $\bar{V}^{(k)} = V^{(k)}(\bar{V})$ .
2.  $\sqrt{M}(\check{V}^{(k)} - \bar{V}^{(k)}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_V^{(k)}(\Sigma_X^{(k)}, \bar{X}^{(k)}, \bar{\rho}^{(k)}))$  if (3.12) is satisfied, where  $\Sigma_V^{(k)}(\Sigma_X^{(k)}, \bar{X}^{(k)}, \bar{\rho}^{(k)})$  is defined in (E.5).

*Proof.* See Appendix E.  $\square$

THEOREM 4.7. *Provided that Assumptions 3.5–3.8 and 4.1 hold, we have the following:*

1.  $\check{V} \xrightarrow{P} \bar{V}$ .
2.  $\sqrt{M}(\check{V} - \bar{V}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_V(\{\Sigma_X^{(k)}\}, \{\bar{X}^{(k)}\}, \{\bar{\rho}^{(k)}\}))$  if (3.12) is satisfied for all  $k$  and if  $K$  simulations are statistically independent, where  $\Sigma_V(\{\Sigma_X^{(k)}\}, \{\bar{X}^{(k)}\}, \{\bar{\rho}^{(k)}\})$  is defined in (F.3).

*Proof.* See Appendix F.  $\square$

**4.3. Error analysis.** According to Theorem 4.7, the estimation error of  $\check{V}$  follows approximately a multivariate normal distribution  $\mathcal{N}(\mathbf{0}, \Sigma_V/M)$  when  $M$  is large enough, and  $\Sigma_V(\{\Sigma_X^{(k)}\}, \{\bar{X}^{(k)}\}, \{\bar{\rho}^{(k)}\})$  can be estimated by replacing  $\bar{X}^{(k)}, \bar{\rho}^{(k)}$  with  $\hat{X}^{(k)}, \hat{\rho}^{(k)}$  if  $\Sigma_X^{(k)}$  is given. Therefore, the remaining key problem is how to estimate  $\Sigma_X^{(k)}$ . In this section we present an algorithm for estimating  $\Sigma_X^{(k)}$  based on the following assumption, which is similar to the assumption proposed in Remark 3.11 and implies that each simulation is driven by an underlying reversible Markov model.

ASSUMPTION 4.8. *For each simulation  $k$ ,  $x_t^{(k)}$  can be expressed as a function of a latent variable  $y_t^{(k)}$  with  $x_t^{(k)} = f^{(k)}(y_t^{(k)})$ , and  $\{y_t^{(k)}\}$  is a stationary, irreducible, and reversible Markov chain.*

Under Assumption 4.8, it can be seen that  $\{\mathcal{V}(\Delta C_t^{(k)})\}$  is also a stationary process, where the definition of  $\Delta C_t^{(k)} = [\Delta C_{t,ij}^{(k)}]$  is the same as in Remark 3.11. To describe the estimation algorithm, we also need some new notation. We denote by  $\kappa^{(k)}(h) = \text{Cov}(\mathcal{V}(\Delta C_t^{(k)}), \mathcal{V}(\Delta C_{t+h}^{(k)}))$  the autocovariance of  $\{\mathcal{V}(\Delta C_t^{(k)})\}$  with lag  $h$ , define  $\Gamma^{(k)}(l) = \kappa^{(k)}(2l+1) + \kappa^{(k)}(2l+2)$ , and let  $\eta^{(k)}(l)$  be a sum of  $n^2$  elements in  $\Gamma^{(k)}(l)$ , which can be represented as

$$(4.15) \quad \eta^{(k)}(l) = \sum_{i,j} \text{Cov}(\Delta C_{t,ji}^{(k)}, \Delta C_{t+2l+1,ij}^{(k)}) + \text{Cov}(\Delta C_{t,ji}^{(k)}, \Delta C_{t+2l+2,ij}^{(k)}).$$

THEOREM 4.9. *If Assumption 4.8 and (3.12) hold, and the series  $\sum_{h=0}^{\infty} \kappa^{(k)}(h)$  is convergent, then*

1.  $\Sigma_X^{(k)} = \kappa^{(k)}(0) + \sum_{l=0}^{\infty} (\Gamma^{(k)}(l) + \Gamma^{(k)}(l)^T)$ ;
2.  $\eta^{(k)}(l) \geq \|\Gamma^{(k)}(l)\|_{\max}$  and  $\eta^{(k)}(l) \leq \eta^{(k)}(l+1)$  for  $l \geq 0$ .

*Proof.* See Appendix G.  $\square$

The above theorem provides an intuitive way to estimate  $\Sigma_X^{(k)}$ :

$$(4.16) \quad \hat{\Sigma}_X^{(k)} = \hat{\kappa}^{(k)}(0) + \sum_{l=0}^{\lfloor \frac{M-3}{2} \rfloor} (\hat{\Gamma}^{(k)}(l) + \hat{\Gamma}^{(k)}(l)^T),$$

where  $\hat{\kappa}^{(k)}(0)$  and  $\hat{\Gamma}^{(k)}(l)$  denote the estimates of  $\kappa^{(k)}(0)$  and  $\Gamma^{(k)}(l)$ . We now inves-

tigate the calculation of  $\hat{\kappa}^{(k)}(0)$  and  $\hat{\Gamma}^{(k)}(l)$ .

**Estimation of  $\kappa^{(k)}(0)$ .** It is easy to verify that the  $((i-1)n+j, (m-1)n+m')$ th element of  $\kappa^{(k)}(0)$  equals  $1_{(i,j)=(m,m')} \bar{X}_{ij}^{(k)} - \bar{X}_{ij}^{(k)} \bar{X}_{mm'}^{(k)}$ ; therefore we can calculate the element in the same position in  $\hat{\kappa}^{(k)}(0)$  by  $1_{(i,j)=(m,m')} \check{X}_{ij}^{(k)} - \check{X}_{ij}^{(k)} \check{X}_{mm'}^{(k)}$  with  $\check{X}_{ij}^{(k)} \propto \exp(-\check{Z}_{ij}^{(k)})$ .

**Estimation of  $\Gamma^{(k)}(l)$ .** For  $h > 0$ ,  $\kappa^{(k)}(h)$  can be estimated by the empirical autocovariance:

(4.17)

$$\hat{\kappa}'^{(k)}(h) = \frac{1}{M-h} \sum_{t=1}^{M-h} \left( \mathcal{V}(\Delta C_t^{(k)}) - \mathcal{V}(\check{X}^{(k)}) \right) \left( \mathcal{V}(\Delta C_{t+h}^{(k)}) - \mathcal{V}(\check{X}^{(k)}) \right)^{\text{T}},$$

where  $\check{X}^{(k)}$  is an estimate of  $\mathbb{E}[\Delta C_t^{(k)}] = \bar{X}^{(k)}$ . Then  $\Gamma^{(k)}(l)$  can be estimated as

$$(4.18) \quad \hat{\Gamma}'^{(k)}(l) = \hat{\kappa}'^{(k)}(2l+1) + \hat{\kappa}'^{(k)}(2l+2).$$

However, the estimation error (4.18) will increase substantially as  $l$  approaches  $\lfloor (M-3)/2 \rfloor$ . So here we modify  $\hat{\Gamma}'^{(k)}(l)$  by correcting the corresponding estimated value of  $\eta^{(k)}(l)$ :

$$(4.19) \quad \hat{\Gamma}^{(k)}(l) = \begin{cases} \hat{\Gamma}'^{(k)}(l), & l = 0, \\ \min \left\{ \frac{\hat{\eta}^{(k)}(l-1)}{\hat{\eta}'^{(k)}(l)}, 1 \right\} \cdot \hat{\Gamma}'^{(k)}(l), & l > 1 \text{ and } \hat{\eta}'^{(k)}(l) > 0, \\ 0, & l > 1 \text{ and } \hat{\eta}'^{(k)}(l) \leq 0, \end{cases}$$

where  $\hat{\eta}'^{(k)}(l)$  and  $\hat{\eta}^{(k)}(l)$  denote the values of  $\eta^{(k)}(l)$  obtained from  $\hat{\Gamma}'^{(k)}(l)$  and  $\hat{\Gamma}^{(k)}(l)$ . It can be seen that  $\hat{\eta}^{(k)}(l)$  is nonnegative and decreasing with  $l$ , which is consistent with the conclusion of Theorem 4.9. Besides, we can show that  $\hat{\Gamma}^{(k)}(l) \equiv 0$  for  $l \geq l'$  if  $\hat{\eta}^{(k)}(l') \leq 0$ . Thus the estimator of  $\Sigma_X^{(k)}$  in this section is, in fact, a time window estimator [9], where the large-lag terms outside the window are set to be zero, and the window size  $l_w = \min\{l | \hat{\eta}^{(k)}(l) \leq 0\}$  implies that the curve of  $\|\Gamma^{(k)}(l)\|_{\max}$  goes below the noise level at  $l = l_w$ .

*Remark 4.10.* From the definition of  $\Sigma_X^{(k)}$  we can deduce that  $\Sigma_X^{(k)} \succeq 0$  and  $\mathbf{1}^{\text{T}} \Sigma_X^{(k)} \mathbf{1} = 0$ , but the  $\hat{\Sigma}_X^{(k)}$  obtained by (4.16) may not satisfy the constraints. For this problem, we can correct the value of  $\hat{\Sigma}_X^{(k)}$  as  $\hat{\Sigma}_X^{(k)} := \mathcal{M}_P \circ \hat{\Sigma}_X^{(k)}$ , where  $\mathcal{M}_P \circ G = (\mathbf{I} - \frac{1}{m} \mathbf{1}^{\text{T}} \mathbf{1})(G - \min\{\lambda_{\min}(G), 0\} \mathbf{I})(\mathbf{I} - \frac{1}{m} \mathbf{1}^{\text{T}} \mathbf{1})$  for a symmetric matrix  $G \in \mathbb{R}^{m \times m}$  with the smallest eigenvalue  $\lambda_{\min}(G)$ . It is easy to see that  $\mathcal{M}_P$  is a mapping from the symmetric matrix set to the set  $\{G | G \succeq 0, \mathbf{1}^{\text{T}} G \mathbf{1} = 0\}$ , and  $\mathcal{M}_P \circ G = G$  if  $G \succeq 0$  and  $\mathbf{1}^{\text{T}} G \mathbf{1} = 0$  hold.

**4.4. Comparison to related work.** It is interesting to compare the approximate MLE with the weighted histogram analysis method (WHAM), which estimates  $\pi_i$  as

$$(4.20) \quad \hat{\pi}_i^{\text{WHAM}} \propto \sum_k c_i^{(k)} \hat{\pi}_i^{(k)},$$

where  $\hat{\pi}^{(k)} = [\hat{\pi}_i^{(k)}]$  is the estimate of  $\pi^{(k)}$  obtained from the histogram of simulation  $k$ , which has the same definition as in (3.8), and  $c_i^{(k)}$  is the weight of  $\hat{\pi}_i^{(k)}$ , which

is selected so as to minimize the statistical error (see [26] for more details). From (4.14) and (4.20), it can be observed that both the approximate maximum likelihood estimator and the WHAM estimator can be expressed as linear combinations of  $K$  “local” estimators. The differences between them are the following: (1) In WHAM, the biased estimate  $\hat{\pi}^{(k)}$  of each simulation  $k$  is a “valid” estimate only if the global equilibrium is reached in the simulation. The approximate MLE overcomes this limitation by using a transition matrix based algorithm to get biased estimates of  $\{V^{(k)}\}$ , where the dynamic information contained in simulation trajectories can be exploited to improve the estimation accuracy and only the local equilibrium assumption is required. (2) The combination weights in the approximate MLE are designed under the Markov assumption, which is more “reasonable” than the independent and identically distributed (i.i.d.) assumption used in WHAM because the biased simulations are autocorrelated in most practical cases. (The statistical error in WHAM is derived under the assumption that all  $\{x_t^{(k)}\}$  are i.i.d. processes.)

Another Markov modeling motivated method for the estimation problem of multiple biased simulations is the multiple Markov transition matrix method (MMMM) developed in [23], which performs the estimation of  $T^{(k)}, \pi^{(k)}$  and the combination of biased estimation results in an approximate Bayesian manner. In contrast to the proposed approximate MLE, MMMM suffers from the following disadvantages: (1) The detailed balance condition is not considered in MMMM, which might lead to relatively poor estimates in local estimations. (The influence of the detailed balance condition on the Bayesian inference of Markov models was studied in [17, 15].) (2) The combination operation in MMMM involves a highly nonlinear and nonconvex optimization problem and is therefore numerically unstable, which is caused by the nonlinear relationship between  $\pi$  and  $\pi^{(k)}$ . The approximate MLE avoids this problem by constructing the optimization model w.r.t.  $Z^{(k)}, V$  instead of  $T^{(k)}, \pi$  (see Remark 4.5). Furthermore, the convergence and the estimation error of MMMM in non-Markovian cases have not been analyzed.

**5. Extension to the general case.** In previous discussions, we have described the MLE of the reference stationary distribution or free energy from multiple biased simulations under the condition that all simulations share the same state space and simulation length. Actually, this restriction can be removed easily by extension of the proposed methods and results.

In this section, we consider a more general case where simulation trajectories are generated in the same way as in section 2, except that here we run a biased simulation from time 0 to time  $M^{(k)}$  on a state space  $\mathcal{S}^{(k)} \subseteq \mathcal{S}$  with  $|\mathcal{S}^{(k)}| = n^{(k)}$  for each  $k$ , and we do not assume that  $\mathcal{S}^{(1)} = \dots = \mathcal{S}^{(K)}$  or  $M^{(1)} = \dots = M^{(K)}$ . Then the  $K$  simulation trajectories can be represented as  $\{x_t^{(1)}\}_{t=0}^{M^{(1)}}, \dots, \{x_t^{(K)}\}_{t=0}^{M^{(K)}}$ , where we use  $x_t^{(k)} = i$  ( $i \in \{1, \dots, n^{(k)}\}$ ) to denote that the  $i$ th state in  $\mathcal{S}^{(k)}$  is observed at time  $t$  in the  $k$ th simulation (all  $\mathcal{S}^{(k)}$  are assumed to be ordered sets), and the relationships between  $V^{(k)}, \pi^{(k)}$  and  $V, \pi$  can be written as

$$(5.1) \quad V_i^{(k)} = V_{j_i} + U_{j_i}^{(k)} - \frac{1}{n^{(k)}} \sum_{m=1}^{n^{(k)}} \left( V_{j_m} + U_{j_m}^{(k)} \right)$$

and  $\pi_i^{(k)} \propto \exp(-U_{j_i}^{(k)})\pi_{j_i}$  if the  $i$ th state in  $\mathcal{S}^{(k)}$  corresponds to the state  $j_i$  in  $\mathcal{S}$ . With the above notation and results, (3.3), (4.8), and (4.14) can be directly used to perform the MLE and approximate MLE of  $\pi$  and  $V$ .

We now investigate the convergence of the results of the MLE and approximate MLE in the general case. For convenience of analysis, we let  $M = (\sum_k M^{(k)})/K$  such that the convergence still can be stated w.r.t.  $M \rightarrow \infty$ . Moreover, here we make a new assumption on simulation lengths.

**ASSUMPTION 5.1.** *The limit  $\hat{w}^{(k)} \xrightarrow{P} \bar{w}^{(k)}$  exists and  $\bar{w}^{(k)} > 0$  for all  $k = 1, \dots, K$ , where  $\hat{w}^{(k)} = M^{(k)}/M$ .*

Based on Assumption 5.1 and assumptions considered in sections 3.2 and 4.2, we can prove the following theorems on the MLE and approximate MLE in the general case.

**THEOREM 5.2.** *If Assumptions 3.5–3.8 and 5.1 hold, then  $\tilde{\theta} \xrightarrow{P} \bar{\theta}$ .*

**THEOREM 5.3.** *If Assumptions 3.5–3.8 and 5.1 hold and the conditions*

1. *for each  $k$ , there exists a  $\Sigma_X^{(k)}$  such that*

$$(5.2) \quad \sqrt{M^{(k)}} \left( \mathcal{V} \left( \hat{X}^{(k)} \right) - \mathcal{V} \left( \bar{X}^{(k)} \right) \right) \xrightarrow{d} \mathcal{N} \left( \mathbf{0}, \Sigma_X^{(k)} \right);$$

2. *conditions 2–5 in Theorem 3.10 hold;*
3.  *$\sqrt{M}(\hat{w}^{(k)} - \bar{w}^{(k)}) \xrightarrow{P} 0$*

*are satisfied, then (3.13) holds with  $\Sigma_X = \text{diag}(\bar{w}^{(1)}\Sigma_X^{(1)}, \dots, \bar{w}^{(K)}\Sigma_X^{(K)})$ .*

**THEOREM 5.4.** *Provided that Assumptions 3.5–3.8, 4.1, and 5.1 hold, we have the following:*

1.  *$\hat{V} \xrightarrow{P} \bar{V}$ .*
2.  *$\sqrt{M}(\hat{V} - \bar{V}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_V(\{\bar{w}^{(k)}\Sigma_X^{(k)}\}, \{\bar{X}^{(k)}\}, \{\bar{\rho}^{(k)}\}))$  if (5.2) is satisfied for all  $k$  and if  $K$  simulations are statistically independent, where  $\bar{\rho}^{(k)}$  is a vector consisting of  $\{\bar{Z}_{ij}^{(k)} | \bar{X}_{ij}^{(k)} > 0, i \leq j, (i, j) \neq (n^{(k)}, n^{(k)})\}$  and  $\Sigma_V(\cdot)$  has the same definition as in Theorem 4.7.*

*Remark 5.5.* Detailed optimization algorithms and proofs of theorems are all omitted in this section because they are similar to those in sections 3 and 4.

**6. Numerical experiments.** In this section, the approximate MLE proposed in this paper will be applied to some numerical examples of multiple biased simulations, and the performance will be compared to that of WHAM and MMMM. For convenience, here we denote a set of multiple biased simulations described in section 2 by MBS( $K, M$ ).

**6.1. Umbrella sampling with Markovian simulations.** Umbrella sampling is a commonly used biased simulation technique, where each biasing potential (also called “umbrella potential”) is designed to confine the system around some region of state space and achieve a more efficient sampling, especially at transition states which the unbiased simulation would visit only rarely. In this example, the umbrella sampling simulations are employed on a reference system with state set  $\mathcal{S} = \{s_i = -5 + 10(i - 1)/99 | i = 1, \dots, 100\}$  and free energy  $V = [V_i] = [0.25s_i^4 - 5s_i^2 - 9.9874]$ . As shown in Figure 6.1, the reference system has two metastable states centered at  $A$  and  $B$ , and the switching between metastable states is blocked by an energy barrier with peak position  $O$ .

For umbrella sampling simulations, we design the following 15 different biased potentials:  $U^{(k)} = [U_i^{(k)}] = [4(s_i + \frac{15}{14}k - \frac{60}{7})^2]$  for  $1 \leq k \leq 15$ . Note that these potentials will be repeatedly used if the simulation number is larger than 15; i.e.,  $U^{(k)} = U^{((k-1) \bmod 15 + 1)}$  if  $k > 15$ . The simulation trajectory  $\{x_t^{(k)}\}_{t=0}^M$  is generated by a Metropolis simulation model, which is a reversible Markov chain with initial distribution  $\text{Pr}(x_0^{(k)} = s_i) \propto \exp(-U_i^{(k)})$ , stationary distribution  $\pi_i^{(k)} \propto \exp(-V_i -$



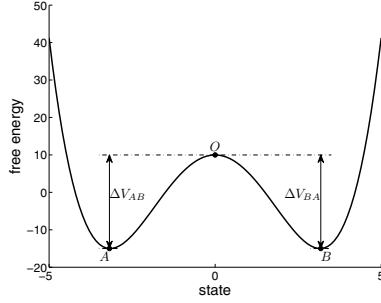


FIG. 6.1. Free energy profile of the reference system, which has two potential wells with minima at  $A$  and  $B$  separated by an energy barrier. The highest energy position  $O$  of the barrier represents the transition state, and the energy barrier heights for transitions  $A \rightarrow B$  and  $B \rightarrow A$  are defined as  $\Delta V_{AB} = |V_O - V_A|$  and  $\Delta V_{BA} = |V_O - V_B|$  with  $V_A$ ,  $V_B$ , and  $V_O$  the potentials of  $A$ ,  $B$ , and  $O$ .

$U_i^{(k)}$ ), and transition probability

$$(6.1) \quad \Pr(x_{t+1}^{(k)} = s_j | x_t^{(k)} = s_i) = \begin{cases} \min \left\{ \frac{\exp(-V_j - U_j^{(k)})}{\exp(-V_i - U_i^{(k)})} q_{ji}, q_{ij} \right\}, & i \neq j, \\ 1 - \sum_{l \neq i} \Pr(x_{t+1}^{(k)} = s_l | x_t^{(k)} = s_i), & i = j, \end{cases}$$

where  $q_{ij}$  satisfies  $q_{ij} \propto 1_{|i-j| \leq 2}$  and  $\sum_j q_{ij} = 1$ .

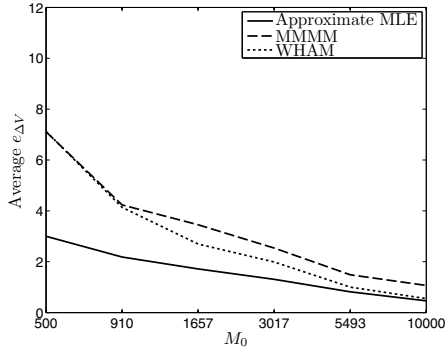
The comparisons between the estimation methods are based on the mean error of approximations of energy barrier heights:

$$(6.2) \quad e_{\Delta V} = \frac{1}{2} (|\Delta V_{AB} - \Delta V_{AB}^{\text{approx}}| + |\Delta V_{BA} - \Delta V_{BA}^{\text{approx}}|),$$

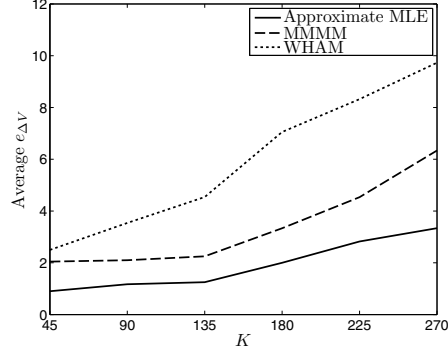
where the definitions of  $\Delta V_{AB}$  and  $\Delta V_{BA}$  are given in Figure 3.1, and the superscript “approx” represents the approximate value obtained from the estimated  $V$ .

We first set  $K = 15$  and  $M = 500, 910, 1657, 3017, 5493, 10000$ , and perform 30 independent runs of MBS ( $K, M$ ) for each value of  $M$ . Figure 6.2(a) displays the average  $e_{\Delta V}$  of the approximate MLE, MMMM, and WHAM for different  $M$ , and Figures 6.2(c) and 6.2(d) show the estimates of  $V$  obtained from a run of MBS (15, 500) and MBS (15, 10000). It can be seen that the estimation errors of all three methods decrease with increasing simulation length, and the proposed approximate MLE performs significantly better than the other two methods. Note that MMMM is also a Markov chain model based method, but its performance turns out to be worse than WHAM in this numerical experiment, especially for large simulation lengths  $M$ .

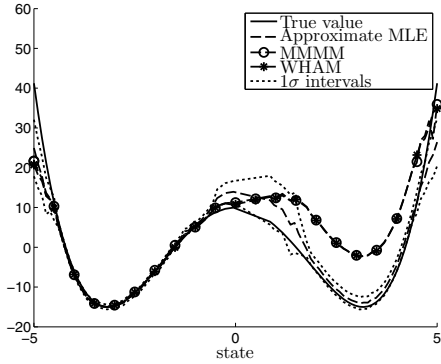
Next, we validate whether the estimation methods can reconstruct the free energy  $V$  from very short simulations. Here we fix the total simulation time  $MK$ , and set  $M = \lceil 22500/K \rceil$  with  $K = 45, 90, 135, 180, 225, 270$ . The estimation results are summarized in Figures 6.2(b), 6.2(e), and 6.2(f). (The  $1\sigma$  confidence intervals in Figure 6.2(f) are provided by using the sample standard deviation of  $\check{V}$  calculated from the 30 independent runs of MBS (270, 83) because the simulation length is too short such that the error analysis approach in section 4.3 is not applicable.) It should be noted that the equilibrium assumption used by WHAM does not hold if  $M$  is too small, because the initial distributions of simulations differ from the biased stationary



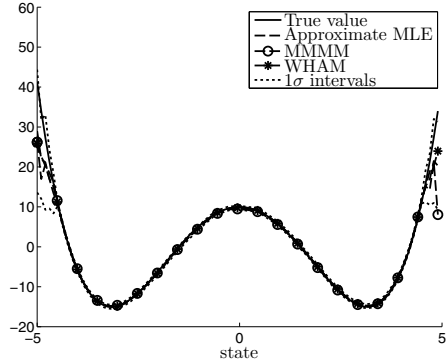
(a) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 15$  and  $M = 500, 910, 1657, 3017, 5493, 10000$ .



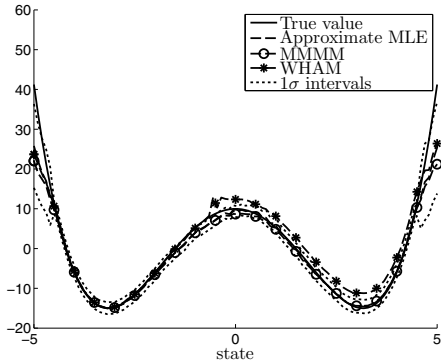
(b) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 45, 90, 135, 180, 225, 270$  and  $M = \lceil 22500/K \rceil$ .



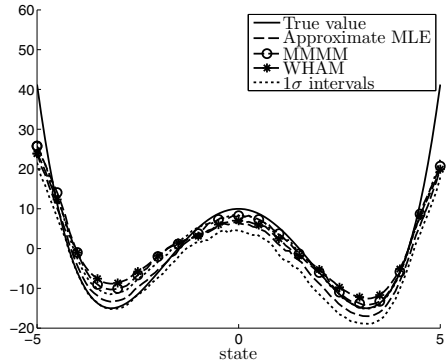
(c) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 5000), where the  $e_{\Delta V}$  of approximate MLE = 3.5015,  $e_{\Delta V}$  of MMMM = 6.3153, and  $e_{\Delta V}$  of WHAM = 6.3500.



(d) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 10000), where the  $e_{\Delta V}$  of approximate MLE = 0.3060,  $e_{\Delta V}$  of MMMM = 0.6002, and  $e_{\Delta V}$  of WHAM = 0.3397.



(e) Estimates of  $V$  generated by the different estimators on a run of MBS(45, 5000), where the  $e_{\Delta V}$  of approximate MLE = 0.4570,  $e_{\Delta V}$  of MMMM = 1.5012, and  $e_{\Delta V}$  of WHAM = 1.8948.



(f) Estimates of  $V$  generated by the different estimators on a run of MBS(270, 83), where the  $e_{\Delta V}$  of approximate MLE = 3.1939,  $e_{\Delta V}$  of MMMM = 4.2559, and  $e_{\Delta V}$  of WHAM = 7.2764.

FIG. 6.2. Estimation results of umbrella sampling with Markovian simulations. The  $1\sigma$  confidence intervals in (c), (d), and (e) are obtained by the approach described in section 4.3, and those in (f) are obtained from the sample standard deviation of  $\bar{V}$  in the 30 independent runs.

distributions. Therefore the estimation accuracy of WHAM is reduced when the individual simulation lengths are shorter, although the total data size stays almost the same. In contrast, the proposed approximate MLE and MMMM are less affected by the change in the length of individual simulations. This is because these methods rely on having local rather than global equilibrium assumptions. Furthermore, the proposed method outperforms both WHAM and MMMM in this numerical experiment.

**6.2. Umbrella sampling with non-Markovian simulations.** We now consider the estimation problem from an umbrella sampling simulation in the case that the Markov assumption does not hold, i.e., the bins used to estimate the free energy do not correspond to the Markov states of the underlying simulation. The simulation model and the other settings in this section are basically the same as in section 6.1 except that the state set is defined as  $\mathcal{S} = \{\bar{s}_1, \dots, \bar{s}_{10}\}$  with  $\bar{s}_1 = \{s_1, \dots, s_{10}\}$ ,  $\bar{s}_2 = \{s_{11}, \dots, s_{15}\}$ ,  $\bar{s}_3 = \{s_{16}, \dots, s_{20}\}$ ,  $\dots$ ,  $\bar{s}_{17} = \{s_{86}, \dots, s_{90}\}$ , and  $\bar{s}_{18} = \{s_{91}, \dots, s_{100}\}$ . It is clear that the observed state sequences in simulations do not satisfy the Markov property with this definition of states.

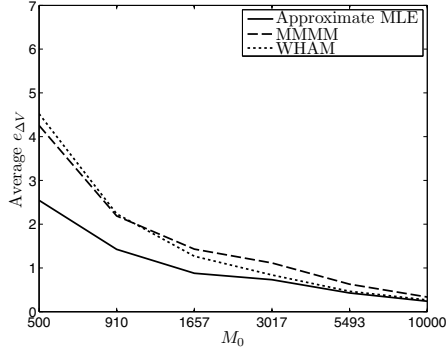
We utilize the three methods to approximate the free energy  $V$  by using the non-Markovian simulation data, and the estimation results with different  $(K, M)$  are shown in Figure 6.3, where  $e_{\Delta V}$  is defined in the same way as in section 6.1 with  $A, B$  and  $O$  the local minimum and peak positions in  $\mathcal{S}$ . As observed from the figures, the estimates obtained from the approximate MLE are more precise than those obtained from the other estimators for various values of  $(K, M)$ .

**6.3. Metadynamics with Markovian simulations.** Metadynamics is another biased simulation technique often employed in computational physics and chemistry, which is able to escape local free energy minima and improve the searching properties of simulations through iteratively modifying the biasing potential. Given  $K, M$  and a reference system as in section 6.1, a metadynamics procedure can also be expressed as a run of MBS( $K, M$ ) with  $U_i^{(k)} = 0$  for  $k = 1$  and  $U_i^{(k)} = U_i^{(k-1)} + u_c(s_i | x_M^{(k-1)})$  for  $k > 1$ , where  $u_c(s|x)$  denote a Gaussian function of  $s$  centered at  $x$ . Thus, for each of the  $K$  simulations in an MBS run, a Gaussian hat is added to the potential at the last point of the previous simulation. This effectively fills up the potential energy basins with increasing  $k$ . Ultimately the effective potential becomes approximately flat. Here we define  $u_c(s|x) = 5 \exp(-(s-x)^2)$ , and the simulation data  $\{x_t^{(k)}\}_{t=0}^M$  is also generated by the Metropolis sampling model with  $x_0^{(k)} = x_M^{(k-1)}$ .

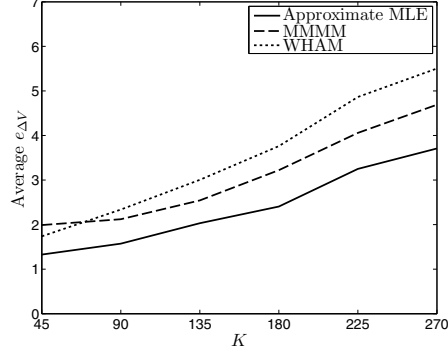
The three estimation methods are applied to reconstruct the free energy of the reference system by data generated by metadynamics with different  $(K, M)$ , and the estimation results are shown in Figure 6.4. The superior performance of the presented method is clearly evident from the figures.

**6.4. Metadynamics with non-Markovian simulations.** In this example, the free energy estimation problem of metadynamics with non-Markovian simulations is investigated. We generate the simulation data as in section 6.3 and convert the state sequences to non-Markovian processes as in section 6.2. Then the three methods can be used to estimate the unbiased free energy of states  $\bar{s}_1, \dots, \bar{s}_{18}$ .

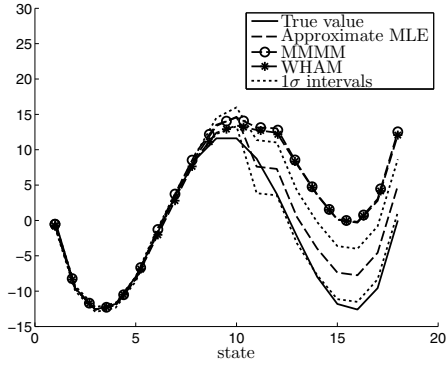
All the estimation results are displayed in Figure 6.5. It is obvious that the approximate MLE does a much better job in the free energy estimation than the other two methods.



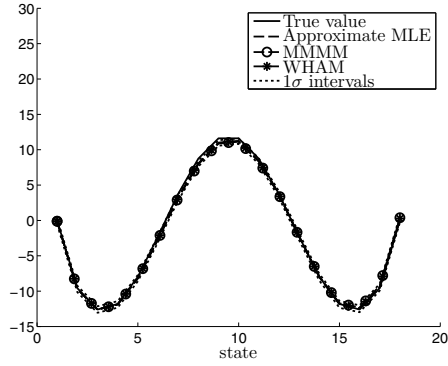
(a)  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 15$  and  $M = 500, 910, 1657, 3017, 5493, 10000$ .



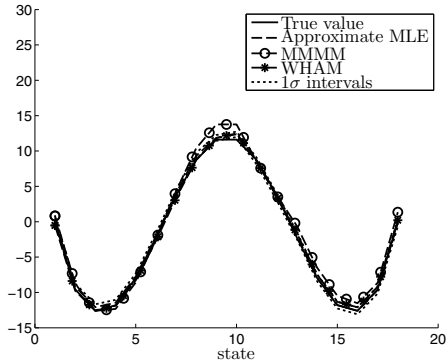
(b) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 45, 90, 135, 180, 225, 270$  and  $M = [22500/K]$ .



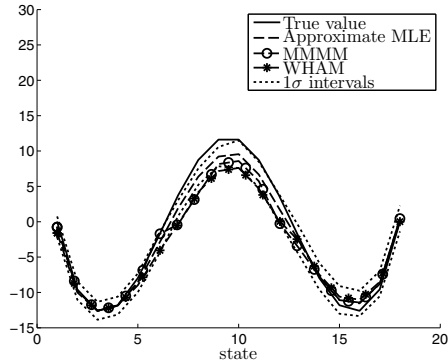
(c) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 500), where the  $e_{\Delta V}$  of approximate MLE = 2.4208,  $e_{\Delta V}$  of MMMM = 6.2008, and  $e_{\Delta V}$  of WHAM = 6.1413.



(d) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 10000), where the  $e_{\Delta V}$  of approximate MLE = 0.4532,  $e_{\Delta V}$  of MMMM = 0.6289, and  $e_{\Delta V}$  of WHAM = 0.4711.

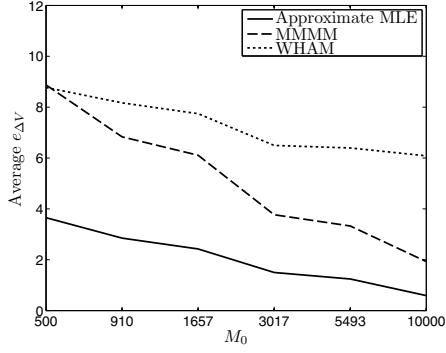


(e) Estimates of  $V$  generated by the different estimators on a run of MBS(45, 500), where the  $e_{\Delta V}$  of approximate MLE = 0.4989,  $e_{\Delta V}$  of MMMM = 1.6107, and  $e_{\Delta V}$  of WHAM = 0.6257.

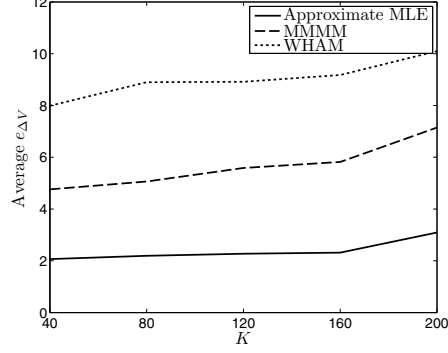


(f) Estimates of  $V$  generated by the different estimators on a run of MBS(270, 83), where the  $e_{\Delta V}$  of approximate MLE = 2.5862,  $e_{\Delta V}$  of MMMM = 3.6059, and  $e_{\Delta V}$  of WHAM = 4.7896.

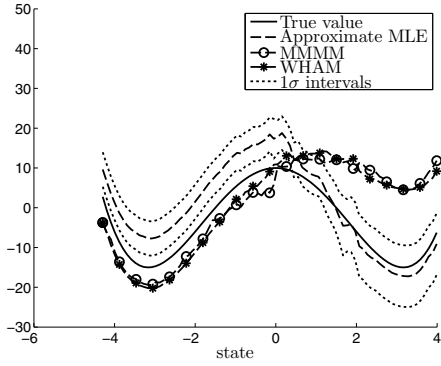
FIG. 6.3. Estimation results of umbrella sampling with non-Markovian simulations. The  $1\sigma$  confidence intervals in (c), (d), and (e) are obtained by the approach described in section 4.3, and those in (f) are obtained from the sample standard deviation of  $\hat{V}$  in the 30 independent runs.



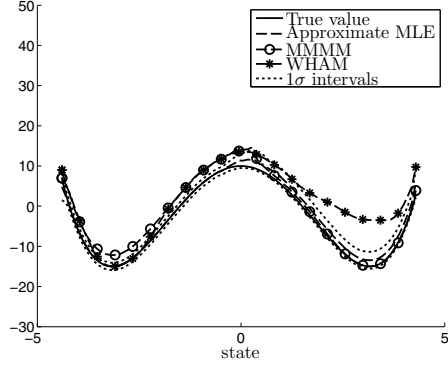
(a) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 15$  and  $M = 500, 910, 1657, 3017, 5493, 10000$ .



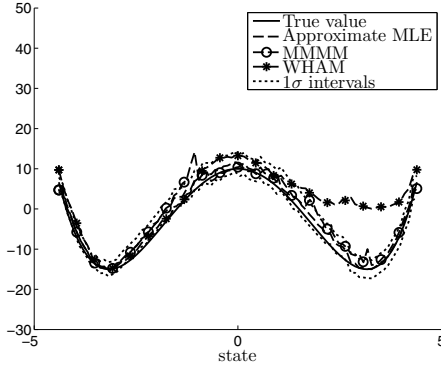
(b) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 40, 80, 120, 160, 200$  and  $M = \lceil 2000/K \rceil$ .



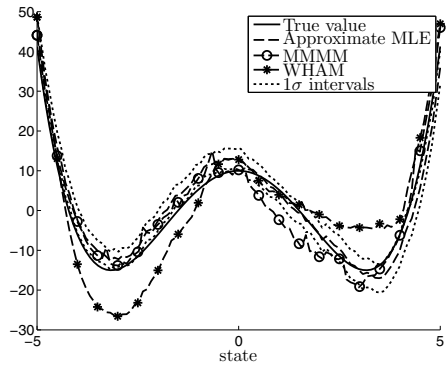
(c) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 500), where the  $e_{\Delta V}$  of approximate MLE = 6.2858,  $e_{\Delta V}$  of MMMM = 11.9022, and  $e_{\Delta V}$  of WHAM = 12.4188.



(d) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 10000), where the  $e_{\Delta V}$  of approximate MLE = 0.8246,  $e_{\Delta V}$  of MMMM = 3.1642, and  $e_{\Delta V}$  of WHAM = 5.7601.

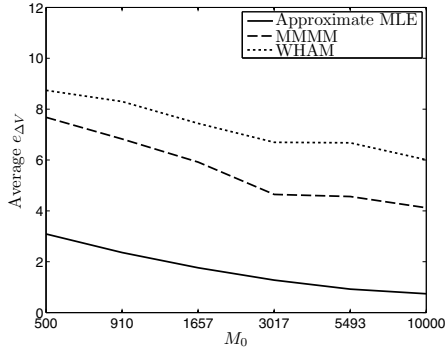


(e) Estimates of  $V$  generated by the different estimators on a run of MBS(40, 50), where the  $e_{\Delta V}$  of approximate MLE = 1.6589,  $e_{\Delta V}$  of MMMM = 3.6707, and  $e_{\Delta V}$  of WHAM = 7.5342.

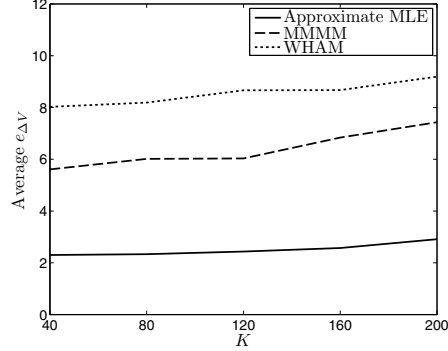


(f) Estimates of  $V$  generated by the different estimators on a run of MBS(200, 10), where the  $e_{\Delta V}$  of approximate MLE = 2.8518,  $e_{\Delta V}$  of MMMM = 6.5913, and  $e_{\Delta V}$  of WHAM = 10.9540.

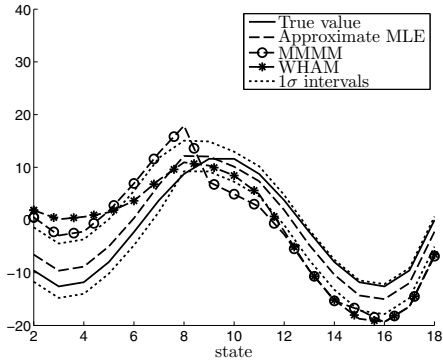
FIG. 6.4. Estimation results of metadynamics with Markovian simulations. The  $1\sigma$  confidence intervals in (c) and (d) are obtained by the approach described in section 4.3, and those in (e) and (f) are obtained from the sample standard deviation of  $\tilde{V}$  in the 30 independent runs.



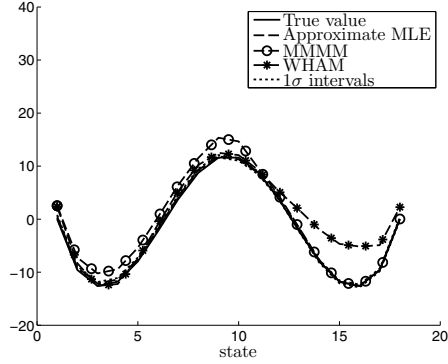
(a) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 15$  and  $M = 500, 910, 1657, 3017, 5493, 10000$ .



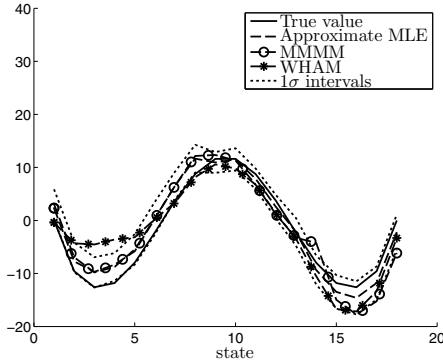
(b) Average  $e_{\Delta V}$  calculated over 30 independent runs of MBS( $K, M$ ) for  $K = 40, 80, 120, 160, 200$  and  $M = \lceil 2000/K \rceil$ .



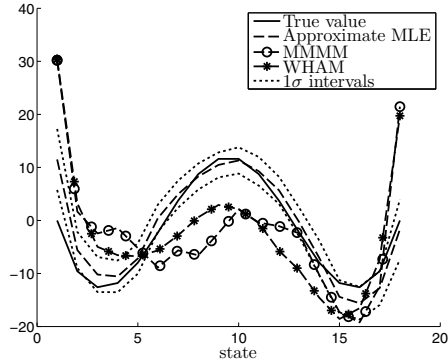
(c) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 500), where the  $e_{\Delta V}$  of approximate MLE = 2.6950,  $e_{\Delta V}$  of MMMM = 8.2297, and  $e_{\Delta V}$  of WHAM = 9.7315.



(d) Estimates of  $V$  generated by the different estimators on a run of MBS(15, 10000), where the  $e_{\Delta V}$  of approximate MLE = 0.2325,  $e_{\Delta V}$  of MMMM = 2.5563, and  $e_{\Delta V}$  of WHAM = 3.7233.



(e) Estimates of  $V$  generated by the different estimators on a run of MBS(40, 50), where the  $e_{\Delta V}$  of approximate MLE = 2.4020,  $e_{\Delta V}$  of MMMM = 4.0281, and  $e_{\Delta V}$  of WHAM = 6.3157.



(f) Estimates of  $V$  generated by the different estimators on a run of MBS(200, 10), where the  $e_{\Delta V}$  of approximate MLE = 2.5255,  $e_{\Delta V}$  of MMMM = 8.0936, and  $e_{\Delta V}$  of WHAM = 8.6291.

FIG. 6.5. Estimation results of metadynamics with Markovian simulations. The  $1\sigma$  confidence intervals in (c) and (d) are obtained by the approach described in section 4.3, and those in (e) and (f) are obtained from the sample standard deviation of  $\hat{V}$  in the 30 independent runs.

**7. Conclusions.** We have presented a transition matrix based estimation method for stationary distributions or free energy profiles using data from biased simulations, such as umbrella sampling or metadynamics. In contrast to existing estimators such as the weighted histogram analysis method (WHAM), the present estimator is not based on absolute counts in histogram bins, but rather is based on the transition counts between an arbitrary state space discretization. This discretization may be in a single or a few order parameters, e.g., those order parameters in which the umbrella sampling or metadynamics simulations are driven, or they may come from the clustering of a higher-dimensional space, such as is frequently used in Markov modeling. The only condition is that the energy bias used in the biased simulations can be associated to the discrete states, suggesting that at least the order parameters used to drive the umbrella sampling/metadynamics simulation should be discretized finely. The stationary probabilities or free energies are then reconstructed on the discrete states used. The estimator presented here has a number of advantages over existing methods such as WHAM. Most importantly, in all scenarios tested here, the estimation error of the transition matrix based estimator was significantly smaller than that of existing estimation methods. The reason for this is that the estimator does not rely on the biased simulation to fully equilibrate within one simulation condition, but asks only for local equilibrium in the discrete states, which is a much weaker requirement. As a consequence, the present method can also be used to estimate free energy profiles and stationary distributions from metadynamics simulations using all simulation data. Previously, metadynamics simulations could only be analyzed using the fraction of the simulation generated after the free energy minima have been filled and the simulation samples from an approximately flat free energy landscape. These advantages may lead to very substantial savings of CPU time for a given system and, in addition, permit the simulation of systems that were otherwise out of reach.

**Appendix A. Proof of Theorem 3.2.** For convenience, here we define  $\Theta$  to be the solution set defined by constraints in (3.3),  $\Theta_1$  to be the set of feasible solutions which satisfies  $T_{ij}^{(k)} = 0$  for  $(i, j, k) \in \{(i, j, k) | C_{ij}^{(k)} + C_{ji}^{(k)} = 0 \text{ and } i \neq j\}$ , and  $\Theta_2$  to be the set of feasible solutions which satisfies  $1_{T_{ij}^{(k)} > 0} = 1_{C_{ij}^{(k)} > 0}$  for all  $i, j, k$ .

*Part (1).* In this part, we will prove the optimal solution existence of (3.3). Suppose that  $(\pi', T'^{(1)}, \dots, T'^{(K)})$  is a feasible solution with objective value  $L' > -\infty$ . We can define a new objective function  $L_+(\pi, T^{(1)}, \dots, T^{(K)}) = \max\{L(\pi, T^{(1)}, \dots, T^{(K)}), L' - a\}$ , where  $a > 0$  is a constant. It is easy to verify that  $L_+$  is a continuous function on  $\Theta$ . Thus, the optimization problem  $\max_{(\pi, T^{(1)}, \dots, T^{(K)}) \in \Theta} L_+(\pi, T^{(1)}, \dots, T^{(K)})$  has a global optimal solution  $(\pi'', T''^{(1)}, \dots, T''^{(K)})$  with  $L_+(\pi'', T''^{(1)}, \dots, T''^{(K)}) = L''$  because  $\Theta$  is a closed set. Noting that  $L'' \geq L' > L' - a$ , we have  $L(\pi'', T''^{(1)}, \dots, T''^{(K)}) = L''$ . Therefore, for any  $(\pi, T^{(1)}, \dots, T^{(K)}) \in \Theta$ , we have  $L(\pi, T^{(1)}, \dots, T^{(K)}) \leq L_+(\pi, T^{(1)}, \dots, T^{(K)}) = L(\pi'', T''^{(1)}, \dots, T''^{(K)})$ .

*Part (2).* In this part, we will prove the first conclusion of the theorem. Suppose that  $(\pi', T'^{(1)}, \dots, T'^{(K)})$  is an optimal solution. We can define a new solution  $(\pi'', T''^{(1)}, \dots, T''^{(K)})$  with

$$(A.1) \quad T''_{ij}^{(k)} = \begin{cases} 1_{C_{ij}^{(k)} + C_{ji}^{(k)} > 0} \cdot T'_{ij}^{(k)}, & i \neq j, \\ 1 - \sum_{l \neq i} T''_{il}^{(k)}, & i = j. \end{cases}$$

Obviously,  $(\pi'', T''^{(1)}, \dots, T''^{(K)})$  is a feasible solution belonging to  $\Theta_1$ , and  $T''_{ii}^{(k)} \geq$

$T''_{ii}{}^{(k)}$ . We have

$$\begin{aligned} L\left(\pi', T''^{(1)}, \dots, T''^{(K)}\right) &= L\left(\pi', T'^{(1)}, \dots, T'^{(K)}\right) + \sum_{i,k} C_{ii}^{(k)} \left(\log T''_{ii}{}^{(k)} - \log T'^{(k)}\right) \\ (A.2) \qquad \qquad \qquad &\geq L\left(\pi', T'^{(1)}, \dots, T'^{(K)}\right). \end{aligned}$$

Therefore,  $(\pi', T''^{(1)}, \dots, T''^{(K)})$  is also an optimal solution.

*Part (3).* We now prove the second conclusion. Suppose there is an optimal solution  $(\pi', T'^{(1)}, \dots, T'^{(K)})$  belonging to  $\Theta_1 \setminus \Theta_2$ . Then there exist  $i, j, k$  such that  $T'_{ij}{}^{(k)} = 0$  and  $C_{ij}^{(k)} > 0$ , and  $L(\pi', T'^{(1)}, \dots, T'^{(K)}) = -\infty$ . This leads to a contradiction with the optimality of  $(\pi', T'^{(1)}, \dots, T'^{(K)})$ . Thus, the optimal solution belonging to  $\Theta_1$  must be an element of  $\Theta_2$  if  $C_{ii}^{(k)} > 0$  and  $1_{C_{ij}^{(k)} > 0} = 1_{C_{ji}^{(k)} > 0}$  for all  $i, j, k$ .

**Appendix B. Proof of Theorem 3.9.** Let  $\Theta$  be the feasible set of  $\theta$  defined by constraints in (3.3),  $\hat{Q}(\theta) = \sum_k L^{(k)}(T^{(k)} | \hat{X}^{(k)})$ , and  $\bar{Q}(\theta) = \sum_k L^{(k)}(T^{(k)} | \bar{X}^{(k)})$ .

From Assumption 3.7, we have  $\hat{X}_{ij}^{(k)} = 0$  if  $\bar{X}_{ij}^{(k)} = 0$ . Then  $\hat{Q}(\bar{\theta}), \bar{Q}(\bar{\theta}) > -\infty$ , and we can define the following new functions:  $\hat{Q}_+(\theta) = \max\{\hat{Q}(\theta), \hat{Q}(\bar{\theta}) - a\}$  and  $\bar{Q}_+(\theta) = \max\{\bar{Q}(\theta), \bar{Q}(\bar{\theta}) - a\}$ , where  $a > 0$  is a constant.

*Part (1).* First, we will prove that  $\bar{\theta}$  is the unique solution of  $\max_{\theta \in \Theta} \bar{Q}(\theta)$ . We note that

$$(B.1) \qquad \bar{Q}(\theta) = - \sum_k \text{KLR}_{\pi^{(k)}} \left( \bar{T}^{(k)} || T^{(k)} \right) + \sum_{i,j,k} \bar{X}_{ij}^{(k)} \log \bar{T}_{ij}^{(k)}.$$

According to the property of the KL divergence rate and Assumption 3.5,  $\bar{Q}(\theta)$  can achieve the maximal value if and only if  $T^{(k)} = \bar{T}^{(k)}$  for all  $k$ . Then we can conclude from Assumption 3.8 that  $\theta = \bar{\theta}$  is the unique solution of  $\max_{\theta \in \Theta} \bar{Q}(\theta)$ .

*Part (2).* It is easy to verify that

$$(B.2) \qquad \theta = \arg \max_{\theta \in \Theta} \hat{Q}(\theta) \Leftrightarrow \theta = \arg \max_{\theta \in \Theta} \hat{Q}_+(\theta)$$

and  $\bar{\theta} = \arg \max_{\theta \in \Theta} \bar{Q}_+(\theta)$ . The proof is omitted because it is trivial.

*Part (3).* In this part, we will prove that  $\sup_{\theta \in \Theta} |\hat{Q}_+(\theta) - \bar{Q}_+(\theta)| \xrightarrow{P} 0$ . Define the event

$$(B.3) \qquad \omega : \hat{X}_{ij}^{(k)} \geq \epsilon \quad \forall (i, j, k) \in S_I \quad \text{and} \quad \hat{Q}(\bar{\theta}) \geq \bar{Q}(\bar{\theta}) - \epsilon,$$

and set

$$(B.4) \qquad \Theta_1 = \left\{ \theta | T_{ij}^{(k)} \geq \exp\left(\frac{\bar{Q}(\bar{\theta}) - \epsilon - a}{\epsilon}\right) \text{ for } (i, j, k) \in S_I \right\} \cap \Theta,$$

where  $S_I = \{(i, j, k) | \bar{X}_{ij}^{(k)} > 0\}$  and  $\epsilon \in (0, \min_{(i,j,k) \in S_I} \bar{X}_{ij}^{(k)})$ . According to the definitions of  $\hat{Q}_+(\theta)$  and  $\bar{Q}_+(\theta)$ , we can get

$$\begin{aligned} (B.5) \qquad 1_\omega \cdot \sup_{\theta \in \Theta} \left| \hat{Q}_+(\theta) - \bar{Q}_+(\theta) \right| &\leq \max \left\{ 1_{\theta \in \Theta_1} \cdot \left| \hat{Q}(\theta) - \bar{Q}(\theta) \right|, \left| \hat{Q}(\bar{\theta}) - \bar{Q}(\bar{\theta}) \right| \right\} \\ &\leq \left| \frac{\bar{Q}(\bar{\theta}) - \epsilon - a}{\epsilon} \right| \sum_{(i,j,k) \in S_I} \left| \hat{X}_{ij}^{(k)} - \bar{X}_{ij}^{(k)} \right|. \end{aligned}$$



Moreover, considering that  $\hat{X}_{ij}^{(k)} \xrightarrow{P} \bar{X}_{ij}^{(k)}$ , we have  $1_\omega \xrightarrow{P} 1$ . Therefore

$$(B.6) \quad \sup_{\theta \in \Theta} \left| \hat{Q}_+(\theta) - \bar{Q}_+(\theta) \right| \xrightarrow{P} 0.$$

According to the definitions of  $\hat{Q}_+(\theta)$  and  $\bar{Q}_+(\theta)$  and the conclusions of Parts (1)–(3), it can be easily verified that  $\hat{Q}_+(\theta)$  satisfies the following conditions: (i)  $\bar{Q}_+(\theta)$  is uniquely maximized at  $\bar{\theta}$ ; (ii)  $\Theta$  is compact; (iii)  $\bar{Q}_+(\theta)$  is continuous; (iv)  $\hat{Q}_+(\theta)$  converges uniformly in probability to  $\bar{Q}_+(\theta)$ ; and (v)  $\bar{\theta} = \arg \max_{\theta \in \Theta} \hat{Q}(\theta)$ . Then we have  $\hat{\theta} \xrightarrow{P} \bar{\theta}$  by using Theorem 2.1 in [16] and (B.2).

**Appendix C. Proof of Theorem 3.10.** Let  $\Theta_r$  be the feasible set of  $\theta_r$  defined by constraints in (3.3), where  $T_{ij}^{(k)}$  and  $\pi_n$  which do not belong to  $\theta_r$  can be treated as functions of  $\theta_r$ . It is easy to see that  $\bar{\theta}_r$  is an interior point of  $\Theta_r$ , and that

$$(C.1) \quad \Theta_1 = \left\{ \theta_r \mid T_{ij}^{(k)}(\theta_r) \geq \epsilon \text{ for } (i, j, k) \in S_I \text{ and } \pi_i(\theta_r) \geq \epsilon \forall i \right\} \cap \Theta_r$$

is a closed neighborhood of  $\bar{\theta}_r$ , where  $\epsilon \in (0, \min\{\min_{(i,j,k) \in S_I} \bar{T}_{ij}^{(k)}, \min_i \bar{\pi}_i\})$  and  $S_I$  has the same definition as in Appendix B.

It is easy to verify that

$$\begin{aligned} \sqrt{M} \left( \nabla_{\theta_r} \hat{Q}(\theta(\bar{\theta}_r)) - \nabla_{\theta_r} \bar{Q}(\theta(\bar{\theta}_r)) \right)^T &= \sqrt{M} \nabla_{\theta_r} \Phi(\theta(\bar{\theta}_r))^T (\hat{V}_X - \bar{V}_X) \\ &\stackrel{d}{\rightarrow} \mathcal{N}(0, \Sigma) \end{aligned}$$

and  $\sup_{\theta_r \in \Theta_1} \|\nabla_{\theta_r, \theta_r} \hat{Q}(\theta(\theta_r)) - \nabla_{\theta_r, \theta_r} \bar{Q}(\theta(\theta_r))\| \xrightarrow{P} 0$ , where  $\hat{Q}(\theta)$  and  $\bar{Q}(\theta)$  have the same definition as in Appendix B, and  $V_X = \mathcal{V}(X^{(1)}, \dots, X^{(K)})$ . Then by Theorem 3.1 in [16], (3.13) holds.

**Appendix D. Bilevel optimization procedure for (4.12).** Here we define  $\underline{\rho}^{(k)}$  as a vector consisting of the elements of  $\{Z_{ij}^{(k)} \mid \underline{C}_{ij}^{(k)} > 0, i \leq j, (i, j) \neq (n, n)\}$ . (Note that  $\underline{\rho}^{(k)}$  is different from the  $\rho^{(k)}$  defined in section 4.2 because  $1_{\underline{C}_{ij}^{(k)} > 0} = 1_{\bar{X}_{ij}^{(k)} > 0}$  may not hold for all  $i, j, k$ .) Then we can eliminate the first two constraints of (4.12) by using the substitution method and regarding  $Z^{(k)}$  and  $Z_i^{(k)}$  as functions of  $\underline{\rho}^{(k)}$ , and express (4.12) as a bilevel optimization problem consisting of an upper-level problem

$$(D.1) \quad \begin{aligned} \max_V \quad & \sum \check{L}_V^{(k)}(V^{(k)}(V)) \\ \text{s.t.} \quad & \mathbf{1}^T V = 0 \end{aligned}$$

and  $K$  lower-level problems

$$(D.2) \quad \begin{aligned} \check{L}_V^{(k)}(V^{(k)}) &= \max_{\underline{\rho}^{(k)}} \check{L}_Z^{(k)}(Z^{(k)}(\underline{\rho}^{(k)}) \mid \underline{C}^{(k)}) \\ \text{s.t.} \quad & V_Z(\check{Z}^{(k)}) + \sum_{i,j} \nabla_{\mathcal{V}(Z^{(k)})} V_Z(\check{Z}^{(k)}) \cdot \mathcal{V}(Z^{(k)}(\underline{\rho}^{(k)}) - \check{Z}^{(k)}) = V^{(k)} \end{aligned}$$

for  $k = 1, \dots, K$ . Note that  $\check{Z}^{(k)} = Z^{(k)}(\check{\rho}^{(k)})$  with

$$(D.3) \quad \check{\rho}^{(k)} = \operatorname{argmax}_{\underline{\rho}^{(k)}} \check{L}_Z^{(k)}(Z^{(k)}(\underline{\rho}^{(k)}) \mid \underline{C}^{(k)})$$

and both sides of the constraint of (D.2) are zero-mean. Then (D.2) can be simplified as

$$(D.4) \quad L_V^{(k)}(V^{(k)}) = \max_{\underline{\rho}^{(k)}} L_Z^{(k)}(\check{Z}^{(k)}|\underline{C}^{(k)}) + \frac{1}{2}(\underline{\rho}^{(k)} - \check{\underline{\rho}}^{(k)})^T H_\rho^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)})(\underline{\rho}^{(k)} - \check{\underline{\rho}}^{(k)})$$

s.t.  $[\mathbf{I} \ \mathbf{0}] \nabla_{\rho^{(k)}} V_Z(\check{\underline{\rho}}^{(k)})(\underline{\rho}^{(k)} - \check{\underline{\rho}}^{(k)}) = [\mathbf{I} \ \mathbf{0}](V^{(k)} - \check{V}^{(k)}),$

where  $H_\rho^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) = \sum_i \underline{C}_i^{(k)} \nabla_{\underline{\rho}^{(k)}} Z_i^k(\check{\underline{\rho}}^{(k)}) \prec 0$  and where we denote  $V_Z(Z^{(k)}(\underline{\rho}^{(k)}))$  by  $V_Z(\underline{\rho}^{(k)})$  for convenience of notation. (The negative-definiteness of  $H_\rho^{(k)}(\underline{C}^{(k)}, \underline{\rho}^{(k)})$  can be easily verified according to its definition and Assumption 4.1.) It is easy to verify that  $[\mathbf{I} \ \mathbf{0}] \nabla_{\underline{\rho}^{(k)}} V_Z(\check{\underline{\rho}}^{(k)})$  is full row rank. Then, using the Lagrange multiplier method, we get

$$(D.5) \quad \check{L}_V^{(k)}(V^{(k)}) = L_Z^{(k)}(\check{Z}^{(k)}|\underline{C}^{(k)}) + \frac{1}{2}(V - V^{(k)})^T \Xi'^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)})(V - V^{(k)}),$$

where

$$(D.6) \quad \Xi'^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) = \left[ \begin{array}{c} \mathbf{I} \\ \mathbf{0}^T \end{array} \right] \left( [\mathbf{I} \ \mathbf{0}] \nabla_{\underline{\rho}^{(k)}} V_Z(\check{\underline{\rho}}^{(k)}) (H_\rho^k(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}))^{-1} \cdot ([\mathbf{I} \ \mathbf{0}] \nabla_{\underline{\rho}^{(k)}} V_Z(\check{\underline{\rho}}^{(k)}))^T \right)^{-1} [\mathbf{I} \ \mathbf{0}]$$

is negative-semidefinite and satisfies  $[\mathbf{I} \ \mathbf{0}] \Xi'^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) [\mathbf{I} \ \mathbf{0}]^T \prec 0$ . Substituting (D.5) into (D.1) and applying the KKT conditions, it is easy to verify that the solution of (4.12) is (4.14) with

$$(D.7) \quad \Xi^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) = \left( \sum_{m=1}^K (\nabla_V V^{(m)}(V)) \Xi'^{(m)}(\underline{C}^{(m)}, \check{\underline{\rho}}^{(m)}) \nabla_V V^{(m)}(V) \right)^+ \cdot (\nabla_V V^{(k)}(V))^T \Xi'^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)})$$

and

$$(D.8) \quad b^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) = \Xi^{(k)}(\underline{C}^{(k)}, \check{\underline{\rho}}^{(k)}) \nabla_V V^{(k)}(V) U^{(k)}.$$

**Appendix E. Proof of Lemma 4.6.** First we define  $\check{Q}^{(k)}(\rho^{(k)}) = L_Z^{(k)}(Z^{(k)}(\rho^{(k)})|\hat{\underline{X}}^{(k)})$  and  $\bar{Q}^{(k)}(\rho^{(k)}) = L_Z^{(k)}(Z^{(k)}(\rho^{(k)})|\bar{X}^{(k)})$  with  $\hat{\underline{X}}^{(k)} = [\hat{X}_{ij}^{(k)}] = \underline{C}^{(k)}/M$ . It is clear that there is a function  $\Phi^{(k)}(\cdot)$  such that  $\check{Q}^{(k)}(\rho^{(k)})$  and  $\bar{Q}^{(k)}(\rho^{(k)})$  can be written as  $\check{Q}^{(k)}(\rho^{(k)}) = \mathcal{V}(\hat{\underline{X}}^{(k)})^T \Phi^{(k)}(\rho^{(k)})$  and  $\bar{Q}^{(k)}(\rho^{(k)}) = \mathcal{V}(\bar{X}^{(k)})^T \Phi^{(k)}(\rho^{(k)})$ .

Under Assumptions 3.5–3.8 and 4.1, it is easy to see that  $\bar{X}^{(k)}$  is irreducible and its diagonal elements are positive, which implies that  $\bar{T}^{(k)}$  is an ergodic transition matrix with unique stationary distribution. Thus  $\bar{V}^{(k)} = V_Z(\bar{\rho}^{(k)})$ . ( $V_Z(\rho^{(k)})$  denotes  $V_Z(Z^{(k)}(\rho^{(k)}))$ .)

Define the event

$$(E.1) \quad \omega^{(k)} : 1_{C_{ij}^{(k)} > 0} = 1_{\bar{X}_{ij}^{(k)} > 0} \quad \forall i, j$$

and

$$(E.2) \quad \check{\rho}^{(k)} = \begin{cases} \check{\underline{\rho}}^{(k)}, & 1_{\omega^{(k)}} = 1, \\ \mathbf{0}, & 1_{\omega^{(k)}} = 0. \end{cases}$$

It is clear that the functions  $1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)})$  and  $\bar{Q}^k(\rho^{(k)})$  satisfy the following: (i)  $\check{\rho}^{(k)} = \arg \max_{\rho^{(k)}} 1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)})$  and  $\bar{Q}^k(\rho^{(k)})$  is uniquely maximized at  $\bar{\rho}^{(k)}$ ; (ii)  $1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)})$  is concave for  $\nabla_{\rho^{(k)} \rho^{(k)}} \check{Q}^k(\rho^{(k)}) = H_{\rho}^{(k)}(\hat{X}^{(k)}, \rho^{(k)}) \prec 0$  if  $\omega^{(k)}$  holds; (iii)  $1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)}) \xrightarrow{P} \bar{Q}^k(\rho^{(k)})$  for any  $\rho^{(k)}$  because  $1_{\omega^{(k)}} \xrightarrow{P} 1$ . Then we have  $\check{\underline{\rho}}^k \xrightarrow{P} \check{\rho}^k \xrightarrow{P} \bar{\rho}^k$  according to Theorem 2.7 in [16], and  $\check{V}^k \xrightarrow{P} \bar{V}^k$  since  $V_Z(\rho^{(k)})$  is a continuous function of  $\rho^{(k)}$ .

We now show the second conclusion of the lemma. Because it holds that  $\bar{\rho}^{(k)} = \arg \max_{\rho^{(k)}} \bar{Q}^k(\rho^{(k)})$ , it follows that  $\nabla_{\rho^{(k)}} \bar{Q}^k(\bar{\rho}^{(k)}) = \mathbf{0}^T$ . Then we have

$$(E.3) \quad \begin{aligned} \sqrt{M}(\nabla_{\rho^{(k)}}(1_{\omega^{(k)}} \cdot \check{Q}^k(\bar{\rho}^{(k)})))^T &\xrightarrow{P} 1_{\omega^{(k)}} \cdot (\nabla_{\rho^{(k)}} \Phi^k(\bar{\rho}^k))^T \cdot (\mathcal{V}(\hat{X}^k) - \mathcal{V}(\bar{X}^k)) \\ &\xrightarrow{d} \mathcal{N}(\mathbf{0}, (\nabla_{\rho^{(k)}} \Phi^k(\bar{\rho}^k))^T \Sigma_X^{(k)} (\nabla_{\rho^{(k)}} \Phi^k(\bar{\rho}^k))). \end{aligned}$$

Furthermore, the Hessian matrices of  $1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)})$  and  $\bar{Q}^k(\rho^{(k)})$  satisfy  $\nabla_{\rho^{(k)} \rho^{(k)}}(1_{\omega^{(k)}} \cdot \check{Q}^k(\rho^{(k)})) \xrightarrow{P} \nabla_{\rho^{(k)} \rho^{(k)}} \bar{Q}^k(\rho^{(k)})$  for any  $\rho^{(k)}$ . Using the mean value theorem and Theorem 3.1 in [16], we can conclude that  $\sqrt{M}(\check{\rho}^{(k)} - \bar{\rho}^{(k)}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_{\rho}^{(k)}(\Sigma_X^{(k)}, \bar{X}^{(k)}, \bar{\rho}^{(k)}))$  and

$$(E.4) \quad \sqrt{M}(\check{V}^{(k)} - \bar{V}^{(k)}) \xrightarrow{d} \mathcal{N}(\mathbf{0}, \Sigma_V^{(k)}(\Sigma_X^{(k)}, \bar{X}^{(k)}, \bar{\rho}^{(k)}))$$

with

$$(E.5) \quad \begin{aligned} \Sigma_V^{(k)}(\Sigma_X^{(k)}, X^{(k)}, \rho^{(k)}) &= \nabla_{\rho^k} V_Z(Z^{(k)}(\rho^{(k)})) \\ &\cdot \Sigma_{\rho}^k(\Sigma_X^{(k)}, X^{(k)}, \rho^{(k)}) (\nabla_{\rho^k} V_Z(Z^{(k)}(\rho^{(k)})))^T, \end{aligned}$$

where

$$(E.6) \quad \begin{aligned} \Sigma_{\rho}^k(\Sigma_X^{(k)}, X^{(k)}, \rho^{(k)}) &= (H_{\rho}^{(k)}(X^{(k)}, \rho^{(k)}))^{-1} (\nabla_{\rho^{(k)}} \Phi^{(k)}(\rho^{(k)}))^T \Sigma_X^{(k)} \\ &\cdot \nabla_{\rho^{(k)}} \Phi^{(k)}(\rho^{(k)}) (H_{\rho}^{(k)}(X^{(k)}, \rho^{(k)}))^{-1} \end{aligned}$$

and  $H_{\rho}^{(k)}(X^{(k)}, \rho^{(k)}) = \sum_{i,j} X_{ij}^{(k)} \nabla_{\rho^{(k)} \rho^{(k)}} Z_i^{(k)}(\rho^{(k)})$ .

**Appendix F. Proof of Theorem 4.7.** Since the value of  $\check{V}$  will not be affected if we replace  $\underline{C}^{(k)}$  with  $\hat{X}^{(k)} = \underline{C}^{(k)}/M$ ,  $\check{V}$  can be expressed as

$$(F.1) \quad \begin{aligned} \check{V} &= \left( \sum_{k=1}^K A^T \Xi^{(k)}(\hat{X}^{(k)}, \check{\underline{\rho}}^{(k)}) A \right)^+ \\ &\cdot \left( \sum_{k=1}^K A^T \Xi^{(k)}(\hat{X}^{(k)}, \check{\underline{\rho}}^{(k)}) (\check{V}^{(k)} - AU^{(k)}) \right) \end{aligned}$$

with  $A = \mathbf{I} - \frac{1}{n}\mathbf{1}\mathbf{1}^\top \in \mathbb{R}^{n \times n}$  and  $\Xi^{(k)}(\cdot)$  defined by (D.6), and  $\check{V} = \arg \max_{\mathbf{1}^\top V=0} Q_q(V; \{\check{V}^{(k)}\})$  with

$$(F.2) \quad Q_q(V; \{\check{V}^{(k)}\}) = \sum_{k=1}^K \frac{1}{2} \left( A(V + U^{(k)}) - \check{V}^{(k)} \right)^\top \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)}) \cdot \left( A(V + U^{(k)}) - \check{V}^{(k)} \right).$$

*Part (1).* We first prove that  $\bar{V}$  is the unique maximum point of  $Q_q(V; \{\bar{V}^{(k)}\})$  under the constraint  $\mathbf{1}^\top V = 0$ . It is clear that  $\bar{V}$  is a maximum point since  $Q_q(\bar{V}; \{\bar{V}^{(k)}\}) = 0$  and  $Q_q(V; \{\bar{V}^{(k)}\}) \leq 0$ . We now show the uniqueness of  $\bar{V}$  by contradiction. Suppose that  $V'$  is another maximum point which satisfies  $Q_q(V'; \{\bar{V}^{(k)}\}) = 0$  and  $V' \neq \bar{V}$ . Then  $[\mathbf{I} \ \mathbf{0}](A(V' + U^{(k)}) - \bar{V}^{(k)}) = \mathbf{0}$  for any  $k$  because  $[\mathbf{I} \ \mathbf{0}] \cdot \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)}) \cdot [\mathbf{I} \ \mathbf{0}]^\top < 0$ , which implies that the first  $n-1$  elements of  $A(V' + U^{(k)}) - \bar{V}^{(k)}$  are zero. Further, considering that  $\mathbf{1}^\top A(V' + U^{(k)}) = \mathbf{1}^\top \bar{V}^{(k)} = 0$ , we can conclude that  $A(V' + U^{(k)}) = \bar{V}^{(k)}$  for each  $k$ . Therefore the probability distribution  $\pi' = [\pi'_i]$  with  $\pi'_i \propto \exp(-V'_i)$  satisfies (3.11), which contradicts Assumption 3.8. Thus  $\bar{V}$  is the unique solution of  $\arg \max_{\mathbf{1}^\top V=0} Q_q(V; \{\bar{V}^{(k)}\})$ .

*Part (2).* In this part we will prove that  $\text{null}(\sum_{k=1}^K A^\top \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)})A) = \text{span}(\mathbf{1})$  for any  $\hat{X}^{(k)}, \check{\rho}^{(k)}$ . It is clear that  $\text{span}(\mathbf{1}) \subseteq \text{null}(\sum_{k=1}^K A^\top \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)})A)$  since  $A\mathbf{1} = \mathbf{0}$ . Suppose that  $v \notin \text{span}(\mathbf{1})$  is another vector which belongs to  $\text{null}(\sum_{k=1}^K A^\top \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)})A)$  and satisfies  $\mathbf{1}^\top v = 0$ ; then  $\bar{V} + v \neq \bar{V}$  is also a maximum point of  $\arg \max_{\mathbf{1}^\top V=0} Q_q(V; \{\bar{V}^{(k)}\})$ . This is a contradiction to the result of Part (1). Therefore  $\text{null}(\sum_{k=1}^K A^\top \Xi^{(k)}(\hat{X}^{(k)}, \check{\rho}^{(k)})A) = \text{span}(\mathbf{1})$ .

*Part (3).* Combining the result of Part (2) and Theorem 5.2 in [27] leads to the conclusion of the theorem with

$$(F.3) \quad \Sigma_V \left( \{\Sigma_X^{(k)}\}, \{X^{(k)}\}, \{\rho^{(k)}\} \right) = \sum_{k=1}^K \Xi^{(k)} \left( X^{(k)}, \rho^{(k)} \right) \Sigma_V^{(k)} \left( \Sigma_X^{(k)}, X^{(k)}, \rho^{(k)} \right) \cdot \left( \Xi^{(k)} \left( X^{(k)}, \rho^{(k)} \right) \right)^\top.$$

**Appendix G. Proof of Theorem 4.9.** Let  $G(\cdot)$  be a function of  $y_t^{(k)}$  with  $G(y_t^{(k)}) = [G_{ij}(y_t^{(k)})] = [1_{f^{(k)}(y_t^{(k)})=i} \Pr(x_{t+1}^{(k)} = j | y_t^{(k)})]$ , and let  $\kappa_G^{(k)}(h) = \text{Cov}(\mathcal{V}(G(y_t^{(k)})), \mathcal{V}(G(y_{t+h}^{(k)})))$  be the  $h$  lag autocovariance of  $\{\mathcal{V}(G(y_t^{(k)}))\}$ . It is easy to verify that  $\kappa_G^{(k)}(h)$  and  $\kappa_G^{(k)}(h+1)$  are composed of the same elements but in different arrangements for  $h \geq 0$ , and that  $\eta^{(k)}(l) = \text{tr}(\kappa_G^{(k)}(2l) + \kappa_G^{(k)}(2l+1))$ .

From the above results, we can conclude that  $\eta^{(k)}(l) = \text{tr}(\kappa_G^{(k)}(2l) + \kappa_G^{(k)}(2l+1))$  is a nonnegative and decreasing function of  $l$  for  $l \geq 0$ , and that  $\kappa_G^{(k)}(2l) + \kappa_G^{(k)}(2l+1) \geq 0$  by using Theorem 3.1 in [9]. Therefore the second conclusion of the theorem can be shown.

We now show the first conclusion. According to Theorem 2.1 in [9] and considering

that  $\sum_{h=0}^{\infty} \kappa_G^{(k)}(h)$  is convergent, we have

$$\begin{aligned} N\text{Var}\left(\frac{1}{N}\sum_{t=1}^N \mathcal{V}\left(\Delta C_t^{(k)}\right)\right) &= \kappa^{(k)}(0) + \sum_{h=1}^{N-1} \frac{N-h}{N} \left(\kappa^{(k)}(h) + \kappa^{(k)}(h)^T\right) \\ \text{(G.1)} \qquad \qquad \qquad &\rightarrow \kappa^{(k)}(0) + \sum_{h=1}^{\infty} \left(\kappa^{(k)}(h) + \kappa^{(k)}(h)^T\right) \end{aligned}$$

as  $N \rightarrow \infty$ , and  $\sum_X^{(k)} = \kappa^{(k)}(0) + \sum_{l=0}^{\infty} (\Gamma^{(k)}(l) + \Gamma^{(k)}(l)^T)$ .

**Acknowledgment.** We have benefited from stimulating discussions with Benjamin Trendelkamp-Schroer (FU Berlin).

#### REFERENCES

- [1] A. BARDUCCI, G. BUSSI, AND M. PARRINELLO, *Well-tempered metadynamics: A smoothly converging and tunable free-energy method*, Phys. Rev. Lett., 100 (2008), 020603.
- [2] K. A. BEAUCHAMP, G. R. BOWMAN, T. J. LANE, L. MAIBAUM, I. S. HAQUE, AND V. S. PANDE, *MSMBuilder2: Modeling conformational dynamics at the picosecond to millisecond scale*, J. Chem. Theory Comput., 7 (2011), pp. 3412–3419.
- [3] K. BINDER AND A. P. YOUNG, *Spin glasses: Experimental facts, theoretical concepts, and open questions*, Rev. Mod. Phys., 58 (1986), pp. 801–976.
- [4] G. R. BOWMAN, K. A. BEAUCHAMP, G. BOXER, AND V. S. PANDE, *Progress and challenges in the automated construction of Markov state models for full protein systems*, J. Chem. Phys., 131 (2009), 124101.
- [5] S. BOYD AND L. VANDENBERGHE, *Convex Optimization*, Cambridge University Press, Cambridge, UK, 2004.
- [6] J. D. CHODERA, K. A. DILL, N. SINGHAL, V. S. PANDE, W. C. SWOPE, AND J. W. PITERA, *Automatic discovery of metastable states for the construction of Markov models of macromolecular conformational dynamics*, J. Chem. Phys., 126 (2007), 155101.
- [7] A. M. FERREBERG AND R. H. SWENDSEN, *Optimized Monte Carlo data analysis*, Phys. Rev. Lett., 63 (1989), pp. 1195–1198.
- [8] H. FRAUENFELDER, S. G. SLIGAR, AND P. G. WOLYNES, *The energy landscapes and motions of proteins*, Science, 254 (1991), pp. 1598–1603.
- [9] C. GEYER, *Practical Markov chain Monte Carlo*, Statist. Sci., 7 (1992), pp. 473–483.
- [10] H. GRUBMÜLLER, *Predicting slow structural transitions in macromolecular systems: Conformational flooding*, Phys. Rev. E, 52 (1995), pp. 2893–2906.
- [11] O. HÄGGSTRÖM, *On the central limit theorem for geometrically ergodic Markov chains*, Probab. Theory Related Fields, 132 (2005), pp. 74–82.
- [12] G. JONES, *On the Markov chain central limit theorem*, Probab. Surv., 1 (2004), pp. 299–320.
- [13] A. LAIO AND M. PARRINELLO, *Escaping free energy minima*, Proc. Natl. Acad. Sci. USA, 99 (2002), 12562.
- [14] D. MARX AND J. HUTTER, *Ab initio molecular dynamics: Theory and implementation*, in Modern Methods and Algorithms of Quantum Chemistry, NIC Ser. 1, J. Grotendorst, ed., John von Neumann Institute for Computing, Jülich, Germany, 2000, pp. 301–449.
- [15] P. METZNER, F. NOÉ, AND C. SCHÜTTE, *Estimating the sampling error: Distribution of transition matrices and functions of transition matrices for given trajectory data*, Phys. Rev. E, 80 (2009), 021106.
- [16] W. K. NEWBY AND D. MCFADDEN, *Large sample estimation and hypothesis testing*, in Handbook of Econometrics, Vol. 4, Handbooks in Econom. 2, North-Holland, Amsterdam, 1994, pp. 2111–2245.
- [17] F. NOÉ, *Probability distributions of molecular observables computed from Markov models*, J. Chem. Phys., 128 (2008), 244103.
- [18] F. NOÉ, I. HORENKO, C. SCHÜTTE, AND J. C. SMITH, *Hierarchical analysis of conformational dynamics in biomolecules: Transition networks of metastable states*, J. Chem. Phys., 126 (2007), 155102.
- [19] S. PIANA AND A. LAIO, *A bias-exchange approach to protein folding*, J. Phys. Chem. B, 111 (2007), pp. 4553–4559.

- [20] J.-H. PRINZ, M. HELD, J. C. SMITH, AND F. NOÉ, *Efficient computation, sensitivity, and error analysis of committor probabilities for complex dynamical processes*, Multiscale Model. Simul., 9 (2011), pp. 545–567.
- [21] J.-H. PRINZ, H. WU, M. SARICH, B. KELLER, M. SENNE, M. HELD, J. CHODERA, C. SCHÜTTE, AND F. NOÉ, *Markov models of molecular kinetics: Generation and validation*, J. Chem. Phys., 134 (2011), 174105.
- [22] Z. RACHED, F. ALAJAJI, AND L. CAMPBELL, *The Kullback-Leibler divergence rate between Markov sources*, IEEE Trans. Inform. Theory, 50 (2004), pp. 917–921.
- [23] S. SAKURABA AND A. KITAO, *Multiple Markov transition matrix method: Obtaining the stationary probability distribution from multiple simulations*, J. Comput. Chem., 30 (2009), pp. 1850–1858.
- [24] D. SHERRINGTON AND S. KIRKPATRICK, *Solvable model of a spin-glass*, Phys. Rev. Lett., 35 (1975), pp. 1792–1796.
- [25] A. SMOLA, S. VISHWANATHAN, AND T. HOFMANN, *Kernel methods for missing variables*, in Proceedings of the Tenth International Workshop on Artificial Intelligence and Statistics, Barbados, The Society for Artificial Intelligence and Statistics, 2005.
- [26] M. SOUAILLE AND B. ROUX, *Extension to the weighted histogram analysis method: Combining umbrella sampling with free energy calculations*, Comput. Phys. Comm., 135 (2001), pp. 40–57.
- [27] G. W. STEWART, *On the continuity of the generalized inverse*, SIAM J. Appl. Math., 17 (1969), pp. 33–45.
- [28] W. C. SWOPE, J. W. PITERA, AND F. SUITS, *Describing protein folding kinetics by molecular dynamics simulations. 1. Theory*, J. Phys. Chem. B, 108 (2004), pp. 6571–6581.
- [29] G. M. TORRIE AND J. P. VALLEAU, *Nonphysical sampling distributions in Monte Carlo free-energy estimation: Umbrella sampling*, J. Comput. Phys., 23 (1977), pp. 187–199.
- [30] G. YUAN, X. LU, AND Z. WEI, *A conjugate gradient method with descent direction for unconstrained optimization*, J. Comput. Appl. Math., 233 (2009), pp. 519–530.