

---

# Spectral learning of dynamic systems from nonequilibrium data\*

---

**Hao Wu and Frank Noé**

Department of Mathematics and Computer Science  
Freie Universität Berlin  
Arnimallee 6, 14195 Berlin  
{hao.wu, frank.noe}@fu-berlin.de

## Abstract

Observable operator models (OOMs) and related models are one of the most important and powerful tools for modeling and analyzing stochastic systems. They can exactly describe dynamics of finite-rank systems, and be efficiently learned from data by moment based algorithms. Almost all OOM learning algorithms are developed based on the assumption of equilibrium data which is very difficult to guarantee in real life, especially for complex processes with large time scales. In this paper, we derive a nonequilibrium learning algorithm for OOMs, which dismisses this assumption and can effectively extract the equilibrium dynamics of a system from nonequilibrium observation data. In addition, we propose binless OOMs for the application of nonequilibrium learning to continuous-valued systems. In comparison with the other OOMs with continuous observations, binless OOMs can achieve consistent estimation from nonequilibrium data with only linear computational complexity.

## 1 Introduction

In the last two decades, a collection of highly related dynamical models including observable operator models (OOMs) [1–3], predictive state representations [4–6] and spectral learning based hidden Markov models [7, 8], have become powerful and increasingly popular tools for analysis of dynamical data. These models are largely similar and can be unified in a general learning framework of multiplicity automata, or equivalently sequential systems [9, 10]. We focus in this paper only on stochastic systems without control inputs. Because all of above mentioned models can be expressed in the form of OOMs for such systems, we will refer to them as OOMs below.

In contrast with the other commonly used models such as Markov models [11], Langevin models [12], traditional hidden Markov models (HMMs) [13] and recurrent neural networks [14], OOMs can exactly characterize the dynamics of a stochastic system without any a priori knowledge except the assumption of finite dynamical rank (i.e., the rank of Hankel matrix) [10], and the parameter estimation can be efficiently performed by the method of moments for discrete-valued systems without solving any intractable inverse or optimization problem.

A major challenge for OOM based dynamical modeling approaches arises from nonequilibrium data. In most literature, the observation data are assumed to be equilibrium so that the expected values of observables associated with OOM learning can be reliably computed by simple averaging. However, the equilibrium assumption can be approximately satisfied only if most of observation data are generated after the system has mixed. In many practical situations, especially where metastable physical or chemical processes are involved, this assumption can be severely violated due to the

---

\*Original title: Learning Observable Operator Models from Nonequilibrium Data

limit of experimental technique or computational capacity. A notable example is the distributed computing project Folding@home [15], which explores protein folding processes that occur on the timescales of microseconds to milliseconds based on molecular dynamics simulations on the order of nanoseconds in length. In such a case, it is still unknown how to obtain promising estimates of OOMs from nonequilibrium data consisting of short trajectories. In [16], a hybrid estimation algorithm was proposed to improve OOM learning of large-time-scale processes by using both dynamic and static data, but it still requires assumption of equilibrium data. One solution to reduce the statistical bias caused by nonequilibrium data is to discard the observation data generated before the system reaches steady state, which is a common trick in applied statistics [17]. Obviously, this way suffers from substantial information loss and is infeasible when observation trajectories are shorter than mixing times. Another possible way would be to learn OOMs by likelihood-based estimation instead of moment-based estimation, but there is no effective maximum likelihood or Bayesian estimator of OOMs until now. The maximum pseudo-likelihood estimator of OOMs proposed in [18] demands high computational cost and its consistency is yet unverified.

Another difficulty for OOM based modeling approaches is learning with continuous data, where density estimation problems are involved. The density estimation can be performed by parametric methods such as the fuzzy interpolation [19] and the kernel density estimation [8]. But these methods would reduce the flexibility of OOMs for dynamical modeling because of their limited expressive capacity. Recently, a kernel embedding based OOM learning algorithm was proposed to cope with continuous data [20], which avoids explicit density estimation and learns OOMs in a nonparametric manner. However, the kernel embedding usually yields a very large computational complexity, which greatly limits practical applications of this algorithm to real-world systems.

The purpose of this paper is to address the challenge of nonequilibrium learning of OOMs due to the requirements of analysis of both discrete- and continuous-valued systems. We provide a modified moment-based method for discrete-valued stochastic systems which allows us to consistently estimate the equilibrium dynamics from nonequilibrium data, and then extend this method to OOM learning with continuous observations in a binless manner. In comparison with the existing learning methods for continuous OOMs, the proposed binless method does not rely on any density estimator, and can achieve consistent estimation with linear computational complexity in data size even if the equilibrium assumption of observations does not hold. Moreover, some numerical experiments are provided to demonstrate the capability of the proposed nonequilibrium learning methods.

## 2 Preliminaries

### 2.1 Notation

In this paper, we use  $\mathbb{P}$  to denote probability distribution for discrete random variables and probability density for continuous random variables. The indicator function of event  $e$  is denoted by  $1_e$  and the dirac delta function centered at  $x$  is denoted by  $\delta_x(\cdot)$ . For a given process  $\{a_t\}$ , we write the subsequence  $(a_k, a_{k+1}, \dots, a_{k'})$  as  $a_{k:k'}$ , and  $\mathbb{E}_\infty[a_t] \triangleq \lim_{t \rightarrow \infty} \mathbb{E}[a_t]$  means the expected value of  $a_t$  in equilibrium if the limit exists. In addition, the convergence in probability is denoted by  $\xrightarrow{p}$ .

### 2.2 Observable operator models

An  $m$ -dimensional observable operator model (OOM) with observation space  $\mathcal{O}$  can be represented by a tuple  $\mathcal{M} = (\boldsymbol{\omega}, \{\Xi(x)\}_{x \in \mathcal{O}}, \boldsymbol{\sigma})$ , which consists of an initial state vector  $\boldsymbol{\omega} \in \mathbb{R}^{1 \times m}$ , an evaluation vector  $\boldsymbol{\sigma} \in \mathbb{R}^{m \times 1}$  and an observable operator matrix  $\Xi(x) \in \mathbb{R}^{m \times m}$  associated to each element  $x \in \mathcal{O}$ .  $\mathcal{M}$  defines a stochastic process  $\{x_t\}$  in  $\mathcal{O}$  as

$$\mathbb{P}(x_{1:t} | \mathcal{M}) = \boldsymbol{\omega} \Xi(x_{1:t}) \boldsymbol{\sigma} \quad (1)$$

with  $\Xi(x_{1:t}) \triangleq \Xi(x_1) \dots \Xi(x_t)$ . It is interesting to note that (1) can also be represented in the form of state space models as

$$\begin{aligned} \boldsymbol{\omega}_t &= (\boldsymbol{\omega}_{t-1} \Xi(x_t) \boldsymbol{\sigma})^{-1} \boldsymbol{\omega}_{t-1} \Xi(x_t) \\ \mathbb{P}(x_{t+1} | \boldsymbol{\omega}_t) &= \boldsymbol{\omega}_t \Xi(x_{t+1}) \boldsymbol{\sigma} \end{aligned} \quad (2)$$

Here the internal state  $\boldsymbol{\omega}_t$  in (2) is a sufficient statistics the process at each time  $t$ , which contains all the information needed to predict the future observations and is initialized by  $\boldsymbol{\omega}_0 = \boldsymbol{\omega}$ . It is clear that two OOMs  $\mathcal{M}$  and  $\mathcal{M}'$  are equivalent if and only if  $\mathbb{P}(x_{1:t} | \mathcal{M}) \equiv \mathbb{P}(x_{1:t} | \mathcal{M}')$ .

---

**Algorithm 1** General procedure for OOM learning
 

---

**INPUT:** Observation trajectories generated by a stochastic process  $\{x_t\}$  in  $\mathcal{O}$ 
**OUTPUT:**  $\hat{\mathcal{M}} = (\hat{\omega}, \{\hat{\Xi}(x)\}_{x \in \mathcal{O}}, \hat{\sigma})$ 
**PARAMETER:**  $m$ : dimension of the OOM.  $D_1, D_2$ : numbers of feature functions.  $L$ : order of feature functions.

- 1: Construct feature functions  $\phi_1 = (\varphi_{1,1}, \dots, \varphi_{1,D_1})^\top$  and  $\phi_2 = (\varphi_{2,1}, \dots, \varphi_{2,D_2})^\top$ , where each  $\varphi_{i,j}$  is a mapping from  $\mathcal{O}^L$  to  $\mathbb{R}$  and  $D_1, D_2 \geq m$ .
- 2: Approximate

$$\bar{\phi}_1 \triangleq \mathbb{E}_\infty [\phi_1(x_{t+1:t+L})], \quad \bar{\phi}_2 \triangleq \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})] \quad (6)$$

$$\mathbf{C}_{1,2} \triangleq \mathbb{E}_\infty [\phi_1(x_{t-L:t-1})\phi_2(x_{t:t+L-1})^\top] \quad (7)$$

$$\mathbf{C}_{1,3}(x) \triangleq \mathbb{E}_\infty [1_{x_t=x} \cdot \phi_1(x_{t-L:t-1})\phi_2(x_{t+1:t+L})^\top], \quad \forall x \in \mathcal{O} \quad (8)$$

 by their empirical means  $\hat{\phi}_1, \hat{\phi}_2, \hat{\mathbf{C}}_{1,2}$  and  $\hat{\mathbf{C}}_{1,3}(x)$  over observation data.

- 3: Choose matrix  $\mathbf{F}_1 \in \mathbb{R}^{D_1 \times m}$ ,  $\mathbf{F}_2 \in \mathbb{R}^{D_2 \times m}$  such that  $\mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2$  is invertible.
- 4: Compute

$$\hat{\sigma} = \left( \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \hat{\phi}_1 \quad (9)$$

$$\hat{\Xi}(x) = \left( \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,3}(x) \mathbf{F}_2, \quad \forall x \in \mathcal{O} \quad (10)$$

$$\hat{\omega} = \hat{\phi}_2^\top \mathbf{F}_2 \quad (11)$$


---

### 3 Learning OOMs using moments

#### 3.1 Algorithm

Here and hereafter, we only consider the case that the observation space  $\mathcal{O}$  is a finite set. (Learning with continuous observations will be discussed in Section 5.) A large number of largely similar methods have been developed to learn OOMs from discrete data, and the generic learning procedure of these methods is summarized in Algorithm 1 by omitting details of algorithm implementation and parameter choice. For convenience of description and analysis, we specify in this paper the formula for calculating  $\hat{\phi}_1, \hat{\phi}_2, \hat{\mathbf{C}}_{1,2}$  and  $\hat{\mathbf{C}}_{1,3}(x)$  in Line 2 of Algorithm 1 as follows:

$$\hat{\phi}_1 \triangleq \frac{1}{N} \sum_{n=1}^N \phi_1(\bar{s}_n^1), \quad \hat{\phi}_2 \triangleq \frac{1}{N} \sum_{n=1}^N \phi_2(\bar{s}_n^2) \quad (3)$$

$$\hat{\mathbf{C}}_{1,2} \triangleq \frac{1}{N} \sum_{n=1}^N \phi_1(\bar{s}_n^1) \phi_2(\bar{s}_n^2)^\top \quad (4)$$

$$\hat{\mathbf{C}}_{1,3}(x) \triangleq \frac{1}{N} \sum_{n=1}^N 1_{s_n^2=x} \phi_1(\bar{s}_n^1) \phi_2(\bar{s}_n^3)^\top, \quad \forall x \in \mathcal{O} \quad (5)$$

Here  $\{(\bar{s}_n^1, s_n^2, \bar{s}_n^3)\}_{n=1}^N$  is the collection of all subsequences of length  $(2L+1)$  appearing in observation data ( $N = T - 2L$  for a single observation trajectory of length  $T$ ). For instance, if an observation subsequence  $x_{t-L:t+L}$  is denoted by  $(\bar{s}_n^1, s_n^2, \bar{s}_n^3)$  with some  $n$ , then  $\bar{s}_n^1 = x_{t-L:t-1}$  and  $\bar{s}_n^3 = x_{t+1:t+L}$  represents the prefix and suffix of  $x_{t-L:t+L}$  of length  $L$ ,  $s_n^2 = x_t$  is the intermediate observation value, and  $\bar{s}_n^2 = x_{t:t+L-1}$  is an ‘‘intermediate part’’ of the subsequence of length  $L$  starting from time  $t$  (see Fig. 1 for a graphical illustration).

Algorithm 1 is much more efficient than the commonly used likelihood-based learning algorithms and does not suffer from local optima issues. In addition, and more importantly, this algorithm can be shown to be consistent if the observation data are equilibrium so that empirical estimates of  $\bar{\phi}_1, \bar{\phi}_2, \mathbf{C}_{1,2}$  and  $\mathbf{C}_{1,3}(x)$  converge to their true values with increasing data size (see, e.g., [3, 8, 10] for

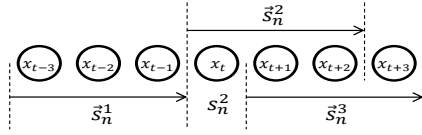


Figure 1: Illustration of variables  $s_n^1, s_n^2, s_n^3$  and  $s_n^2$  used in Eqs. (3)-(5) with  $(s_n^1, s_n^2, s_n^3) = x_{t-L:t+L}$  and  $L = 3$ .

related works). However, the consistent estimation of OOMs with nonequilibrium data is still an unsolved problem.

### 3.2 Theoretical analysis

We now analyze statistical properties of the OOM learning algorithm without the assumption of equilibrium observations. Before stating our main result, some assumptions on observation data are listed as follows:

**Assumption 1.** *The observation data consists of  $I$  independent trajectories of length  $T$  produced by a stochastic process  $\{x_t\}$ , and the data size tends to infinity with (i)  $I \rightarrow \infty$  and  $T = T_0$  or (ii)  $T \rightarrow \infty$  and  $I = I_0$ .*

**Assumption 2.**  *$\{x_t\}$  is driven by an  $m$ -dimensional OOM  $\mathcal{M} = (\omega, \{\Xi(x)\}_{x \in \mathcal{O}}, \sigma)$ , and satisfies*

$$\frac{1}{T'} \sum_{t=1}^{T'} f_t \xrightarrow{P} \mathbb{E}_\infty [f_t] = \mathbb{E}_\infty [f_t | x_{1:k}] \quad (12)$$

as  $T' \rightarrow \infty$  for all  $k, l, x_{1:k}$  and  $f_t = f(x_{t:t+l-1})$ .

**Assumption 3.** *The limit of  $\mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \in \mathbb{R}^{m \times m}$  is invertible.*

Notice that we do not assume stationarity of processes as previously done in the literature, and Assumption 2 only states the asymptotic stationarity of  $\{x_t\}$ . Therefore, estimates of  $\hat{\phi}_1, \hat{\phi}_2, \hat{\mathbf{C}}_{1,2}$  and  $\hat{\mathbf{C}}_{1,3}(x)$  obtained from empirical means may not be consistent if lengths of observation trajectories are kept at finite values (i.e., Case (i) in Assumption 1). Assumption 3 ensures that the limit of  $\hat{\mathcal{M}}$  given by Algorithm 1 is well defined.

Based on the above assumptions, we have the following theorem concerning the consistency of the OOM learning algorithm (see Appendix A.1 for proof):

**Theorem 1.** *Under Assumptions 1-3, the estimated OOM  $\hat{\mathcal{M}} = (\hat{\omega}, \{\hat{\Xi}(x)\}_{x \in \mathcal{O}}, \hat{\sigma})$  given by Algorithm 1 satisfies*

$$\hat{\sigma} \xrightarrow{P} \sigma_{\text{eq}}, \quad \hat{\Xi}(x) \xrightarrow{P} \Xi_{\text{eq}}(x), \quad \forall x \in \mathcal{O} \quad (13)$$

where  $\mathcal{M}_{\text{eq}} = (\omega_{\text{eq}}, \{\Xi_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \sigma_{\text{eq}})$  is an  $m$ -dimensional OOM equivalent to  $\mathcal{M}$ .

This theorem is central in this paper, and implies that the moment based learning algorithm can achieve consistent estimation of all parameters of OOMs except initial state vectors even for nonequilibrium data. ( $\hat{\omega} \xrightarrow{P} \omega_{\text{eq}}$  does not hold in most cases except when  $\{x_t\}$  is stationary. See Appendix A.1 for details). It can be further generalized according to requirements in more complicated situations where, for example, the data set consists of both several long trajectories and many short trajectories or trajectories are not independent from each other. The following two generalizations are particularly worth mentioning due to their importance for practical applications:

1. The  $i$ -th observation trajectories is generated by OOM  $\mathcal{M} = (\omega^i, \{\Xi(x)\}_{x \in \mathcal{O}}, \sigma)$  for  $i = 1, \dots, I$  (i.e., observations are generated with multiple different initial conditions), and the mean value of  $\{\omega^i\}_{i=1}^I$  tends to a constant in probability for  $I \rightarrow \infty$ .
2. Matrices  $\mathbf{F}_1$  and  $\mathbf{F}_2$  are not constant but given by the singular value decomposition of  $\hat{\mathbf{C}}_{1,2}$  as in the spectral learning algorithm [21, 7, 22].

We show in Appendix A that the above two generalizations do not affect the consistency of  $\hat{\Xi}(x)$  and  $\hat{\sigma}$ . In fact, it can be proved by similar proofs that all theoretical conclusions in this paper hold for the two generalizations.

## 4 Nonequilibrium learning of OOMs

According to the discussion in the previous section, the only remaining problem for learning OOMs from nonequilibrium data is how to estimate initial state vectors. Considering that the purpose of dynamical modeling is to predict properties of the system in equilibrium in many situations, here we only approximate equilibrium values of internal states of OOMs (see below) rather than actual initial state vectors, because the latter depend on initial conditions of data generation and the former are more physically interesting for analysis of equilibrium dynamics.

Given parameters of  $\mathcal{M}_{\text{eq}}$  in Theorem 1, the equilibrium value of the internal state is defined as

$$\bar{\omega}_{\text{eq}} = \lim_{t \rightarrow \infty} \omega_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O})^t \quad (14)$$

if the limit exists, where  $\Xi_{\text{eq}}(\mathcal{O}) = \sum_{x \in \mathcal{O}} \Xi_{\text{eq}}(x)$ . Then the equilibrium dynamics of  $\{x_t\}$  can be characterized as

$$\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+k} = z_{1:k}) = \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(z_{1:k}) \sigma_{\text{eq}} \quad (15)$$

From (14) and (15), we have

$$\begin{cases} \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O}) = \lim_{t \rightarrow \infty} \omega_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O})^{t+1} = \bar{\omega}_{\text{eq}} \\ \bar{\omega}_{\text{eq}} \sigma_{\text{eq}} = \lim_{t \rightarrow \infty} \sum_{x \in \mathcal{O}} \mathbb{P}(x_{t+1} = x) = 1 \end{cases} \quad (16)$$

This motivates the following nonequilibrium learning algorithm for OOMs: *Perform Algorithm 1 to get  $\hat{\Xi}(x)$  and  $\hat{\sigma}$  and calculate  $\hat{\omega}$  by a quadratic programming problem*

$$\hat{\omega} = \arg \min_{\mathbf{w} \in \{\mathbf{w} | \mathbf{w} \hat{\sigma} = 1\}} \left\| \mathbf{w} \hat{\Xi}(\mathcal{O}) - \mathbf{w} \right\|^2 \quad (17)$$

(See Appendix A.4 for a closed-form expression of the solution to (17).)

The existence and uniqueness of  $\bar{\omega}_{\text{eq}}$  are shown in Appendix A.4, which yield the following theorem:

**Theorem 2.** *Under Assumptions 1-3, the estimated OOM  $\hat{\mathcal{M}}$  provided by the nonequilibrium learning algorithm satisfies*

$$\mathbb{P}(x_{1:l} = z_{1:l} | \hat{\mathcal{M}}) \xrightarrow{P} \lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+l} = z_{1:l}) \quad (18)$$

for all  $l$  and  $z_{1:l}$ .

*Remark 1.* Some OOM learning algorithms for equilibrium data [23] also calculate  $\hat{\omega}$  based on (16), where feature functions  $\phi_1, \phi_2$  and matrices  $\mathbf{F}_1, \mathbf{F}_2$  are specifically constructed so that  $\hat{\omega} \hat{\Xi}(\mathcal{O}) = \hat{\omega}, \hat{\omega} \hat{\sigma} = 1$  can be exactly satisfied even if the statistical noise is considered. In comparison with these algorithms, the nonequilibrium learning algorithm does not require such a restriction, and is shown to be applicable to nonequilibrium data.

## 5 Binless learning of OOMs

We now consider how to learn OOMs from continuous data. In the case of a real observation space  $\mathcal{O} \subset \mathbb{R}^d$ ,  $\mathcal{M}$  defines probability densities of paths of  $\{x_t\}$  as in (1), and  $\mathbf{C}_{1,3}(x)$  becomes a matrix-valued density function  $\mathbf{C}_{1,3}(x) = \frac{1}{dx} \mathbb{E}_{\infty} [1_{x_t \in dx} \cdot \phi_1(x_{t-L:t-1}) \phi_2(x_{t+1:t+L})^\top]$  with general feature functions  $\phi_1, \phi_2$  on  $\mathbb{R}^d$ , which is difficult to approximate for each  $x \in \mathcal{O}$ . The existing continuous learning algorithms overcome this problem by using parametric methods [19, 8] or kernel embeddings [20], but none of them can achieve consistent estimation with a low computational complexity like discrete learning algorithms even for equilibrium data.

Here we present a binless strategy to perform dynamical modeling with continuous and nonequilibrium data, which simply views each available observation as a discrete probability atom in the observation space and approximates  $\mathbf{C}_{1,3}(x)$  by

$$\hat{\mathbf{C}}_{1,3}(x) = \frac{1}{N} \sum_{n=1}^N \delta_{s_n^2}(x) \phi_1(s_n^1) \phi_2(s_n^3)^\top \quad (19)$$

instead of (5). Using this strategy, a binless OOM  $\hat{\mathcal{M}} = (\hat{\omega}, \{\hat{\Xi}(x)\}_{x \in \mathcal{O}}, \hat{\sigma})$  with the observable operator matrix  $\hat{\Xi}(x) = \sum_{z \in \mathcal{X}} \hat{\mathbf{W}}_z \delta_z(x)$  supported on  $\mathcal{X} = \{s_n^2\}_{n=1}^N$  can be constructed by

---

**Algorithm 2** Nonequilibrium learning procedure of Binless OOMs
 

---

**INPUT:** Observation trajectories generated by a stochastic process  $\{x_t\}$  in  $\mathcal{O} \subset \mathbb{R}^d$

**OUTPUT:** Binless OOM  $\hat{\mathcal{M}} = (\hat{\omega}, \{\hat{\Xi}(x)\}_{x \in \mathcal{O}}, \hat{\sigma})$

- 1: Construct feature functions  $\phi_1 : \mathbb{R}^{Ld} \mapsto \mathbb{R}^{D_1}$  and  $\phi_2 : \mathbb{R}^{Ld} \mapsto \mathbb{R}^{D_2}$  with  $D_1, D_2 \geq m$ .
- 2: Calculate  $\bar{\phi}_1, \bar{\phi}_2, \mathbf{C}_{1,2}, \mathbf{C}_{1,3}(x)$  by (3), (4) and (19).
- 3: Choose matrix  $\mathbf{F}_1 \in \mathbb{R}^{D_1 \times m}, \mathbf{F}_2 \in \mathbb{R}^{D_2 \times m}$  such that  $\mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2$  is invertible.
- 4: Compute  $\hat{\sigma}, \hat{\omega}$  and  $\hat{\Xi}(x) = \sum_{z \in \mathcal{X}} \hat{\mathbf{W}}_z \delta_z(x)$  by (9), (17) and

$$\hat{\mathbf{W}}_{s_n^2} = \frac{1}{N} \left( \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \phi_1(\bar{s}_n^1) \phi_2(\bar{s}_n^3)^\top \mathbf{F}_2 \quad (21)$$

where  $\hat{\Xi}(\mathcal{O}) = \int dx \hat{\Xi}(x) = \sum_{z \in \mathcal{X}} \hat{\mathbf{W}}_z$ .

---

nonequilibrium learning with computational complexity  $O(N)$  as in Algorithm 2, where feature functions can be selected as splines, radial basis functions or other commonly used activation functions for single-layer neural networks in practice in order to digest adequate dynamical information from observation data. Note the binless strategy can be applied to more general cases where observations are strings, graphs or other structured variables, and is very similar to that used in Monte Carlo integration or nonparametric maximum likelihood estimation [24]. Although we cannot use the binless OOM to evaluate path probability densities of  $\{x_t\}$  as in (18), the equilibrium expectation of any observable  $g_t = g(x_{t+1:t+r})$  of  $\{x_t\}$  can be approximated as

$$\begin{aligned} \mathbb{E}_\infty [g_t] &\approx \mathbb{E} [g_t | \hat{\mathcal{M}}] \\ &= \sum_{x_{1:r} \in \mathcal{X}^r} g(x_{1:r}) \hat{\omega} \hat{\mathbf{W}}_{z_1} \dots \hat{\mathbf{W}}_{z_r} \hat{\sigma} \end{aligned} \quad (20)$$

By adding a technical assumption, our previous result on consistency of nonequilibrium learning of OOMs can be extended to the binless case as follows (see Appendix A.4 for proof):

**Assumption 4.** *The observation space  $\mathcal{O}$  is a closed set in  $\mathbb{R}^d$  and feature functions  $\phi_1, \phi_2$  are bounded on  $\mathcal{O}^L$ .*

**Theorem 3.** *Under Assumptions 1-4, the binless OOM provided by Algorithm 2 satisfies*

$$\mathbb{E} [g(x_{1:r}) | \hat{\mathcal{M}}] \xrightarrow{P} \mathbb{E}_\infty [g(x_{t+1:t+r})] \quad (22)$$

(i) *for all continuous functions  $g : \mathcal{O}^r \mapsto \mathbb{R}$ .*

(ii) *for all bounded and Borel measurable functions  $g : \mathcal{O}^r \mapsto \mathbb{R}$ , if there exist constants  $\bar{\xi}$  and  $\underline{\xi}$  so that  $\|\Xi(x)\| \leq \bar{\xi}$  and  $\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+r} = z_{1:r}) \geq \underline{\xi}$  for all  $x \in \mathcal{O}$  and  $z_{1:r} \in \mathcal{O}^r$ .*

Note that we do not assume the observed dynamics coincides with a parametric model defined by feature functions in Theorem 3. This theorem shows that binless OOMs allow us to consistently and efficiently extract equilibrium histograms, principle components, time-cross correlations, etc., of a dynamical systems from nonequilibrium data, which is important especially for thermodynamic and kinetic analysis in computational physics and chemistry.

*Remark 2.* The computational complexity of (20) is  $O(N^r)$ , which is unaffordable for large data sets if  $r > 1$ . In this paper, we focus on estimation of specific observables in the forms of  $\mathbb{E}_\infty [a(x_t)]$  and  $\mathbb{E}_\infty [a(x_t) b(x_{t+k})]$  by binless OOMs, which only require  $O(N)$  time. The efficient estimation of  $\mathbb{E}_\infty [g(x_{t+1:t+r})]$  for general  $g$  with is outside the scope of this paper and will be dealt with separately in another paper.

## 6 Applications

In this section, we evaluate our algorithms on two stochastic systems driven by Brownian dynamics and the molecular dynamics of alanine dipeptide, and compare them to several alternatives. The detailed settings of simulations and algorithms are provided in Appendix B.

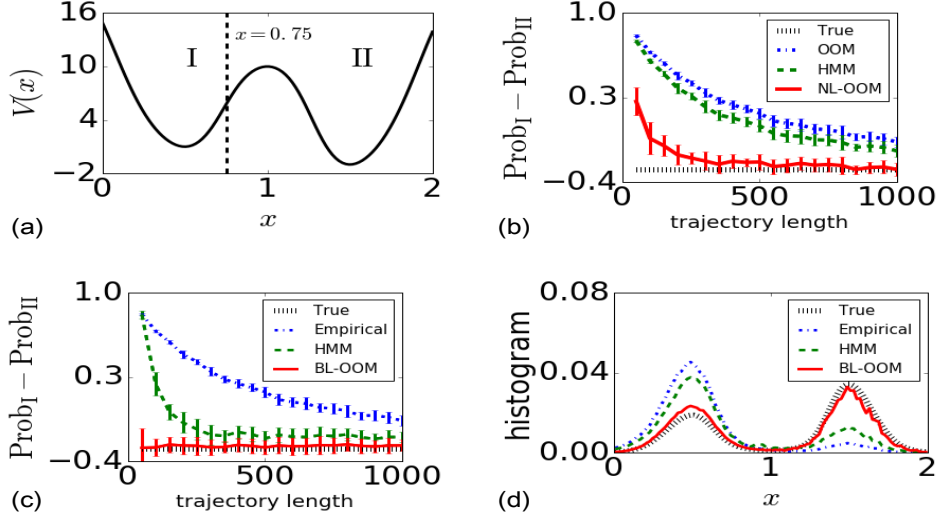


Figure 2: Comparison of modeling methods for a one-dimensional diffusion process. (a) Potential function. (b) Estimates of the difference between equilibrium probabilities of I and II given by the traditional OOM, HMM and OOM using nonequilibrium learning (NL-OOM) with  $\mathcal{O} = \{I, II\}$ . (c) Estimates of the probability difference given by the empirical estimator, HMM and Binless OOM (BL-OOM) using nonequilibrium learning with  $\mathcal{O} = [0, 2]$ . (d) Equilibrium histograms of  $\{x_t\}$  with 100 uniform bins estimated from trajectories with length 50. The initial  $x_0$  are uniformly drawn from  $[0, 0.5]$ , length of each trajectory is  $T = 50 \sim 1000$  and the number of trajectories is  $[10^5/T]$ . Error bars are standard deviations over 30 independent experiments.

**Brownian dynamics** Fig. 2(a) shows the potential function of a one-dimensional diffusion process  $\{x_t\}$  on  $[0, 2]$  driven by Brownian dynamics, where the state space is discretized into two clusters I, II. It is obvious that the equilibrium probability of finding  $x_t$  in I is smaller than that of  $x_t \in II$ , because the potential well contained in II is deeper than the other one. In this example, all simulations are performed by starting from a uniform distribution on  $[0, 0.2]$ , which implies that simulations are highly nonequilibrium and it is difficult to accurately estimate the equilibrium probabilities  $\text{Prob}_I = \mathbb{E}_\infty [1_{x_t \in I}]$  and  $\text{Prob}_{II} = \mathbb{E}_\infty [1_{x_t \in II}]$  of I and II from the simulation data. We first utilize the traditional OOM learning, expectation-maximization based HMM learning and the proposed nonequilibrium learning algorithm of OOMs to estimate  $\text{Prob}_I$  and  $\text{Prob}_{II}$  by assuming that we only know which cluster the  $x_t$  is in for each time  $t$ , i.e., the observation space  $\mathcal{O} = \{I, II\}$ . Fig. 2(b) summarizes the estimation results with different simulation lengths. It can be seen that estimates given by the traditional OOM and the HMM are far away from true values even for the largest simulation length  $T = 1000$ . In addition, it is worth pointing out that estimates given by the traditional OOM are very similar to empirical means of  $1_{x_t \in I}$  and  $1_{x_t \in II}$  because the OOM learning algorithm is essentially a moment matching algorithm and the estimated moments cannot be corrected in the traditional learning algorithm. (See Fig. 2(c). Note that the empirical estimates of  $\text{Prob}_I$  and  $\text{Prob}_{II}$  are the same for discrete and continuous observations.) In contrast to previous methods, the nonequilibrium learning based OOM effectively reduce the statistical bias in the nonequilibrium data, and achieves statistically correct estimation at  $T = 300$ .

Figs. 2(c) and 2(d) plot estimates of the equilibrium state distribution given by the empirical estimator, HMM and binless OOM using nonequilibrium learning under the condition that the value of  $x_t$  is exactly known and  $\mathcal{O} = [0, 2]$ , where the empirical estimator calculates statistics through averaging over all observations. The observation model of the HMM is constructed based on 100 uniform bins on the state space, where samples within the same bin are assumed to be independent. With such a fine discretization, the performance of the HMM is improved, but estimation errors of the HMM for short trajectory lengths are still large. Here, the proposed binless OOM significantly outperform the other methods, and its estimates are very close to true values even for extremely small short trajectories.

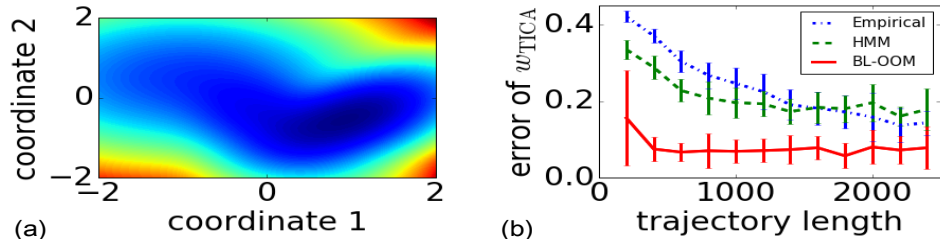


Figure 3: Comparison of modeling methods for a two-dimensional diffusion process. (a) Potential function. (b) Estimates of the coefficient vector  $w_{\text{TICA}} \in \mathbb{R}^2$  of the first TIC with lag time 100, which depends on  $\mathbb{E}_\infty [x_t]$ ,  $\mathbb{E}_\infty [x_t x_t^\top]$ , and  $\mathbb{E}_\infty [x_t x_{t+\tau}^\top]$ . The initial  $x_0$  are uniformly drawn from  $[-2, 0] \times [-2, 0]$ , length of each trajectory is  $T = 200 \sim 2500$  and the number of trajectories is  $[10^5/T]$ . Error bars are standard deviations over 30 independent experiments.

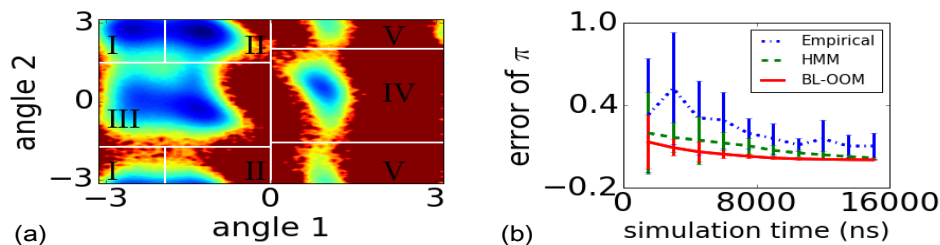


Figure 4: Comparison of modeling methods for molecular dynamics of alanine dipeptide. (a) Reduced free energy. (b) Estimates of  $\pi$ , the vector of equilibrium probabilities of metastable states I  $\sim$  V, where the horizontal axis denotes the total simulation time  $T \times I$ . Length of each trajectory is  $T = 10\text{ns}$  and the number of trajectories is  $I = 150 \sim 1500$ . Error bars are standard deviations over 30 independent experiments.

Fig. 3 provides an example of applying binless OOMs to kinetic analysis. The goal of this experiment is to perform the time-structure based independent component (TIC) analysis [25] of a two-dimensional Brownian dynamics based on nonequilibrium observation data. Fig. 3(b) displays the estimation errors of the coefficient vector of the first TIC obtained from different learning models, which also demonstrates the superiority of the proposed binless OOM method.

**Alanine dipeptide** Alanine dipeptide is a small molecule which consists of two alanine amino acid units, and its configuration can be described by two backbone dihedral angles. Fig. 4(a) shows the potential profile of the alanine dipeptide with respect to the two angles, which contains five metastable states. We perform multiple short molecular dynamics simulations starting from the metastable state IV, where each simulation length is 10ns, and utilizes different methods to approximate the stationary distribution of the five metastable states. As shown in Fig. 4(b), the proposed binless OOM yields lower estimation error compared to each of the alternatives.

## 7 Conclusion

In this paper, we investigated the statistical properties of the general OOM learning procedure for nonequilibrium data, and developed a general framework for learning dynamical models from nonequilibrium data. Under this framework, the existing learning approaches of OOMs and the other related models can be conveniently and efficiently applied to nonequilibrium (discrete or continuous) data by using the nonequilibrium learning technique and the binless learning technique. The main ideas of the two techniques are to correct the model parameters by the algebraic constraints under the equilibrium condition and to handle continuous observations in a binless manner. Interesting directions of future research include approximation error of nonequilibrium learning with finite data size and applications of nonequilibrium learning to controlled systems.



## References

- [1] H. Jaeger, “Observable operator models for discrete stochastic time series,” *Neural Comput.*, vol. 12, no. 6, pp. 1371–1398, 2000.
- [2] M.-J. Zhao, H. Jaeger, and M. Thon, “A bound on modeling error in observable operator models and an associated learning algorithm,” *Neural Comput.*, vol. 21, no. 9, pp. 2687–2712, 2009.
- [3] H. Jaeger, “Discrete-time, discrete-valued observable operator models: a tutorial,” tech. rep., International University Bremen, 2012.
- [4] M. L. Littman, R. S. Sutton, and S. Singh, “Predictive representations of state,” in *Adv. Neural. Inf. Process. Syst. 14 (NIPS 2001)*, pp. 1555–1561, 2001.
- [5] S. Singh, M. James, and M. Rudary, “Predictive state representations: A new theory for modeling dynamical systems,” in *Proc. 20th Conf. Uncertainty Artif. Intell. (UAI 2004)*, pp. 512–519, 2004.
- [6] E. Wiewiora, “Learning predictive representations from a history,” in *Proc. 22nd Intl. Conf. on Mach. Learn. (ICML 2005)*, pp. 964–971, 2005.
- [7] D. Hsu, S. M. Kakade, and T. Zhang, “A spectral algorithm for learning hidden Markov models,” in *Proc. 22nd Conf. Learning Theory (COLT 2009)*, pp. 964–971, 2005.
- [8] S. Siddiqi, B. Boots, and G. Gordon, “Reduced-rank hidden Markov models,” in *Proc. 13th Intl. Conf. Artif. Intell. Stat. (AISTATS 2010)*, vol. 9, pp. 741–748, 2010.
- [9] A. Beigel, F. Bergadano, N. H. Bshouty, E. Kushilevitz, and S. Varricchio, “Learning functions represented as multiplicity automata,” *J. ACM*, vol. 47, no. 3, pp. 506–530, 2000.
- [10] M. Thon and H. Jaeger, “Links between multiplicity automata, observable operator, models and predictive state representations — a unified learning framework,” *J. Mach. Learn. Res.*, vol. 16, pp. 103–147, 2015.
- [11] G. R. Bowman, V. S. Pande, and F. Noé, *An introduction to Markov state models and their application to long timescale molecular simulation*. Springer, 2013.
- [12] N. Schaudinnus, B. Bastian, R. Hegger, and G. Stock, “Multidimensional langevin modeling of nonoverdamped dynamics,” *Phys. Rev. Lett.*, vol. 115, no. 5, p. 050602, 2015.
- [13] L. R. Rabiner, “A tutorial on hidden markov models and selected applications in speech recognition,” *Proc. IEEE*, vol. 77, no. 2, pp. 257–286, 1989.
- [14] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Comp.*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [15] M. Shirts and V. S. Pande, “Screen savers of the world unite,” *Science*, vol. 290, pp. 1903–1904, 2000.
- [16] T.-K. Huang and J. Schneider, “Spectral learning of hidden Markov models from dynamic and static data,” in *Proc. 30th Intl. Conf. on Mach. Learn. (ICML 2013)*, pp. 630–638, 2013.
- [17] M. K. Cowles and B. P. Carlin, “Markov chain monte carlo convergence diagnostics: a comparative review,” *J. Am. Stat. Assoc.*, vol. 91, no. 434, pp. 883–904, 1996.
- [18] N. Jiang, A. Kulesza, and S. Singh, “Improving predictive state representations via gradient descent,” in *Proc. 30th AAAI Conf. Artif. Intell. (AAAI 2016)*, 2016.
- [19] H. Jaeger, “Modeling and learning continuous-valued stochastic processes with OOMs,” Tech. Rep. GMD-102, German National Research Center for Information Technology (GMD), 2001.
- [20] B. Boots, S. M. Siddiqi, G. Gordon, and A. Smola, “Hilbert space embeddings of hidden markov models,” in *Proc. 27th Intl. Conf. on Mach. Learn. (ICML 2010)*, 2010.
- [21] M. Rosencrantz, G. Gordon, and S. Thrun, “Learning low dimensional predictive representations,” in *Proc. 22nd Intl. Conf. on Mach. Learn. (ICML 2004)*, pp. 88–95, ACM, 2004.
- [22] B. Boots, *Spectral Approaches to Learning Predictive Representations*. PhD thesis, Carnegie Mellon University, 2012.
- [23] M.-J. Zhao, H. Jaeger, and M. Thon, “Making the error-controlling algorithm of observable operator models constructive,” *Neural Comput.*, vol. 21, no. 12, pp. 3460–3486, 2009.
- [24] A. Kong, P. McCullagh, X.-L. Meng, D. Nicolae, and Z. Tan, “A theory of statistical models for monte carlo integration,” *J. R. Stat. Soc. B*, vol. 65, no. 3, pp. 585–604, 2003.
- [25] G. Perez-Hernandez, F. Paul, T. Giorgino, G. De Fabritiis, and F. Noé, “Identification of slow molecular order parameters for markov model construction,” *J. Chem. Phys.*, vol. 139, no. 1, p. 015102, 2013.

# Supplementary Information

## A Proofs

### A.1 Proof of Theorem 1

For convenience, here we define

$$\boldsymbol{\omega}'(T) = \frac{1}{T-2L} \boldsymbol{\omega} \sum_{t=1}^{T-2L} \boldsymbol{\Xi}(\mathcal{O})^{t-1}$$

and

$$\mathbf{G}_\sigma = \sum_{z_{1:L}} \phi_2(z_{1:L}) \boldsymbol{\sigma}^\top \boldsymbol{\Xi}(z_{1:L})^\top \quad (\text{A.1})$$

**Part (1)** We first show the theorem in the case of  $T = T_0$  and  $I \rightarrow \infty$ .

Let

$$\mathbf{G}_\omega = \sum_{z_{1:L}} \phi_1(z_{1:L}) \boldsymbol{\omega}'(T_0) \boldsymbol{\Xi}(z_{1:L}) \quad (\text{A.2})$$

Since  $I \rightarrow \infty$ , we have

$$\begin{aligned} \hat{\boldsymbol{\phi}}_1 &\xrightarrow{p} \mathbb{E} [\hat{\boldsymbol{\phi}}_1] = \mathbf{G}_\omega \boldsymbol{\sigma} \\ \hat{\boldsymbol{\phi}}_2^\top &\xrightarrow{p} \mathbb{E} [\hat{\boldsymbol{\phi}}_2^\top] = \boldsymbol{\omega}'(T_0) \mathbf{G}_\sigma^\top \\ \hat{\mathbf{C}}_{1,2} &\xrightarrow{p} \mathbb{E} [\hat{\mathbf{C}}_{1,2}] = \mathbf{G}_\omega \mathbf{G}_\sigma^\top \\ \hat{\mathbf{C}}_{1,3}(x) &\xrightarrow{p} \mathbb{E} [\hat{\mathbf{C}}_{1,3}(x)] = \mathbf{G}_\omega \boldsymbol{\Xi}(x) \mathbf{G}_\sigma^\top \end{aligned}$$

In addition, we can obtain from Assumption 3 that

$$\text{rank}(\mathbf{G}_\omega) = \text{rank}(\mathbf{F}_1^\top \mathbf{G}_\omega) = \text{rank}(\mathbf{G}_\sigma) = \text{rank}(\mathbf{G}_\sigma^\top \mathbf{F}_2) = m$$

Therefore,  $\hat{\mathcal{M}}$  satisfies

$$\begin{aligned} \hat{\boldsymbol{\omega}} &= \hat{\boldsymbol{\phi}}_2^\top \mathbf{F}_2 \\ &\xrightarrow{p} \boldsymbol{\omega}'(T_0) \mathbf{G}_\sigma \mathbf{F}_2 \\ \hat{\boldsymbol{\Xi}}(x) &= \left( \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,3}(x) \mathbf{F}_2 \\ &\xrightarrow{p} \left( \mathbf{F}_1^\top \mathbf{G}_\omega \mathbf{G}_\sigma^\top \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \mathbf{G}_\omega \boldsymbol{\Xi}(x) \mathbf{G}_\sigma^\top \mathbf{F}_2 \\ &= \boldsymbol{\Xi}_{\text{eq}}(x) \\ \hat{\boldsymbol{\sigma}} &= \left( \mathbf{F}_1^\top \hat{\mathbf{C}}_{1,2} \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \hat{\boldsymbol{\phi}}_1 \\ &\xrightarrow{p} \left( \mathbf{F}_1^\top \mathbf{G}_\omega \mathbf{G}_\sigma^\top \mathbf{F}_2 \right)^{-1} \mathbf{F}_1^\top \mathbf{G}_\omega \boldsymbol{\sigma} \\ &= \boldsymbol{\sigma}_{\text{eq}} \end{aligned}$$

where  $\mathcal{M}_{\text{eq}} = (\boldsymbol{\omega}_{\text{eq}}, \{\boldsymbol{\Xi}_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \boldsymbol{\sigma}_{\text{eq}})$  is an OOM which equivalent to  $\mathcal{M}$  as

$$\begin{aligned} \boldsymbol{\omega}_{\text{eq}} &= \boldsymbol{\omega} \mathbf{G}_\sigma^\top \mathbf{F}_2 \\ \boldsymbol{\Xi}_{\text{eq}}(x) &= \left( \mathbf{G}_\sigma^\top \mathbf{F}_2 \right)^{-1} \boldsymbol{\Xi}(x) \left( \mathbf{G}_\sigma^\top \mathbf{F}_2 \right) \\ \boldsymbol{\sigma}_{\text{eq}} &= \left( \mathbf{G}_\sigma^\top \mathbf{F}_2 \right)^{-1} \boldsymbol{\sigma} \end{aligned} \quad (\text{A.3})$$

Note  $\hat{\boldsymbol{\omega}} \xrightarrow{p} \boldsymbol{\omega}_{\text{eq}}$  does not hold in general cases.

**Part (2)** We now consider the case of  $I = I_0$  and  $T \rightarrow \infty$ .

According to Assumption 2, the limit

$$\begin{aligned}\hat{\mathbf{C}}_{1,2} &\xrightarrow{p} \mathbb{E}_\infty [\phi_1(x_{t-L:t-1})\phi_2(x_{t:t+L-1})^\top] \\ &= \lim_{k \rightarrow \infty} \sum_{z_{1:L}} \phi_1(z_{1:L})\boldsymbol{\omega}\Xi(\mathcal{O})^k \Xi(z_{1:L})\mathbf{G}_\sigma^\top\end{aligned}$$

exists. Then

$$\begin{aligned}\hat{\phi}_1 &\xrightarrow{p} \mathbb{E}_\infty [\phi_1(x_{t-L:t-1})] = \mathbf{G}_\omega \boldsymbol{\sigma} \\ \hat{\phi}_2^\top &\xrightarrow{p} \mathbb{E}_\infty [\phi_2(x_{t:t+L-1})^\top] = \lim_{k \rightarrow \infty} \boldsymbol{\omega}\Xi(\mathcal{O})^k \mathbf{G}_\sigma^\top \mathbf{F}_2 \\ \hat{\mathbf{C}}_{1,2} &\xrightarrow{p} \mathbb{E}_\infty [\hat{\mathbf{C}}_{1,2}] = \mathbf{G}_\omega \mathbf{G}_\sigma^\top \\ \hat{\mathbf{C}}_{1,3}(x) &\xrightarrow{p} \mathbb{E}_\infty [\hat{\mathbf{C}}_{1,2}(x)] = \mathbf{G}_\omega \Xi(x) \mathbf{G}_\sigma^\top\end{aligned}$$

with

$$\mathbf{G}_\omega = \lim_{k \rightarrow \infty} \sum_{z_{1:L}} \phi_1(z_{1:L})\boldsymbol{\omega}\Xi(\mathcal{O})^k \Xi(z_{1:L}) \quad (\text{A.4})$$

and we can get

$$\text{rank}(\mathbf{G}_\omega) = \text{rank}(\mathbf{F}_1^\top \mathbf{G}_\omega) = \text{rank}(\mathbf{G}_\sigma) = \text{rank}(\mathbf{G}_\sigma^\top \mathbf{F}_2) = m$$

according to Assumption 3.

Therefore,

$$\begin{aligned}\hat{\boldsymbol{\omega}} &\xrightarrow{p} \lim_{k \rightarrow \infty} \boldsymbol{\omega}\Xi(\mathcal{O})^k \mathbf{G}_\sigma^\top \mathbf{F}_2 \\ \hat{\Xi}(x) &\xrightarrow{p} \Xi_{\text{eq}}(x) \\ \hat{\boldsymbol{\sigma}} &\xrightarrow{p} \boldsymbol{\sigma}_{\text{eq}}\end{aligned}$$

where  $\mathcal{M}_{\text{eq}} = (\boldsymbol{\omega}_{\text{eq}}, \{\Xi_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \boldsymbol{\sigma}_{\text{eq}})$  has the same definition as in (A.3). Note  $\hat{\boldsymbol{\omega}} \xrightarrow{p} \boldsymbol{\omega}_{\text{eq}}$  does not hold in general cases where  $\boldsymbol{\omega}\Xi(\mathcal{O}) \neq \boldsymbol{\omega}$ .

## A.2 Asymptotic correctness of nonequilibrium learning with different initial states

If the  $i$ -th observation trajectories is generated by OOM  $\mathcal{M} = (\boldsymbol{\omega}^i, \{\Xi(x)\}_{x \in \mathcal{O}}, \boldsymbol{\sigma})$  for  $i = 1, \dots, I$ , and

$$\boldsymbol{\omega}'' = \begin{cases} \frac{1}{I} \sum_{i=1}^I \boldsymbol{\omega}^i, & \text{for } T \rightarrow \infty \\ \text{plim}_{I \rightarrow \infty} \frac{1}{I} \sum_{i=1}^I \boldsymbol{\omega}^i, & \text{for } I \rightarrow \infty \end{cases}$$

the asymptotic correctness can also shown as in Appendix A.1 by setting

$$\mathbf{G}_\omega = \sum_{z_{1:L}} \phi_1(z_{1:L})\boldsymbol{\omega}'(T_0)\Xi(z_{1:L})$$

with

$$\boldsymbol{\omega}'(T) = \frac{1}{T-2L} \boldsymbol{\omega}'' \sum_{t=1}^{T-2L} \Xi(\mathcal{O})^{t-1}$$

for  $I \rightarrow \infty$ , and

$$\mathbf{G}_\omega = \lim_{k \rightarrow \infty} \sum_{z_{1:L}} \phi_1(z_{1:L})\boldsymbol{\omega}''\Xi(\mathcal{O})^k \Xi(z_{1:L})$$

for  $T \rightarrow \infty$ .

### A.3 Asymptotic correctness of nonequilibrium learning based on the spectral learning algorithm

In the spectral learning algorithm, matrices  $\mathbf{F}_1$  and  $\mathbf{F}_2$  are given by singular value decomposition (SVD) as

$$\mathbf{F}_1 = \mathbf{U}, \quad \mathbf{F}_2 = \mathbf{V}\mathbf{\Sigma}^{-1} \quad (\text{A.5})$$

where  $\mathbf{\Sigma} \in \mathbb{R}^{m \times m}$  is a diagonal matrix contains the top  $m$  singular values of  $\hat{\mathbf{C}}_{1,2}$ , and  $\mathbf{U}$  and  $\mathbf{V}$  consist of the corresponding  $m$  left and right singular vectors of  $\hat{\mathbf{C}}_{1,2}$ . In this appendix, we will prove the following theorem:

**Theorem 4.** *Under Assumptions 1 and 2, for the estimated OOM  $\hat{\mathcal{M}} = (\hat{\omega}, \{\hat{\Xi}(x)\}_{x \in \mathcal{O}}, \hat{\sigma})$  given by Algorithm 1 with  $\mathbf{F}_1, \mathbf{F}_2$  defined by (A.5), there exists an OOM  $\mathcal{M}' = (\omega', \{\Xi'(x)\}_{x \in \mathcal{O}}, \sigma')$  which is equivalent to  $\hat{\mathcal{M}}$  and satisfies*

$$\Xi'(x) \xrightarrow{p} \Xi_{\text{eq}}(x), \quad \forall x \in \mathcal{O} \quad (\text{A.6})$$

$$\sigma' \xrightarrow{p} \sigma_{\text{eq}} \quad (\text{A.7})$$

where  $\mathcal{M}_{\text{eq}} = (\omega_{\text{eq}}, \{\Xi_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \sigma_{\text{eq}})$  is an  $m$ -dimensional OOM equivalent to  $\mathcal{M}$ , if the rank of the limit of  $\hat{\mathbf{C}}_{1,2}$  is not less than  $m$ .

*Proof.* According to Appendix A.1, the limit of  $\hat{\mathbf{C}}_{1,2}$  can be expressed as

$$\hat{\mathbf{C}}_{1,2} \xrightarrow{p} \mathbf{G}_\omega \mathbf{G}_\sigma^\top$$

where  $\mathbf{G}_\omega$  and  $\mathbf{G}_\sigma$  have the same definitions as in Appendix A.1. So the limit of  $\hat{\mathbf{C}}_{1,2}$  has rank  $m$ . By the Eckart-Young-Mirsky Theorem,  $\hat{\mathbf{C}}_{1,2}^{\text{trun}} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$  is the best rank  $m$  approximation to  $\hat{\mathbf{C}}_{1,2}$  and therefore

$$\hat{\mathbf{C}}_{1,2}^{\text{trun}} \xrightarrow{p} \mathbf{G}_\omega \mathbf{G}_\sigma^\top$$

Let

$$\mathbf{G}_\omega \mathbf{G}_\sigma^\top = \tilde{\mathbf{U}} \tilde{\mathbf{\Sigma}} \tilde{\mathbf{V}}^\top$$

be the SVD of  $\mathbf{G}_\omega \mathbf{G}_\sigma^\top$ ,

$$\tilde{\mathbf{F}}_1 = \tilde{\mathbf{U}}, \quad \tilde{\mathbf{F}}_2 = \tilde{\mathbf{V}} \tilde{\mathbf{\Sigma}}^{-1}$$

and  $\mathcal{M}_{\text{eq}} = (\omega_{\text{eq}}, \{\Xi_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \sigma_{\text{eq}})$  be an OOM which is equivalent to  $\mathcal{M}$  with

$$\begin{aligned} \omega_{\text{eq}} &= \omega \mathbf{G}_\sigma^\top \tilde{\mathbf{F}}_2 \\ \Xi_{\text{eq}}(x) &= \left( \mathbf{G}_\sigma^\top \tilde{\mathbf{F}}_2 \right)^{-1} \Xi(x) \left( \mathbf{G}_\sigma^\top \tilde{\mathbf{F}}_2 \right) \\ \sigma_{\text{eq}} &= \left( \mathbf{G}_\sigma^\top \tilde{\mathbf{F}}_2 \right)^{-1} \sigma \end{aligned}$$

We can obtain from the Wedin Theorem and the continuity of singular values of matrix that

$$\begin{aligned} \min_{\mathbf{R}} \left\| \mathbf{U}\mathbf{R} - \tilde{\mathbf{U}} \right\| &= \left\| \mathbf{U}\mathbf{U}^\top \tilde{\mathbf{U}} - \tilde{\mathbf{U}} \right\| \xrightarrow{p} 0 \\ \min_{\mathbf{R}} \left\| \mathbf{V}\mathbf{R} - \tilde{\mathbf{V}} \right\| &= \left\| \mathbf{V}\mathbf{V}^\top \tilde{\mathbf{V}} - \tilde{\mathbf{V}} \right\| \xrightarrow{p} 0 \\ &\quad \mathbf{\Sigma} \xrightarrow{p} \tilde{\mathbf{\Sigma}} \end{aligned}$$

Therefore, we can construct an OOM  $\mathcal{M}' = (\omega', \{\Xi'(x)\}_{x \in \mathcal{O}}, \sigma')$  with

$$\begin{aligned}
\omega' &= \hat{\omega} \left( \Sigma \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right) \\
&= \hat{\phi}_2^\top \mathbf{V} \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \\
\Xi'(x) &= \left( \Sigma \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \hat{\Xi}(x) \left( \Sigma \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right) \\
&= \left( \tilde{\mathbf{U}}^\top \mathbf{U} \mathbf{U}^\top \hat{\mathbf{C}}_{1,2} \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \tilde{\mathbf{U}}^\top \mathbf{U} \mathbf{U}^\top \hat{\mathbf{C}}_{1,3}(x) \mathbf{V} \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \\
\sigma' &= \left( \Sigma \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \hat{\sigma} \\
&= \left( \tilde{\mathbf{U}}^\top \mathbf{U} \mathbf{U}^\top \hat{\mathbf{C}}_{1,2} \mathbf{V} \mathbf{V}^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \tilde{\mathbf{U}}^\top \mathbf{U} \mathbf{U}^\top \hat{\phi}_1
\end{aligned}$$

which is equivalent to  $\hat{\mathcal{M}}$  and satisfies

$$\begin{aligned}
\Xi'(x) &\xrightarrow{p} \left( \tilde{\mathbf{U}}^\top \mathbf{G}_\omega \mathbf{G}_\sigma^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \tilde{\mathbf{U}}^\top \mathbf{G}_\omega \Xi(x) \mathbf{G}_\sigma^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \\
&= \Xi_{\text{eq}}(x) \\
\sigma' &\xrightarrow{p} \left( \tilde{\mathbf{U}}^\top \mathbf{G}_\omega \mathbf{G}_\sigma^\top \tilde{\mathbf{V}} \tilde{\Sigma}^{-1} \right)^{-1} \tilde{\mathbf{U}}^\top \mathbf{G}_\omega \sigma \\
&= \sigma_{\text{eq}}
\end{aligned}$$

□

It is worth pointing out that we can also show conclusions of Theorems 2 and 3 with (A.5) by using similar proofs. The details proofs are omitted because they are trivial.

#### A.4 Proof of Theorem 2

**Part (1)** We first show that there is an OOM  $\bar{\mathcal{M}}_{\text{eq}} = (\bar{\omega}_{\text{eq}}, \{\bar{\Xi}_{\text{eq}}(x)\}_{x \in \mathcal{O}}, \sigma_{\text{eq}})$  which can describe the equilibrium dynamics of  $\{x_t\}$ , where  $\bar{\Xi}_{\text{eq}}(x)$  and  $\sigma_{\text{eq}}$  are defined in (A.3).

In the case of  $T = T_0$  and  $I \rightarrow \infty$ , we can obtain from Assumptions 2 and 3 that

$$\begin{aligned}
\lim_{k \rightarrow \infty} \mathbf{G}_\omega \Xi(\mathcal{O})^k \mathbf{G}_\sigma^\top &= \lim_{k \rightarrow \infty} \frac{1}{T_0 - 2L} \sum_{t=0}^{T_0 - 2L - 1} \mathbb{E} \left[ \phi_1(x_{t+1:t+L}) \phi_2(x_{t+L+k+1:t+2L+k})^\top \right] \\
&= \left( \frac{1}{T_0 - 2L} \sum_{t=0}^{T_0 - 2L - 1} \mathbb{E} [\phi_1(x_{t+1:t+L})] \right) \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \\
&= \mathbf{G}_\omega \sigma \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \\
\Rightarrow \lim_{k \rightarrow \infty} \Xi(\mathcal{O})^k &= \sigma \bar{\omega}
\end{aligned} \tag{A.8}$$

with

$$\bar{\omega} = \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \mathbf{G}_\sigma^{+\top} \tag{A.9}$$

where  $\mathbf{G}_\omega$  and  $\mathbf{G}_\sigma$  are defined by (A.2) and (A.1), and  $\mathbf{G}_\sigma^+$  denotes the Moore-Penrose pseudoinverse of  $\mathbf{G}_\sigma$ . Then

$$\begin{aligned}
\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+l} = z_{1:l}) &= \lim_{t \rightarrow \infty} \omega \Xi(\mathcal{O})^t \Xi(z_{1:l}) \sigma \\
&= \omega \Xi(\mathcal{O}) \sigma \bar{\omega} \Xi(z_{1:l}) \sigma \\
&= \bar{\omega} \Xi(z_{1:l}) \sigma \\
&= \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(z_{1:l}) \sigma_{\text{eq}}
\end{aligned}$$

with

$$\bar{\omega}_{\text{eq}} = \bar{\omega} \mathbf{G}_\sigma^\top \mathbf{F}_2 \tag{A.10}$$

In the case of  $T = T_0$  and  $I \rightarrow \infty$ , because  $\text{rank}(\mathbf{G}_\omega) = m$  for  $\mathbf{G}_\omega$  defined by (A.4), there is a sufficiently large but finite  $T'$  so that  $\text{rank}(\mathbf{G}'_\omega) = m$  with

$$\mathbf{G}'_\omega = \sum_{z_{1:L}} \phi_1(z_{1:L}) \omega \Xi(\mathcal{O})^{T'} \Xi(z_{1:L})$$

Considering

$$\begin{aligned} \lim_{k \rightarrow \infty} \mathbf{G}'_\omega \Xi(\mathcal{O})^k \mathbf{G}'_\omega^\top &= \lim_{k \rightarrow \infty} \mathbb{E} \left[ \phi_1(x_{T'+1:T'+L}) \phi_2(x_{T'+L+k+1:T'+2L+k})^\top \right] \\ &= \mathbf{G}'_\omega \sigma \left( \mathbb{E}_\infty \left[ \phi_2(x_{t+1:t+L})^\top \right] \right) \\ \Rightarrow \lim_{k \rightarrow \infty} \Xi(\mathcal{O})^k &= \sigma \bar{\omega} \end{aligned} \quad (\text{A.11})$$

with  $\bar{\omega}$  defined by (A.9), we can also conclude that

$$\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+l} = z_{1:l}) = \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(z_{1:l}) \sigma_{\text{eq}}$$

with  $\bar{\omega}_{\text{eq}}$  defined by (A.10).

Note in both cases,  $\bar{\omega}_{\text{eq}}$  satisfies  $\omega_{\text{eq}} \lim_{k \rightarrow \infty} \Xi(\mathcal{O})^k = \bar{\omega}_{\text{eq}}$  and

$$\begin{aligned} \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O}) &= \lim_{t \rightarrow \infty} \omega_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O})^{t+1} \\ &= \bar{\omega}_{\text{eq}} \\ \bar{\omega}_{\text{eq}} \sigma_{\text{eq}} &= \bar{\omega}_{\text{eq}} \Xi_{\text{eq}}(\mathcal{O}) \sigma_{\text{eq}} \\ &= \lim_{t \rightarrow \infty} \sum_{x \in \mathcal{O}} \mathbb{P}(x_t = x) = 1 \end{aligned}$$

**Part (2)** In this part, we show that

$$\mathbf{w} \Xi_{\text{eq}}(\mathcal{O}) = \mathbf{w}, \quad \mathbf{w} \sigma_{\text{eq}} = 1$$

has a unique solution  $\mathbf{w} = \bar{\omega}_{\text{eq}}$ .

According to Appendix A.1 and (A.8), (A.11), if  $\mathbf{w} \Xi_{\text{eq}}(\mathcal{O}) = \mathbf{w}$  and  $\mathbf{w} \sigma_{\text{eq}} = 1$ , we have

$$\begin{aligned} \mathbf{w} &= \lim_{k \rightarrow \infty} \mathbf{w} \Xi_{\text{eq}}(\mathcal{O})^k \\ &= \lim_{k \rightarrow \infty} \mathbf{w} (\mathbf{G}_\sigma^\top \mathbf{F}_2)^{-1} \Xi(\mathcal{O})^k (\mathbf{G}_\sigma^\top \mathbf{F}_2) \\ &= \mathbf{w} (\mathbf{G}_\sigma^\top \mathbf{F}_2)^{-1} \sigma \bar{\omega} (\mathbf{G}_\sigma^\top \mathbf{F}_2) \\ &= \mathbf{w} \sigma_{\text{eq}} \bar{\omega}_{\text{eq}} \\ &= \bar{\omega}_{\text{eq}} \end{aligned}$$

**Part (3)** We now show Theorem 2.

The optimization problem (17) can be equivalently transformed into an unconstrained problem

$$\hat{\omega} = \min_{\mathbf{w}} \left\| \mathbf{w}^{\text{proj}} \hat{\Xi}(\mathcal{O}) - \mathbf{w}^{\text{proj}} \right\|^2 + \left\| \mathbf{w}^{\text{proj}} - \mathbf{w} \right\|^2$$

where

$$\mathbf{w}^{\text{proj}} = \mathbf{w} (\mathbf{I} - \hat{\sigma} \hat{\sigma}^+) + \hat{\sigma}^+ \quad (\text{A.12})$$

denotes the projection of  $\mathbf{w}$  on the space  $\{\mathbf{w} | \mathbf{w} \hat{\sigma} = 1\}$ ,  $\mathbf{I}$  denotes the identity matrix of appropriate dimension, and  $\bar{\omega}_{\text{eq}}$  is the unique solution if  $\hat{\Xi}(\mathcal{O}) = \Xi_{\text{eq}}(\mathcal{O})$  and  $\hat{\sigma} = \sigma_{\text{eq}}$  according to the conclusion in Part (2). Then we can obtain that  $\hat{\omega} \xrightarrow{P} \bar{\omega}_{\text{eq}}$  according to Theorem 2.7 in [1], which yields the conclusion of Theorem 2.

**Part (4)** We derive in this part the closed-form solution to (17).

Since the projection of  $\mathbf{w}$  on the space  $\{\mathbf{w} | \mathbf{w}\hat{\sigma} = 1\}$  is  $\mathbf{w}^{\text{proj}}$  defined by (A.12), (17) can be equivalent transformed into

$$\min_{\mathbf{w}} \left\| \mathbf{w} (\mathbf{I} - \hat{\sigma}\hat{\sigma}^+) \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) + \hat{\sigma}^+ \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) \right\|^2$$

The solution to this problem is

$$\mathbf{w}^* = -\hat{\sigma}^+ \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) \left( (\mathbf{I} - \hat{\sigma}\hat{\sigma}^+) \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) \right)^+$$

which provides the optimal value of  $\hat{\omega}$  as

$$\begin{aligned} \hat{\omega} &= \mathbf{w}^* (\mathbf{I} - \hat{\sigma}\hat{\sigma}^+) + \hat{\sigma}^+ \\ &= \hat{\sigma}^+ - \hat{\sigma}^+ \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) \left( (\mathbf{I} - \hat{\sigma}\hat{\sigma}^+) \left( \hat{\Xi}(\mathcal{O}) - \mathbf{I} \right) \right)^+ (\mathbf{I} - \hat{\sigma}\hat{\sigma}^+) \end{aligned} \quad (\text{A.13})$$

### A.5 Proof of Theorem 3

Here we only consider the consistency of the binless OOM as  $I \rightarrow \infty$ . The proof can be easily extended to the case of  $T \rightarrow \infty$ . In addition, we denote  $\mathbb{E}_\infty[g(x_{t+1:t+r})]$  and  $\mathbb{E}[g(x_{1:r})|\hat{\mathcal{M}}]$  by  $\mathbb{E}_\infty[g]$  and  $\mathbb{E}_{\hat{\mathcal{M}}}[g]$  for convenience of notation.

**Part (1)** We first show that Theorem 3 holds for  $g(x_{t+1:t+r}) = 1_{x_{t+1:t+r} \in \mathcal{B}_{i_1} \times \mathcal{B}_{i_2} \times \dots \times \mathcal{B}_{i_r}}$ , where  $\mathcal{B}_1, \dots, \mathcal{B}_K$  is a partition of  $\mathcal{O}$ ,  $i_{1:r} \in \{1, \dots, K\}^r$ , and  $1_e$  denotes the indicator function of event  $e$ . In this case, we can construct a discrete OOM with observation space  $\{\mathcal{B}_1, \dots, \mathcal{B}_K\}$  by the nonequilibrium learning algorithm, which can provide the same estimate of  $\mathbb{E}_\infty[g(x_{t+1:t+r})]$  as  $\hat{\mathcal{M}}$ . Therefore, we can show  $\mathbb{E}_{\hat{\mathcal{M}}}[g] \xrightarrow{P} \mathbb{E}_\infty[g]$  by using the similar proof of Theorem 2.

**Part (2)** We now consider the case that  $g$  is a continuous function. According to the Heine-Cantor theorem,  $g$  is also uniformly continuous. Then, for an arbitrary  $\epsilon > 0$ , we can construct a simple function

$$\hat{g}(x_{t+1:t+r}) = \sum_{i_1, \dots, i_r} c_{i_1 i_2 \dots i_r} 1_{x_{t+1:t+r} \in \mathcal{B}_{i_1} \times \dots \times \mathcal{B}_{i_r}}$$

so that

$$|g(z_{1:r}) - \hat{g}(z_{1:r})| \leq \epsilon, \quad \forall z_{1:r} \in \mathcal{O}^r$$

where  $\{\mathcal{B}_1, \dots, \mathcal{B}_K\}$  is a partition of  $\mathcal{O}$ . Then, we have

$$|\mathbb{E}_\infty[g] - \mathbb{E}_\infty[\hat{g}]| \leq \mathbb{E}_\infty[|g - \hat{g}|] \leq \epsilon$$

and

$$|\mathbb{E}_\infty[\hat{g}] - \mathbb{E}_{\hat{\mathcal{M}}}[\hat{g}]| \xrightarrow{P} 0$$

as  $I \rightarrow \infty$  according to the conclusion of Part (1), where  $\mathbb{E}_\infty[g] = \mathbb{E}_\infty[g(x_{t+1:t+r})]$  and  $\mathbb{E}_{\hat{\mathcal{M}}}[g] = \mathbb{E}[g(x_{1:r})|\hat{\mathcal{M}}]$ .

It can be known from Assumption 4, there exists a constant  $\xi$  such that

$$1_{\max_{x \in \mathcal{X}} \|\hat{\Xi}(x)\| < \xi/|\mathcal{X}|} \xrightarrow{P} 1 \quad (\text{A.14})$$

Under the condition that  $\max_{x \in \mathcal{X}} \|\hat{\Xi}(x)\| < \xi/|\mathcal{X}|$ , we have

$$\begin{aligned} |\mathbb{E}_{\hat{\mathcal{M}}}[\hat{g}] - \mathbb{E}_{\hat{\mathcal{M}}}[g]| &= \hat{\omega}_0 \left( \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \hat{\Xi}(z_{1:r}) \right) \hat{\sigma} \\ &\leq \|\hat{\omega}_0\| \|\hat{\sigma}\| \left( \sum_{z_{1:r} \in \mathcal{X}^r} \frac{\xi^r \epsilon}{|\mathcal{X}|^r} \right) \\ &= \|\hat{\omega}_0\| \|\hat{\sigma}\| \xi^r \epsilon \end{aligned}$$

In addition, considering that we can show as in Appendix A.1 that  $\hat{\omega} \xrightarrow{P} \bar{\omega}_{\text{eq}}$  and  $\hat{\sigma} \xrightarrow{P} \sigma_{\text{eq}}$ , we can obtain

$$1_{\|\hat{\omega}\| \|\hat{\sigma}\| \leq \xi_0} \xrightarrow{P} 1 \quad (\text{A.15})$$

and

$$1_{|\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \leq \xi_0 \xi^r \epsilon} \xrightarrow{P} 1$$

where  $\xi_0$  is a constant larger than  $\|\bar{\omega}_{\text{eq}}\| \cdot \|\sigma_{\text{eq}}\|$ .

Based on the above analysis and the fact that

$$\begin{aligned} |\mathbb{E}_{\infty}[g] - \mathbb{E}_{\mathcal{M}}[g]| &= |\mathbb{E}_{\infty}[g] - \mathbb{E}_{\infty}[\hat{g}] + \mathbb{E}_{\infty}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}] + \mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \\ &\leq |\mathbb{E}_{\infty}[g] - \mathbb{E}_{\infty}[\hat{g}]| + |\mathbb{E}_{\infty}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| + |\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \end{aligned}$$

we can get

$$\begin{aligned} \Pr(|\mathbb{E}_{\infty}[g] - \mathbb{E}_{\mathcal{M}}[g]| \leq (\xi_0 \xi^r + 2)\epsilon) &\geq \Pr(|\mathbb{E}_{\infty}[g] - \mathbb{E}_{\infty}[\hat{g}]| \leq \epsilon, |\mathbb{E}_{\infty}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| \leq \epsilon, \\ &\quad |\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \leq \xi_0 \xi^r \epsilon) \\ &\rightarrow 1 \end{aligned}$$

Because this equation holds for all  $\epsilon > 0$ , we can conclude that  $\mathbb{E}_{\mathcal{M}}[g] \xrightarrow{P} \mathbb{E}_{\infty}[g]$ .

**Part (3)** In this part, we prove the conclusion of the theorem in the case where  $g$  is a Borel measurable function and bounded with  $|g(z_{1:r})| < \xi_g$  for all  $z_{1:r} \in \mathcal{O}^r$ , and there exist constants  $\bar{\xi}$  and  $\underline{\xi}$  so that  $\|\Xi(x)\| \leq \bar{\xi}$  and  $\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+r} = z_{1:r}) \geq \underline{\xi}$  for all  $x \in \mathcal{O}$  and  $z_{1:r} \in \mathcal{O}^r$ .

According to Theorem 2.2 in [2], for an arbitrary  $\epsilon > 0$ , there is a continuous function  $\hat{g}'$  satisfies  $\mathbb{E}_{\infty}[1_{x_{t+1:t+r} \in \mathcal{K}_{\epsilon}(\hat{g}')}] < \epsilon$ , where  $\mathcal{K}_{\epsilon}(\hat{g}') = \{z_{1:r} | z_{1:r} \in \mathcal{O}^r, |\hat{g}'(z_{1:r}) - g(z_{1:r})| > \epsilon\}$ . Define

$$\hat{g}(z_{1:r}) = \begin{cases} \hat{g}'(z_{1:r}), & |\hat{g}'(z_{1:r})| \leq \xi_g \\ -\xi_g, & \hat{g}'(z_{1:r}) < -\xi_g \\ \xi_g, & \hat{g}'(z_{1:r}) > \xi_g \end{cases}$$

It can be seen that  $\hat{g}$  is a continuous function which is also satisfies  $\mathbb{E}_{\infty}[1_{x_{t+1:t+r} \in \mathcal{K}_{\epsilon}(\hat{g})}] < \epsilon$  and bounded with  $|\hat{g}(z_{1:r})| < \xi_g$ . So the difference between  $\mathbb{E}_{\infty}[g]$  and  $\mathbb{E}_{\infty}[\hat{g}]$  satisfies

$$\begin{aligned} |\mathbb{E}_{\infty}[g] - \mathbb{E}_{\infty}[\hat{g}]| &\leq \mathbb{E}_{\infty}[|g(x_{t+1:t+r}) - \hat{g}(x_{t+1:t+r})|] \\ &= \mathbb{E}_{\infty}[1_{x_{t+1:t+r} \in \mathcal{K}_{\epsilon}(\hat{g})}] \mathbb{E}_{\infty}[|g(x_{t+1:t+r}) - \hat{g}(x_{t+1:t+r})| | x_{t+1:t+r} \in \mathcal{K}_{\epsilon}(\hat{g})] \\ &\quad + \mathbb{E}_{\infty}[1_{x_{t+1:t+r} \notin \mathcal{K}_{\epsilon}(\hat{g})}] \mathbb{E}_{\infty}[|g(x_{t+1:t+r}) - \hat{g}(x_{t+1:t+r})| | x_{t+1:t+r} \notin \mathcal{K}_{\epsilon}(\hat{g})] \\ &\leq \epsilon \cdot 2\xi_g + \epsilon = (2\xi_g + 1)\epsilon \end{aligned}$$

For the difference between  $\mathbb{E}_{\infty}[\hat{g}]$  and  $\mathbb{E}_{\mathcal{M}}[\hat{g}]$ , we can obtain from the above that  $|\mathbb{E}_{\infty}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| \xrightarrow{P} 0$  as  $I \rightarrow \infty$  by considering that  $\hat{g}$  is continuous, which implies that there is an  $I_0$  such that

$$\Pr(|\mathbb{E}_{\infty}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| > \epsilon) < \epsilon, \quad \forall I > I_0$$

Next, let us consider the value of  $|\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]|$ . Note that

$$\begin{aligned} |\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| &\leq \|\hat{\omega}_0\| \|\hat{\sigma}\| \left\| \sum_{z_{1:n} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \hat{\Xi}(z_{1:r}) \right\| \\ &< \frac{\xi_0 \xi^r}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \end{aligned}$$

under the condition that  $\|\hat{\Xi}(x)\| < \xi/|\mathcal{X}|$  and  $\|\hat{\omega}\| \|\hat{\sigma}\| \leq \xi_0$ . Therefore, there exists an  $I_1$  such that

$$\Pr\left(|\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \geq \frac{\xi_0 \xi^r}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right|\right) < \epsilon, \quad \forall I > I_1 \quad (\text{A.16})$$



due to (A.14) and (A.15). Let  $x'_{1:r}$  denotes a random sample taken uniformly from  $\mathcal{X}^r$ . We can obtain that

$$\begin{aligned}\mathbb{P}(x'_{1:r}) &= \mathbb{P}(x'_1) \dots \mathbb{P}(x'_r) \\ &\leq (\|\omega\| \|\sigma\| \xi_O \bar{\xi})^r\end{aligned}$$

where  $\xi_O \geq \|\Xi(\mathcal{O})^k\|$  for any  $k \geq 0$ . Note  $\xi_O < \infty$  because we can show the existing of the limit of  $\{\|\Xi(\mathcal{O})^0\|, \|\Xi(\mathcal{O})^1\|, \dots\}$  by similar steps in Appendix A.4. Thus

$$\begin{aligned}\mathbb{E} \left[ \frac{1}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \right] &\leq \mathbb{E} [\mathbb{E} [|\hat{g}(x'_{1:r}) - g(x'_{1:r})| | \mathcal{X}]] \\ &= \mathbb{E} [|\hat{g}(x'_{1:r}) - g(x'_{1:r})|] \\ &= \mathbb{E} [1_{x'_{1:r} \in \mathcal{K}_\epsilon(\hat{g})} \mathbb{E} [|\hat{g}(x'_{1:r}) - g(x'_{1:r})| | x'_{1:r} \in \mathcal{K}_\epsilon(\hat{g})] \\ &\quad + \mathbb{E} [1_{x'_{1:r} \notin \mathcal{K}_\epsilon(\hat{g})} \mathbb{E} [|\hat{g}(x'_{1:r}) - g(x'_{1:r})| | x'_{1:r} \notin \mathcal{K}_\epsilon(\hat{g})]] \\ &\leq \xi_\mu \epsilon \cdot 2\xi_g + \epsilon = (2\xi_g \xi_\mu + 1) \epsilon\end{aligned}$$

where  $\xi_\mu = (\|\omega\| \|\sigma\| \xi_O \bar{\xi})^r / \xi$ . By the Markov's inequality, we have

$$\Pr \left[ \frac{1}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \geq \sqrt{\epsilon} \right] \leq (2\xi_g \xi_\mu + 1) \sqrt{\epsilon} \quad (\text{A.17})$$

Combining (A.16) and (A.17) leads to

$$\begin{aligned}\Pr (|\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \geq \xi_0 \xi^r \sqrt{\epsilon}) &\leq \Pr \left( |\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \geq \frac{\xi_0 \xi^r}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \right) \\ &\quad + \Pr \left( \frac{1}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \geq \sqrt{\epsilon} \right) \\ &\leq \epsilon + (2\xi_g \xi_\mu + 1) \sqrt{\epsilon}\end{aligned}$$

for all  $I > I_1$ .

From all the above, we have

$$\begin{aligned}\Pr (|\mathbb{E}_\infty[g] - \mathbb{E}_{\mathcal{M}}[g]| \leq 2(\xi_g + 1)\epsilon + \xi_0 \xi^r \sqrt{\epsilon}) \\ &\geq \Pr (|\mathbb{E}_\infty[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| \leq \epsilon, |\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| \leq \xi_0 \xi^r \sqrt{\epsilon}) \\ &\geq 1 - \Pr (|\mathbb{E}_\infty[\hat{g}] - \mathbb{E}_{\mathcal{M}}[\hat{g}]| > \epsilon) - \Pr (|\mathbb{E}_{\mathcal{M}}[\hat{g}] - \mathbb{E}_{\mathcal{M}}[g]| > \xi_0 \xi^r \sqrt{\epsilon}) \\ &\geq 1 - 2\epsilon - (2\xi_g \xi_\mu + 1) \sqrt{\epsilon}\end{aligned}$$

for all  $I > \max\{I_0, I_1\}$ , which yields  $\mathbb{E}_{\mathcal{M}}[g] \xrightarrow{P} \mathbb{E}_\infty[g]$  due to the arbitrariness of  $\epsilon$ .

## B Settings in applications

### B.1 Models

The diffusion processes in Section 6 are driven by the Brownian dynamics

$$dx_t = -\nabla V(x_t) dt + \sqrt{2\beta^{-1}} dW_t$$

with  $\beta = 0.3$ , sample interval 0.002s,

$$V(x) = \frac{\sum_{i=1}^5 (|x - c_i| + 0.001)^{-2} u_i}{\sum_{i=1}^5 (|x - c_i| + 0.001)^{-2}}$$

for the one-dimensional process, and  $\beta = 2$ , sample interval 0.01s,

$$V(x) = -\log \left( \sum_{i=1}^3 p_i \mathcal{N}(x | \mu_i, \Sigma_i) \right)$$

for the two-dimensional process, where  $c_{1:5} = (-0.3, 0.5, 1, 1.5, 2.3)$ ,  $u_{1:5} = (21, 4, 8, -1, 20)$ ,  $p_{1:3} = (0.25, 0.25, 0.5)$ ,  $\mu_1 = (0, -0.5)$ ,  $\mu_2 = (-1, 0.5)$ ,  $\mu_3 = (1, -0.5)$ . The simulation details of alanine dipeptide is given in [3].

## B.2 Algorithms

The parameters of discrete OOMs are chosen as:  $L = 3$ ,  $m = 10$ ,  $\mathbf{F}_1, \mathbf{F}_2$  are given by the truncated SVD and  $\phi_1 = \phi_2$  are indicator functions of all  $\mathcal{O}^L$  observation subsequences with length  $L$ .

The parameters of binless OOMs are almost the same as discrete ones, except  $\phi_1 = \phi_2$  are Gaussian activation functions with random weights of functional link neural networks with  $D_1 = D_2 = 100$ .

The number of hidden states of HMMs is 10. For continuous data, we partition the state space into 100 discrete bins  $k$ -mean clustering (except for the one-dimensional process), and then learn HMMs by the EM algorithm under the assumption that samples within the same bin are drawn independently.

## References

- [1] W. K. Newey and D. McFadden, “Large sample estimation and hypothesis testing,” *Handbook of Econometrics*, vol. 4, pp. 2111–2245, 1994.
- [2] K. Hornik, M. Stinchcombe, and H. White, “Multilayer feedforward networks are universal approximators,” *Neural Netw.*, vol. 2, no. 5, pp. 359–366, 1989.
- [3] B. Trendelkamp-Schroer and F. Noé, “Efficient estimation of rare-event kinetics,” *Phys. Rev. X*, vol. 6, pp. 011009, 2016.