

# Supplementary Information

## A Proofs

### A.1 Proof of Theorem 1

For convenience, here we define

$$\begin{aligned}\omega^*(T) &= \frac{1}{T-2L} \omega \sum_{t=1}^{T-2L} \Xi(\mathcal{O})^{t-1} \\ \mathbf{G}_\sigma &= \sum_{z_{1:L}} \phi_2(z_{1:L}) \sigma^\top \Xi(z_{1:L})^\top\end{aligned}\tag{A.1}$$

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

Let

$$\mathbf{G}_\omega = \sum_{z_{1:L}} \phi_1(z_{1:L}) \omega^*(T_0) \Xi(z_{1:L})\tag{A.2}$$

Since  $I \rightarrow \infty$ , we have

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

According to Assumption 3 and the Eckart-Young-Mirsky Theorem, we can conclude that

$$\text{rank}(\mathbf{G}_\omega) = \text{rank}(\mathbf{G}_\sigma) = \text{rank}(\hat{\mathbf{C}}_{1,2}) = m$$

and

$$\hat{\mathbf{C}}_{1,2}^{\text{trun}} = \mathbf{U} \Sigma \mathbf{V}^\top \xrightarrow{P} \mathbf{G}_\omega \mathbf{G}_\sigma^\top$$

By using the SVD of  $\mathbf{G}_\omega \mathbf{G}_\sigma^\top$

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

with  $\text{rank}(\tilde{\mathbf{U}}) = \text{rank}(\tilde{\mathbf{V}}) = \text{rank}(\tilde{\Sigma})$ , we can construct an OOM  $\mathcal{M}' = (\omega', \{\Xi'(x)\}_{x \in \mathcal{O}}, \sigma')$  with

$$\omega' = \hat{\omega} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1}\tag{A.3}$$

$$\Xi'(x) = (\mathbf{G}_\sigma^\top \mathbf{V}) \hat{\Xi}(x) (\mathbf{G}_\sigma^\top \mathbf{V})^{-1}\tag{A.4}$$

$$\sigma' = (\mathbf{G}_\sigma^\top \mathbf{V}) \hat{\sigma}\tag{A.5}$$

which is obviously equivalent to  $\hat{\mathcal{M}}$ .

We can obtain from  $\text{rank}(\mathbf{U} \Sigma \mathbf{V}^\top) = \text{rank}(\mathbf{G}_\omega \mathbf{G}_\sigma^\top) = m$  that

$$(\mathbf{U} \Sigma \mathbf{V}^\top)^+ = \mathbf{V} \Sigma^{-1} \mathbf{U}^\top \xrightarrow{P} (\mathbf{G}_\omega \mathbf{G}_\sigma^\top)^+$$

where  $\mathbf{A}^+$  denotes the Moore-Penrose pseudoinverse of  $\mathbf{A}$ , so

$$\begin{aligned}
\boldsymbol{\omega}' &= \hat{\phi}_2^\top \mathbf{V} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1} \\
&\xrightarrow{p} \boldsymbol{\omega}^*(T_0) \\
\boldsymbol{\Xi}'(x) &= (\mathbf{G}_\sigma^\top \mathbf{V}) \boldsymbol{\Sigma}^{-1} \mathbf{U}^\top \hat{\mathbf{C}}_{1,3}(x) \mathbf{V} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1} \\
&\xrightarrow{p} \mathbf{G}_\sigma^\top \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^\top \mathbf{G}_\omega \boldsymbol{\Xi}(x) \\
&\xrightarrow{p} \mathbf{G}_\sigma^\top (\mathbf{G}_\omega \mathbf{G}_\sigma^\top)^+ \mathbf{G}_\omega \boldsymbol{\Xi}(x) \\
&= \mathbf{G}_\omega^+ \mathbf{G}_\omega \mathbf{G}_\sigma^\top (\mathbf{G}_\omega \mathbf{G}_\sigma^\top)^+ \mathbf{G}_\omega \mathbf{G}_\sigma^\top \mathbf{G}_\sigma^{+\top} \boldsymbol{\Xi}(x) \\
&= \boldsymbol{\Xi}(x) \\
\boldsymbol{\sigma}' &= \mathbf{G}_\sigma^\top \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^\top \hat{\phi}_1 \\
&\xrightarrow{p} \boldsymbol{\sigma}
\end{aligned}$$

Note  $\boldsymbol{\omega}' \xrightarrow{p} \boldsymbol{\omega}$  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

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

exists. Then

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

with

$$\mathbf{G}_\omega = \lim_{k \rightarrow \infty} \sum_{z_{1:L}} \phi_1(z_{1:L}) \boldsymbol{\omega} \boldsymbol{\Xi}(\mathcal{O})^k \boldsymbol{\Xi}(z_{1:L}) \quad (\text{A.6})$$

The remaining part of the proof is omitted because it is the same as in Part (1).

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

If the  $i$ -th observation trajectories is generated by OOM  $\mathcal{M} = (\boldsymbol{\omega}^i, \{\boldsymbol{\Xi}(x)\}_{x \in \mathcal{O}}, \boldsymbol{\sigma})$  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 be shown as in Appendix A.1 by setting

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

with

$$\boldsymbol{\omega}^*(T) = \frac{1}{T-2L} \boldsymbol{\omega}^{**} \sum_{t=1}^{T-2L} \boldsymbol{\Xi}(\mathcal{O})^{t-1}$$

for  $I \rightarrow \infty$ , and

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

for  $T \rightarrow \infty$ .

### A.3 Proof of Theorem 2

**Part (1)** We first show that there is an OOM  $\mathcal{M}_{\text{eq}} = (\omega_{\text{eq}}, \{\Xi(x)\}_{x \in \mathcal{O}}, \sigma)$  which can describe the equilibrium dynamics of  $\{x_t\}$ .

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

$$\begin{aligned}
\lim_{k \rightarrow \infty} \mathbf{G}_\omega \Xi(\mathcal{O})^k \mathbf{G}_\sigma^\top &= \lim_{k \rightarrow \infty} \frac{1}{T_0 - 2L} \sum_{t=0}^{T_0-2L-1} \mathbb{E} \left[ \phi_1(x_{t+1:t+L}) \phi_2(x_{t+L+k+1:t+2L+k})^\top \right] \\
&= \left( \frac{1}{T_0 - 2L} \sum_{t=0}^{T_0-2L-1} \mathbb{E} [\phi_1(x_{t+1:t+L})] \right) \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \\
&= \mathbf{G}_\omega \sigma \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \\
\Rightarrow \lim_{k \rightarrow \infty} \Xi(\mathcal{O})^k &= \sigma \omega_{\text{eq}}
\end{aligned} \tag{A.7}$$

with

$$\omega_{\text{eq}} = \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \mathbf{G}_\sigma^{+\top} \tag{A.8}$$

where  $\mathbf{G}_\omega$  and  $\mathbf{G}_\sigma$  are defined by (A.2) and (A.1). Then

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

In the case of  $I = I_0$  and  $T \rightarrow \infty$ , because  $\text{rank}(\mathbf{G}_\omega) = m$  for  $\mathbf{G}_\omega$  defined by (A.6), 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

$$\begin{aligned}
\lim_{k \rightarrow \infty} \mathbf{G}'_\omega \Xi(\mathcal{O})^k \mathbf{G}_\sigma^\top &= \lim_{k \rightarrow \infty} \mathbb{E} \left[ \phi_1(x_{T'+1:T'+L}) \phi_2(x_{T'+L+k+1:T'+2L+k})^\top \right] \\
&= \mathbf{G}'_\omega \sigma \left( \mathbb{E}_\infty [\phi_2(x_{t+1:t+L})^\top] \right) \\
\Rightarrow \lim_{k \rightarrow \infty} \Xi(\mathcal{O})^k &= \sigma \omega_{\text{eq}}
\end{aligned} \tag{A.9}$$

with  $\omega_{\text{eq}}$  defined by (A.8), we can also conclude that

$$\lim_{t \rightarrow \infty} \mathbb{P}(x_{t+1:t+l} = z_{1:l}) = \omega_{\text{eq}} \Xi(z_{1:l}) \sigma$$

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

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

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

$$\mathbf{w} \Xi(\mathcal{O}) = \mathbf{w}, \quad \mathbf{w} \sigma = 1$$

has a unique solution  $\mathbf{w} = \omega_{\text{eq}}$ .

According to Appendix A.1 and (A.7), (A.9), if  $\mathbf{w} \Xi(\mathcal{O}) = \mathbf{w}$  and  $\mathbf{w} \sigma = 1$ , we have

$$\begin{aligned}
\mathbf{w} &= \lim_{k \rightarrow \infty} \mathbf{w} \Xi(\mathcal{O})^k \\
&= \mathbf{w} \sigma \omega_{\text{eq}} \\
&= \omega_{\text{eq}}
\end{aligned}$$

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

The problem (16) is equivalent to

$$\begin{aligned} \min_{\mathbf{w}'} E(\mathbf{w}') &= (\mathbf{w}' \Xi'(\mathcal{O}) - \mathbf{w}') (\mathbf{G}_\sigma^\top \mathbf{V}) \\ &\quad \cdot (\mathbf{G}_\sigma^\top \mathbf{V})^\top (\mathbf{w}' \Xi'(\mathcal{O}) - \mathbf{w}')^\top \\ \text{s.t.} \quad &\mathbf{w}' \sigma' = 1 \end{aligned}$$

where  $\Xi'(\mathcal{O}) = \sum_{x \in \mathcal{O}} \Xi'(x)$ ,  $\Xi'(x)$  and  $\sigma'$  are given by (A.4) and (A.5), and  $\mathbf{w}'$  is related to  $\mathbf{w}$  with  $\mathbf{w}' = \mathbf{w} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1}$ . This problem can be further transformed into an unconstrained one

$$\min_{\mathbf{w}'} E(\mathbf{w}' (\mathbf{I} - \sigma' \sigma'^+) + \sigma'^+) + \|\mathbf{w}' (\mathbf{I} - \sigma' \sigma'^+) + \sigma'^+ - \mathbf{w}'\|^2 \quad (\text{A.10})$$

where  $\mathbf{w}' (\mathbf{I} - \sigma' \sigma'^+) + \sigma'^+$  is the projection of  $\mathbf{w}'$  on the space  $\{\mathbf{w}' | \mathbf{w}' \sigma' = 1\}$  and  $\mathbf{I}$  denotes the identity matrix of appropriate dimension. Considering that  $\Xi'(x) \xrightarrow{P} \Xi(x)$ ,  $\sigma' \xrightarrow{P} \sigma$ ,

$$\begin{aligned} (\mathbf{G}_\sigma^\top \mathbf{V}) (\mathbf{G}_\sigma^\top \mathbf{V})^\top &= \mathbf{G}_\sigma^\top \mathbf{V} \Sigma^{-1} \mathbf{U}^\top \mathbf{U} \Sigma \mathbf{V}^\top \mathbf{G}_\sigma \\ &\xrightarrow{P} \mathbf{G}_\sigma^\top (\mathbf{G}_\omega^\top \mathbf{G}_\sigma)^+ \mathbf{G}_\omega \mathbf{G}_\sigma^\top \mathbf{G}_\sigma \\ &= \mathbf{G}_\sigma^\top \mathbf{G}_\sigma \end{aligned}$$

and the conclusion in Part (2), we can obtain that the optimal solution of (A.10) converges to  $\omega_{\text{eq}}$  in probability and  $\hat{\omega}_{\text{eq}} \xrightarrow{P} \omega_{\text{eq}} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1}$  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 (16).

Since the projection of  $\mathbf{w}''$  on the space  $\{\mathbf{w}'' | \mathbf{w}'' \hat{\sigma} = 1\}$  is  $\mathbf{w}'' (\mathbf{I} - \hat{\sigma} \hat{\sigma}^+) + \hat{\sigma}^+$ , (16) can be equivalent transformed into

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

The solution to this problem is

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

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

$$\begin{aligned} \hat{\omega}_{\text{eq}} &= \mathbf{w}^* (\mathbf{I} - \hat{\sigma} \hat{\sigma}^+) + \hat{\sigma}^+ \\ &= \hat{\sigma}^+ - \hat{\sigma}^+ (\hat{\Xi}(\mathcal{O}) - \mathbf{I}) \left( (\mathbf{I} - \hat{\sigma} \hat{\sigma}^+) (\hat{\Xi}(\mathcal{O}) - \mathbf{I}) \right)^+ (\mathbf{I} - \hat{\sigma} \hat{\sigma}^+) \quad (\text{A.11}) \end{aligned}$$

#### A.4 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  $\mathbb{E}[g(x_{1:r}) | \hat{\mathcal{M}}_{\text{eq}}]$  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  $\mathcal{O}$  and  $i_{1:r} \in \{1, \dots, K\}^r$ . 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  $\hat{\mathcal{M}}_{\text{eq}}$ . 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

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

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  $\mathbb{E}_{\mathcal{M}}[g] = \mathbb{E}[g(x_{1:r})|\mathcal{M}_{\text{eq}}]$ .

It can be known from the boundness of feature functions, there exists a constant  $\xi$  such that

$$1_{\max_{x \in \mathcal{X}} \|\hat{\mathbf{W}}_x\| < \xi/|\mathcal{X}|} \xrightarrow{P} 1 \quad (\text{A.12})$$

Under the condition that  $\max_{x \in \mathcal{X}} \|\hat{\mathbf{W}}_x\| < \xi/|\mathcal{X}|$ , we have

$$\begin{aligned} |\mathbb{E}_{\mathcal{M}_{\text{eq}}}[\hat{g}] - \mathbb{E}_{\mathcal{M}_{\text{eq}}}[g]| &= \hat{\omega}_{\text{eq}} \left( \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \mathbf{W}_{z_1} \dots \mathbf{W}_{z_r} \right) \hat{\sigma} \\ &\leq \|\hat{\omega}_{\text{eq}}\| \|\hat{\sigma}\| \left( \sum_{z_{1:r} \in \mathcal{X}^r} \frac{\xi^r \epsilon}{|\mathcal{X}|^r} \right) \\ &= \|\hat{\omega}_{\text{eq}}\| \|\hat{\sigma}\| \xi^r \epsilon \end{aligned}$$

In addition, considering that we can show as in Appendix A.1 that

$$\begin{aligned} \hat{\omega}_{\text{eq}} &\xrightarrow{P} \omega_{\text{eq}} \mathbf{G}_\sigma^\top \mathbf{V} \\ \hat{\sigma} &\xrightarrow{P} (\mathbf{G}_\sigma^\top \mathbf{V})^{-1} \sigma \end{aligned}$$

we can obtain

$$1_{\|\hat{\omega}_{\text{eq}}\| \|\hat{\sigma}\| \leq \xi_0} \xrightarrow{P} 1 \quad (\text{A.13})$$

and

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

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

Based on the above analysis and the fact that

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

we can get

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

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

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

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

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

$$\begin{aligned} |\mathbb{E}_{\hat{\mathcal{M}}}[\hat{g}] - \mathbb{E}_{\hat{\mathcal{M}}}[g]| &\leq \|\hat{\omega}_0\| \|\hat{\sigma}\| \left\| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \hat{\mathbf{W}}_{z_1} \dots \hat{\mathbf{W}}_{z_r} \right\| \\ &< \frac{\xi_0 \xi^r}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \end{aligned}$$

under the condition that  $\|\hat{\mathbf{W}}_x\| < \xi/|\mathcal{X}|$  and  $\|\hat{\omega}_{\text{eq}}\| \|\hat{\sigma}\| \leq \xi_0$ . Therefore, there exists an  $I_1$  such that

$$\Pr\left(|\mathbb{E}_{\hat{\mathcal{M}}_{\text{eq}}}[\hat{g}] - \mathbb{E}_{\hat{\mathcal{M}}_{\text{eq}}}[g]| \geq \frac{\xi_0 \xi^r}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right|\right) < \epsilon, \quad \forall I > I_1 \quad (\text{A.14})$$

due to (A.12) and (A.13). 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.3. 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

$$\Pr \left[ \frac{1}{|\mathcal{X}|^r} \left| \sum_{z_{1:r} \in \mathcal{X}^r} (\hat{g}(z_{1:r}) - g(z_{1:r})) \right| \geq \sqrt{\epsilon} \right] \leq (2\xi_g \xi_\mu + 1) \sqrt{\epsilon} \quad (\text{A.15})$$

Combining (A.14) and (A.15) leads to

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

for all  $I > I_1$ .

From all the above, we have

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

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

## B Settings in applications

### B.1 Models

The one-dimensional diffusion processes in Section 5 are driven by the Brownian dynamics with  $\beta = 0.3$ ,

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

and the sample interval is 0.002. For the two-dimensional process,  $\beta = 2$ ,

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

and the sample interval is 0.01, 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 spectral learning are chosen as:  $L = 3$ ,  $m = 10$ , and  $\phi_1 = \phi_2$  are indicator functions of all  $\mathcal{O}^L$  observation subsequences with length  $L$ .

The parameters of binless spectral learning 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, and then learn HMMs by the EM algorithm, where the HMM package in PyEMMA [4] is used. All observation samples within the same bin are assumed to be independent for quantitative analysis.

## 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.
- [4] M. K. Scherer, B. Trendelkamp-Schroer, F. Paul, G. Pérez-Hernández, M. Hoffmann, N. Plattner, C. Wehmeyer, J. -H. Prinz, and F. Noé, "PyEMMA 2: A software package for estimation, validation, and analysis of Markov models," *J. Chem. Theory Comput.*, vol. 11, no. 11, pp. 5525–5542, 2015.