

Learning Good State and Action Representations for Markov Decision Process via Tensor Decomposition

Chengzhuo Ni

*Department of Electrical & Computer Engineering
Princeton University*

CHENGZHUO.NI@PRINCETON.EDU

Yaqi Duan

*LABORATORY FOR INFORMATION & DECISION SYSTEMS
MASSACHUSETTS INSTITUTE OF TECHNOLOGY*

YAQID@MIT.EDU

Munther Dahleh

*Electrical Engineering & Computer Science Department
Massachusetts Institute of Technology*

DAHLEH@MIT.EDU

Mengdi Wang¹

*Department of Electrical & Computer Engineering
Princeton University*

MENGDIW@PRINCETON.EDU

Anru R. Zhang¹

*Department of Biostatistics & Bioinformatics
Duke University*

ANRU.ZHANG@DUKE.EDU

Editor: Animashree Anandkumar

Abstract

The transition kernel of a continuous-state-action Markov decision process (MDP) admits a natural tensor structure. This paper proposes a tensor-inspired unsupervised learning method to identify meaningful low-dimensional state and action representations from empirical trajectories. The method exploits the MDP's tensor structure by kernelization, importance sampling and low-Tucker-rank approximation. This method can be further used to cluster states and actions respectively and find the best discrete MDP abstraction. We provide sharp statistical error bounds for tensor concentration and the preservation of diffusion distance after embedding. We further prove that the learned state/action abstractions provide accurate approximations to latent block structures if they exist, enabling function approximation in downstream tasks such as policy evaluation.

1. Introduction

State abstraction is a core problem at the heart of control and reinforcement learning (RL). In high-dimension RL, a naive grid discretization of the continuous state space often leads to exponentially many discrete states - an open challenge known as the curse of dimensionality. Having good state representations will significantly improve the efficiency of RL, by enabling the use of function approximation to better generalize knowledge from seen states to unseen states.

1. To whom the correspondence should be addressed to.

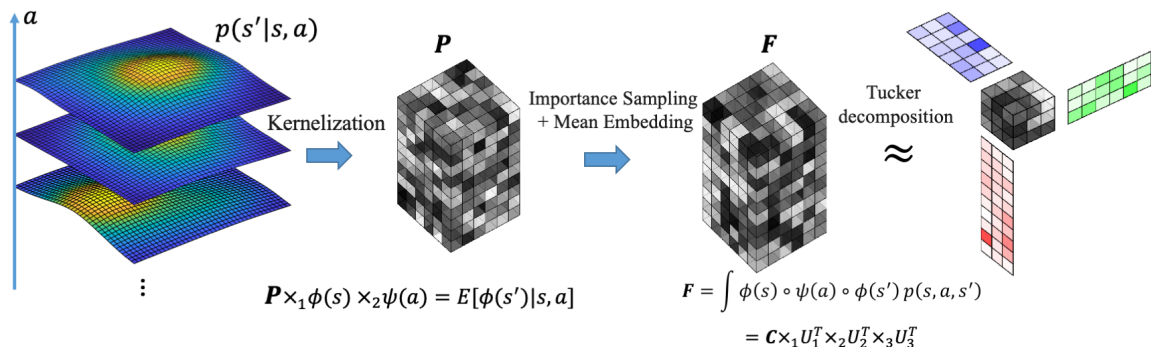


Figure 1: An illustration of our tensor-inspired state and action embedding method

We say a state/action representation is “good”, if it enables the use of function approximation to extrapolate and predict future value of unseen states. Suppose there is a representation allowing exact linear parametrization of the transition and value functions, then the sample complexity of RL reduces to depend linearly on d - the representation’s dimension (Lagoudakis and Parr, 2003; Zanette et al., 2019; Yang and Wang, 2019; Jin et al., 2019). Even if exact parametrization is not possible, a good representation can be still useful for solving RL with approximation error guarantee (see discussions in Du et al. (2019a); Lattimore and Szepesvari (2019)). An important related problem is to strategically explore in online RL while learning state abstractions (Du et al., 2019b; Misra et al., 2019). Motivated by these advances, we desire methods that can learn good representations, for RL with high-dimensional state and action spaces, automatically from empirical data.

What further complicates the problem is the large action space. An action can be either a one-step decision or a sequence of multi-step decisions (known as *option*). States under different actions lead to very different dynamics. Although states and actions may admit separate low-dimensional structures, they are entangled with each other in sample trajectories. This necessitates the tensor approach to decouple actions from states, so that we can learn their abstractions respectively.

1.1 Our Approach

In this paper, we study the state and action abstraction of Markov decision processes (MDP) from a tensor decomposition view. We focus on the batch data setting. The Tucker decomposition structure of a transition kernel p provides natural abstractions of the state and action spaces. We illustrate the low-Tucker-rank property in a number of reduced-order MDP models, including the block MDP (i.e., hard aggregation), latent-state model (i.e., soft aggregation).

Suppose we are given state-action-state transition samples $\mathcal{D} = \{(s, a, s')\}$ from a long sample path generated by a behavior policy. Our objective is to identify a state embedding map and an action embedding map, which map the original state and action spaces (maybe continuous and high-dimensional) into low-dimensional representations, respectively. The embedding maps are desired to be maximally “predicative”, by preserving a notion of

kernelized diffusion distance that measures similarity between states in terms of their future dynamics.

To handle continuous state and action spaces, we use nonparametric function approximation with known kernel functions over the state and action spaces. By approximately decomposing the kernel into finitely many features, we are able to handle the continuous problem by estimating a transition tensor of finite dimensions. Next, we leverage importance sampling and low-rank tensor approximation to identify the desired state and action embedding maps. They yield “good” representations of states and actions that are useful for linear function approximation in RL. Further, these representations can be used to find the best discrete approximation to the MDP, and in particular, recover the latent structures of block MDP with high accuracy. To the best of knowledge, this paper makes the first attempt to learn low-rank representations for high-dimensional continuous Markov decision, with statistical guarantee. Figure 1 illustrates the main idea of our approach. **Contributions** of this paper include:

- A tensor-inspired kernelized embedding method to learn low-dimensional state and action representations from empirical trajectories. The method exploits the MDP’s tensor structure by importance sampling, mean embedding and low-rank approximation.
- Theoretical guarantee that the embedding maps largely preserve the “predictability” of states and actions in terms of a kernelized diffusion distance, which is proved using a novel tensor concentration analysis.
- A discrete state/action abstraction method that provably recovers latent block structures of aggregable MDP. Theoretical guarantee that the learned abstractions are “good” representations for approximating transition/value functions within a small error tolerance.
- The numerical studies to corroborate our theoretical findings. The simulation results show the advantage of the proposed method over the baselines of vanilla and top r kernel PCA methods.

1.2 Related Literature

Spectral and low-rank methods for dimension reduction have a long history. Our approach traces back to the diffusion map approach for manifold learning and graph analysis (Lafon and Lee, 2006), which comes with a notion of diffusion distance that quantifies similarity between two nodes in a random walk. Coifman et al. (2008) extended the idea to systems driven by stochastic differential equations. Schütte et al. (2011) and Klus et al. (2016, 2020) studied how to infer dynamics of a system from leading spectrum of transition operator and find coresets of the state space.

The statistical theory of low-rank Markov model estimation received attention in recent years. Zhang and Wang (2020); Zhu et al. (2022) studied the low-rank estimation of finite-state Markov chains. Löffler and Picard (2021) studied the nonparametric estimation of transition kernel for continuous-state reversible Markov processes with exponentially decaying eigenvalues. Sun et al. (2019) studied kernelized state embedding and statistical estimation of metastable clusters. These results only apply to Markov processes.

In control theory and RL, state aggregation is a long known approach for reducing the complexity of the state space; see e.g., Moore (1991); Bertsekas and Tsitsiklis (1996); Singh et al. (1995); Tsitsiklis and Van Roy (1996); Ren and Krogh (2002). Representation learning methods were proposed that uses diagonalization or dilation of some Laplacian operator as a surrogate of the transition operator; see e.g. Johns and Mahadevan (2007); Mahadevan (2005); Parr et al. (2007); Petrik (2007). See Mahadevan et al. (2009) for a review. For online RL problems, representation learning approaches have been proposed to find good state-action representations while maintaining a sub-linear regret (Modi et al., 2021; Agarwal et al., 2020; Uehara et al., 2021; Zhang et al., 2022). (Ni et al., 2023) recently applied the representation learning approach to the multi-agent setting. These methods typically require prior knowledge about structures of the problem such as the transition function, or assume access to a finite feature class that covers the ground-truth feature. For tensor-based methods, Mahajan et al. (2021) uses low-rank tensor approximations to model agent interactions in the multi-agent setting. The approach views the Q-function as a tensor whose modes correspond to the action spaces of different agents. Van Der Vaart et al. (2021) considers model-based multi-agent RL and applies low-rank tensor approximation to estimate the transition probabilities and rewards. These approaches only apply to finite state-action MDPs with a low CP rank.

General methods for tensor decomposition and low-rank approximation have been studied in the applied math, statistics, and computer science literature, including the high-order singular value decomposition (HOSVD) (De Lathauwer et al., 2000b), high-order orthogonal iteration (HOOI) (De Lathauwer et al., 2000a), best low-rank approximation (Richard and Montanari, 2014; Zhang and Xia, 2018), sketched-based algorithms (Song et al., 2016), power iteration, k -means power iteration (Anandkumar et al., 2014; Sun et al., 2017), sparse high-order SVD (Zhang and Han, 2019), generalized tensor decomposition (Hong et al., 2020; Han et al., 2022), etc. The readers are also referred to surveys on tensor decomposition (Kolda and Bader, 2009; Cichocki et al., 2015) and their applications in machine learning (Sidiropoulos et al., 2017; Janzamin et al., 2019; Panagakis et al., 2021).

1.3 Markov Decision Process

An instance of a Markov decision process can be specified by a tuple $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$, where \mathcal{S} and \mathcal{A} are state and action spaces, p is the transition probability kernel, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the one-step reward function. At each step t , suppose the current state is s_t . If the agent chooses an action a_t , she will receive an instant reward $r(s_t, a_t) \in [0, 1]$ and system's state will transit to s_{t+1} according to the probability distribution $p(\cdot|s_t, a_t)$. A policy π is a rule for choosing actions based on states, where $\pi(\cdot|s)$ is a probability distribution over \mathcal{A} conditioned on $s \in \mathcal{S}$. Under a given policy, the transition of the MDP will reduce to a Markov chain, whose transition kernel is denoted by p^π where $p^\pi(s'|s) = p^{1,\pi}(s'|s) = \int_{\mathcal{A}} \pi(a|s)p(s'|s, a)da$. Based on that, we define the t -step transition kernel $p^{t,\pi}(\cdot|s)$ inductively by $p^{t,\pi}(\cdot|s) = \int p^{t-1,\pi}(s'|s)p^\pi(\cdot|s')ds'$. And we further use ν^π to denote the invariant distribution of that Markov chain. Define the worst-case mixing time (Levin et al., 2009, Page 55) as

$$t_{mix} = \max_{\pi} \min \{t \mid \|p^{t,\pi}(\cdot|s_0) - \nu^\pi\|_{TV} \leq 1/4, \forall s_0 \in \mathcal{S}, t' \geq t\},$$

where $\|\cdot\|_{TV}$ denotes the total variation distance. Throughout the paper, we use C to denote generic constants, while the actual values of C may vary from line to line.

1.4 Tensor and Tucker Decomposition

For a general tensor $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times \dots \times p_N}$, we denote $\mathbf{X} \times_n \mathbf{U}$ as the product between \mathbf{X} and a matrix $\mathbf{U} \in \mathbb{R}^{q \times p_n}$ on the n th mode, which is of size $p_1 \times \dots \times p_{n-1} \times q \times p_{n+1} \times \dots \times p_N$. Each element of $\mathbf{X} \times_n \mathbf{U}$ is defined as $(\mathbf{X} \times_n \mathbf{U})_{i_1 \dots i_{n-1} j i_{n+1} \dots i_N} = \sum_{i_n=1}^{p_n} \mathbf{X}_{i_1 \dots i_{n-1} i_n i_{n+1} \dots i_N} \mathbf{U}_{j i_n}$. We denote by $\mathcal{M}_k(\mathbf{X}) \in \mathbb{R}^{p_k \times \prod_{i \neq k} p_i}$ the factor- k matricization (or flattening) of \mathbf{X} . The Tucker decomposition of \mathbf{X} is of the form $\mathbf{X} = \mathbf{G} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N$, where $\mathbf{G} \in \mathbb{R}^{q_1 \times \dots \times q_N}$ is a smaller core tensor. In particular, we call the smallest size of \mathbf{G} the Tucker-rank of \mathbf{X} . Rigorously, we define $\text{Tucker-Rank}(\mathbf{X}) = (R_1, R_2, \dots, R_N)$, where $R_k = \text{Rank}(\mathcal{M}_k(\mathbf{X}))$. The inner product between two tensors $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{p_1 \times p_2 \times \dots \times p_N}$ is defined as

$$\langle \mathbf{X}, \mathbf{Y} \rangle = \sum_{i_1=1}^{p_1} \sum_{i_2=1}^{p_2} \dots \sum_{i_N=1}^{p_N} \mathbf{X}_{i_1 i_2 \dots i_N} \mathbf{Y}_{i_1 i_2 \dots i_N}.$$

The spectral norm and Frobenius norm of a tensor $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times \dots \times p_N}$ are defined as

$$\|\mathbf{X}\|_\sigma = \sup_{\|u_i\|=1, 1 \leq i \leq N} \langle \mathbf{X}, u_1 \circ u_2 \circ \dots \circ u_N \rangle, \quad \|\mathbf{X}\|_F = \sqrt{\langle \mathbf{X}, \mathbf{X} \rangle}.$$

Suppose S, A are reproducing kernel Hilbert space. We define the Tucker-rank of an operator $\mathbb{P} : S \times A \rightarrow S$ as an analogue of Tucker decomposition of tabular tensors: suppose there exist $c_{ijk} \in \mathbb{R}$ and functions $u_i, w_k \in \mathcal{H}_S, v_j \in \mathcal{H}_A, i \in [r], j \in [l], k \in [m]$, such that $(\mathbb{P}f)(s, a) = \sum_{i=1}^r \sum_{j=1}^l \sum_{k=1}^m c_{ijk} u_i(s) v_j(a) \langle f, w_k \rangle_{\mathcal{H}_S}$. Then write $\text{Tucker-Rank}(\mathbb{P})$ as the minimum (r, l, m) that ensure this equation holds.

2. A Tensor View of Markov Decision Process

Consider a continuous-state MDP with the transition kernel p , where each $p(\cdot|s, a)$ is a conditional transition density function. We adopt a tensor view to exploit structures of p for abstractions of state and action spaces. The Tucker rank of p turns out related to commonly used reduced-order models such as state aggregation and latent models. We handle the continuous state and action spaces using kernel function approximation. Suppose we have a Reproducing Kernel Hilbert Space (RKHS) \mathcal{H}_S for functions over states and a RKHS \mathcal{H}_A for functions over actions. We make the assumption that the MDP's transition kernel p can be represented in these function spaces.

Assumption 1 *Let \mathbb{P} be the transition operator of p , i.e., $(\mathbb{P}f)(s, a) = \int p(s'|s, a) f(s') ds'$. Assume that $\text{Tucker-Rank}(\mathbb{P}) \leq (r, l, m)$ ¹, and $\mathbb{P}f \in \mathcal{H}_S \times \mathcal{H}_A, \forall f \in \mathcal{H}_S$.*

Here, the low-Tucker rankness assumption captures the structure that state/action space can be compressed into a lower-dimensional space while preserving the dynamics. This assumption naturally holds in many well-known reinforcement learning models, such as soft state aggregation (Singh et al., 1995; Bertsekas, 2007; Sutton and Barto, 1998), rich-observation MDP (Azizzadenesheli et al., 2016; Du et al., 2019b), contextual MDP (Jiang

et al., 2017), linear/factor MDP (Jin et al., 2019), kernel MDP (Ormoneit and Glynn, 2002; Chowdhury and Gopalan, 2019).

In the remainder of the paper, we assume without loss of generality that the state and action kernel spaces admit finitely many known basis functions, which we refer to as state features $\phi(s) \in \mathbb{R}^{d_S}$ and action features $\psi(a) \in \mathbb{R}^{d_A}$. This is a rather mild assumption: Even if we do not know the basis function but are only given kernel functions K_S and K_A for \mathcal{H}_S and \mathcal{H}_A . According to Rahimi and Recht (2008), one can generate finitely many random features to approximately span these kernel spaces such that $K_S(s, s') \approx \sum_{i=1}^{d_S} \phi_i(s)^\top \phi_i(s')$ and $K_A(a, a') \approx \sum_{i=1}^{d_A} \psi_i(a)^\top \psi_i(a')$. Also note that our approach applies to *arbitrary* state and action spaces, as long as they come with appropriate kernel functions. Although p is infinitely dimensional, we use the given kernel spaces and represent p with a finite-dimensional tensor. In particular, Assumption 1 implies the following tensor linear model:

Lemma 1 (Conditional transition tensor and linear model) *Suppose Assumption 1 holds. There exists a tensor $\mathbf{P} \in \mathbb{R}^{d_S \times d_A \times d_S}$ such that Tucker-Rank(\mathbf{P}) $\leq (r, l, m)$ and*

$$\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top = \mathbb{E}[\phi(s')|s, a], \quad \forall s \in S, a \in A.$$

Tucker decomposition is one of the most general low-rank structure for tensors. Remarkably, the low-Tucker-rank property (Assumption 1) turns out to be universal in a number of reduced-order MDP models. Typical examples include block MDP (Du et al., 2019b) and soft MDP aggregation (Singh et al., 1995), whose detailed descriptions are placed in Appendix B. The low-Tucker-rank property also holds in MDPs with rich observations (Azizzadenesheli et al., 2016), and is related to the Bellman rank (Jiang et al., 2017). We remark that the tensor rank is determined solely by the transition model p (i.e., the environment), regardless of the reward r .

3. Tensor-Inspired State and Action Embedding Learning

In this section, we develop a tensor-inspired representation learning method, which embeds states and actions into decoupled low-dimensional spaces. Next, we will develop the method step by step, and provide theoretical guarantees.

3.1 Tensor MDP Mean Embedding by Importance Sampling

Suppose we have a batch dataset of state-action samples.

Assumption 2 *The data $\mathcal{D} = \{(s, a, s')\}$ consists of state-action-state transitions from a single sample path generated by a known behavior policy $\bar{\pi}$.*

Let ξ be the stationary state distribution of the sample path under policy $\bar{\pi}$. Let η be a positive probability measure over the action space. Consider the tensor mean embedding

$$\mathbf{F} = \int \phi(s) \circ \psi(a) \circ \phi(s') p(s, a, s') ds da ds' \in \mathbb{R}^{d_S \times d_A \times d_S},$$

where $p(s, a, s') = p(s'|s, a)\xi(s)\eta(a)$.

Lemma 2 *Assumption 1 implies Tucker-Rank(\mathbf{F}) $\leq (r, l, m)$.*

We estimate the mean embedding tensor \mathbf{F} by importance sampling:

$$\bar{\mathbf{F}} = n^{-1} \sum_{i=1}^n \frac{\eta(a_i)}{\bar{\pi}(a_i|s_i)} \cdot \phi(s_i) \circ \psi(a_i) \circ \phi(s'_i). \quad (1)$$

The mean embedding tensor \mathbf{F} is related to the transition tensor \mathbf{P} through a simple relation.

Lemma 3 (Relation between \mathbf{P} and \mathbf{F}) *When $\{\psi_i(\cdot)\}_{i=1}^{d_A}$ forms a set of orthogonal basis with respect to $L^2(\eta)$, we have $\mathbf{P} = \mathbf{F} \times_1 \mathbf{\Sigma}^{-1}$, where $\mathbf{\Sigma} = \int \xi(s)\phi(s)\phi(s)^\top ds$.*

Necessity of importance sampling. The importance sampling step (1) is necessary to decouple states from actions. Without importance sampling, the naive mean embedding tensor

$$\mathbf{W} := \int \phi(s) \circ \psi(a) \circ \phi(s') \xi(s) \bar{\pi}(a|s) p(s'|a, s) ds da ds'$$

may have large ranks on the first two dimensions. This is due to that the behavior policy $\bar{\pi}$ couples the state and action spaces together, therefore their independent low-dimensional structures are lost in the mean embedding tensor \mathbf{W} . Without using importance sampling, if we replace \mathbf{F} with plain mean \mathbf{W} , Lemma 2 and Lemma 3 no longer hold. As a result, one cannot learn the best low-dimensional structure of p from \mathbf{W} .

3.2 Low-Rank Estimation of Transition Tensor

We estimate a low-rank approximation to \mathbf{F} by solving:

$$\hat{\mathbf{F}} = \operatorname{argmin} \|\mathbf{Q} - \bar{\mathbf{F}}\|_\sigma, \text{ subject to Tucker-Rank}(\mathbf{Q}) \leq (r, l, m) \quad (2)$$

and estimate the transition operator \mathbf{P} by $\hat{\mathbf{P}} = \hat{\mathbf{F}} \times_1 \hat{\mathbf{\Sigma}}^{-1}$, where $\hat{\mathbf{\Sigma}} = \frac{1}{n} \sum_{i=1}^n \phi(s_i)\phi(s_i)^\top$. Define

$$\begin{aligned} K_{max} &= \max \left\{ \sup_s K_S(s, s), \sup_a K_A(a, a) \right\}, \\ \bar{\mu} &= \|\mathbb{E}[K_1(S, S)\phi(S)\phi(S)^\top]\|_\sigma, \\ \kappa &= \sup_{s \in \mathcal{S}, a \in \mathcal{A}} \frac{\eta(a)}{\pi(a|s)}, \\ \bar{\lambda} &= \sup_{\mathbf{u}, \mathbf{v}, \mathbf{w}} \mathbb{E}_{\xi \circ \eta \circ p(\cdot, \cdot)} [[(\mathbf{u}^\top \phi(S))(\mathbf{v}^\top \psi(A))(\mathbf{w}^\top \phi(S'))]^2], \text{ where } \mathbf{u}, \mathbf{w} \in S^{d_S-1}, \mathbf{v} \in S^{d_A-1}. \end{aligned}$$

Here, K_S, K_A are the kernels associated with the state RKHS space \mathcal{H}_S and action RKHS space \mathcal{H}_A , respectively.

Theorem 4 (Low-rank estimation of the transition tensor \mathbf{P}) *Suppose Assumptions 1-2 hold. Suppose ψ is orthonormal with respect to $L^2(\eta)$, and*

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\|\mathbf{\Sigma}^{-1}\|_\sigma^2 \bar{\mu} + \frac{K_{max}^2}{\bar{\mu}} + \frac{\kappa K_{max}^3}{\bar{\lambda}} \right) \left(\log \frac{2t_{mix}}{\delta} + 8(d_S + d_A) \right),$$

then with probability $1 - \delta$, we have

$$\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma \leq 256 \|\mathbf{\Sigma}^{-1}\|_\sigma \sqrt{\frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\mathbf{\Sigma}^{-1}\|_\sigma^2)}{(n/t_{mix}) \log^{-2}(n/t_{mix})}}.$$

The derivation of $\hat{\mathbf{P}}$ also provides a tractable way to estimate $\mathbb{E}[\phi(s')|s, a]$ by

$$\hat{\mathbb{E}}[\phi(s')|s, a] := \hat{\mathbf{P}} \times_1 \phi(s)^\top \times_2 \psi(a)^\top.$$

And we have the following guarantee on the estimation error:

$$\|\hat{\mathbb{E}}[\phi(s')|s, a] - \mathbb{E}[\phi(s')|s, a]\| \leq K_{max} \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma.$$

Rank selection Our theory assumes the prior knowledge of tensor rank. In practice, it is common to tune the rank parameters by checking the elbow in the scree plot and using cross validation (see discussions in the classical literature on PCA, e.g., Jolliffe (1986)). In theory, rank estimation is hard unless one makes additional strong assumptions, like that the eigengap is bounded from below.

Computation Finding the exact optimum of (2) can be computationally intense in general De Silva and Lim (2008). In practice, we can apply classic tensor decomposition algorithms, such as higher-order orthogonal iteration (HOOI) (De Lathauwer et al., 2000a), high-order SVD (De Lathauwer et al., 2000b), sequential-HOSVD (Vannieuwenhoven et al., 2012), gradient descent (Han et al., 2022), to find an approximate solution to (2). In particular, the statistical optimality of tensor power iterations, e.g., HOOI and HOSVD (Appendix A), have been justified in some special cases Zhang and Xia (2018). We expect these approximations also work for our problems, which is later validated in our experiment.

3.3 Learning State and Action Embeddings

Next, we show how to embed states and actions to low-dimensional representations to be maximally “predictive.” Consider a kernelized diffusion distance of the MDP, which measures similarity in terms of future dynamics restricted to a function class:

$$\text{dist}[(s_1, a_1), (s_2, a_2)] = \sup_{\|f\|_{\mathcal{H}_S} \leq 1} |\mathbb{E}[f(s')|s_1, a_1] - \mathbb{E}[f(s')|s_2, a_2]|.$$

This distance quantifies how well one can generalize the predicted value at a seen state-action pair (s, a) to a new (s', a') . Under the low-tensor-rank assumption, we have $\mathbf{P} = \mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3$, where $\mathbf{U}_1, \mathbf{U}_2, \mathbf{U}_3$ are columnwisely orthonormal matrices. Then we can define the *kernelized state diffusion map*, *kernelized action diffusion map* and their joint map as

$$f(\cdot) := \mathbf{U}_1^\top \phi(\cdot), \quad g(\cdot) := \mathbf{U}_2^\top \psi(\cdot), \quad \Phi(s, a) := \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top,$$

respectively. It follows that $\text{dist}[(s, a), (s', a')] = \|\Phi(s, a) - \Phi(s', a')\|$, if ϕ is a collection of orthonormal basis functions of \mathcal{H}_S . Motivated by the preceding analysis, we propose to estimate state and action embedding maps based on the tensor estimator. After we obtain $\hat{\mathbf{P}}$, we can simply find the corresponding state and action embedding maps from factors of its Tucker decomposition

$$\hat{\mathbf{P}} = \hat{\mathbf{C}} \times_1 \hat{\mathbf{U}}_1 \times_2 \hat{\mathbf{U}}_2 \times_3 \hat{\mathbf{U}}_3,$$

where we require that $\hat{\mathbf{U}}_k, k = 1, 2, 3$ are column-wisely orthonormal. The full procedure is given in Algorithm 1. Now we have obtained the state embedding map \hat{f} and the action

Algorithm 1 Learning State and Action Embedding Maps

 1: **Input:** $\{(s_i, a_i, s'_i)\}_{i=1}^n, (r, l, m)$

2: Calculate

$$\bar{\mathbf{F}} = \frac{1}{n} \sum_{i=1}^n \frac{\eta(a_i)}{\pi(a_i|s_i)} \phi(s_i) \circ \psi(a_i) \circ \phi(s'_i),$$

 and get $\hat{\mathbf{F}}$ as the low-rank approximation of $\bar{\mathbf{F}}$ using (2)

 3: Calculate $\hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n \phi(s_i) \phi^\top(s_i)$, $\hat{\mathbf{P}} = \hat{\mathbf{F}} \times_1 \hat{\Sigma}^{-1}$

 4: Let $\mathbf{P}_1 = \hat{\mathbf{P}}$. For $k = 1, 2, 3$, derive $\hat{\mathbf{U}}_k$ from the SVD

$$\mathcal{M}_k(\mathbf{P}_k) = \hat{\mathbf{U}}_k \mathbf{\Lambda}_k \mathbf{V}_k^\top,$$

 and let $\mathbf{P}_{k+1} = \mathbf{P}_k \times_k \hat{\mathbf{U}}_k$.

 5: **Output:**

 State and action embedding maps $\hat{f} : s \mapsto \hat{\mathbf{U}}_1^\top \phi(s)$, $\hat{g} : a \mapsto \hat{\mathbf{U}}_2^\top \psi(a)$;

 Core transition tensor $\hat{\mathbf{C}} = \mathbf{P}_4$.

embedding map \hat{g} . Accordingly, we define the joint state-action embedding and the empirical embedding distance as

$$\hat{\Phi}(s, a) = \hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top, \quad \widehat{\text{dist}}[(s, a), (s', a')] = \|\hat{\Phi}(s, a) - \hat{\Phi}(s', a')\|.$$

Theorem 5 (Embedding error bound) *Let Assumptions 1-2 hold. Suppose ϕ is an orthogonal basis of \mathcal{H}_S , and ψ is orthogonal w.r.t $L^2(\eta)$, then we can find an orthogonal matrix \mathbf{O} , such that*

$$\begin{aligned} \|\hat{\Phi}(s, a) - \mathbf{O}\Phi(s, a)\| &\leq \epsilon, \\ |\widehat{\text{dist}}[(s, a), (s', a')] - \text{dist}[(s, a), (s', a')]| &\leq 2\epsilon, \forall s, a, s', a', \end{aligned}$$

where ϵ is controlled by the low-rank estimation error, $\epsilon := K_{\max} \left(1 + \frac{\sqrt{2}\|\mathbf{P}\|_\sigma}{\sigma}\right) \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma$, and $\sigma := \sup_{\|w\| \leq 1} \sigma_m(\mathbf{P} \times_1 w^\top)$, where σ_m denotes the m -th singular value of a matrix.

Advantage of tensor method. As an alternative, one could ignore the tensor structure and treat the state and action jointly, yielding a low-dimensional representation for the pair (s, a) directly. This approach may be favorable if the (s, a) has a very simple joint structure. However, the tensor approach may be significantly more sample efficient if s and a admit *separate* low-dimensional structures. To see this, suppose the state and action features have dimensions d_S and d_A before embedding. Also assume the Tucker rank is $r = l = m$ for simplicity. By treating (s, a) jointly and ignoring the tensor structure, one would need $\tilde{\Omega}(d_A d_S r)$ samples to reliably recover the low-dimensional structure. In comparison, our tensor-based approach requires only $\tilde{\Omega}((d_X + d_A)r)$ samples.

4. Estimating the Optimal Discrete MDP Abstraction

Next, we study how to provably reduce a continuous-state continuous-action MDP into a discrete one, by an application of the learned kernelized diffusion distance to partition the state and action spaces.

4.1 Optimal Partition of State and Action Spaces

Our goal is to learn an optimal discretization of a continuous MDP. Specifically, we want to find a partition of \mathcal{S} and \mathcal{A} , denoted as blocks $A_i, B_j, i \in [n_s], j \in [n_a]$ and a collection of probability transition distributions $\{q_{ij}(\cdot)\}$ on the blocks. For each state-action pair $(s, a) \in A_i \times B_j$, and some function $f \in \mathcal{H}_S$, we want to approximate the one-step expected value $(\mathbb{P}f)(s, a) = \int p(s'|s, a)f(s')ds'$ by

$$(\mathbb{P}f)(s, a) = \int p(s'|s, a)f(s')ds' \approx \int q_{ij}(s')f(s')ds'.$$

We formalize the *optimal state-action partition problem* as:

$$\min_{\{A_i, B_j, q_{ij}\}} L(\{A_i, B_j, q_{ij}\}) := \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|p(\cdot|s, a) - q_{ij}(\cdot)\|_{\mathcal{H}_S}^2 dsda, \quad (3)$$

whose solution is denoted as $\{A_i^*, B_j^*, q_{ij}^*\}$, and the corresponding optimal value is denoted by L^* .

In particular, if $|\mathcal{A}| = 1$, the MDP reduces to a Markov process and the optimization problem reduces to $\min_{A_i} \min_{\{q_i\}} \sum_{i=1}^{n_s} \int_{A_i} \xi(s) \|p(\cdot|s) - q_i(\cdot)\|_{\mathcal{H}_S}^2 dsda$, which becomes equivalent to the metastable state partition problem for random walk and dynamic systems E et al. (2008).

4.2 Decoupled State and Action Clustering

Next, consider the RL setting where we wish to learn $\{A_i^*, B_j^*\}$ when p is unknown. Observe that the optimal partition is determined solely by the kernelized diffusion distance equipped by the state-action space. This allows the approximation of (3) by the *empirical state-action clustering problem*:

$$\min_{\{A_i, B_j, z_{ij}\}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \cdot \|\hat{\Phi}(s, a) - z_{ij}\|^2 dsda, \quad (4)$$

whose solution is denoted by $\{\hat{A}_i, \hat{B}_j, \hat{z}_{ij}\}$. Then the corresponding discrete transition distribution from state abstraction i and action abstraction j takes the form $\hat{q}_{ij}(\cdot) = \hat{z}_{ij}^\top \hat{U}_3^\top \phi(\cdot)$.

To facilitate computation, we provide a relaxation of problem (4) that can be solved using k -means-type algorithms. By taking $z_{ij} = \hat{C} \times_1 f_i^\top \times_2 g_j^\top$ for some f_i, g_j , the partition problem becomes

$$\min_{\{A_i, B_j\}} \min_{\{f_i, g_j\}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \cdot \|\hat{C} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \hat{C} \times_1 f_i^\top \times_2 g_j^\top\|^2 dsda.$$

Using the relation

$$\|\hat{C} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \hat{C} \times_1 f_i^\top \times_2 g_j^\top\|^2 \leq 2\|\hat{C}\|_\sigma^2 K_{max} (\|\hat{f}(s) - f_i\|^2 + \|\hat{g}(a) - g_j\|^2),$$

we can relax the problem (4) into two simpler subproblems:

$$\min_{A_i} \min_{f_i} \sum_i \int_{A_i} \xi(s) \|\hat{f}(s) - f_i\|^2 ds; \quad \min_{B_j} \min_{g_j} \sum_j \int_{B_j} \eta(a) \|\hat{g}(a) - g_j\|^2 da.$$

In short, one can efficiently compute the decoupled state and action clusters using the learned representations from the tensor method. The full procedure is given in Alg. 2.

Algorithm 2 Learning Optimal State-Action Abstractions

- 1: **Input:** $\{(s_i, a_i, s'_i)\}_{i=1}^n, (r, l, m)$
- 2: Estimate the state embedding and action embedding maps $\hat{f}, \hat{g}, \hat{U}_3$ using Algorithm 1.
- 3: Apply k-means to solve the following two problems, respectively:

$$\min_{\{A_i\}} \min_{\{f_i\}} \sum_i \int_{A_i} \xi(s) \|\hat{f}(s) - f_i\|^2 ds,$$

$$\min_{\{B_j\}} \min_{\{g_j\}} \sum_j \int_{B_j} \eta(a) \|\hat{g}(a) - g_j\|^2 da.$$

- 4: **Output:** The state partition $\{\hat{A}_j\}$ and action partition $\{\hat{B}_j\}$
-

4.3 Theoretical Guarantee

The following theorem guarantees that the empirical discretion is not far from the groundtruth.

Theorem 6 (Mean squared clustering error) *Let Assumptions 1-2 hold. Suppose ψ is a orthonormal basis with respect to $L^2(\eta)$ and n sufficiently large, then with probability at least $1 - \delta$, we have*

$$L(\{\hat{A}_i, \hat{B}_j, q_{ij}\}) \leq C \frac{\|\Sigma^{-1}\|_{\sigma} r l m (1 + \frac{2\bar{\lambda}\|\Sigma^{-1}\|_{\sigma}^2}{\sigma^2})}{\max\{r, l, m\}} \cdot \frac{\bar{\lambda}(\log(2t_{mix}/\delta) + d_S + d_A)(\kappa + \bar{\mu}\|\Sigma^{-1}\|_{\sigma}^2)}{(n/t_{mix}) \log^{-2}(n/t_{mix})} + 4L^*,$$

where L^* is the optimal value of problem (3), σ is defined as in Theorem 5, C is an absolute constant.

Next, we focus on the case where the true MDP has latent block structures.

Assumption 3 *Let there be blocks on the state and action spaces $A_i, B_j, i \in [n_s], j \in [n_a]$, i.e.,*

$$p(\cdot|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} q_{ij}^*(\cdot) \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j},$$

for some probability density functions q_{ij} .

Suppose we have applied Algorithm 2 to recover the latent blocks. Let $\{\hat{A}_i\}_{i=1}^{n_s}, \{\hat{B}_j\}_{j=1}^{n_a}$ be the estimated state and action clusters. Define the misclassification error as

$$M(\{\hat{A}_i\}, \{\hat{B}_j\}) = \min_{\sigma_1, \sigma_2} \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_{\sigma_1(i)} \times \hat{B}_{\sigma_2(j)}))}{(\xi \times \eta)(A_i \times B_j)},$$

where σ_1 and σ_2 are permutations over the state and action blocks, respectively. We prove the following clustering error bound:

Theorem 7 (Misclassification error for block MDP) *Let Assumptions 1, 2, 3 hold. For n sufficiently large, with probability $1 - \delta$, we have*

$$M(\{\hat{A}_i\}, \{\hat{B}_j\}) \leq C \|\Sigma^{-1}\|_\sigma \frac{rlm(1 + \frac{2\bar{\lambda}\|\Sigma^{-1}\|_\sigma^2}{\sigma^2})}{\max\{r, l, m\}} \cdot \frac{\bar{\lambda}(\log(2t_{mix}/\delta) + d_S + d_A)(\kappa + \bar{\mu}\|\Sigma^{-1}\|_\sigma^2)}{\Delta_1^2(n/t_{mix})(\log \frac{n}{t_{mix}})^{-2}},$$

where $\Delta_1^2 := \min_{i,j} \min_{(k,l) \neq (i,j)} \xi(A_i)\eta(B_j) \|q_{ij}^*(\cdot) - q_{kl}^*(\cdot)\|_{\mathcal{H}_S}^2$ and C is an absolute constant.

Remark 1 *The bounds in Theorems 6 and 7 grow proportionally to rlm (i.e., the product of Tucker ranks of \mathbf{P}), which can be large even when r, l, m are individually small. However, the term is essential as it represents the degree of freedom of a Tucker rank (r, l, m) tensor, which is given by $p(r + l + m) + rlm$ (Zhang, 2019, Proposition 1).*

Next we investigate how the statistical inaccuracy of state abstraction would affect downstream RL tasks. We consider block-structured MDP whose transition kernel p , reward r and policy π are defined on state and action blocks $\{A_i\}, \{B_j\}$. We use $M = (p, \{r_h\}_{h=1}^H)$ to denote such an MDP instance, and use \mathcal{M} to denote the collection of all such M . One may wonder if the state abstraction error would blow up, particularly if we want to evaluate a multi-step cumulative return. Define the H -step state abstraction error as the worst-case policy evaluation error over horizon H , given by

$$D(\{\hat{A}_i\}, \{\hat{B}_j\}) = \sup_{M \in \mathcal{M}, \pi} \inf_{\hat{p}} \left| \mathbb{E}_p^\pi \left[\sum_{h=1}^H r_h(S_h, A_h) \right] - \mathbb{E}_{\hat{p}}^\pi \left[\sum_{h=1}^H r_h(S_h, A_h) \right] \right|,$$

where the supremum is taken over all block-structured MDP instances and policies on $\{A_i, B_j\}$, and the infimum is to find a best-fit transition model on the estimated clusters $\{\hat{A}_i, \hat{B}_j\}$, of the form

$$\hat{p}(s'|s, a) = \sum_{i,k=1}^{n_s} \sum_{j=1}^{n_a} \frac{\hat{q}(k|i,j)}{\xi(\hat{A}_k)} \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \mathbf{1}_{s' \in \hat{A}_k},$$

where \hat{q} is a set of discrete transition probabilities.

Theorem 8 (Policy evaluation error due to inaccurate state abstraction) *Let Assumptions 1, 2, 3 hold. Then for n sufficiently large, with probability $1 - \delta$, we have*

$$D(\{\hat{A}_i\}, \{\hat{B}_j\}) \leq 4\bar{c}\underline{c}^{-1}H^2M(\{\hat{A}_i\}, \{\hat{B}_j\}),$$

where $\underline{c} = \min_{i,j}(\xi \times \eta)(A_i \times B_j)$, $\bar{c} = \max_{i,j}(\xi \times \eta)(A_i \times B_j)$.

Theorem 8 shows that the H -step state abstraction error grows at most quadratically with H , not exponentially. In other words, inaccuracy in state abstraction does not suffer from the curse of horizon. Thus the learned state and action abstractions are useful for approximate policy evaluation.

5. Numerical Experiment

We test our approach on a particular MDP derived from a controlled stochastic process. Let the state and action spaces be both \mathbb{R}^2 . Suppose the state-action pair at step k is (s_k, a_k) . Then the next state s_{k+1} is set to be $X_{\tau(k+1)}$ for some $\tau > 0$, where X_t is the solution of the SDE:

$$dX_t = -[\nabla V(X_t) + F(a_k)]dt + \sqrt{2}dB_t, k\tau \leq t \leq (k+1)\tau,$$

where $V(\cdot)$ is a wavy potential function, $F(\cdot)$ is a block-wise constant function (Figure 2), B_t is the standard Brownian motion. Let the behavior policy be always choosing a from a standard normal distribution. We use the Gaussian kernels and obtain state/action features by generating N random Fourier features $h = [h_1, h_2, \dots, h_N]$ such that $K(x, y) \approx \sum_{i=1}^N h_i(x)h_i(y)$.

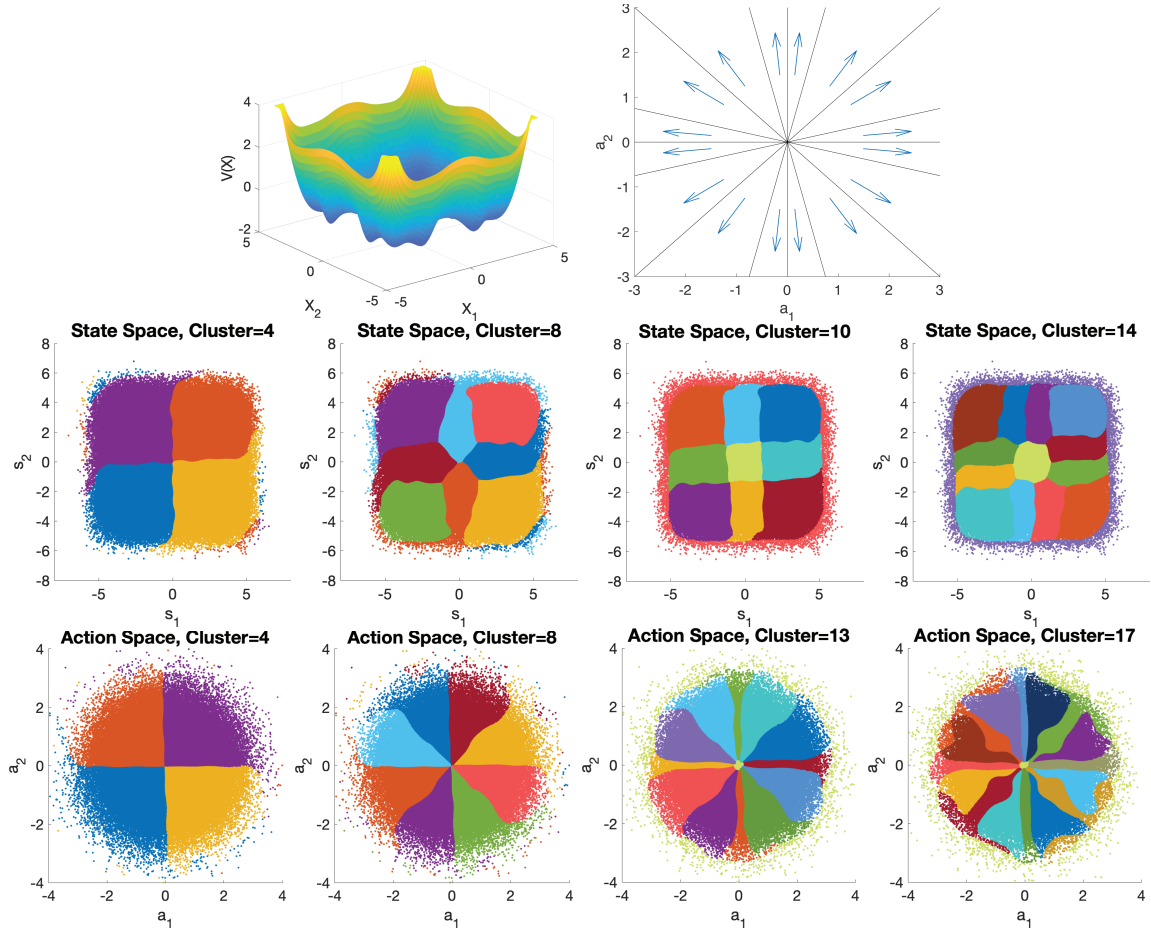


Figure 2: Row 1: Left: Potential function $V(\cdot)$; Right: Block-wise control function $F(\cdot)$. The action space has 16 blocks, and in each block $F(\cdot)$ is a constant drift vector (see the arrows); Row 2: Learned state abstractions with varying clustering sizes; Row 3: Learned action abstractions with varying clustering sizes.

State-action Clustering We first apply Algorithm 2 to estimate state and action clusters. The results are shown at the right side of Figure 2. Comparing them with the ground truth, we can validate that our method indeed reveal the latent state and action blocks.

Low-Rank Estimation of the Transition Tensor We then investigate the efficiency of estimating \mathbf{P} via our tensor method. We compare our method with two baselines: (1) The vanilla method, which directly estimates the transition tensor by $\hat{\mathbf{P}} = \hat{\mathbf{F}} \times_1 \hat{\Sigma}^{-1}$ without any low-rank approximation; (2) The “top r ” method, whose the procedure is: i) calculate the top r (or l, m) principle components of the sample covariance per mode; ii) project features onto the subspace spanned by the top principle components; iii) estimate the transition tensor via the vanilla method (discussed above) in the space of projected features. Fig. 3 visualizes the estimation errors of these methods with different choices of (r, l, m) , where errors are averaged over five independent runs. We observe that, for most of the time, our method consistently outperforms the baselines. Note that the top r method performs slightly better when n is very small, because in this case data is too small to get meaningful estimate of \mathbf{P} . The three approaches have similar performance when the rank constraint is set to be $(60, 30, 60)$ or higher. This is because the rank constraint is already close to the dimensions of the original state-action features, which reduces the impact of the rank-constrained estimator and introduces additional noise due to computational limitations. In practice, small rank constraints are preferred for both statistical and computational reasons.

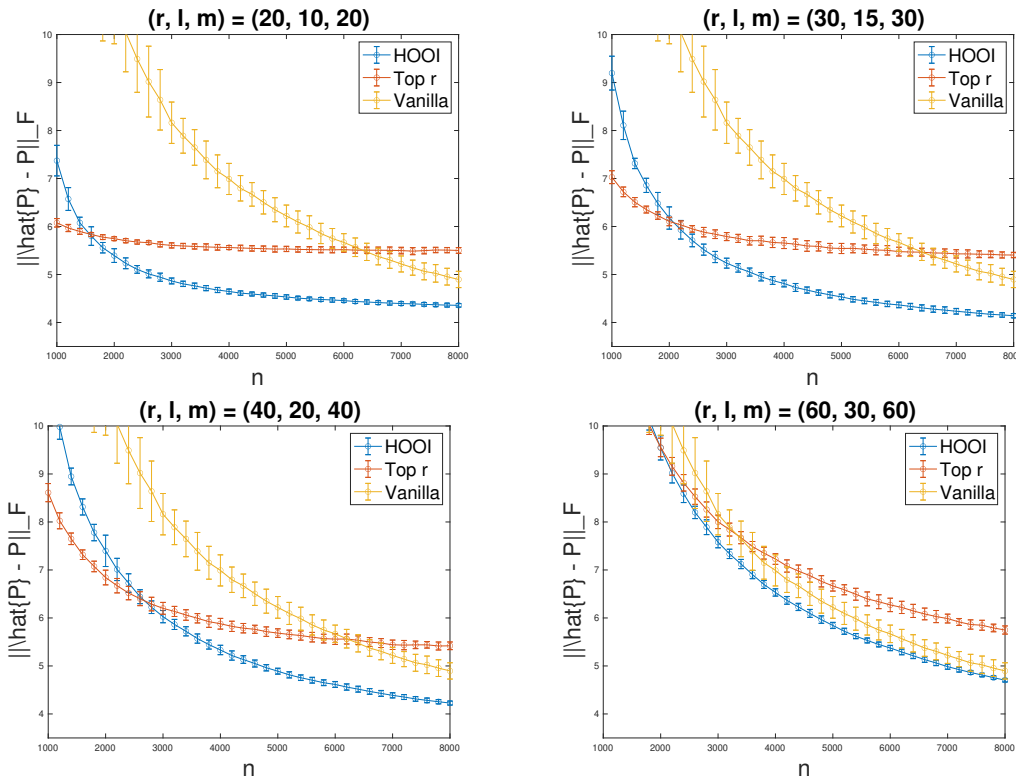


Figure 3: Low-tensor-rank estimation of \mathbf{P} , compared with baseline methods.

References

- Alekh Agarwal, Sham Kakade, Akshay Krishnamurthy, and Wen Sun. Flambe: Structural complexity and representation learning of low rank mdps. *Advances in neural information processing systems*, 33:20095–20107, 2020.
- Animashree Anandkumar, Rong Ge, and Majid Janzamin. Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates. *arXiv preprint arXiv:1402.5180*, 2014.
- Kamyar Azizzadenesheli, Alessandro Lazaric, and Animashree Anandkumar. Reinforcement learning in rich-observation mdps using spectral methods. *arXiv preprint arXiv:1611.03907*, 2016.
- Dimitri P Bertsekas. *Dynamic programming and optimal control*. Athena scientific Belmont, MA, 2007.
- Dimitri P Bertsekas and John N Tsitsiklis. *Neuro-dynamic programming*. Athena Scientific, Belmont, MA, 1996.
- T Tony Cai and Anru Zhang. Rate-optimal perturbation bounds for singular subspaces with applications to high-dimensional statistics. *The Annals of Statistics*, 46(1):60–89, 2018.
- Sayak Ray Chowdhury and Aditya Gopalan. Online learning in kernelized markov decision processes. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 3197–3205. PMLR, 2019.
- Andrzej Cichocki, Danilo Mandic, Lieven De Lathauwer, Guoxu Zhou, Qibin Zhao, Cesar Caiafa, and Huy Anh Phan. Tensor decompositions for signal processing applications: From two-way to multiway component analysis. *IEEE Signal Processing Magazine*, 32(2): 145–163, 2015.
- Ronald R. Coifman, Ioannis G. Kevrekidis, Stéphane Lafon, Mauro Maggioni, and Boaz Nadler. Diffusion maps, reduction coordinates, and low dimensional representation of stochastic systems. *SIAM Journal on Multiscale Modeling and Simulation*, 7(2):852–864, 2008.
- Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. On the best rank-1 and rank-(r_1, r_2, \dots, r_n) approximation of higher-order tensors. *SIAM journal on Matrix Analysis and Applications*, 21(4):1324–1342, 2000a.
- Lieven De Lathauwer, Bart De Moor, and Joos Vandewalle. A multilinear singular value decomposition. *SIAM journal on Matrix Analysis and Applications*, 21(4):1253–1278, 2000b.
- Vin De Silva and Lek-Heng Lim. Tensor rank and the ill-posedness of the best low-rank approximation problem. *SIAM Journal on Matrix Analysis and Applications*, 30(3): 1084–1127, 2008.

- Simon S Du, Sham M Kakade, Ruosong Wang, and Lin F Yang. Is a good representation sufficient for sample efficient reinforcement learning? *arXiv preprint arXiv:1910.03016*, 2019a.
- Simon S Du, Akshay Krishnamurthy, Nan Jiang, Alekh Agarwal, Miroslav Dudík, and John Langford. Provably efficient rl with rich observations via latent state decoding. *arXiv preprint arXiv:1901.09018*, 2019b.
- Weinan E, Tiejun Li, and Eric Vanden-Eijnden. Optimal partition and effective dynamics of complex networks. *Proceedings of the National Academy of Sciences*, 105(23):7907–7912, 2008.
- Rungang Han, Rebecca Willett, and Anru R Zhang. An optimal statistical and computational framework for generalized tensor estimation. *The Annals of Statistics*, 50(1):1–29, 2022.
- David Hong, Tamara G Kolda, and Jed A Duersch. Generalized canonical polyadic tensor decomposition. *SIAM Review*, 62(1):133–163, 2020.
- Majid Janzamin, Rong Ge, Jean Kossaifi, and Anima Anandkumar. Spectral learning on matrices and tensors. *Foundations and Trends® in Machine Learning*, 12(5-6):393–536, 2019.
- Nan Jiang, Akshay Krishnamurthy, Alekh Agarwal, John Langford, and Robert E Schapire. Contextual decision processes with low bellman rank are pac-learnable. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1704–1713. JMLR. org, 2017.
- Chi Jin, Zhuoran Yang, Zhaoran Wang, and Michael I Jordan. Provably efficient reinforcement learning with linear function approximation. *arXiv preprint arXiv:1907.05388*, 2019.
- Jeff Johns and Sridhar Mahadevan. Constructing basis functions from directed graphs for value function approximation. In *Proceedings of the 24th international conference on Machine learning*, pages 385–392. ACM, 2007.
- Ian T Jolliffe. Principal components in regression analysis. In *Principal component analysis*, pages 129–155. Springer, 1986.
- Stefan Klus, Péter Koltai, and Christof Schütte. On the numerical approximation of the perron–frobenius and koopman operator. *Journal of Computational Dynamics*, 3(1):51–79, 2016.
- Stefan Klus, Ingmar Schuster, and Krikamol Muandet. Eigendecompositions of transfer operators in reproducing kernel hilbert spaces. *Journal of Nonlinear Science*, 30(1): 283–315, 2020.
- Tamara G Kolda and Brett W Bader. Tensor decompositions and applications. *SIAM review*, 51(3):455–500, 2009.

- Stéphane Lafon and Ann Lee. Diffusion maps and coarse-graining: A unified framework for dimensionality reduction, graph partitioning, and data set parameterization. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 29(9):1393–1403, 2006.
- Michail G Lagoudakis and Ronald Parr. Least-squares policy iteration. *Journal of machine learning research*, 4(Dec):1107–1149, 2003.
- Tor Lattimore and Csaba Szepesvari. Learning with good feature representations in bandits and in rl with a generative model. *arXiv preprint arXiv:1911.07676*, 2019.
- David Asher Levin, Yuval Peres, and Elizabeth Lee Wilmer. *Markov chains and mixing times*. American Mathematical Soc., 2009.
- Matthias Löffler and Antoine Picard. Spectral thresholding for the estimation of markov chain transition operators. *Electronic Journal of Statistics*, 15(2):6281–6310, 2021.
- Sridhar Mahadevan. Proto-value functions: Developmental reinforcement learning. In *Proceedings of the 22nd international conference on Machine learning*, pages 553–560. ACM, 2005.
- Sridhar Mahadevan et al. Learning representation and control in markov decision processes: New frontiers. *Foundations and Trends® in Machine Learning*, 1(4):403–565, 2009.
- Anuj Mahajan, Mikayel Samvelyan, Lei Mao, Viktor Makoviyuchuk, Animesh Garg, Jean Kossai, Shimon Whiteson, Yuke Zhu, and Animashree Anandkumar. Tesseract: Tensorised actors for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 7301–7312. PMLR, 2021.
- Dipendra Misra, Mikael Henaff, Akshay Krishnamurthy, and John Langford. Kinematic state abstraction and provably efficient rich-observation reinforcement learning. *arXiv preprint arXiv:1911.05815*, 2019.
- Aditya Modi, Jinglin Chen, Akshay Krishnamurthy, Nan Jiang, and Alekh Agarwal. Model-free representation learning and exploration in low-rank mdps. *arXiv preprint arXiv:2102.07035*, 2021.
- Andrew W Moore. Variable resolution dynamic programming: Efficiently learning action maps in multivariate real-valued state-spaces. In *Machine Learning Proceedings 1991*, pages 333–337. Elsevier, 1991.
- Chengzhuo Ni, Yuda Song, Xuezhou Zhang, Chi Jin, and Mengdi Wang. Representation learning for general-sum low-rank markov games. *International Conference on Learning Representations*, 2023.
- Dirk Ormoneit and Peter Glynn. Kernel-based reinforcement learning in average-cost problems. *IEEE Transactions on Automatic Control*, 47(10):1624–1636, 2002.
- Yannis Panagakis, Jean Kossai, Grigorios G Chrysos, James Oldfield, Mihalis A Nicolaou, Anima Anandkumar, and Stefanos Zafeiriou. Tensor methods in computer vision and deep learning. *Proceedings of the IEEE*, 109(5):863–890, 2021.

- Ronald Parr, Christopher Painter-Wakefield, Lihong Li, and Michael Littman. Analyzing feature generation for value-function approximation. In *Proceedings of the 24th international conference on Machine learning*, pages 737–744. ACM, 2007.
- Marek Petrik. An analysis of laplacian methods for value function approximation in mdps. In *IJCAI*, pages 2574–2579, 2007.
- Ali Rahimi and Benjamin Recht. Random features for large-scale kernel machines. In *Advances in neural information processing systems*, pages 1177–1184, 2008.
- Zhiyuan Ren and Bruce H Krogh. State aggregation in markov decision processes. In *Decision and Control, 2002, Proceedings of the 41st IEEE Conference on*, volume 4, pages 3819–3824. IEEE, 2002.
- Emile Richard and Andrea Montanari. A statistical model for tensor pca. In *Advances in Neural Information Processing Systems*, pages 2897–2905, 2014.
- Christof Schütte, Frank Noe, Jianfeng Lu, Macro Sarich, and Eric Vanden-Eijnden. Markov state models based on milestoning. *The Journal of Chemical Physics*, 134(20):204105, 2011.
- Nicholas D Sidiropoulos, Lieven De Lathauwer, Xiao Fu, Kejun Huang, Evangelos E Papalexakis, and Christos Faloutsos. Tensor decomposition for signal processing and machine learning. *IEEE Transactions on Signal Processing*, 65(13):3551–3582, 2017.
- Satinder P Singh, Tommi Jaakkola, and Michael I Jordan. Reinforcement learning with soft state aggregation. In *Advances in neural information processing systems*, pages 361–368, 1995.
- Zhao Song, David Woodruff, and Huan Zhang. Sublinear time orthogonal tensor decomposition. In *Advances in Neural Information Processing Systems*, pages 793–801, 2016.
- Will Wei Sun, Junwei Lu, Han Liu, and Guang Cheng. Provable sparse tensor decomposition. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 79(3):899–916, 2017.
- Yifan Sun, Yaqi Duan, Hao Gong, and Mengdi Wang. Learning low-dimensional state embeddings and metastable clusters from time series data. In *Advances in Neural Information Processing Systems*, pages 4563–4572, 2019.
- Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- Joel A Tropp. Freedman’s inequality for matrix martingales. *Electron. Commun. Probab*, 16:262–270, 2011.
- John N Tsitsiklis and Benjamin Van Roy. Feature-based methods for large scale dynamic programming. *Machine Learning*, 22(1-3):59–94, 1996.
- Masatoshi Uehara, Xuezhou Zhang, and Wen Sun. Representation learning for online and offline rl in low-rank mdps. *arXiv preprint arXiv:2110.04652*, 2021.

- Pascal Van Der Vaart, Anuj Mahajan, and Shimon Whiteson. Model based multi-agent reinforcement learning with tensor decompositions. *arXiv preprint arXiv:2110.14524*, 2021.
- Nick Vannieuwenhoven, Raf Vandebril, and Karl Meerbergen. A new truncation strategy for the higher-order singular value decomposition. *SIAM Journal on Scientific Computing*, 34(2):A1027–A1052, 2012.
- Roman Vershynin. *High-Dimensional Probability*. Cambridge University Press (to appear), 2017.
- Per-Åke Wedin. Perturbation bounds in connection with singular value decomposition. *BIT Numerical Mathematics*, 12(1):99–111, 1972.
- Lin Yang and Mengdi Wang. Sample-optimal parametric q-learning using linearly additive features. In *International Conference on Machine Learning*, pages 6995–7004, 2019.
- Andrea Zanette, Alessandro Lazaric, Mykel J Kochenderfer, and Emma Brunskill. Limiting extrapolation in linear approximate value iteration. In *Advances in Neural Information Processing Systems*, pages 5616–5625, 2019.
- Anru Zhang. Cross: Efficient low-rank tensor completion. *The Annals of Statistics*, 47(2): 936–964, 2019.
- Anru Zhang and Rungang Han. Optimal sparse singular value decomposition for high-dimensional high-order data. *Journal of the American Statistical Association*, pages 1708–1725, 2019.
- Anru Zhang and Mengdi Wang. Spectral state compression of markov processes. *IEEE transactions on information theory*, 66(5):3202–3231, 2020.
- Anru Zhang and Dong Xia. Tensor svd: Statistical and computational limits. *IEEE Transactions on Information Theory*, 64(11):7311–7338, 2018.
- Xuezhou Zhang, Yuda Song, Masatoshi Uehara, Mengdi Wang, Alekh Agarwal, and Wen Sun. Efficient reinforcement learning in block mdps: A model-free representation learning approach. In *International Conference on Machine Learning*, pages 26517–26547. PMLR, 2022.
- Ziwei Zhu, Xudong Li, Mengdi Wang, and Anru Zhang. Learning markov models via low-rank optimization. *Operations Research*, 70(4):2384–2398, 2022.

Appendix

A. The HOOI Algorithm

Algorithm 3 HOOI for MDP Tensor Decomposition

- 1: **Input:** tensor mean embedding $\bar{\mathbf{F}}$, $(r, l, m), t_{\max}$
- 2: Initialization:

$$\bar{\mathbf{U}}_1^{(0)} = \text{SVD}_r(\mathcal{M}_1(\bar{\mathbf{F}})), \bar{\mathbf{U}}_2^{(0)} = \text{SVD}_l(\mathcal{M}_2(\bar{\mathbf{F}} \times_1 \bar{\mathbf{U}}_1^{(0)\top})),$$

$$\bar{\mathbf{U}}_3^{(0)} = \text{SVD}_m(\mathcal{M}_3(\bar{\mathbf{F}} \times_1 \bar{\mathbf{U}}_1^{(0)\top} \times_2 \bar{\mathbf{U}}_2^{(0)\top})),$$

where $\text{SVD}_r(\cdot)$ is the operation that returns the leading r singular vector of matrix \cdot .

- 3: **for** $t = 1, \dots, t_{\max}$ **do**
 - 4: $\bar{\mathbf{U}}_1^{(t)} = \text{SVD}_r(\mathcal{M}_1(\bar{\mathbf{F}} \times_2 \bar{\mathbf{U}}_2^{(t-1)\top} \times_3 \bar{\mathbf{U}}_3^{(t-1)\top}))$,
 - $\bar{\mathbf{U}}_2^{(t)} = \text{SVD}_l(\mathcal{M}_2(\bar{\mathbf{F}} \times_1 \bar{\mathbf{U}}_1^{(t)\top} \times_3 \bar{\mathbf{U}}_3^{(t-1)\top}))$,
 - $\bar{\mathbf{U}}_3^{(t)} = \text{SVD}_m(\mathcal{M}_3(\bar{\mathbf{F}} \times_1 \bar{\mathbf{U}}_1^{(t)\top} \times_2 \bar{\mathbf{U}}_2^{(t)\top}))$.
 - 5: **end for**
 - 6: **Output:** $\hat{\mathbf{F}} = \bar{\mathbf{F}} \times_1 (\bar{\mathbf{U}}_1^{(t_{\max})\top})^\top \bar{\mathbf{U}}_1^{(t_{\max})\top} \times_2 (\bar{\mathbf{U}}_2^{(t_{\max})\top})^\top \bar{\mathbf{U}}_2^{(t_{\max})\top} \times_3 (\bar{\mathbf{U}}_3^{(t_{\max})\top})^\top \bar{\mathbf{U}}_3^{(t_{\max})\top}$
-

B. Examples of Low Rank MDPs

We give two basic examples.

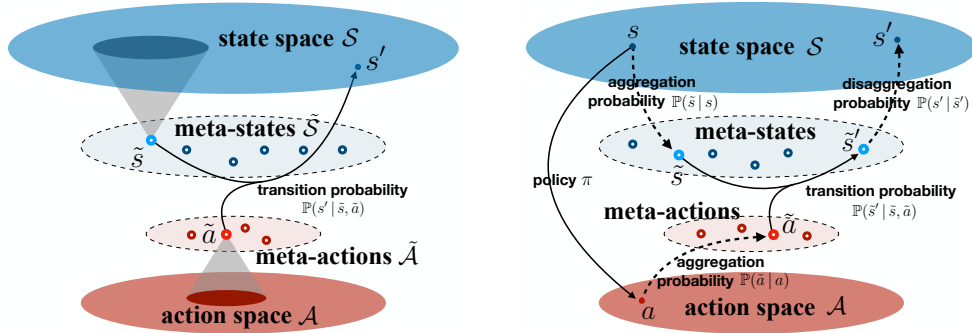


Figure 4: Left: Block MDP (aka hard aggregation); Right: Latent-state-action MDP (aka soft aggregation).

Example 1 (Block MDP (Hard Aggregation)) Let $\tilde{\mathcal{S}}$ and $\tilde{\mathcal{A}}$ be finite sets. Suppose there exists state and action abstractions $f : \mathcal{S} \mapsto \tilde{\mathcal{S}}$ and $g : \mathcal{A} \mapsto \tilde{\mathcal{A}}$ such that

$$p(\cdot|s, a) = p(\cdot|s', a') \text{ if } f(s) = f(s'), g(a) = g(a').$$

Then p has Tucker rank at most $(|\tilde{\mathcal{S}}|, |\tilde{\mathcal{A}}|, |\mathcal{S}|)$.

Example 2 (Latent-State-Action MDP (Soft Aggregation)) Given an MDP $\mathcal{M} = (\mathcal{S}, \mathcal{A}, p, r)$, we say \mathcal{M} has an (r, l, m) -latent variable model if there exist a latent state-action-state stochastic process $\{\tilde{s}_t, \tilde{a}_t, \tilde{s}'_t\} \subseteq \tilde{\mathcal{S}} \times \tilde{\mathcal{A}} \times \tilde{\mathcal{S}}'$, with $|\tilde{\mathcal{S}}| = r, |\tilde{\mathcal{A}}| = l, |\tilde{\mathcal{S}}'| = m$, such that

$$\begin{aligned} \mathbb{P}(\tilde{s}_t, \tilde{a}_t | s_1, a_1, \dots, s_t, a_t) &= \mathbb{P}(\tilde{s}_t | s_t) \mathbb{P}(\tilde{a}_t | a_t), \mathbb{P}(\tilde{s}'_t | s_1, a_1, \dots, s_t, a_t, \tilde{s}_t, \tilde{a}_t) = \mathbb{P}(\tilde{s}'_t | \tilde{s}_t, \tilde{a}_t), \\ \mathbb{P}(s_{t+1} | s_1, a_1, \dots, s_t, a_t, \tilde{s}_t, \tilde{a}_t, \tilde{s}'_t) &= \mathbb{P}(s_{t+1} | \tilde{s}'_t). \end{aligned}$$

In this case, one can verify that p has Tucker rank (r, l, m) .

We give an illustrative example to show the advantage of utilizing tensor MDP formulation as opposed to the matrix ones.

Example 3 Consider $\mathcal{A} = \{1, 2\}$, $\mathcal{S} = \{1, 2, 3, 4\}$. Construct the MDP transition tensor \mathbf{P} as

$$\begin{aligned} \mathbf{P}_{\cdot 1} &= \begin{bmatrix} 1/6 & 1/6 & 1/3 & 1/3 \\ 1/6 & 1/6 & 1/3 & 1/3 \\ 1/3 & 1/3 & 1/6 & 1/6 \\ 1/3 & 1/3 & 1/6 & 1/6 \end{bmatrix}, \\ \mathbf{P}_{\cdot 2} &= \begin{bmatrix} 1/3 & 1/3 & 1/6 & 1/6 \\ 1/3 & 1/3 & 1/6 & 1/6 \\ 1/6 & 1/6 & 1/3 & 1/3 \\ 1/6 & 1/6 & 1/3 & 1/3 \end{bmatrix}. \end{aligned}$$

Then, $\mathbf{P} = \mathbf{C} \times_1 \mathbf{U} \times_3 \mathbf{U}$ for

$$\mathbf{U} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \end{bmatrix}, \mathbf{C}_{\cdot 1} = \begin{bmatrix} 1/6 & 1/3 \\ 1/3 & 1/6 \end{bmatrix}, \mathbf{C}_{\cdot 2} = \begin{bmatrix} 1/3 & 1/6 \\ 1/6 & 1/3 \end{bmatrix}$$

and the state-space is aggregatable into two meta-states: $\{1, 2\}$ and $\{3, 4\}$. Consider a random policy: $\pi(a|s) = 1/2$ for $a = 1, 2$. Without taking into account the tensor structure induced by the policy, one can check that the state transitions $\{s_0, s_1, \dots\}$ form a Markov process with the following transition matrix

$$\tilde{\mathbf{P}} = \frac{1}{2} \mathbf{P}_{\cdot 1} + \frac{1}{2} \mathbf{P}_{\cdot 2} = \begin{bmatrix} 1/4 & 1/4 & 1/4 & 1/4 \\ 1/4 & 1/4 & 1/4 & 1/4 \\ 1/4 & 1/4 & 1/4 & 1/4 \\ 1/4 & 1/4 & 1/4 & 1/4 \end{bmatrix}.$$

Clearly, the meta-state partition is “averaged out” using any matrix methods and there is no hope to extract the meta-state information merely from the state transitions $\{s_0, s_1, \dots\}$. On the other hand, the tensor formulation, which preserves the original state-action-state, allows a reliable state aggregation efficiently.

C. Derivation of Optimization Problem (4)

The original optimization objective is

$$\begin{aligned} & \min_{A_i, B_j} \min_{\{q_{ij}\}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|p(\cdot|s, a) - q_{ij}(\cdot)\|_{\mathcal{H}_S}^2 dsda \\ &= \min_{A_i, B_j} \min_{\{q_{ij}\}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle q_{ij}(\cdot), \phi(\cdot) \rangle\|^2 dsda. \end{aligned}$$

Simple calculations show that given fixed A_i, B_j , the best choice of $\langle q_{ij}(\cdot), \phi(\cdot) \rangle$ is

$$\begin{aligned} \langle q_{ij}(\cdot), \phi(\cdot) \rangle &= \frac{1}{\xi(A_i)\eta(B_j)} \int_{A_i \times B_j} \xi(s)\eta(a) \langle p(\cdot|s, a), \phi(\cdot) \rangle dsda \\ &= \frac{1}{\xi(A_i)\eta(B_j)} \int_{A_i \times B_j} \xi(s)\eta(a) (\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top) dsda \\ &= \frac{1}{\xi(A_i)\eta(B_j)} \int_{A_i \times B_j} \xi(s)\eta(a) (\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{U}_3) dsda \\ &= \mathbf{U}_3 \mathbf{z}_{ij}, \end{aligned}$$

where $\mathbf{z}_{ij} = \frac{1}{\xi(A_i)\eta(B_j)} \int_{A_i \times B_j} \xi(s)\eta(a) (\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top) dsda$. Note that

$$\begin{aligned} & \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \mathbf{U}_3 \mathbf{z}_{ij}\|^2 dsda \\ &= \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{U}_3 - \mathbf{U}_3 \mathbf{z}_{ij}\|^2 dsda \\ &= \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \mathbf{z}_{ij}\|^2 dsda. \end{aligned}$$

Therefore, our problem can be further formalized as

$$\min_{A_i, B_j} \min_{\mathbf{z}_{ij}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \mathbf{z}_{ij}\|^2 dsda.$$

When we only have empirical data, the above problem can be approximated by

$$\min_{A_i, B_j} \min_{\mathbf{z}_{ij}} \sum_{i,j} \int_{A_i \times B_j} \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{z}_{ij}\|^2 dsda,$$

which is exactly (4).

D. Experiment Details

In the experiment, we use the Gaussian kernels $K_1(x, y) = K_A(x, y) = \frac{1}{2\pi\sigma^2} \exp\{-\frac{\|x-y\|^2}{2\sigma^2}\}$. And the features are obtained by generating N_s (or N_a) random Fourier features $h =$

$[h_1, h_2, \dots, h_{N_s}]$ such that $K(x, y) \approx \sum_{i=1}^{N_s(N_a)} h_i(x)h_i(y)$. And the action features are then orthogonalized with respect to $L^2(\eta)$. In the experiment, we choose $\tau = 0.1, \sigma = 0.5, N_s = 100, N_a = 50$.

For the clustering problem, we choose sample size $n = 10^6$ and $(r, l, m) = (3, 3, 3)$, and the state features are further orthogonalized with respect to $L^2(\xi)$. For the estimation problem, the ground-truth is approximately obtained from the vanilla method with sample size $n = 10^6$. The following figure shows the clustering result of the top- r method, which does the clustering on the subspace spanned by the top- r (or l, m) eigenvectors of the covariance matrix. From the figure we can see that the top- r method does not capture the correct clustering information of the transition kernel compared with our method.

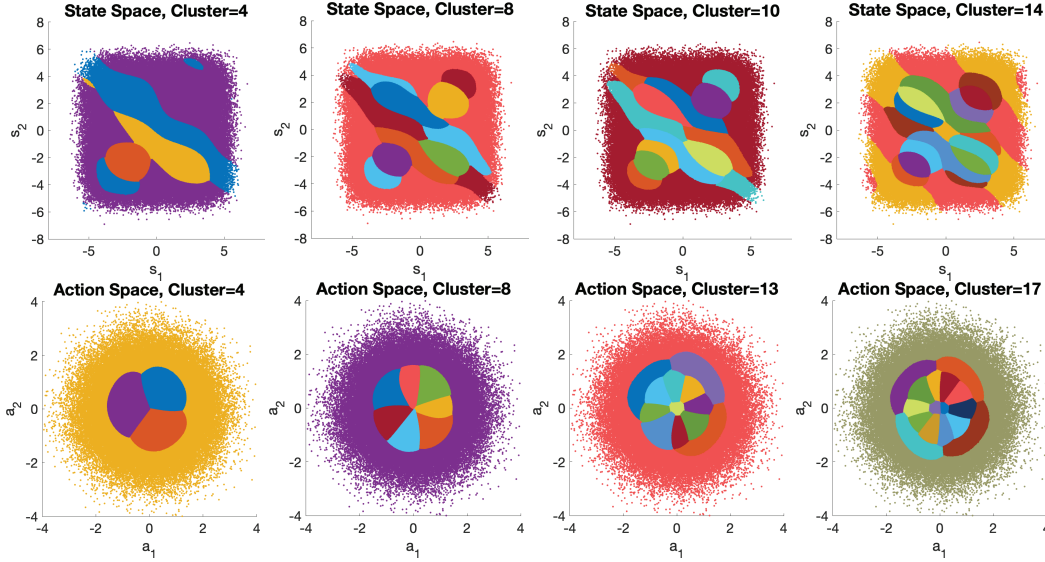


Figure 5: Row 1: Learned state abstractions with varying clustering sizes by top- r method; Row 2: Learned action abstractions with varying clustering sizes by top- r method.

E. Technical Lemmas

Lemma 9 *Suppose*

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\bar{\mu} \|\Sigma^{-1}\|_{\sigma}^2 + \frac{K_{max}^2}{\bar{\mu}} \right) \log \frac{dst_{mix}}{\delta}.$$

Then with probability $1 - \delta$, we have

$$\|\hat{\Sigma}^{-1} - \Sigma^{-1}\|_{\sigma} \leq 32 \|\Sigma^{-1}\|_{\sigma}^2 \sqrt{\frac{\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Lemma 10 *For any tensor $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times \dots \times p_N}$, such that $\text{Tucker-Rank}(\mathbf{X}) \leq (r_1, r_2, \dots, r_N)$, we can always find column-wise orthonormal matrices $\mathbf{U}_1 \in \mathbb{R}^{p_1 \times r_1}, \dots, \mathbf{U}_N \in \mathbb{R}^{p_N \times r_N}$ and*

a core tensor $\mathbf{C} \in \mathbb{R}^{r_1 \times \dots \times r_N}$, such that

$$\mathbf{X} = \mathbf{C} \times_1 \mathbf{U}_1 \times_2 \dots \times_N \mathbf{U}_N.$$

Lemma 11 *Suppose the worst-case mixing time of the MDP is t_{mix} , then for any $\varepsilon > 0$ and policy π , suppose ν^π is the invariant distribution of π , then for any initial distribution μ , we have*

$$\left\| \int p^{t,\pi}(\cdot|s_0)\mu(s_0)ds_0 - \nu^\pi(\cdot) \right\|_{TV} \leq \varepsilon, \forall t \geq 2t_{mix} \log \frac{1}{\varepsilon}.$$

Lemma 12 *For any given tensor $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ such that $\text{Tucker-Rank}(\mathbf{X}) \leq (r_1, r_2, r_3)$, we have*

$$\|\mathbf{X}\|_F \leq \sqrt{\frac{r_1 r_2 r_3}{\max\{r_1, r_2, r_3\}}} \|\mathbf{X}\|_\sigma.$$

Lemma 13 *Given $p \in \mathbb{N}, \varepsilon \in \mathbb{R}$, there always exists an ε -net of the sphere S^{p-1} whose size is no more than $(1 + 2/\varepsilon)^p$.*

Lemma 14 *Given a tensor $\mathbf{X} \in \mathbb{R}^{p_1 \times p_2 \times p_3}$ and three ε -nets of the unit sphere $\mathcal{N}_1 \subset S^{p_1-1}, \mathcal{N}_2 \subset S^{p_2-1}, \mathcal{N}_3 \subset S^{p_3-1}$, we have*

$$\|\mathbf{X}\|_\sigma \leq \frac{\max_{x \in \mathcal{N}_1, y \in \mathcal{N}_2, z \in \mathcal{N}_3} |\langle \mathbf{X}, x \circ y \circ z \rangle|}{1 - 3\varepsilon - 3\varepsilon^2 - \varepsilon^3}.$$

Lemma 15 *Suppose*

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\bar{\mu} \|\Sigma^{-1}\|_\sigma^2 + \frac{K_{max}^2}{\bar{\mu}} \right) \log \frac{d_S t_{mix}}{\delta}.$$

Then with probability $1 - \delta$, we have

$$\|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_\sigma \leq 32 \|\Sigma^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Lemma 16 (Concentration in tensor spectral norm) *Let Assumptions 1-2 hold. Suppose*

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \frac{\kappa K_{max}^3}{\bar{\lambda}} \left(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A) \right),$$

then with probability $1 - \delta$, we have

$$\|\hat{\mathbf{F}} - \mathbf{F}\|_\sigma \leq 64 \sqrt{\frac{\kappa \bar{\lambda} (\log \frac{t_{mix}}{\delta} + 8(d_S + d_A)) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Lemma 17 *Suppose*

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\|\Sigma^{-1}\|_{\sigma}^2 \bar{\mu} + \frac{K_{max}^2}{\bar{\mu}} + \frac{\kappa K_{max}^3}{\lambda} \right) \left(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A) \right).$$

Then with probability $1 - \delta$, we have

$$\begin{aligned} & \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{1/2}) - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{1/2})\|_{\sigma} \\ & \leq 256 \|\Sigma^{-1}\|_{\sigma}^{\frac{1}{2}} \sqrt{\frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\Sigma^{-1}\|_{\sigma}^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}. \end{aligned}$$

Lemma 18 *Suppose $\mathbf{P}, \hat{\mathbf{P}}$ are two order-3 tensors of the same dimension. Suppose $\mathbf{P} = \mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3$ and $\hat{\mathbf{P}} = \hat{\mathbf{C}} \times_1 \hat{\mathbf{U}}_1 \times_2 \hat{\mathbf{U}}_2 \times_3 \hat{\mathbf{U}}_3$, where $\mathbf{C}, \hat{\mathbf{C}} \in \mathbb{R}^{r \times l \times m}$, $\mathbf{U}_3^{\top} \mathbf{U}_3 = \hat{\mathbf{U}}_3^{\top} \hat{\mathbf{U}}_3 = \mathbf{I}$. We have*

$$\|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|_{\sigma} \leq \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_{\sigma}}{\sigma},$$

where

$$\sigma = \sup_{w \in \mathbb{R}^p, \|w\|=1} \sigma_m(\mathbf{P} \times_1 w).$$

F. Proofs

Proof of Lemma 1 Recall that under Assumption 1, there exist $c_{ijk} \in \mathbb{R}$, $u_i, w_k \in \mathcal{H}_S$, $v_j \in \mathcal{H}_A$, $i \in [r]$, $j \in [l]$, $k \in [m]$ such that

$$(\mathbb{P}f)(s, a) = \sum_{i=1}^r \sum_{j=1}^l \sum_{k=1}^m c_{ijk} u_i(s) v_j(a) \langle f, w_k \rangle_{\mathcal{H}_S}, \quad \forall f \in \mathcal{H}_S.$$

Let $\mathbf{C} \in \mathbb{R}^{r \times l \times m}$ be defined as $\mathbf{C}_{ijk} = c_{ijk}$. Then we can rewrite

$$\sum_{i=1}^r \sum_{j=1}^l \sum_{k=1}^m c_{ijk} u_i(s) v_j(a) \langle f, w_k \rangle = \mathbf{C} \times_1 \mathbf{u}(s)^{\top} \times_2 \mathbf{v}(a)^{\top} \times_3 \langle f, \mathbf{w} \rangle^{\top}.$$

Now, because we have $u_i, w_k \in \mathcal{H}_S$, $v_j \in \mathcal{H}_A$, we can find three matrices $\mathbf{U}_1 \in \mathbb{R}^{d_S \times r}$, $\mathbf{U}_2 \in \mathbb{R}^{d_A \times l}$, $\mathbf{U}_3 \in \mathbb{R}^{d_S \times m}$, such that

$$\mathbf{u} = \mathbf{U}_1^{\top} \phi, \mathbf{v} = \mathbf{U}_2^{\top} \psi, \mathbf{w} = \mathbf{U}_3^{\top} \phi.$$

Then we have by the association law (Kolda and Bader, 2009, Section 2.5) that

$$\begin{aligned} \mathbf{C} \times_1 \mathbf{u}(s)^{\top} \times_2 \mathbf{v}(a)^{\top} \times_3 \langle f, \mathbf{w} \rangle^{\top} &= \mathbf{C} \times_1 (\mathbf{U}_1^{\top} \phi(s))^{\top} \times_2 (\mathbf{U}_2^{\top} \psi(a))^{\top} \times_3 (\mathbf{U}_3^{\top} \langle f, \phi \rangle)^{\top} \\ &= (\mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3) \times_1 \phi(s)^{\top} \times_2 \psi(a)^{\top} \times_3 \langle f, \phi \rangle^{\top}. \end{aligned}$$

In particular, we take $f = \phi_i, i = 1, 2, \dots, d_S$, and define \mathbf{V} by

$$\mathbf{V}_{ij} = \langle \phi_i, \phi_j \rangle,$$

then we have

$$\begin{aligned} \mathbb{E}[\phi(s')|s, a] &= (\mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 \mathbf{U}_3) \times_1 \phi(s) \times_2 \psi(a)^\top \times_3 \mathbf{V}^\top \\ &= (\mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 (\mathbf{U}_3^\top \mathbf{V})^\top) \times_1 \phi(s)^\top \times_2 \psi(a)^\top. \end{aligned}$$

Now, we define $\mathbf{P} \in \mathbb{R}^{d_S \times d_A \times d_S}$ by

$$\mathbf{P} = \mathbf{C} \times_1 \mathbf{U}_1 \times_2 \mathbf{U}_2 \times_3 (\mathbf{U}_3^\top \mathbf{V})^\top.$$

Then the Tucker-rank of \mathbf{P} is no larger than the size of \mathbf{C} , i.e.,

$$\text{Tucker-Rank}(\mathbf{P}) \leq (r, l, m),$$

and we have

$$\mathbb{E}[\phi(s')|s, a] = \mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top,$$

which finishes the proof.

Proof of Lemma 2 To show this, simply notice

$$\begin{aligned} \mathbf{F} &= \int \phi(s) \circ \psi(a) \circ \phi(s') p(s'|s, a) \xi(s) \eta(a) ds da ds' \\ &= \int \phi(s) \circ \psi(a) \circ \left(\int \phi(s') p(s'|s, a) ds' \right) \xi(s) \eta(a) ds da \\ &= \int \phi(s) \circ \psi(a) \circ \mathbb{E}[\phi(s')|s, a] \xi(s) \eta(a) ds da. \end{aligned}$$

Now we use the notation in Lemma 1 to write $\mathbb{E}[\phi(s')|s, a] = \mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top$ for some $\mathbf{P} \in \mathbb{R}^{d_S \times d_A \times d_S}$, and get

$$\begin{aligned} \mathbf{F} &= \int \phi(s) \circ \psi(a) \circ (\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top) \xi(s) \eta(a) ds da \\ &= \int \mathbf{P} \times_1 \left(\xi(s) \phi(s) \phi(s)^\top \right)^\top \times_2 \left(\eta(a) \psi(a) \psi(a)^\top \right)^\top ds da \quad (5) \\ &= \mathbf{P} \times_1 \left(\int \xi(s) \phi(s) \phi(s)^\top ds \right)^\top \times_2 \left(\int \eta(a) \psi(a) \psi(a)^\top da \right)^\top. \end{aligned}$$

The result of Lemma 1 shows that $\text{Tucker-Rank}(\mathbf{P}) \leq (r, l, m)$, which implies that

$$\text{Tucker-Rank}(\mathbf{F}) \leq \text{Tucker-Rank}(\mathbf{P}) \leq (r, l, m).$$

Proof of Lemma 3 Notice that when ψ is orthogonal with respect to $L^2(\eta)$, (5) reduces to

$$\begin{aligned} \mathbf{F} &= \int \xi(s)\eta(a)p(s'|s, a)\phi(s) \circ \psi(a) \circ \phi(s')dsdad s' \\ &= \mathbf{P} \times_1 \left(\int \xi(s)\phi(s)\phi(s)^\top ds \right)^\top \times_2 \left(\int \eta(a)\psi(a)\psi(a)^\top da \right)^\top \\ &= \mathbf{P} \times_1 \boldsymbol{\Sigma}, \end{aligned}$$

which implies

$$\mathbf{P} = \mathbf{F} \times_1 \boldsymbol{\Sigma}^{-1}.$$

Proof of Lemma 9 According to the result of Lemma 15, we know that when

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\bar{\mu} \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2 + \frac{K_{max}^2}{\bar{\mu}} \right) \log \frac{dst_{mix}}{\delta}.$$

Then with probability $1 - \delta$, we have

$$\left\| (\hat{\boldsymbol{\Sigma}}^{-1} - \boldsymbol{\Sigma}^{-1}) \boldsymbol{\Sigma}^{1/2} \right\|_\sigma \leq 32 \|\boldsymbol{\Sigma}^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Therefore, we can directly get

$$\begin{aligned} \|\hat{\boldsymbol{\Sigma}}^{-1} - \boldsymbol{\Sigma}^{-1}\|_\sigma &= \|(\hat{\boldsymbol{\Sigma}}^{-1} - \boldsymbol{\Sigma}^{-1}) \boldsymbol{\Sigma}^{\frac{1}{2}} \boldsymbol{\Sigma}^{-\frac{1}{2}}\|_\sigma \\ &\leq \|(\hat{\boldsymbol{\Sigma}}^{-1} - \boldsymbol{\Sigma}^{-1}) \boldsymbol{\Sigma}^{\frac{1}{2}}\|_\sigma \|\boldsymbol{\Sigma}^{-\frac{1}{2}}\|_\sigma \\ &\leq 32 \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2 \sqrt{\frac{\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}, \end{aligned}$$

which finishes the proof.

Proof of Lemma 10 Suppose the SVD of $\mathcal{M}_1(\mathbf{X})$ is

$$\mathcal{M}_1(\mathbf{X}) = \mathbf{U}_1 \boldsymbol{\Sigma} \mathbf{V}^\top,$$

where $\mathbf{U}_1 \in \mathbb{R}^{p_1 \times r_1}$ is a column-wise orthonormal matrix. Therefore,

$$\mathcal{M}_1(\mathbf{X}) = \mathbf{U}_1 (\mathbf{U}_1^\top \mathcal{M}_1(\mathbf{X})),$$

which is equivalent to

$$\mathbf{X} = (\mathbf{X} \times_1 \mathbf{U}_1^\top) \times_1 \mathbf{U}_1.$$

Let $\mathbf{X}_1 = \mathbf{X} \times_1 \mathbf{U}_1^\top$, with a similar procedure we can find some column-wise orthogonal matrix \mathbf{U}_2 , such that

$$\mathbf{X}_1 = (\mathbf{X}_1 \times_2 \mathbf{U}_2^\top) \times_2 \mathbf{U}_2.$$

Repeating this process for N times, we can find a series of column orthogonal matrix $\mathbf{U}_1, \mathbf{U}_2, \dots, \mathbf{U}_N$ and a core tensor $\mathbf{C} = \mathbf{X}_N$, such that

$$\mathbf{X} = \mathbf{C} \times_1 \mathbf{U}_1 \cdots \times_N \mathbf{U}_N,$$

which has finished the proof.

Proof of Lemma 11 Let $\alpha = t_{mix}$, for any initial distribution μ , we use the notation

$$p^{t,\pi}(\cdot|\mu) := \int p^{t,\pi}(\cdot|s)\mu(s)ds$$

to denote the state distribution after t steps starting from initial state distribution μ . One direct fact is

$$\begin{aligned} \|p^{\alpha,\pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} &= \frac{1}{2} \int \left| \int p^{\alpha,\pi}(s|s_0)\mu(s_0)ds_0 - \nu^\pi(s) \right| ds \\ &= \frac{1}{2} \int \left| \int (p^{\alpha,\pi}(s|s_0) - \nu^\pi(s))\mu(s_0)ds_0 \right| ds \\ &\leq \int \mu(s_0) \left(\frac{1}{2} \int |p^{\alpha,\pi}(s|s_0) - \nu^\pi(s)| ds \right) ds_0 \\ &\leq \frac{1}{4} \int \mu(s_0) ds_0 = \frac{1}{4}. \end{aligned}$$

Now, for any initial distribution μ and any $n \geq 2\alpha$, we have

$$\begin{aligned} \|p^{n,\pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} &= \frac{1}{2} \int \left| \int p^{n,\pi}(s|s_0)\mu(s_0)ds_0 - \nu^\pi(s) \right| ds \\ &= \frac{1}{2} \int \left| \int p^{n-\alpha,\pi}(s|s_1) \left[\int p^{\alpha,\pi}(s_1|s_0)\mu(s_0)ds_0 - \nu^\pi(s_1) \right] ds_1 \right| ds \\ &= \frac{1}{2} \int \left| \int p^{n-\alpha,\pi}(s|s_1) [(p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_+ - (p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_-] ds_1 \right| ds \\ &= \frac{1}{2} \int \left| \int p^{n-\alpha,\pi}(s|s_1) \left[(p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_+ - \left(\int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \right) \nu^\pi(s_1) \right. \right. \\ &\quad \left. \left. + \frac{1}{2} \left(\int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s} \right) \nu^\pi(s_1) - (p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_- \right] ds_1 \right| ds, \end{aligned}$$

where we use the relation

$$\int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} = \int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s},$$

because we have

$$0 = \int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s})) d\tilde{s} = \int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} - \int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s}.$$

Therefore, we get

$$\begin{aligned} &\|p^{n,\pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} \\ &\leq \frac{1}{2} \int \left| \int p^{n-\alpha,\pi}(s|s_1) \left[(p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_+ - \left(\int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \right) \nu^\pi(s_1) \right] ds_1 \right| ds \\ &\quad + \frac{1}{2} \int \left| \int p^{n-\alpha,\pi}(s|s_1) \left[\left(\int (p^{\alpha,\pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s} \right) \nu^\pi(s_1) - (p^{\alpha,\pi}(s_1|\mu) - \nu^\pi(s_1))_- \right] ds_1 \right| ds. \end{aligned}$$

For the first term, note that

$$\begin{aligned}
 & \int \left| \int p^{n-\alpha, \pi}(s|s_1) \left[(p^{\alpha, \pi}(s_1|\mu) - \nu^\pi(s_1))_+ - \left(\int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \right) \nu^\pi(s_1) \right] ds_1 \right| ds \\
 &= \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \cdot \int \left| \int p^{n-\alpha, \pi}(s|s_1) \left[\frac{(p^{\alpha, \pi}(s_1|\mu) - \nu^\pi(s_1))_+}{\int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s}} - \nu^\pi(s_1) \right] ds_1 \right| ds \\
 &= \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \cdot \int \left| \int p^{n-\alpha, \pi}(s|s_1) \frac{(p^{\alpha, \pi}(s_1|\mu) - \nu^\pi(s_1))_+}{\int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s}} ds_1 - \nu^\pi(s) \right| ds.
 \end{aligned}$$

Note that

$$\frac{(p^{\alpha, \pi}(\cdot|\mu) - \nu^\pi(\cdot))_+}{\int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s}}$$

is also a probability density of some initial distribution, and $n - \alpha \geq \alpha$, so we have

$$\begin{aligned}
 & \int \left| \int p^{\alpha, \pi}(s|s_1) \left[(p^{\alpha, \pi}(s_1|\mu) - \nu^\pi(s_1))_+ - \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} \right] \nu^\pi(s_1) \right| ds_1 \Big| ds \\
 & \leq \frac{1}{2} \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s}.
 \end{aligned}$$

Similarly, we also get

$$\begin{aligned}
 & \int \left| \int p^{\alpha, \pi}(s|s_1) \left[(p^{\alpha, \pi}(s_1|\mu) - \nu^\pi(s_1))_- - \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s} \right] \nu^\pi(s_1) \right| ds_1 \Big| ds \\
 & \leq \frac{1}{2} \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s},
 \end{aligned}$$

which implies that

$$\begin{aligned}
 \|p^{n, \pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} & \leq \frac{1}{4} \left(\int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_+ d\tilde{s} + \int (p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s}))_- d\tilde{s} \right) \\
 & = \frac{1}{4} \left(\int |p^{\alpha, \pi}(\tilde{s}|\mu) - \nu^\pi(\tilde{s})| d\tilde{s} \right) \leq \frac{1}{4} \cdot \frac{1}{2}.
 \end{aligned}$$

By induction, we can prove that for any $n \geq k\alpha$, we have

$$\|p^{n, \pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} \leq \frac{1}{4} \left(\frac{1}{2} \right)^{k-1}.$$

Therefore, for any $0 < \varepsilon < e^{-1}$, let $k = \lceil 2 \log \frac{1}{\varepsilon} \rceil + 1$, then for any $n \geq 2\alpha \log \frac{1}{\varepsilon} \geq (k-1)\alpha$, we have

$$\|p^{n, \pi}(\cdot|\mu) - \nu^\pi(\cdot)\|_{TV} \leq \frac{1}{4} \left(\frac{1}{2} \right)^{k-2} \leq \varepsilon,$$

which has finished the proof.

Proof of Lemma 12 Without loss of generality, assume that $r_3 = \max\{r_1, r_2, r_3\}$. Since $\text{Tucker-Rank}(\mathbf{X}) \leq (r_1, r_2, r_3)$, according to the result of Lemma 10, there exists a decomposition of \mathbf{X}

$$\mathbf{X} = \mathbf{G} \times_1 \mathbf{U}, \mathbf{G} \in \mathbb{R}^{r_1 \times p_2 \times p_3}, \mathbf{U} \in \mathbb{R}^{p_1 \times r_1},$$

where \mathbf{U} is a column-wise orthonormal matrix, i.e., $\mathbf{U}^\top \mathbf{U} = \mathbf{I}_{r_1}$. This formulation implies $\text{Tucker-Rank}(\mathbf{G}) \leq (r_1, r_2, r_3)$, therefore for each $1 \leq i \leq r_1$, we have $\text{rank}(\mathbf{G}_{i::}) \leq \min\{r_2, r_3\} = r_2$. We then consider the SVD of $\mathbf{G}_{i::}$,

$$\mathbf{G}_{i::} = \mathbf{V}_i \mathbf{\Lambda}_i \mathbf{W}_i^\top,$$

where $\mathbf{\Lambda}_i = \text{diag}(\lambda_{i1}, \lambda_{i2}, \dots, \lambda_{ir_2})$, $\mathbf{V}_i \in \mathbb{R}^{p_2 \times r_2}$, $\mathbf{W}_i \in \mathbb{R}^{p_3 \times r_2}$, $\mathbf{V}_i^\top \mathbf{V}_i = \mathbf{W}_i^\top \mathbf{W}_i = \mathbf{I}_{r_2}$. The above formulation is equivalent to

$$\mathbf{X} = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \lambda_{ij} \mathbf{u}_i \circ \mathbf{v}_{ij} \circ \mathbf{w}_{ij},$$

where \mathbf{v}_{ij} is the j th column of \mathbf{V}_i , \mathbf{w}_{ij} is the j th column of \mathbf{W}_i , \mathbf{u}_i is the i th column of \mathbf{U} . According to the definition of $\|\cdot\|_\sigma$, we have

$$\begin{aligned} \lambda_{ij} &= \left\langle \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \lambda_{ij} \mathbf{u}_i \circ \mathbf{v}_{ij} \circ \mathbf{w}_{ij}, \mathbf{u}_i \circ \mathbf{v}_{ij} \circ \mathbf{w}_{ij} \right\rangle \leq \|\mathbf{X}\|_\sigma \\ -\lambda_{ij} &= \left\langle \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \lambda_{ij} \mathbf{u}_i \circ \mathbf{v}_{ij} \circ \mathbf{w}_{ij}, (-\mathbf{u}_i) \circ \mathbf{v}_{ij} \circ \mathbf{w}_{ij} \right\rangle \leq \|\mathbf{X}\|_\sigma. \end{aligned}$$

So we get

$$|\lambda_{ij}|^2 \leq \|\mathbf{X}\|_\sigma^2, \forall i, j.$$

On the other hand,

$$\begin{aligned} \|\mathbf{X}\|_F^2 &= \sum_{a=1}^{p_1} \sum_{b=1}^{p_2} \sum_{c=1}^{p_3} \left(\sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \lambda_{ij} \mathbf{u}_{ia} \mathbf{v}_{ijb} \mathbf{w}_{ijc} \right)^2 = \sum_{a=1}^{p_1} \sum_{b=1}^{p_2} \sum_{c=1}^{p_3} \sum_{i=1}^{r_1} \mathbf{u}_{ia}^2 \left(\sum_{j=1}^{r_2} \lambda_{ij} \mathbf{v}_{ijb} \mathbf{w}_{ijc} \right)^2 \\ &= \sum_{i=1}^{r_1} \sum_{b=1}^{p_2} \sum_{c=1}^{p_3} \left(\sum_{j=1}^{r_2} \lambda_{ij} \mathbf{v}_{ijb} \mathbf{w}_{ijc} \right)^2 = \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \sum_{b=1}^{p_2} \sum_{c=1}^{p_3} \mathbf{v}_{ijb}^2 \mathbf{w}_{ijc}^2 \lambda_{ij}^2 \\ &= \sum_{i=1}^{r_1} \sum_{j=1}^{r_2} \lambda_{ij}^2 \leq r_1 r_2 \|\mathbf{X}\|_\sigma^2, \end{aligned}$$

which implies

$$\|\mathbf{X}\|_F \leq \sqrt{r_1 r_2} \|\mathbf{X}\|_\sigma = \sqrt{\frac{r_1 r_2 r_3}{\max\{r_1, r_2, r_3\}}} \|\mathbf{X}\|_\sigma.$$

Proof of Lemma 13 The conclusion can be directly derived from Corollary 4.2.13 in Vershynin (2017).

Proof of Lemma 14 According to the definition of $\|\cdot\|_\sigma$, we can always find $x_0 \in S^{p_1-1}, y_0 \in S^{p_2-1}, z_0 \in S^{p_3-1}$ such that

$$\langle \mathbf{X}, x_0 \circ y_0 \circ z_0 \rangle = \|\mathbf{X}\|_\sigma$$

Then according to the definition of ε -net, we can always find x, y, z from these ε -nets such that $\|x - x_0\|_2 \leq \varepsilon, \|y - y_0\|_2 \leq \varepsilon, \|z - z_0\|_2 \leq \varepsilon$, then

$$\begin{aligned} & |\langle \mathbf{X}, x_0 \circ y_0 \circ z_0 \rangle - \langle \mathbf{X}, x \circ y \circ z \rangle| \\ & \leq |\langle \mathbf{X}, (x_0 - x) \circ y_0 \circ z_0 \rangle| + |\langle \mathbf{X}, x \circ (y_0 - y) \circ z_0 \rangle| + |\langle \mathbf{X}, x_0 \circ y_0 \circ (z_0 - z) \rangle| \\ & + |\langle \mathbf{X}, (x_0 - x) \circ (y_0 - y) \circ z_0 \rangle| + |\langle \mathbf{X}, (x_0 - x) \circ y_0 \circ (z_0 - z) \rangle| \\ & + |\langle \mathbf{X}, x_0 \circ (y_0 - y) \circ (z_0 - z) \rangle| + |\langle \mathbf{X}, (x_0 - x) \circ (y_0 - y) \circ (z_0 - z) \rangle| \\ & \leq \|\mathbf{X}\|_\sigma (\|x_0 - x\|_2 + \|y_0 - y\|_2 + \|z_0 - z\|_2) \\ & + \|x_0 - x\|_2 \|y_0 - y\|_2 + \|x_0 - x\|_2 \|z_0 - z\|_2 + \|y_0 - y\|_2 \|z_0 - z\|_2 \\ & + \|x_0 - x\|_2 \|y_0 - y\|_2 \|z_0 - z\|_2 \\ & \leq \|\mathbf{X}\|_\sigma (3\varepsilon + 3\varepsilon^2 + \varepsilon^3), \end{aligned}$$

which implies

$$\begin{aligned} \|\mathbf{X}\|_\sigma & \leq \max_{x \in \mathcal{N}_1, y \in \mathcal{N}_2, z \in \mathcal{N}_3} |\langle \mathbf{X}, x \circ y \circ z \rangle| + \|\mathbf{X}\|_\sigma (3\varepsilon + 3\varepsilon^2 + \varepsilon^3) \\ \Rightarrow \|\mathbf{X}\|_\sigma & \leq \frac{\max_{x \in \mathcal{N}_1, y \in \mathcal{N}_2, z \in \mathcal{N}_3} |\langle \mathbf{X}, x \circ y \circ z \rangle|}{1 - 3\varepsilon - 3\varepsilon^2 - \varepsilon^3}, \end{aligned}$$

which has finished the proof.

Proof of Lemma 15 Step 1:

Let $\mathbf{H}_i = \phi(s_i)\phi(s_i)^\top$. We introduce some sufficiently large integer α such that from any initial distribution μ , one always has

$$\|p^{\alpha, \pi}(s|\mu) - \xi(s)\|_{TV} \leq \frac{\bar{\mu}}{2K_{max}^2} \wedge \frac{t}{2K_{max}}.$$

By Lemma 11, we can simply choose $\alpha = \lceil 2t_{mix} \log \frac{4K_{max}^3}{\bar{\mu}t} \rceil + 1$ to satisfy this condition. For each $0 \leq l \leq \alpha - 1$ and $1 \leq k \leq n_l = \lceil \frac{n-l}{\alpha} \rceil$, we define \mathbf{H}_k^l as $\mathbf{H}_{k\alpha+l}$. We also denote \mathbf{H} as a random matrix independent with $\mathbf{H}_i, 1 \leq i \leq n$, which is defined as

$$\mathbf{H} = \mathbb{E}_{S \sim \xi} [\phi(S)\phi(S)^\top].$$

Denote \mathcal{F}_i as the σ -algebra generated by the history up to step i . Then we have

$$\begin{aligned} \|\mathbf{H}_k^l - \mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma & \leq \|\mathbf{H}_k^l\|_\sigma + \|\mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma \\ & \leq \|\phi(s_{k\alpha+l})\|^2 + \mathbb{E}[\|\phi(s_{k\alpha+l})\|^2 | \mathcal{F}_{(k-1)\alpha+l}] \\ & \leq 2K_{max} \end{aligned}$$

and

$$\begin{aligned}
 \|\mathbb{E}[\mathbf{H}_k^l \mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma &\leq \|\mathbb{E}[\mathbf{H}_k^l \mathbf{H}_k^l - \mathbb{E}[\mathbf{H}\mathbf{H}] | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma + \|\mathbb{E}[\mathbf{H}\mathbf{H}]\|_\sigma \\
 &= \|\mathbb{E}[\mathbf{H}_k^l \mathbf{H}_k^l - \mathbb{E}[\mathbf{H}\mathbf{H}] | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma + \bar{\mu} \\
 &= \left\| \int \|\phi(s)\|^2 \phi(s) \phi(s)^\top (p^{\alpha,\pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)) ds \right\|_\sigma + \bar{\mu} \\
 &\leq 2K_{max}^2 \|p^{\alpha,\pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)\|_{TV} + \bar{\mu} \leq 2\bar{\mu},
 \end{aligned}$$

where μ_i is the state distribution at step i . We then have

$$\|\mathbb{E}[(\mathbf{H}_k^l - \mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}])(\mathbf{H}_k^l - \mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}] | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma \leq \|\mathbb{E}[\mathbf{H}_k^l \mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}]\|_\sigma \leq 2\bar{\mu}.$$

So according to the martingale version of matrix Bernstein's inequality (See e.g. Tropp (2011)), we have

$$\mathbb{P}\left(\left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}])\right\|_\sigma > t\right) \leq 2d_S e^{-\frac{1}{2} \frac{n_l t^2}{2\bar{\mu} + 2K_{max} t/3}}.$$

Step 2:

We have

$$\begin{aligned}
 \|\mathbb{E}[\mathbf{H}_k^l | \mathcal{F}_{(k-1)\alpha+l}] - \mathbf{H}\|_\sigma &= \left\| \int \phi(s) \phi(s)^\top (p^{\alpha,\pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)) ds \right\|_\sigma \\
 &\leq 2K_{max} \|p^{\alpha,\pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)\|_{TV} \leq t,
 \end{aligned}$$

which implies

$$\mathbb{P}\left(\left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbf{H})\right\|_\sigma > 2t\right) \leq 2d_S e^{-\frac{1}{2} \frac{n_l t^2}{2\bar{\mu} + 2K_{max} t/3}}.$$

We then use a union bound to get

$$\begin{aligned}
 \mathbb{P}\left(\exists l, \left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbf{H})\right\|_\sigma > 2t\right) &= \mathbb{P}\left(\bigcup_{l=0}^{\alpha-1} \left\{\left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbf{H})\right\|_\sigma > 2t\right\}\right) \\
 &\leq \sum_{l=0}^{\alpha-1} \mathbb{P}\left(\left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbf{H})\right\|_\sigma > 2t\right) \\
 &\leq \sum_{l=0}^{\alpha-1} 2d_S e^{-\frac{1}{2} \frac{n_l t^2}{2\bar{\mu} + 2K_{max} t/3}} \\
 &\leq 2\alpha d_S e^{-\frac{1}{2} \frac{(n-2\alpha)t^2}{2\alpha(\bar{\mu} + K_{max} t/3)}},
 \end{aligned}$$

where we use the fact that $n_l = \lceil \frac{n-l}{\alpha} \rceil \geq \frac{n}{\alpha} - 2$. Therefore, we get

$$\begin{aligned} \mathbb{P}(\|\hat{\Sigma} - \Sigma\|_\sigma > 2t) &= \mathbb{P}\left(\left\|\frac{1}{n} \sum_{k=1}^n \mathbf{H}_k - \mathbf{H}\right\|_\sigma > 2t\right) \\ &\leq \mathbb{P}\left(\exists l, \left\|\frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{H}_k^l - \mathbf{H})\right\|_\sigma > 2t\right) \\ &\leq 2\alpha d_S e^{-\frac{1}{2} \frac{(n-2\alpha)t^2}{2\alpha(\bar{\mu} + K_{max}t/3)}}. \end{aligned}$$

Replacing $2t$ by t , we get

$$\mathbb{P}(\|\hat{\Sigma} - \Sigma\|_\sigma > t) \leq 2\alpha d_S e^{-\frac{1}{16} \frac{(n-2\alpha)t^2}{\alpha(\bar{\mu} + K_{max}t/6)}}.$$

Now we assume

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 256 \frac{K_{max}^2}{\bar{\mu}} \log \frac{d_S t_{mix}}{\delta}$$

and take

$$t = \sqrt{\frac{256 \bar{\mu} \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Then we have

$$\begin{aligned} \alpha &\leq 4t_{mix} \log \frac{4K_{max}^3}{\bar{\mu}t} = 4t_{mix} \log \left(\frac{K_{max}^3}{\bar{\mu}^{3/2}} \sqrt{\frac{n/t_{mix}}{16 \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}} \right) \\ &\leq 4t_{mix} \log \left(\frac{K_{max}^3}{\bar{\mu}^{3/2}} \sqrt{\frac{n/t_{mix}}{16 \log \frac{n}{t_{mix}}}} \right) \leq 8t_{mix} \log \frac{n/t_{mix}}{\log \frac{n}{t_{mix}}} \leq 8t_{mix} \log \frac{n}{t_{mix}} \leq \frac{1}{4}n. \end{aligned}$$

Meanwhile,

$$\frac{1}{6}K_{max}t = \frac{1}{6}K_{max} \sqrt{\frac{256 \bar{\mu} \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \leq \bar{\mu},$$

so we have

$$\begin{aligned} \mathbb{P}(\|\hat{\Sigma} - \Sigma\|_\sigma > t) &\leq 2\alpha d_S e^{-\frac{1}{16} \frac{(n-2\alpha)t^2}{\alpha(\bar{\mu} + K_{max}t/6)}} \leq 16t_{mix} d_S \log \frac{n}{t_{mix}} e^{-\frac{1}{256} \frac{nt^2}{t_{mix}\bar{\mu} \log \frac{n}{t_{mix}}}} \\ &= 16t_{mix} d_S \log \frac{n}{t_{mix}} e^{-\log \frac{d_S t_{mix}}{\delta} \log \frac{n}{t_{mix}}} \\ &\leq 16t_{mix} d_S \log \frac{n}{t_{mix}} e^{-(\log \frac{d_S t_{mix}}{\delta} + \log \frac{n}{t_{mix}})} \\ &= \frac{16\delta \log \frac{n}{t_{mix}}}{n/t_{mix}} \leq \delta, \end{aligned}$$

i.e., with probability at least $1 - \delta$, we have

$$\|\hat{\Sigma} - \Sigma\|_{\sigma} \leq \sqrt{\frac{256\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Step 3:

Notice the relation

$$\begin{aligned} \|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_{\sigma} &\leq \|\Sigma^{-1}\|_{\sigma} \|\hat{\Sigma} - \Sigma\|_{\sigma} \|\hat{\Sigma}^{-1}\Sigma^{\frac{1}{2}}\|_{\sigma} \\ &\leq (\|(\Sigma^{-1} - \hat{\Sigma}^{-1})\Sigma^{\frac{1}{2}}\|_{\sigma} + \|\Sigma^{-\frac{1}{2}}\|_{\sigma}) \|\Sigma^{-1}\|_{\sigma} \|\hat{\Sigma} - \Sigma\|_{\sigma}, \end{aligned}$$

i.e.,

$$\|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_{\sigma} \leq \frac{\|\Sigma^{-1}\|_{\sigma}^{\frac{3}{2}} \|\hat{\Sigma} - \Sigma\|_{\sigma}}{1 - \|\Sigma^{-1}\|_{\sigma} \|\hat{\Sigma} - \Sigma\|_{\sigma}}.$$

Therefore, if

$$\frac{n/t_{mix}}{(\log n/t_{mix})^2} \geq 1024\bar{\mu} \|\Sigma^{-1}\|_{\sigma}^2 \log \frac{dst_{mix}}{\delta},$$

then we have

$$\|\Sigma^{-1}\|_{\sigma} \|\hat{\Sigma} - \Sigma\|_{\sigma} \leq \frac{1}{2}$$

and

$$\|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_{\sigma} \leq 2\|\Sigma^{-1}\|_{\sigma}^{\frac{3}{2}} \|\hat{\Sigma} - \Sigma\|_{\sigma} \leq 32\|\Sigma^{-1}\|_{\sigma}^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

In summary, if we have

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\bar{\mu} \|\Sigma^{-1}\|_{\sigma}^2 + \frac{K_{max}^2}{\bar{\mu}} \right) \log \frac{dst_{mix}}{\delta},$$

then with probability $1 - \delta$, we have

$$\|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_{\sigma} \leq 32\|\Sigma^{-1}\|_{\sigma}^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{dst_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Proof of Lemma 16 Step 1:

Denote $\mathbf{G}_i = \frac{\eta(a_i)}{\pi(a_i|s_i)} \phi(s_i) \circ \psi(a_i) \circ \phi(s'_i)$ and let $\alpha = \lceil 2t_{mix} \log \frac{4K_{max}^{\frac{9}{2}}}{\lambda t} \rceil + 1$. Then according to the result of Lemma 11, we have that for arbitrary initial state distribution μ ,

$$\|p^{\alpha, \pi}(s|\mu) - \xi(s)\|_{TV} \leq \frac{\bar{\lambda}}{2K_{max}^3} \wedge \frac{t}{2K_{max}^{3/2}}.$$

For each $0 \leq l \leq \alpha - 1$ and $1 \leq k \leq n_l = \lceil \frac{n-l}{\alpha} \rceil$, we define $\mathbf{G}_k^l = \mathbf{G}_{k\alpha+l}$. We also denote \mathbf{G} as a random tensor independent with our data which is defined as

$$\mathbf{G} = \frac{\eta(A)}{\pi(A|S)} \phi(S) \circ \psi(A) \circ \phi(S'), S \sim \xi(\cdot), A \sim \pi(\cdot|S), S' \sim p(\cdot|S, A).$$

Then for any $u, w \in \mathbb{R}^{d_S}, v \in \mathbb{R}^{d_A}$, we have

$$|\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top| \leq \left\| \frac{\eta(a_{k\alpha+l})}{\pi(a_{k\alpha+l}|s_{k\alpha+l})} \phi(s_{k\alpha+l}) \right\| \|\psi(a_{k\alpha+l})\| \|\phi(s'_{k\alpha+l})\| \leq K_{max}^{\frac{3}{2}} \kappa.$$

Therefore,

$$|\mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}]| \leq K_{max}^{\frac{3}{2}} \kappa,$$

which implies

$$|\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top - \mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}]| \leq 2K_{max}^{\frac{3}{2}} \kappa.$$

Meanwhile, we have

$$\begin{aligned} & \mathbb{E}[(\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2 | \mathcal{F}_{(k-1)\alpha+l}] \\ &= \mathbb{E}[(\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2 \\ & \quad - \mathbb{E}[(\mathbf{G} \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2 | \mathcal{F}_{(k-1)\alpha+l}] + \mathbb{E}[(\mathbf{G} \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2] \\ &= \int (\phi(s)^\top u)^2 (\psi(a)^\top v)^2 (\phi(s')^\top w)^2 (p^{\alpha, \pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)) \frac{\eta^2(a)}{\pi(a|s)} p(s'|s, a) ds da ds' \\ & \quad + \mathbb{E}[(\mathbf{G} \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2] \\ &\leq \kappa \int (\phi(s)^\top u)^2 (\psi(a)^\top v)^2 (\phi(s')^\top w)^2 |p^{\alpha, \pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)| \eta(a) p(s'|s, a) ds da ds' + \kappa \bar{\lambda} \\ &\leq 2\kappa \bar{\lambda}, \end{aligned}$$

where μ_i is the state distribution at step i . Therefore, we have

$$\begin{aligned} & \mathbb{E}[(\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top - \mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}])^2 | \mathcal{F}_{(k-1)\alpha+l}] \\ &\leq \mathbb{E}[(\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top)^2 | \mathcal{F}_{(k-1)\alpha+l}] \leq 2\kappa \bar{\lambda}. \end{aligned}$$

Again we apply the matrix Bernstein's inequality Tropp (2011) on $\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top$ (Note that $\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top$ is a scalar, which can be viewed as a 1×1 matrix.), and get

$$\begin{aligned} & \mathbb{P} \left(\left| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top - \mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}]) \right| > t \right) \\ &\leq 2e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa \bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t}}. \end{aligned}$$

Step 2:

Now consider three $\frac{1}{4}$ -nets over $S^{d_S-1}, S^{d_A-1}, S^{d_S-1}$, denoted as $\mathcal{N}_1, \mathcal{N}_2, \mathcal{N}_3$. By Lemma 13, we know that $|\mathcal{N}_1| \leq 9^{d_S}, |\mathcal{N}_2| \leq 9^{d_A}, |\mathcal{N}_3| \leq 9^{d_S}$. Then we can get a union bound by

$$\begin{aligned} & \mathbb{P} \left(\exists u \in \mathcal{N}_1, v \in \mathcal{N}_2, w \in \mathcal{N}_3, \right. \\ & \quad \left. \left| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top - \mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}]) \right| > t \right) \\ & \leq \sum_{u \in \mathcal{N}_1} \sum_{v \in \mathcal{N}_2} \sum_{w \in \mathcal{N}_3} \\ & \quad \mathbb{P} \left(\left| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top - \mathbb{E}[\mathbf{G}_k^l \times_1 u^\top \times_2 v^\top \times_3 w^\top | \mathcal{F}_{(k-1)\alpha+l}]) \right| > t \right) \\ & = 2 \cdot 9^{2d_S+d_A} e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa\bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t}} \leq 2e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa\bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t} + 4(d_S+d_A) \log 3}. \end{aligned}$$

Then according to Lemma 14, we know that

$$\begin{aligned} & \left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}]) \right\|_\sigma \\ & \leq \frac{64}{3} \max_{u \in \mathcal{N}_1, v \in \mathcal{N}_2, w \in \mathcal{N}_3} \left| \left\langle \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}]), u \circ v \circ w \right\rangle \right|, \end{aligned}$$

which implies that

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}]) \right\|_\sigma > \frac{3}{64} t \right) & \leq 2e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa\bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t} + 4(d_S+d_A) \log 3} \\ & \leq 2e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa\bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t} + 8(d_S+d_A)}. \end{aligned}$$

Step 3:

Note that

$$\mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}] - \mathbf{F} = \int \phi(s) \circ \psi(a) \circ \phi(s') (p^{\alpha, \pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)) \eta(a) p(s' | s, a) ds da ds'.$$

So we get

$$\|\mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}] - \mathbf{F}\|_\sigma \leq K_{max}^{\frac{3}{2}} \int |p^{\alpha, \pi}(s | \mu_{(k-1)\alpha+l}) - \xi(s)| ds da ds' \leq t,$$

which implies that

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbf{F}) \right\|_\sigma > 2t \right) & \leq \mathbb{P} \left(\left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbb{E}[\mathbf{G}_k^l | \mathcal{F}_{(k-1)\alpha+l}]) \right\|_\sigma > \frac{3}{64} t \right) \\ & \leq 2e^{-\frac{1}{2} \frac{n_l t^2}{2\kappa\bar{\lambda} + \frac{2}{3} K_{max}^{3/2} \kappa t} + 8(d_S+d_A)}. \end{aligned}$$

Based on the fact that $n_l = \lceil \frac{n-l}{\alpha} \rceil \geq \frac{n}{\alpha} - 2$, and replace $2t$ by t , we further get

$$\mathbb{P} \left(\left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbf{F}) \right\|_{\sigma} > t \right) \leq 2e^{-\frac{1}{16} \frac{(n-2\alpha)t^2}{\alpha(\kappa\bar{\lambda} + \frac{1}{6}K_{max}^3\kappa t)} + 8(d_S + d_A)}.$$

Step 4:

Now, we get a union bound over l , and get

$$\begin{aligned} \mathbb{P} \left(\exists l, \left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbf{F}) \right\|_{\sigma} > t \right) &\leq \sum_{l=0}^{\alpha-1} \mathbb{P} \left(\left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbf{F}) \right\|_{\sigma} > 2t \right) \\ &\leq 2\alpha e^{-\frac{1}{16} \frac{(n-2\alpha)t^2}{\alpha(\kappa\bar{\lambda} + \frac{1}{6}K_{max}^3\kappa t)} + 8(d_S + d_A)}. \end{aligned}$$

Such the result implies that

$$\begin{aligned} \mathbb{P} \left(\left\| \frac{1}{n} \sum_{k=1}^n (\mathbf{G}_k - \mathbf{F}) \right\|_{\sigma} > t \right) &\leq \mathbb{P} \left(\exists l, \left\| \frac{1}{n_l} \sum_{k=1}^{n_l} (\mathbf{G}_k^l - \mathbf{F}) \right\|_{\sigma} > t \right) \\ &\leq 2\alpha e^{-\frac{1}{16} \frac{(n-2\alpha)t^2}{\alpha(\kappa\bar{\lambda} + \frac{1}{6}K_{max}^3\kappa t)} + 8(d_S + d_A)}. \end{aligned}$$

Now we assume

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \frac{\kappa K_{max}^3}{\bar{\lambda}} \left(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A) \right)$$

and take

$$t = \sqrt{\frac{1024\kappa\bar{\lambda}(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A))(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Then we have

$$\begin{aligned} \alpha &\leq 4t_{mix} \log \frac{4K_{max}^{\frac{9}{2}}}{\bar{\lambda}t} = 4t_{mix} \log \left(\frac{K_{max}^{\frac{9}{2}}}{\bar{\lambda}^{3/2}} \sqrt{\frac{n/t_{mix}}{64\kappa(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A))(\log \frac{n}{t_{mix}})^2}} \right) \\ &\leq 4t_{mix} \log \left(\frac{K_{max}^{\frac{9}{2}}}{\bar{\lambda}^{3/2}} \sqrt{\frac{n/t_{mix}}{64 \log \frac{n}{t_{mix}}}} \right) \leq 8t_{mix} \log \frac{n/t_{mix}}{\log \frac{n}{t_{mix}}} \leq 8t_{mix} \log \frac{n}{t_{mix}} \leq \frac{1}{4}n \end{aligned}$$

and

$$\frac{1}{6}K_{max}^{\frac{3}{2}}t = \frac{1}{6}K_{max}^{\frac{3}{2}} \sqrt{\frac{1024\bar{\lambda}\kappa(\log \frac{t_{mix}}{\delta} + d_S + d_A)(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \leq \bar{\lambda}.$$

So we have

$$\begin{aligned}
 \mathbb{P}\left(\left\|\frac{1}{n}\sum_{k=1}^n(\mathbf{G}_k - \mathbf{F})\right\|_\sigma > t\right) &\leq 2\alpha e^{-\frac{1}{16}\frac{(n-2\alpha)t^2}{\alpha(\kappa\bar{\lambda} + \frac{1}{6}K_{max}^{3/2}\kappa t)} + 8(d_S + d_A)} \\
 &\leq 16t_{mix} \log \frac{n}{t_{mix}} e^{-\frac{1}{1024}\frac{nt^2}{t_{mix}\kappa\bar{\lambda}\log\frac{n}{t_{mix}}} + 8(d_S + d_A)} \\
 &= 16t_{mix} \log \frac{n}{t_{mix}} e^{-(\log\frac{t_{mix}}{\delta} + 8(d_S + d_A))\log\frac{n}{t_{mix}} + 8(d_S + d_A)} \\
 &\leq 16t_{mix} \log \frac{n}{t_{mix}} e^{-\log\frac{t_{mix}}{\delta}\log\frac{n}{t_{mix}}} \\
 &\leq 16t_{mix} \log \frac{n}{t_{mix}} e^{-\log\frac{t_{mix}}{\delta} - \log\frac{n}{t_{mix}}} \\
 &= \frac{16\delta \log\frac{n}{t_{mix}}}{n/t_{mix}} \leq \delta,
 \end{aligned}$$

i.e., with probability at least $1 - \delta$, we have

$$\left\|\frac{1}{n}\sum_{k=1}^n(\mathbf{G}_k - \mathbf{F})\right\|_\sigma = \|\bar{\mathbf{F}} - \mathbf{F}\|_\sigma \leq \sqrt{\frac{1024\kappa\bar{\lambda}(\log\frac{t_{mix}}{\delta} + 8(d_S + d_A))(\log\frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Step 5: According to our programming, we have

$$\|\hat{\mathbf{F}} - \mathbf{F}\|_\sigma \leq \|\hat{\mathbf{F}} - \bar{\mathbf{F}}\|_\sigma + \|\bar{\mathbf{F}} - \mathbf{F}\|_\sigma \leq 2\|\bar{\mathbf{F}} - \mathbf{F}\|_\sigma.$$

So with probability at least $1 - \delta$, we have

$$\|\hat{\mathbf{F}} - \mathbf{F}\|_\sigma \leq \sqrt{\frac{4096\kappa\bar{\lambda}(\log\frac{t_{mix}}{\delta} + 8(d_S + d_A))(\log\frac{n}{t_{mix}})^2}{n/t_{mix}}}.$$

Proof of Lemma 17 Given that

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\|\Sigma^{-1}\|_\sigma^2 \bar{\mu} + \frac{K_{max}^2}{\bar{\mu}} + \frac{\kappa K_{max}^3}{\bar{\lambda}} \right) \left(\log\frac{t_{mix}}{\delta} + 8(d_S + d_A) \right),$$

then the assumptions in Lemma 16 and Lemma 15 are satisfied simultaneously, and with probability at least $1 - 2\delta$, we have the following relations hold simultaneously,

$$\begin{aligned}
 \|\hat{\mathbf{F}} - \mathbf{F}\|_\sigma &\leq 64 \sqrt{\frac{\kappa\bar{\lambda}(\log\frac{t_{mix}}{\delta} + 8(d_S + d_A))(\log\frac{n}{t_{mix}})^2}{n/t_{mix}}}, \\
 \|(\hat{\Sigma}^{-1} - \Sigma^{-1})\Sigma^{\frac{1}{2}}\|_\sigma &\leq 32 \|\Sigma^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log\frac{d_S t_{mix}}{\delta} (\log\frac{n}{t_{mix}})^2}{n/t_{mix}}}.
 \end{aligned}$$

So we have

$$\begin{aligned}
 & \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma \\
 &= \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top + \mathbf{F} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma \\
 &\leq \|(\hat{\mathbf{F}} - \mathbf{F}) \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma + \|\mathbf{F} \times_1 ((\hat{\Sigma}^{-1} - \Sigma^{-1})^\top \Sigma^{\frac{1}{2}})\|_\sigma \\
 &\leq \|\mathbf{F} - \hat{\mathbf{F}}\|_\sigma \|\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}}\|_\sigma + \|\mathbf{F}\|_\sigma \|(\hat{\Sigma}^{-1} - \Sigma^{-1}) \Sigma^{\frac{1}{2}}\|_\sigma \\
 &\leq \|\mathbf{F} - \hat{\mathbf{F}}\|_\sigma \|\Sigma^{-\frac{1}{2}}\|_\sigma + \|\mathbf{F}\|_\sigma \|(\hat{\Sigma}^{-1} - \Sigma^{-1}) \Sigma^{\frac{1}{2}}\|_\sigma + \|\mathbf{F} - \hat{\mathbf{F}}\|_\sigma \|(\hat{\Sigma}^{-1} - \Sigma^{-1}) \Sigma^{\frac{1}{2}}\|_\sigma.
 \end{aligned}$$

Note that under our assumptions, we have

$$\begin{aligned}
 \|(\hat{\Sigma}^{-1} - \Sigma^{-1}) \Sigma^{\frac{1}{2}}\|_\sigma &\leq 32 \|\Sigma^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \\
 &\leq 32 \|\Sigma^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{d_S t_{mix}}{\delta}}{1024 \|\Sigma^{-1}\|_\sigma^2 \bar{\mu} (\log \frac{t_{mix}}{\delta} + 8(d_S + d_A))}} \\
 &\leq \|\Sigma^{-\frac{1}{2}}\|_\sigma.
 \end{aligned}$$

Meanwhile, we have

$$\begin{aligned}
 \bar{\lambda} &= \sup_{\|u\| \leq 1, \|v\| \leq 1, \|w\| \leq 1} \mathbb{E}[\langle \phi(S) \circ \psi(A) \circ \phi(S'), u \circ v \circ w \rangle^2] \\
 &\geq \sup_{\|u\| \leq 1, \|v\| \leq 1, \|w\| \leq 1} (\mathbb{E}[\langle \phi(S) \circ \psi(A) \circ \phi(S'), u \circ v \circ w \rangle])^2 \\
 &= \left(\sup_{\|u\| \leq 1, \|v\| \leq 1, \|w\| \leq 1} \mathbb{E}[\langle \phi(S) \circ \psi(A) \circ \phi(S'), u \circ v \circ w \rangle] \right)^2 \\
 &= \|\mathbf{F}\|_\sigma^2.
 \end{aligned}$$

So we get

$$\begin{aligned}
 & \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma \\
 &\leq 2 \|\mathbf{F} - \hat{\mathbf{F}}\|_\sigma \|\Sigma^{-\frac{1}{2}}\|_\sigma + \|\mathbf{F}\|_\sigma \|(\hat{\Sigma}^{-1} - \Sigma^{-1}) \Sigma^{\frac{1}{2}}\|_\sigma \\
 &\leq 128 \|\Sigma^{-1}\|_\sigma^{\frac{1}{2}} \sqrt{\frac{\kappa \bar{\lambda} (\log \frac{t_{mix}}{\delta} + 8(d_S + d_A)) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \\
 &\quad + 32 \|\mathbf{F}\|_\sigma \|\Sigma^{-1}\|_\sigma^{\frac{3}{2}} \sqrt{\frac{\bar{\mu} \log \frac{d_S t_{mix}}{\delta} (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \\
 &\leq 256 \|\Sigma^{-1}\|_\sigma^{\frac{1}{2}} \sqrt{\frac{\bar{\lambda} (\log \frac{t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}.
 \end{aligned}$$

Replacing δ by $\frac{1}{2}\delta$, then with probability at least $1 - \delta$, we get

$$\begin{aligned}
 & \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma \\
 &\leq 256 \|\Sigma^{-1}\|_\sigma^{\frac{1}{2}} \sqrt{\frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}},
 \end{aligned}$$

which has finished the proof.

Proof of Lemma 18 Note that for any vector $w \in \mathbb{R}^{d_S}$ such that $\|w\| = 1$, the columns of \mathbf{U}_3 span the row space of matrix $\mathbf{P} \times_1 w^\top$ and the columns of $\hat{\mathbf{U}}_3$ span the row space of matrix $\hat{\mathbf{P}} \times_1 w^\top$. So according to Wedin's lemma Wedin (1972), we have

$$\|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\| \leq \frac{\|\mathbf{P} \times_1 w^\top - \hat{\mathbf{P}} \times_1 w^\top\|}{\sigma_m(\mathbf{P} \times_1 w^\top)} \leq \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma}{\sigma_m(\mathbf{P} \times_1 w^\top)}.$$

Taking infimum over w , we get

$$\|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\| \leq \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma}{\sup_{\|w\|=1} \sigma_m(\mathbf{P} \times_1 w^\top)} = \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma}{\sigma},$$

which has finished the proof.

Proof of Theorem 4 According to the result of Lemma 17, we know that with probability at least $1 - \delta$, we have

$$\begin{aligned} & \|\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top\|_\sigma \\ & \leq 256 \|\Sigma^{-1}\|_\sigma^{\frac{1}{2}} \sqrt{\frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2)(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \end{aligned}$$

and our result follows directly by noticing

$$\begin{aligned} \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma &= \|\hat{\mathbf{F}} \times_1 \hat{\Sigma}^{-1} - \mathbf{F} \times_1 \Sigma^{-1}\|_\sigma \\ &= \|(\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top) \times_1 \Sigma^{-\frac{1}{2}}\|_\sigma \\ &\leq \|(\hat{\mathbf{F}} \times_1 (\hat{\Sigma}^{-1} \Sigma^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\Sigma^{-1} \Sigma^{\frac{1}{2}})^\top)\|_\sigma \|\Sigma^{-1}\|_\sigma^{\frac{1}{2}} \\ &\leq 256 \|\Sigma^{-1}\|_\sigma \sqrt{\frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2)(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}, \end{aligned}$$

which has finished the proof.

Proof of Theorem 5 We first prove

$$\text{dist}[(s, a), (s', a')] = \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \mathbf{C} \times_1 f(s')^\top \times_2 g(a')^\top\|.$$

For any $f \in \mathcal{H}_S$ such that $\|f\|_{\mathcal{H}_S} \leq 1$, we can always find some weight \mathbf{w} such that $\|\mathbf{w}\| \leq 1$, $f = \mathbf{w}^\top \phi$, and

$$\begin{aligned} & |\langle p(\cdot|s, a), f(\cdot) \rangle - \langle p(\cdot|s', a'), f(\cdot) \rangle| \\ &= |\mathbf{w}^\top \langle p(\cdot|s, a), \phi(\cdot) \rangle - \mathbf{w}^\top \langle p(\cdot|s', a'), \phi(\cdot) \rangle| \\ &\leq \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle p(\cdot|s', a'), \phi(\cdot) \rangle\| \\ &= \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top - \mathbf{P} \times_1 \phi(s')^\top \times_2 \psi(a')^\top\| \\ &= \|\mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s))^\top \times_2 (\mathbf{U}_2^\top \psi(a))^\top \times_3 \mathbf{U}_3 - \mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s'))^\top \times_2 (\mathbf{U}_2^\top \psi(a'))^\top \times_3 \mathbf{U}_3\| \\ &= \|\mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s))^\top \times_2 (\mathbf{U}_2^\top \psi(a))^\top - \mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s'))^\top \times_2 (\mathbf{U}_2^\top \psi(a'))^\top\| \\ &= \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \mathbf{C} \times_1 f(s')^\top \times_2 g(a')^\top\|. \end{aligned}$$

Now we are ready to prove the main result. Notice that

$$\begin{aligned}
 \widehat{\text{dist}}[(s, a), (s', a')] &= \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \hat{\mathbf{C}} \times_1 \hat{f}(s')^\top \times_2 \hat{g}(a')^\top\| \\
 &= \|\hat{\mathbf{C}} \times_1 (\hat{\mathbf{U}}_1^\top \phi(s))^\top \times_2 (\hat{\mathbf{U}}_2^\top \psi(a))^\top - \hat{\mathbf{C}} \times_1 (\hat{\mathbf{U}}_1^\top \phi(s'))^\top \times_2 (\hat{\mathbf{U}}_2^\top \psi(a'))^\top\| \\
 &= \|\hat{\mathbf{P}} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 \hat{\mathbf{U}}_3^\top - \hat{\mathbf{P}} \times_1 \phi(s')^\top \times_2 \psi(a')^\top \times_3 \hat{\mathbf{U}}_3^\top\| \\
 &= \|\hat{\Phi}(s, a) - \hat{\Phi}(s', a')\|
 \end{aligned}$$

and for any orthogonal matrix $\mathbf{O} \in \mathbb{R}^{m \times m}$, we have

$$\begin{aligned}
 \text{dist}[(s, a), (s', a')] &= \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \mathbf{C} \times_1 f(s')^\top \times_2 g(a')^\top\| \\
 &= \|\mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s))^\top \times_2 (\mathbf{U}_2^\top \psi(a))^\top - \mathbf{C} \times_1 (\mathbf{U}_1^\top \phi(s'))^\top \times_2 (\mathbf{U}_2^\top \psi(a'))^\top\| \\
 &= \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 \mathbf{U}_3^\top - \mathbf{P} \times_1 \phi(s')^\top \times_2 \psi(a')^\top \times_3 \mathbf{U}_3^\top\| \\
 &= \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\mathbf{U}_3 \mathbf{O})^\top - \mathbf{P} \times_1 \phi(s')^\top \times_2 \psi(a')^\top \times_3 (\mathbf{U}_3 \mathbf{O})^\top\| \\
 &= \|\mathbf{O}^\top \Phi(s, a) - \mathbf{O}^\top \Phi(s', a')\|.
 \end{aligned}$$

So we have

$$|\widehat{\text{dist}}[(s, a), (s', a')] - \text{dist}[(s, a), (s', a')]| \leq \|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\| + \|\hat{\Phi}(s', a') - \mathbf{O}^\top \Phi(s', a')\|.$$

It suffices to bound $\|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\|$, and $\|\hat{\Phi}(s', a') - \mathbf{O}^\top \Phi(s', a')\|$ can be bounded in the exactly same way. Notice that

$$\begin{aligned}
 &\|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\| \\
 &\leq \|(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 \hat{\mathbf{U}}_3^\top + \mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O})^\top\| \\
 &\leq \|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma K_{max} + \|\mathbf{P}\|_\sigma K_{max} \|\mathbf{U}_3^\top - \hat{\mathbf{U}}_3^\top \mathbf{O}\|.
 \end{aligned}$$

In particular, By Part 3, Lemma 1 in Cai and Zhang (2018), one can choose \mathbf{O} such that

$$\|\mathbf{U}_3^\top - \hat{\mathbf{U}}_3^\top \mathbf{O}\| \leq \sqrt{2} \|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|.$$

Then we have

$$\|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\| \leq \|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma K_{max} + \sqrt{2} \|\mathbf{P}\|_\sigma K_{max} \|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|.$$

Now according to the result of Lemma 18, we know that

$$\|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\| \leq \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma}{\sigma}.$$

So we get

$$\|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\| \leq K_{max} \left(1 + \sqrt{2} \frac{\|\mathbf{P}\|_\sigma}{\sigma}\right) \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma.$$

It follows that

$$\begin{aligned}
 |\widehat{\text{dist}}[(s, a), (s', a')] - \text{dist}[(s, a), (s', a')]| &\leq \|\hat{\Phi}(s, a) - \mathbf{O}^\top \Phi(s, a)\| + \|\hat{\Phi}(s', a') - \mathbf{O}^\top \Phi(s', a')\| \\
 &\leq 2K_{max} \left(1 + \sqrt{2} \frac{\|\mathbf{P}\|_\sigma}{\sigma}\right) \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma,
 \end{aligned}$$

which has finished the proof.

Proof of Theorem 6 Define $\hat{p}_d(\cdot|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \hat{q}_{ij}(\cdot) \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j}$, then we have

$$\begin{aligned} & \sum_{i,j} \int_{\hat{A}_i \times \hat{B}_j} \xi(s) \eta(a) \|p(\cdot|s, a) - \hat{q}_{ij}(\cdot)\|_{\mathcal{H}_S}^2 ds da = \int \xi(s) \eta(a) \|p(\cdot|s, a) - \hat{p}_d(\cdot|s, a)\|_{\mathcal{H}_S}^2 ds da \\ & = \int \xi(s) \eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da. \end{aligned}$$

Note that

$$\begin{aligned} & \int \xi(s) \eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & = \int \xi(s) \eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle + \langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & \leq 2 \int \xi(s) \eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & \quad + 2 \int \xi(s) \eta(a) \|\langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & = 2 \int \xi(s) \eta(a) \|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \phi(s)^\top \times_2 \psi(a)^\top\|^2 ds da \\ & \quad + 2 \int \xi(s) \eta(a) \|\langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da. \end{aligned}$$

We have

$$\begin{aligned} & \int \xi(s) \eta(a) \|\langle \hat{p}(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & = \int \xi(s) \eta(a) \left\| \hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top \times_3 \hat{\mathbf{U}}_3 - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \hat{\mathbf{U}}_3 \hat{\mathbf{z}}_{ij} \right\|^2 ds da \\ & = \int \xi(s) \eta(a) \left\| \hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \hat{\mathbf{z}}_{ij} \right\|^2 ds da. \end{aligned}$$

For any orthogonal matrix $\mathbf{O} \in \mathbb{R}^{m \times m}$, because $\hat{A}_i, \hat{B}_j, \hat{z}_{ij}$ is the minimizer of the above term, we have

$$\begin{aligned}
 & \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \hat{z}_{ij}\|^2 ds da \\
 & \leq \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j} \mathbf{O} z_{ij}\|^2 ds da \\
 & \leq 2 \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O}\|^2 ds da \\
 & \quad + 2 \int \xi(s)\eta(a) \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O} - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j} \mathbf{O} z_{ij}\|^2 ds da \\
 & = 2 \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O}\|^2 ds da \\
 & \quad + 2 \int \xi(s)\eta(a) \|\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top - \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j} z_{ij}\|^2 ds da \\
 & = 2 \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O}\|^2 ds da \\
 & \quad + 2 \int \xi(s)\eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle p_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\
 & = 2 \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O}\|^2 ds da + 2L^*
 \end{aligned}$$

where we denote

$$z_{ij} = \mathbf{U}_3^\top \langle q_{ij}^*(\cdot), \phi(\cdot) \rangle, \quad p_d(\cdot|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} q_{ij}^*(\cdot) \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j}.$$

Furthermore, we have

$$\begin{aligned}
 & \int \xi(s)\eta(a) \|\hat{\mathbf{C}} \times_1 \hat{f}(s)^\top \times_2 \hat{g}(a)^\top - \mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top \times_3 \mathbf{O}\|^2 ds da \\
 & = \int \xi(s)\eta(a) \|\hat{\mathbf{P}} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 \hat{\mathbf{U}}_3^\top - \mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da \\
 & \leq 2 \int \xi(s)\eta(a) \|(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 \hat{\mathbf{U}}_3^\top\|^2 ds da \\
 & \quad + 2 \int \xi(s)\eta(a) \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da \\
 & \leq 2 \int \xi(s)\eta(a) \|(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \phi(s)^\top \times_2 \psi(a)^\top\|^2 ds da \\
 & \quad + 2 \int \xi(s)\eta(a) \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da.
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 & \int \xi(s)\eta(a)\|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\
 & \leq 4L^* + 10 \int \xi(s)\eta(a)\|(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \phi(s)^\top \times_2 \psi(a)^\top\|^2 ds da \\
 & \quad + 8 \int \xi(s)\eta(a)\|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da \\
 & = 10 \sum_{k=1}^{d_S} \int \xi(s)\eta(a)\|(\hat{\mathbf{P}} - \mathbf{P})_{::k} \times_1 \phi(s)^\top \times_2 \psi(a)^\top\|^2 ds da \\
 & \quad + 8 \int \xi(s)\eta(a)\|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da + 4L^* \\
 & = 10 \sum_{k=1}^{d_S} \int \xi(s)\eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \phi(s)\phi(s)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a) ds da \\
 & \quad + 8 \int \xi(s)\eta(a)\|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da + 4L^*.
 \end{aligned}$$

Note that

$$\int \xi(s)\phi(s)\phi(s)^\top ds = \mathbf{\Sigma}, \quad \int \eta(a)\psi(a)\psi(a)^\top da = I_{d_A}$$

So we have

$$\begin{aligned}
 & \int \xi(s)\eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \phi(s)\phi(s)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a) ds da \\
 & = \int \eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \left(\int \xi(s)\phi(s)\phi(s)^\top ds \right) (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a) da \\
 & = \int \eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \mathbf{\Sigma} (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a) da \\
 & = \int \text{tr}[\eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \mathbf{\Sigma} (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a)] da \\
 & = \text{tr} \left[(\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \mathbf{\Sigma} (\hat{\mathbf{P}} - \mathbf{P})_{::k} \int \eta(a)\psi(a)\psi(a)^\top da \right] \\
 & = \text{tr} \left[(\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \mathbf{\Sigma} (\hat{\mathbf{P}} - \mathbf{P})_{::k} \right] \\
 & = \|\mathbf{\Sigma}^{\frac{1}{2}} (\hat{\mathbf{P}} - \mathbf{P})_{::k}\|_F^2,
 \end{aligned}$$

which implies

$$\begin{aligned}
 & \sum_{k=1}^{d_S} \int \xi(s)\eta(a)\psi(a)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k}^\top \phi(s)\phi(s)^\top (\hat{\mathbf{P}} - \mathbf{P})_{::k} \psi(a) ds da \\
 & \leq \sum_{k=1}^{d_S} \|\mathbf{\Sigma}^{\frac{1}{2}} (\hat{\mathbf{P}} - \mathbf{P})_{::k}\|_F^2 = \|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \mathbf{\Sigma}^{\frac{1}{2}}\|_F^2.
 \end{aligned}$$

With a similar argument, we also have

$$\int \xi(s)\eta(a) \|\mathbf{P} \times_1 \phi(s)^\top \times_2 \psi(a)^\top \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|^2 ds da = \|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|_F^2.$$

Therefore,

$$\begin{aligned} & \int \xi(s)\eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & \leq 10 \|(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}}\|_F^2 + 8 \|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|_F^2 + 4L^*. \end{aligned}$$

Note that

$$(\hat{\mathbf{P}} - \mathbf{P}) \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} = \hat{\mathbf{F}} \times_1 (\hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top.$$

According to the result of Lemma 12, and note that Tucker-Rank($\hat{\mathbf{F}} \times_1 (\hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top$) $\leq (2r, 2l, 2m)$, we have

$$\begin{aligned} & \|\hat{\mathbf{F}} \times_1 (\hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top\|_F \\ & \leq 2 \sqrt{\frac{rlm}{\max\{r, l, m\}}} \|\hat{\mathbf{F}} \times_1 (\hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top\|_\sigma. \end{aligned}$$

For the second term, similarly we have

$$\|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|_F \leq 2 \sqrt{\frac{rlm}{\max\{r, l, m\}}} \|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|_\sigma.$$

Note that

$$\|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top)^\top\|_\sigma \leq \|\mathbf{F}\|_\sigma \|\boldsymbol{\Sigma}^{-\frac{1}{2}}\|_\sigma \|\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top\|_\sigma \leq \sqrt{\bar{\lambda}} \|\boldsymbol{\Sigma}^{-\frac{1}{2}}\|_\sigma \|\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top\|_\sigma.$$

In particular, by Part 3, Lemma 1 in Cai and Zhang (2018), we can choose \mathbf{O} such that $\|\hat{\mathbf{U}}_3 - \mathbf{U}_3 \mathbf{O}^\top\|_\sigma \leq \sqrt{2} \|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|_\sigma$ and according to the result of Lemma 18, we know that

$$\|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|_\sigma \leq \frac{\|\mathbf{P} - \hat{\mathbf{P}}\|_\sigma}{\sigma}.$$

By Theorem 4 and Lemma 17, we know that with probability at least $1 - \delta$, we have

$$\begin{aligned} & \|\hat{\mathbf{F}} \times_1 (\hat{\boldsymbol{\Sigma}}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top - \mathbf{F} \times_1 (\boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma}^{\frac{1}{2}})^\top\|_\sigma \\ & \leq 256 \|\boldsymbol{\Sigma}^{-1}\|_\sigma^{\frac{1}{2}} \sqrt{\frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}} \\ & \quad \|\hat{\mathbf{P}} - \mathbf{P}\|_\sigma \\ & \leq 256 \|\boldsymbol{\Sigma}^{-1}\|_\sigma \sqrt{\frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}}}, \end{aligned}$$

which implies that with probability at least $1 - \delta$,

$$\begin{aligned} & \int \xi(s)\eta(a) \|\langle p(\cdot|s, a), \phi(\cdot) \rangle - \langle \hat{p}_d(\cdot|s, a), \phi(\cdot) \rangle\|^2 ds da \\ & \leq 2^{22} \|\boldsymbol{\Sigma}^{-1}\|_\sigma \frac{rlm}{\max\{r, l, m\}} \left(1 + \frac{2\bar{\lambda} \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2}{\sigma^2}\right) \frac{\bar{\lambda} (\log \frac{2t_{mix}}{\delta} + d_S + d_A) (\kappa + \bar{\mu} \|\boldsymbol{\Sigma}^{-1}\|_\sigma^2) (\log \frac{n}{t_{mix}})^2}{n/t_{mix}} + 4L^*. \end{aligned}$$

Then we get the desired result.

Proof of Theorem 7 For each $k \in [n_s], l \in [n_a]$, define

$$\delta_{kl}^2 = \min_{(i,j) \neq (k,l)} \|q_{kl}(\cdot) - q_{ij}(\cdot)\|_{\mathcal{H}_S}^2 = \min_{(i,j) \neq (k,l)} \|\mathbf{z}_{kl} - \mathbf{z}_{ij}\|^2, \quad \bar{\Psi}(s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \hat{\mathbf{z}}_{ij} \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j}.$$

where we still use the notation $\mathbf{z}_{ij} = \mathbf{U}_3^\top \langle q_{ij}(\cdot), \phi(\cdot) \rangle$. For an arbitrary orthogonal matrix $\mathbf{O} \in \mathbb{R}^{m \times m}$, let

$$S_{kl} = \{(s, a) \in A_k \times B_l \mid \|\mathbf{O} \bar{\Psi}(s, a) - \mathbf{z}_{kl}\| \geq \frac{\delta_{kl}}{2}\},$$

then we have

$$\begin{aligned} \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} (\xi \times \eta)(S_{kl}) \delta_{kl}^2 & \leq 4 \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \int_{S_{kl}} \xi(s)\eta(a) \|\mathbf{O} \bar{\Psi}(s, a) - \mathbf{z}_{kl}\|^2 ds da \\ & \leq 4 \int \xi(s)\eta(a) \left\| \mathbf{O} \bar{\Psi}(s, a) - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{s \in A_k} \mathbf{1}_{a \in B_l} \right\|^2 ds da \\ & = 4 \left\| \mathbf{O} \bar{\Psi}(\cdot) - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{A_k \times B_l}(\cdot) \right\|_{L^2(\xi \times \eta)}^2. \end{aligned}$$

Note that

$$\begin{aligned} & \left\| \mathbf{O} \bar{\Psi}(\cdot) - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{A_k \times B_l}(\cdot) \right\|_{L^2(\xi \times \eta)}^2 \\ & \leq \left\| \mathbf{O} \bar{\Psi}(\cdot) - \hat{\mathbf{C}} \times_1 \hat{\mathbf{f}}^\top(\cdot) \times_2 \hat{\mathbf{g}}^\top(\cdot) \times_3 \mathbf{O} \right\|_{L^2(\xi \times \eta)} \\ & \quad + \left\| \hat{\mathbf{C}} \times_1 \hat{\mathbf{f}}(\cdot)^\top \times_2 \hat{\mathbf{g}}(\cdot)^\top \times_3 \mathbf{O} - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{A_k \times B_l}(\cdot) \right\|_{L^2(\xi \times \eta)} \\ & = \left\| \bar{\Psi}(\cdot) - \hat{\mathbf{C}} \times_1 \hat{\mathbf{f}}(\cdot)^\top \times_2 \hat{\mathbf{g}}(\cdot)^\top \right\|_{L^2(\xi \times \eta)} \\ & \quad + \left\| \hat{\mathbf{C}} \times_1 \hat{\mathbf{f}}(\cdot)^\top \times_2 \hat{\mathbf{g}}(\cdot)^\top \times_3 \mathbf{O} - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{A_k \times B_l}(\cdot) \right\|_{L^2(\xi \times \eta)}. \end{aligned}$$

We use the fact that $\hat{A}_i, \hat{B}_j, \hat{z}_{ij}$ is the solution to (4), and derive

$$\begin{aligned}
 & \|\bar{\Psi}(\cdot) - \hat{\mathbf{C}} \times_1 \hat{f}(\cdot)^\top \times_2 \hat{g}(\cdot)^\top\|_{L^2(\xi \times \eta)} \\
 &= \left\| \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \hat{z}_{ij} \mathbf{1}_{\hat{A}_i \times \hat{B}_j}(\cdot) - \hat{\mathbf{C}} \times_1 \hat{f}(\cdot)^\top \times_2 \hat{g}(\cdot)^\top \right\|_{L^2(\xi \times \eta)} \\
 &\leq \left\| \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{O}^\top \mathbf{z}_{ij} \mathbf{1}_{A_i \times B_j}(\cdot) - \hat{\mathbf{C}} \times_1 \hat{f}(\cdot)^\top \times_2 \hat{g}(\cdot)^\top \right\|_{L^2(\xi \times \eta)} \\
 &= \|\mathbf{C} \times_1 f(\cdot)^\top \times_2 g(\cdot)^\top \times_3 \mathbf{O}^\top - \hat{\mathbf{C}} \times_1 \hat{f}(\cdot)^\top \times_2 \hat{g}(\cdot)^\top\|_{L^2(\xi \times \eta)} \\
 &= \|\mathbf{P} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 (\mathbf{U}_3 \mathbf{O})^\top - \hat{\mathbf{P}} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 \hat{\mathbf{U}}_3^\top\|_{L^2(\xi \times \eta)}.
 \end{aligned}$$

Here we use the fact that

$$\mathbf{C} \times_1 f(s)^\top \times_2 g(a)^\top = \mathbf{U}_3^\top \langle p(\cdot | s, a), \phi(\cdot) \rangle = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbf{z}_{ij} \mathbf{1}_{A_i \times B_j}(s, a)$$

because of Assumption 3. Therefore, we have

$$\begin{aligned}
 & \left\| \mathbf{O} \bar{\Psi}(\cdot) - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbf{z}_{kl} \mathbf{1}_{A_k \times B_l}(\cdot) \right\|_{L^2(\xi \times \eta)} \\
 &\leq 2 \|\mathbf{P} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 (\mathbf{U}_3 \mathbf{O})^\top - \hat{\mathbf{P}} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 \hat{\mathbf{U}}_3^\top\|_{L^2(\xi \times \eta)}
 \end{aligned}$$

and

$$\begin{aligned}
 & \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} (\xi \times \eta) (S_{kl}) \delta_{kl}^2 \\
 &\leq 16 \|\mathbf{P} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 (\mathbf{U}_3 \mathbf{O})^\top - \hat{\mathbf{P}} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 \hat{\mathbf{U}}_3^\top\|_{L^2(\xi \times \eta)}^2 \\
 &\leq 32 \|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 \hat{\mathbf{U}}_3^\top\|_{L^2(\xi \times \eta)}^2 \\
 &\quad + 32 \|\mathbf{P} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 (\mathbf{U}_3 \mathbf{O} - \hat{\mathbf{U}}_3)^\top\|_{L^2(\xi \times \eta)}^2 \\
 &\leq 32 \|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top\|_{L^2(\xi \times \eta)}^2 \\
 &\quad + 32 \|\mathbf{P} \times_1 \phi(\cdot)^\top \times_2 \psi(\cdot)^\top \times_3 (\mathbf{U}_3 \mathbf{O} - \hat{\mathbf{U}}_3)^\top\|_{L^2(\xi \times \eta)}^2 \\
 &= 32 \|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \Sigma^{\frac{1}{2}}\|_F^2 + 32 \|\mathbf{P} \times_1 \Sigma^{\frac{1}{2}} \times_3 (\mathbf{U}_3 \mathbf{O} - \hat{\mathbf{U}}_3)^\top\|_F^2.
 \end{aligned}$$

Again by Part 3, Lemma 1 in Cai and Zhang (2018), we can choose \mathbf{O} such that

$$\|\mathbf{U}_3 \mathbf{O}^\top - \hat{\mathbf{U}}_3\| \leq \sqrt{2} \|\sin \Theta(\mathbf{U}_3, \hat{\mathbf{U}}_3)\|.$$

With a similar argument in the proof of Theorem 6, we know when

$$\frac{n/t_{mix}}{(\log(n/t_{mix}))^2} \geq 1024 \left(\|\Sigma^{-1}\|_{\sigma \bar{\mu}}^2 + \frac{K_{max}^2}{\bar{\mu}} + \frac{\kappa K_{max}^3}{\bar{\lambda}} \right) \left(\log \frac{t_{mix}}{\delta} + 8(d_S + d_A) \right),$$

with probability at least $1 - \delta$ we have

$$\begin{aligned} & 32\|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}}\|_F^2 + 32\|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\mathbf{U}_3 \mathbf{O}^\top - \hat{\mathbf{U}}_3)^\top\|_F^2 \\ & \leq 2^{24} \|\boldsymbol{\Sigma}^{-1}\|_\sigma \frac{rlm}{\max\{r, l, m\}} \left(1 + \frac{2\bar{\lambda}\|\boldsymbol{\Sigma}^{-1}\|_\sigma^2}{\sigma^2}\right) \frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu}\|\boldsymbol{\Sigma}^{-1}\|_\sigma^2)(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}. \end{aligned}$$

Now we choose sufficiently large n such that

$$32\|(\mathbf{P} - \hat{\mathbf{P}}) \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}}\|_F^2 + 32\|\mathbf{P} \times_1 \boldsymbol{\Sigma}^{\frac{1}{2}} \times_3 (\mathbf{U}_3 \mathbf{O}^\top - \hat{\mathbf{U}}_3)^\top\|_F^2 < \Delta_1^2.$$

Then for any $1 \leq k \leq n_s, 1 \leq l \leq n_a$ we always have

$$(\xi \times \eta)(S_{kl})\delta_{kl}^2 \leq \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} (\xi \times \eta)(S_{kl})\delta_{kl}^2 < \Delta_1^2 \leq (\xi \times \eta)(A_k \times B_l)\delta_{kl}^2,$$

i.e., $(\xi \times \eta)(S_{kl}) \leq (\xi \times \eta)(A_k \times B_l)$, and $(\xi \times \eta)((A_k \times B_l) \setminus S_{kl}) \neq \emptyset$. Now, for any $(i, j) \neq (k, l)$, we are going to show that $\forall (s, a) \in (A_i \times B_j) \setminus S_{ij}, (s', a') \in (A_k \times B_l) \setminus S_{kl}$, we always have $\bar{\Psi}(s, a) \neq \bar{\Psi}(s', a')$. Suppose the claim is not true, then we can find the corresponding $(s, a), (s', a')$, such that $\bar{\Psi}(s, a) = \bar{\Psi}(s', a')$. However, then we have

$$\max\{\delta_{ij}, \delta_{kl}\} \leq \|\mathbf{z}_{ij} - \mathbf{z}_{kl}\| \leq \|\mathbf{z}_{ij} - \mathbf{O}\bar{\Psi}(s, a)\| + \|\mathbf{O}\bar{\Psi}(s', a') - \mathbf{z}_{kl}\| \leq \frac{\delta_{ij}}{2} + \frac{\delta_{kl}}{2},$$

which leads to a contradiction.

Furthermore, notice that for any $(s, a), (s', a') \in (A_i \times B_j) \setminus S_{ij}$, we always have $\bar{\Psi}(s, a) = \bar{\Psi}(s', a')$, otherwise $\bar{\Psi}(s, a)$ will take more than $n_s n_a$ values.

The above two claims show that we can find two one-to-one mappings $\sigma_1 : [n_s] \rightarrow [n_s], \sigma_2 : [n_a] \rightarrow [n_a]$, such that for any $1 \leq k \leq n_s, 1 \leq l \leq n_a$, we have

$$(A_k \times B_l) \setminus S_{kl} \subseteq \hat{A}_{\sigma_1(k)} \times \hat{B}_{\sigma_2(l)}.$$

Without loss of generality we can assume that $\sigma_1(k) = k, \sigma_2(l) = l$, which is always possible after we rearrange the indexes of \hat{A}_k, \hat{B}_l . Then we have

$$(A_k \times B_l) \setminus S_{kl} \subseteq (\hat{A}_k \times \hat{B}_l),$$

which implies

$$(A_k \times B_l) \setminus (\hat{A}_k \times \hat{B}_l) \subseteq S_{kl}.$$

Therefore, for n sufficiently large, we have

$$\begin{aligned} & \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \\ & \leq \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)(S_{ij})\delta_{ij}^2}{\Delta_1^2} \\ & \leq 2^{24} \frac{1}{\Delta_1^2} \|\boldsymbol{\Sigma}^{-1}\|_\sigma \frac{rlm}{\max\{r, l, m\}} \left(1 + \frac{2\bar{\lambda}\|\boldsymbol{\Sigma}^{-1}\|_\sigma^2}{\sigma^2}\right) \frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu}\|\boldsymbol{\Sigma}^{-1}\|_\sigma^2)(\log \frac{n}{t_{mix}})^2}{n/t_{mix}}, \end{aligned}$$

which has finished the proof.

Proof of Theorem 8 According to the result of Theorem 7, we can find some mapping σ_1, σ_2 , such that

$$\begin{aligned} & \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_{\sigma_1(i)} \times \hat{B}_{\sigma_2(j)}))}{(\xi \times \eta)(A_i \times B_j)} \\ & \leq 2^{24} \frac{1}{\Delta_1^2} \|\Sigma^{-1}\|_\sigma \frac{rlm}{\max\{r, l, m\}} \left(1 + \frac{2\bar{\lambda} \|\Sigma^{-1}\|_\sigma^2}{\sigma^2}\right) \frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2)}{(n/t_{mix})(\log \frac{n}{t_{mix}})^{-2}}. \end{aligned}$$

Without loss of generality we assume that $\sigma_1(k) = k, \sigma_2(l) = l, \forall k \in [n_s], l \in [n_a]$. In this case, the result changes into

$$\begin{aligned} & \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \\ & \leq 2^{24} \frac{1}{\Delta_1^2} \|\Sigma^{-1}\|_\sigma \frac{rlm}{\max\{r, l, m\}} \left(1 + \frac{2\bar{\lambda} \|\Sigma^{-1}\|_\sigma^2}{\sigma^2}\right) \frac{\bar{\lambda}(\log \frac{2t_{mix}}{\delta} + d_S + d_A)(\kappa + \bar{\mu} \|\Sigma^{-1}\|_\sigma^2)}{(n/t_{mix})(\log \frac{n}{t_{mix}})^{-2}} =: \varepsilon. \end{aligned}$$

Based on that, we want to further bound $\sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)}$. Notice that

$$\begin{aligned} \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} (\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j)) &= \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} (\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j)) \\ &= \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} (\xi \times \eta)(A_i \times B_j) \\ &\leq \bar{c} \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \leq \bar{c}\varepsilon \end{aligned}$$

and

$$(\xi \times \eta)(A_i \times B_j) \leq (\xi \times \eta)(\hat{A}_i \times \hat{B}_j) + (\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j)) \leq (\xi \times \eta)(\hat{A}_i \times \hat{B}_j) + \bar{c}\varepsilon.$$

So we have

$$\begin{aligned} \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} &\leq \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(A_i \times B_j) - \bar{c}\varepsilon} \\ &\leq \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{\underline{c} - \bar{c}\varepsilon} \\ &= \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{\underline{c} - \bar{c}\varepsilon} \\ &\leq \frac{\bar{c}}{\underline{c} - \bar{c}\varepsilon} \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \\ &\leq \frac{\bar{c}}{\underline{c} - \bar{c}\varepsilon} M(\{\hat{A}_i\}_{i=1}^{n_s}, \{\hat{B}_j\}_{j=1}^{n_a}). \end{aligned}$$

Therefore, when n is sufficiently large such that $\bar{c}\varepsilon \leq \frac{1}{2}\underline{c}$, we have

$$\sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \leq 2 \frac{\bar{c}}{\underline{c}} M(\{\hat{A}_i\}_{i=1}^{n_s}, \{\hat{B}_j\}_{j=1}^{n_a}).$$

Now we are ready to prove the main result. By definition, the transition dynamic can be written as

$$\hat{p}(s'|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \sum_{k=1}^{n_s} \frac{1}{\xi(\hat{A}_k)} \hat{q}(k|i, j) \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \mathbf{1}_{s' \in \hat{A}_k}.$$

Now, without loss of generality, we can assume that the groundtruth transition dynamic also takes the form

$$p(s'|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \sum_{k=1}^{n_s} \frac{1}{\xi(A_k)} q(k|i, j) \mathbf{1}_{s \in A_i} \mathbf{1}_{a \in B_j} \mathbf{1}_{s' \in A_k}.$$

which is due to the fact that the policy and reward are the same in the same block, so for a general transition dynamic satisfying assumption 3, we can simply set $q(k|i, j) = \int_{A_k} q_{ij}^*(s') ds$, then the transition dynamic provided by the above formula will lead to exactly the same H -step value.

Because the infimum is taken over all possible \hat{q} , we can choose \hat{q} exactly to be q , and get

$$\hat{p}^\pi(s', a'|s, a) = \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \frac{q^\pi(k, l|i, j) \mathbf{1}_{s \in \hat{A}_i} \mathbf{1}_{a \in \hat{B}_j} \mathbf{1}_{s' \in \hat{A}_k} \mathbf{1}_{a' \in \hat{B}_l}}{\xi(\hat{A}_k) \eta(\hat{B}_l)},$$

where

$$q^\pi(k, l|i, j) = \pi(l|k) q(k|i, j).$$

For each single path $(s_1, a_1), (s_2, a_2), \dots, (s_H, a_H)$, consider a series of auxiliary random variables T_1, T_2, \dots, T_H , which are defined inductively by

$$T_0 = 0, T_h = \max \left\{ T_{h-1}, \mathbf{1}_{(s_h, a_h) \notin \bigcup_{i=1}^{n_s} (A_i \cap \hat{A}_i) \times \bigcup_{j=1}^{n_a} (B_j \cap \hat{B}_j)} \right\}.$$

Then for each reward function $r_h, h \in [H]$ and any initial distribution μ , we have

$$\begin{aligned} \mathbb{E}_\mu^\pi[r_h(s_h, a_h)] &= \mathbb{E}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=0}] + \mathbb{E}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=1}], \\ \hat{\mathbb{E}}_\mu^\pi[r_h(s_h, a_h)] &= \hat{\mathbb{E}}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=0}] + \hat{\mathbb{E}}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=1}] \end{aligned}$$

and

$$\begin{aligned} &|\mathbb{E}_\mu^\pi[r_h(s_h, a_h)] - \hat{\mathbb{E}}_\mu^\pi[r_h(s_h, a_h)]| \\ &\leq |\mathbb{E}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=0}] - \hat{\mathbb{E}}_\mu^\pi[r_h(s_h, a_h) \mathbf{1}_{T_h=0}]| + (\mathbb{P}_\mu^\pi(T_h = 1) + \hat{\mathbb{P}}_\mu^\pi(T_h = 1)). \end{aligned}$$

For the second term, note that

$$\begin{aligned}
 & \mathbb{P}_\mu^\pi(T_{h+1} = 1) \\
 &= \mathbb{P}_\mu^\pi(T_{h+1} = 1|T_h = 1)\mathbb{P}_\mu^\pi(T_h = 1) + \mathbb{P}_\mu^\pi(T_{h+1} = 1|T_h = 0)\mathbb{P}_\mu^\pi(T_h = 0) \\
 &= \mathbb{P}_\mu^\pi(T_h = 1) + \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbb{P}_\mu^\pi(s_{h+1} \in A_i, a_{h+1} \in B_j, T_{h+1} = 1|T_h = 0) \\
 &= \mathbb{P}_\mu^\pi(T_h = 1) + \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \mathbb{P}_\mu^\pi((s_{h+1}, a_{h+1}) \in (A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j) | s_{h+1} \in A_i, a_{h+1} \in B_j) \\
 &\quad \cdot \mathbb{P}_\mu^\pi(s_{h+1} \in A_i, a_{h+1} \in B_j | T_h = 0) \\
 &= \mathbb{P}_\mu^\pi(T_h = 1) + \sum_{i=1}^{n_s} \sum_{j=1}^{n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \mathbb{P}_\mu^\pi(s_{h+1} \in A_i, a_{h+1} \in B_j | T_h = 0) \\
 &\leq \mathbb{P}_\mu^\pi(T_h = 1) + \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)},
 \end{aligned}$$

which implies

$$\mathbb{P}_\mu^\pi(T_h = 1) \leq H \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)}.$$

Similarly, we have

$$\hat{\mathbb{P}}_\mu^\pi(T_h = 1) \leq H \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)}.$$

For the first term, we consider a discrete MDP over $(A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)$, $1 \leq i \leq n_s$, $1 \leq j \leq n_a$ plus an absorbing state C_0 , and two corresponding transition probabilities

$$\begin{aligned}
 \mathbb{P}((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l) | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)) &= q^\pi(k, l | i, j) \frac{(\xi \times \eta)((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l))}{(\xi \times \eta)(A_k \times B_l)}, \\
 \mathbb{P}(C_0 | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)) &= 1 - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \mathbb{P}((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l) | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)) \\
 \hat{\mathbb{P}}((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l) | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)) &= q^\pi(k, l | i, j) \frac{(\xi \times \eta)((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l))}{(\xi \times \eta)(\hat{A}_k \times \hat{B}_l)}, \\
 \hat{\mathbb{P}}(C_0 | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)) &= 1 - \sum_{k=1}^{n_s} \sum_{l=1}^{n_a} \hat{\mathbb{P}}((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l) | (A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)).
 \end{aligned}$$

We use $\mathbf{P}, \hat{\mathbf{P}}$ to denote the two corresponding transition matrices restricted on $(A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)$, $1 \leq i \leq n_s$, $1 \leq j \leq n_a$ (i.e., the entries related with C_0 are not included, so $\mathbf{P}, \hat{\mathbf{P}}$ are sub-matrices of two stochastic matrices). For two different distributions $\mu, \tilde{\mu}$ over

$(A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j)$, $1 \leq i \leq n_s$, $1 \leq j \leq n_a$, we have

$$\begin{aligned} \|\mu \mathbf{P} - \tilde{\mu} \hat{\mathbf{P}}\|_1 &\leq \|\mu \mathbf{P} - \mu \hat{\mathbf{P}}\|_1 + \|\mu \hat{\mathbf{P}} - \tilde{\mu} \hat{\mathbf{P}}\|_1 \\ &\leq \|\mu\|_1 \max_{(i,j),(k,l)} |(\mathbf{P} - \hat{\mathbf{P}})_{(i,j),(k,l)}| + \|\mu - \tilde{\mu}\|_1 \max_{(i,j),(k,l)} |\hat{\mathbf{P}}_{(i,j),(k,l)}| \\ &\leq \max_{(i,j),(k,l)} |(\mathbf{P} - \hat{\mathbf{P}})_{(i,j),(k,l)}| + \|\mu - \tilde{\mu}\|_1. \end{aligned}$$

Note that

$$(\mathbf{P} - \hat{\mathbf{P}})_{(i,j),(k,l)} = \left(\frac{(\xi \times \eta)((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l))}{(\xi \times \eta)(A_k \times B_l)} - \frac{(\xi \times \eta)((A_k \cap \hat{A}_k) \times (B_l \cap \hat{B}_l))}{(\xi \times \eta)(\hat{A}_k \times \hat{B}_l)} \right) q^\pi(k, l | i, j),$$

so we have

$$\begin{aligned} &\max_{(i,j),(k,l)} |(\mathbf{P} - \hat{\mathbf{P}})_{(i,j),(k,l)}| \\ &\leq \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \left| \frac{(\xi \times \eta)((A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} - \frac{(\xi \times \eta)((A_i \cap \hat{A}_i) \times (B_j \cap \hat{B}_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \right| \\ &\leq \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} + \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)}. \end{aligned}$$

So we get

$$\begin{aligned} &\|\mu \mathbf{P}^h - \mu \hat{\mathbf{P}}^h\|_1 \\ &\leq \max_{\substack{1 \leq i \leq n_s \\ 1 \leq j \leq n_a}} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} + \max_{\substack{1 \leq i \leq n_s \\ 1 \leq j \leq n_a}} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \\ &\quad + \|\mu \mathbf{P}^{h-1} - \mu \hat{\mathbf{P}}^{h-1}\|_1 \\ &\leq \dots \\ &\leq H \left(\max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} \right. \\ &\quad \left. + \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \right), \end{aligned}$$

Note that the transition dynamic of the original state-action space can be embedded into the discrete MDP mentioned above according to the following mapping: Denote \tilde{s}_h as the state of the discrete MDP at step h , then we let $\tilde{s}_h = C_0$ if $T_h = 1$, and let $\tilde{s}_h = (A_i \times B_j) \cap (\hat{A}_i \times \hat{B}_j)$ if $T_h = 0$ and $(s_h, a_h) \in (A_i \times B_j) \cap (\hat{A}_i \times \hat{B}_j)$. Because the reward is only related with blocks and is bounded between $[0, 1]$, so we have

$$\begin{aligned} &|\mathbb{E}_\mu^\pi[r_h(S_h, A_h) \mathbf{1}_{T_h=0}] - \hat{\mathbb{E}}_\mu^\pi[r_h(S_h, A_h) \mathbf{1}_{T_h=0}]| \\ &\leq \|\mu \mathbf{P}^h - \mu \hat{\mathbf{P}}^h\|_1 \\ &\leq H \left(\max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} + \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \right). \end{aligned}$$

In summary, we get

$$\begin{aligned}
 & |\mathbb{E}^\pi[r_h(S_h, A_h)] - \hat{\mathbb{E}}^\pi[r_h(S_h, A_h)]| \\
 & \leq 2H \left(\max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((A_i \times B_j) \setminus (\hat{A}_i \times \hat{B}_j))}{(\xi \times \eta)(A_i \times B_j)} + \max_{1 \leq i \leq n_s, 1 \leq j \leq n_a} \frac{(\xi \times \eta)((\hat{A}_i \times \hat{B}_j) \setminus (A_i \times B_j))}{(\xi \times \eta)(\hat{A}_i \times \hat{B}_j)} \right) \\
 & \leq 4H \frac{\bar{c}}{\underline{c}} M(\{\hat{A}_i\}_{i=1}^{n_s}, \{\hat{B}_j\}_{j=1}^{n_a})
 \end{aligned}$$

and

$$\begin{aligned}
 & \left| \sum_{h=1}^H \mathbb{E}^\pi[r_h(S_h, A_h)] - \sum_{h=1}^H \hat{\mathbb{E}}^\pi[r_h(S_h, A_h)] \right| \\
 & \leq \sum_{h=1}^H |\mathbb{E}^\pi[r_h(S_h, A_h)] - \hat{\mathbb{E}}^\pi[r_h(S_h, A_h)]| \\
 & \leq 4H^2 \frac{\bar{c}}{\underline{c}} M(\{\hat{A}_i\}_{i=1}^{n_s}, \{\hat{B}_j\}_{j=1}^{n_a})
 \end{aligned}$$

which finishes the proof.