

---

# Learning Mixtures of Linear Dynamical Systems

---

Yanxi Chen<sup>1</sup> H. Vincent Poor<sup>1</sup>

## Abstract

We study the problem of learning a mixture of multiple linear dynamical systems (LDSs) from unlabeled short sample trajectories, each generated by one of the LDS models. Despite the wide applicability of mixture models for time-series data, learning algorithms that come with end-to-end performance guarantees are largely absent from existing literature. There are multiple sources of technical challenges, including but not limited to (1) the presence of latent variables (i.e. the unknown labels of trajectories); (2) the possibility that the sample trajectories might have lengths much smaller than the dimension  $d$  of the LDS models; and (3) the complicated temporal dependence inherent to time-series data. To tackle these challenges, we develop a two-stage meta-algorithm, which is guaranteed to efficiently recover each ground-truth LDS model up to error  $\tilde{O}(\sqrt{d/T})$ , where  $T$  is the total sample size. We validate our theoretical studies with numerical experiments, confirming the efficacy of the proposed algorithm.

## 1. Introduction

Imagine that we are asked to learn multiple linear dynamical systems (LDSs) from a mixture of unlabeled sample trajectories — namely, each sample trajectory is generated by one of the LDSs of interest, but we have no idea which system it is. To set the stage and facilitate discussion, recall that in a classical LDS, one might observe a sample trajectory  $\{\mathbf{x}_t\}_{0 \leq t \leq T}$  generated by an LDS obeying  $\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t$ , where  $\mathbf{A} \in \mathbb{R}^{d \times d}$  determines the system dynamics in the noiseless case, and  $\{\mathbf{w}_t\}_{t \geq 0}$  denote independent zero-mean noise vectors with covariance  $\text{cov}(\mathbf{w}_t) = \mathbf{W} \succ \mathbf{0}$ . The mixed LDSs setting considered herein extends classical LDSs by allowing for mixed

measurements as described below; see Figure 1 for a visualization of the scenario.

- *Multiple linear systems.* Suppose that there are  $K$  different LDSs as represented by  $\{(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})\}_{1 \leq k \leq K}$ , where  $\mathbf{A}^{(k)} \in \mathbb{R}^{d \times d}$  and  $\mathbf{W}^{(k)} \in \mathbb{R}^{d \times d}$  represent the state transition matrix and noise covariance matrix of the  $k$ -th LDS, respectively. Here and throughout, we shall refer to  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  as the system matrix for the  $k$ -th LDS. We only assume that  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  for any  $k \neq \ell$ , whereas  $\mathbf{A}^{(k)} = \mathbf{A}^{(\ell)}$  or  $\mathbf{W}^{(k)} = \mathbf{W}^{(\ell)}$  is allowed.
- *Mixed sample trajectories.* We collect a total number of  $M$  unlabeled trajectories from these LDSs. More specifically, the  $m$ -th sample trajectory is drawn from one of the LDSs in the following manner: set  $(\mathbf{A}, \mathbf{W})$  to be  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  for some  $1 \leq k \leq K$ , and generate a trajectory of length  $T_m$  obeying

$$\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t, \quad \text{where the } \mathbf{w}_t \text{'s are i.i.d.,}$$
$$\mathbb{E}[\mathbf{w}_t] = \mathbf{0}, \quad \text{cov}(\mathbf{w}_t) = \mathbf{W} \succ \mathbf{0}.$$

Note, however, that the label  $k$  associated with each sample trajectory is a latent variable not revealed to us, resulting in a mixture of unlabeled trajectories. The current paper focuses on the case where the length of each trajectory is somewhat short, making it infeasible to estimate the system matrix from a single trajectory.

- *Goal.* The aim is to jointly learn the system matrices  $\{(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})\}_{1 \leq k \leq K}$  from the mixture of sample trajectories. In particular, we seek to accomplish this task in a sample-efficient manner, where the total sample size is defined to be the aggregate trajectory length  $\sum_{m=1}^M T_m$ .

**Motivations.** The mixed LDSs setting described above is motivated by many real-world scenarios where a single time-series model is insufficient to capture the complex and heterogeneous patterns in temporal data. For instance, in psychology (Bulteel et al., 2016), researchers collect multiple time-series trajectories (e.g. depression-related symptoms over a period of time) from different patients. Fitting this data with multi-modal LDSs (instead of a single model) not only achieves better fitting performance, but also helps

---

<sup>1</sup>Department of Electrical and Computer Engineering, Princeton University, Princeton, NJ 08544, USA. Correspondence to: Yanxi Chen <yanxic@princeton.edu>.

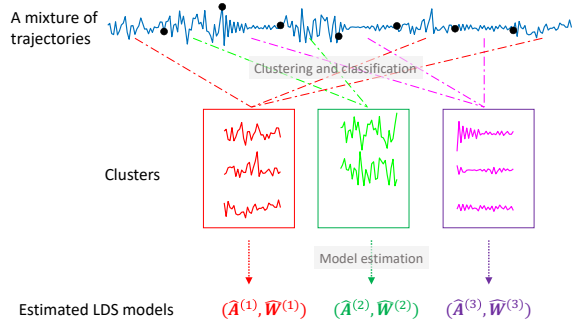


Figure 1. A high-level visualization of the mixed LDSs formulation, and the algorithmic idea of combining clustering, classification, and model estimation. Here, we consider the special case where the multiple short trajectories come from the segments of a single continuous trajectory. The black dots within the continuous trajectory represent the time steps when the latent variable (i.e. the unknown label) changes.

identify subgroups of the persons, which further inspires interpretations of the results and tailored treatments for patients from different subgroups. Another example concerns an automobile sensor dataset, which consists of a single continuous (but possibly time-varying) trajectory of measurements from the sensors (Hallac et al., 2017). Through segmentation of the trajectory, clustering of the short pieces, and learning within each cluster, one can discover, and obtain meaningful interpretations for, a small number of key driving modes (such as “driving straight”, “slowing down”, “turning”). See Section 4 for a longer list of applications.

**Challenges.** While there is no shortage of potential applications, a mixture of LDSs is far more challenging to learn compared to the classical setting with a single LDS. In particular, the presence of the *latent variables*, i.e. the unknown labels of the sample trajectories, significantly complicates matters. One straightforward idea is to learn a coarse model for each trajectory, followed by proper clustering of these coarse models (to be utilized for refined model estimation); however, this idea becomes infeasible in the high-dimensional setting unless all trajectories are sufficiently long. Another popular approach is to alternate between model estimation and clustering of trajectories, based on, say, the expectation-maximization (EM) algorithm; unfortunately, there is no theoretical support for such EM-type algorithms, and we cannot preclude the possibilities that the algorithms get stuck at undesired local optima. The lack of theoretical guarantees in prior literature motivates us to come up with algorithms that enjoy provable performance guarantees. Finally, the present paper is also inspired by the recent progress in *meta-learning for mixed linear regression* (Kong et al., 2020b;a), where the goal is to learn multiple linear models from a mixture of independent samples; note, however, that the temporal dependence underlying time-

series data in our case poses substantial challenges and calls for the development of new algorithmic ideas.

**Main contributions.** In this work, we take an important step towards guaranteed learning of mixed LDSs, focusing on algorithm design that comes with end-to-end theoretical guarantees. In particular, we propose a two-stage meta-algorithm to tackle the challenge of mixed LDSs:

1. *Coarse estimation*: perform a coarse-level clustering of the unlabeled sample trajectories (assisted by dimension reduction), and compute initial model estimation for each cluster;
2. *Refinement*: classify additional trajectories (and add each of them to the corresponding cluster) based on the above coarse model estimates, followed by refined model estimation with the updated clusters.

This two-stage meta approach, as well as the specific methods for individual steps, will be elucidated in Section 2.

Encouragingly, assuming that the noise vectors  $\{w_{m,t}\}$  are independent Gaussian random vectors, the proposed two-stage algorithm is not only computationally efficient, but also guaranteed to succeed in the presence of a polynomial sample size. Informally, our algorithm achieves exact clustering/classification of the sample trajectories as well as faithful model estimation, under the following conditions (with a focus on the dependency on the dimension  $d$ ): (1) each short trajectory length  $T_m$  is allowed to be much smaller than  $d$ ; (2) the total trajectory length of each stage is linear in  $d$  (up to logarithmic factors); (3) to achieve a final model estimation error  $\epsilon \rightarrow 0$  (in the spectral norm), it suffices to have a total trajectory length of order  $d/\epsilon^2$  for each LDS model. See Section 3 (in particular, Corollary 3.7) for the precise statements of our main results, which will also be validated numerically.

It is worth noting that, although we focus on mixed LDSs for concreteness, we will make clear that the proposed modular algorithm is fairly flexible and can be adapted to learning mixtures of other time-series models, as long as certain technical conditions are satisfied; see Remark 2.3 at the end of Section 2 for a detailed discussion.

**Notation.** Throughout this paper, vectors and matrices are represented by boldface letters. For a vector  $\mathbf{x} \in \mathbb{R}^d$ , we define  $(\mathbf{x})_i \in \mathbb{R}$  as the  $i$ -th entry of  $\mathbf{x}$ , and  $\|\mathbf{x}\|_2$  as its  $\ell_2$  norm; for a matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$ , we define  $(\mathbf{X})_i \in \mathbb{R}^n$  as the transpose of the  $i$ -th row of  $\mathbf{X}$ , and  $\|\mathbf{X}\|$  (resp.  $\|\mathbf{X}\|_F$ ) as its spectral (resp. Frobenius) norm. For a symmetric matrix  $\mathbf{X} \in \mathbb{R}^{d \times d}$ , we denote its maximal (resp. minimal) eigenvalue as  $\lambda_{\max}(\mathbf{X})$  (resp.  $\lambda_{\min}(\mathbf{X})$ ); if in addition  $\mathbf{X}$  is positive definite, we denote its condition number as  $\kappa(\mathbf{X}) :=$

**Algorithm 1** A two-stage algorithm for mixed LDSs

- 1: **Input:**  $M$  short trajectories  $\{\mathbf{X}_m\}_{1 \leq m \leq M}$  (where  $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}$ ); parameters  $\tau, G$ .
- 2: // Stage 1: coarse estimation.
- 3: Run Algorithm 2 with  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{subspace}}}$  to obtain subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$ .
- 4: Run Algorithm 3 with  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{clustering}}}$ ,  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$ ,  $\tau, G$ , to obtain clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .
- 5: Run Algorithm 4 with  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  to obtain coarse models  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ .
- 6: // Stage 2: refinement.
- 7: Run Algorithm 5 with  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{classification}}}$ ,  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ ,  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ , to update clusters  $\{\mathcal{C}_k\}$ .
- 8: Run Algorithm 4 with  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  to obtain refined models  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ .
- 9: **Output:** final model estimation  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$  and clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .

$\lambda_{\max}(\mathbf{X})/\lambda_{\min}(\mathbf{X}) \geq 1$ . If  $\mathbf{A} = [A_{ij}]_{1 \leq i \leq m, 1 \leq j \leq n}$  and  $\mathbf{B} = [B_{ij}]_{1 \leq i \leq m, 1 \leq j \leq n}$  are matrices of the same dimension, we denote by  $\langle \mathbf{A}, \mathbf{B} \rangle := \sum_{i=1}^m \sum_{j=1}^n A_{ij} B_{ij}$  their inner product. Given vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$  where  $n < d$ , let  $\text{span}\{\mathbf{x}_i, 1 \leq i \leq n\} \in \mathbb{R}^{d \times n}$  represent the subspace spanned by these vectors. Let  $\mathbf{I}_d$  be the  $d \times d$  identity matrix.

We always use the superscript “ $(k)$ ” to indicate “the  $k$ -th model”, as in  $\mathbf{A}^{(k)}$  and  $\mathbf{W}^{(k)}$ ; this is to be distinguished from the superscript “ $k$ ” without the parentheses, which simply means the power of  $k$ . For a discrete set  $\Omega$ , we denote by  $|\Omega|$  its cardinality. Define  $\mathbb{1}(\mathcal{E})$  to be the indicator function, which takes value 1 if the event  $\mathcal{E}$  happens, and 0 otherwise. Let  $a_n \lesssim b_n$  indicate that  $a_n \leq C_0 b_n$  for all  $n = 1, 2, \dots$ , where  $C_0 > 0$  is some universal constant; moreover,  $a_n \gtrsim b_n$  is equivalent to  $b_n \lesssim a_n$ , and  $a_n \asymp b_n$  indicates that  $a_n \lesssim b_n$  and  $a_n \gtrsim b_n$  hold simultaneously.

## 2. Algorithms

We propose a two-stage paradigm for solving the mixed LDSs problem, as summarized in Algorithm 1. It consists of several subroutines as described in Algorithms 2–5; due to space limitation, in this section we will only illustrate the key ideas behind these subroutines, with full details deferred to Appendix A. Note that Algorithm 1 is stated in a modular manner, and one might replace certain subroutines by alternative schemes in order to handle different settings and model assumptions.

Let us first introduce some additional notation and assumptions that will be useful for our presentation, without much loss of generality. To begin with, we augment the notation for each sample trajectory with its trajectory index; that is, for each  $1 \leq m \leq M$ , the  $m$ -th trajectory — denoted by  $\mathbf{X}_m := \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}$  — starts with some initial state  $\mathbf{x}_{m,0} \in \mathbb{R}^d$ , and evolves according to the  $k_m$ -th LDS for some *unknown* label  $1 \leq k_m \leq K$  such that

$$\mathbf{x}_{m,t+1} = \mathbf{A}^{(k_m)} \mathbf{x}_{m,t} + \mathbf{w}_{m,t}, \text{ where the } \mathbf{w}_{m,t} \text{'s are i.i.d.,}$$

$$\mathbb{E}[\mathbf{w}_{m,t}] = \mathbf{0}, \text{ cov}(\mathbf{w}_{m,t}) = \mathbf{W}^{(k_m)} \succ \mathbf{0} \quad (1)$$

for all  $0 \leq t \leq T_m - 1$ . Next, we divide the  $M$  sample tra-

jectories  $\{\mathbf{X}_m\}_{1 \leq m \leq M}$  in hand into three disjoint subsets  $\mathcal{M}_{\text{subspace}}, \mathcal{M}_{\text{clustering}}, \mathcal{M}_{\text{classification}}$  satisfying

$$\mathcal{M}_{\text{subspace}} \cup \mathcal{M}_{\text{clustering}} \cup \mathcal{M}_{\text{classification}} = \{1, 2, \dots, M\},$$

where each subset of samples will be employed to perform one subroutine. We assume that all trajectories within each subset have the same length, namely,

$$T_m = \begin{cases} T_{\text{subspace}} & \text{if } m \in \mathcal{M}_{\text{subspace}}, \\ T_{\text{clustering}} & \text{if } m \in \mathcal{M}_{\text{clustering}}, \\ T_{\text{classification}} & \text{if } m \in \mathcal{M}_{\text{classification}}. \end{cases} \quad (2)$$

Finally, we assume that  $K \leq d$ , so that performing subspace estimation in Algorithm 1 will be helpful (otherwise one might simply eliminate this step). The interested readers are referred to Appendix E for discussions of some potential extensions of our algorithms (e.g. adapting to the case where the trajectories within a subset have different lengths, or the case where certain parameters are *a priori* unknown).

### 2.1. Preliminary facts

We first introduce some preliminary background on the autocovariance structures and mixing property of linear dynamical systems, which form the basis of our algorithms for subspace estimation and clustering of trajectories.

**Stationary covariance matrices.** Consider first a single LDS model  $(\mathbf{A}, \mathbf{W})$ , with  $\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t$ . If  $\mathbb{E}[\mathbf{x}_t] = \mathbf{0}$ ,  $\text{cov}(\mathbf{x}_t) = \mathbf{\Gamma}$  and  $\mathbb{E}[\mathbf{w}_t] = \mathbf{0}$ ,  $\text{cov}(\mathbf{w}_t) = \mathbf{W}$ , then it follows that  $\mathbb{E}[\mathbf{x}_{t+1}] = \mathbf{0}$ , and  $\text{cov}(\mathbf{x}_{t+1}) = \mathbf{A} \cdot \text{cov}(\mathbf{x}_t) \cdot \mathbf{A}^\top + \text{cov}(\mathbf{w}_t) = \mathbf{A}\mathbf{\Gamma}\mathbf{A}^\top + \mathbf{W}$ . Under certain assumption on stability, this leads to the *order-0 stationary autocovariance matrix* defined as (Kailath et al., 2000)

$$\mathbf{\Gamma}(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}]$$

$$= \mathbf{A} \cdot \mathbf{\Gamma}(\mathbf{A}, \mathbf{W}) \cdot \mathbf{A}^\top + \mathbf{W} = \sum_{t=0}^{\infty} \mathbf{A}^t \mathbf{W} (\mathbf{A}^t)^\top, \quad (3)$$

and the *order-1 stationary autocovariance matrix*

$$\mathbf{Y}(\mathbf{A}, \mathbf{W}) := \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^\top | \mathbf{A}, \mathbf{W}] = \mathbf{A} \cdot \mathbf{\Gamma}(\mathbf{A}, \mathbf{W}). \quad (4)$$

For the mixed LDSs case (1) with  $K$  models, we abbreviate

$$\mathbf{\Gamma}^{(k)} := \mathbf{\Gamma}(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}), \mathbf{Y}^{(k)} := \mathbf{Y}(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}). \quad (5)$$

**Mixing.** Expanding the LDS recursion  $\mathbf{x}_{t+1} = \mathbf{A}\mathbf{x}_t + \mathbf{w}_t$  with  $\text{cov}(\mathbf{w}_t) = \mathbf{W}$ , we have

$$\mathbf{x}_{t+s} = \mathbf{A}^s \mathbf{x}_t + \sum_{i=1}^s \mathbf{A}^{i-1} \mathbf{w}_{t+s-i}, \quad s = 1, 2, \dots \quad (6)$$

If  $\mathbf{A}$  satisfies certain assumptions regarding stability and if  $s$  is larger than certain mixing time of the LDS, then the first term on the right-hand side of (6) approaches zero, while the second term is independent of the history up to time  $t$ , with covariance close to the stationary autocovariance  $\mathbf{\Gamma}(\mathbf{A}, \mathbf{W})$ . This suggests that, for two samples within the same trajectory that are sufficiently far apart, we can treat them as being (almost) independent of each other; this simple fact shall inspire our algorithmic design and streamline statistical analysis later on. It is noteworthy that the proposed algorithms do not require prior knowledge about the mixing times of the LDS models.

## 2.2. Subspace estimation

Recall that the notation  $(\mathbf{\Gamma}^{(k)})_i \in \mathbb{R}^d$  (resp.  $(\mathbf{Y}^{(k)})_i$ ) represents the transpose of the  $i$ -th row of  $\mathbf{\Gamma}^{(k)}$  (resp.  $\mathbf{Y}^{(k)}$ ) defined in (5). Consider the following subspaces:

$$\begin{aligned} \mathbf{V}_i^* &:= \text{span}\left\{(\mathbf{\Gamma}^{(k)})_i, 1 \leq k \leq K\right\}, \\ \mathbf{U}_i^* &:= \text{span}\left\{(\mathbf{Y}^{(k)})_i, 1 \leq k \leq K\right\}, \quad 1 \leq i \leq d, \end{aligned}$$

each of which has rank at most  $K$ . We develop a spectral method to estimate these subspaces, which will in turn allow for proper dimension reduction in subsequent steps. Here, we only introduce the idea for estimating  $\{\mathbf{U}_i^*\}$ , since the idea for  $\{\mathbf{V}_i^*\}$  is very similar.

We first divide  $\{0, 1, \dots, T_{\text{subspace}}\}$  into four segments of the same size, and denote the 2nd (resp. 4th) segment as  $\Omega_1$  (resp.  $\Omega_2$ ). According to the mixing property of LDS, if  $T_{\text{subspace}}$  is larger than some appropriately defined mixing time, then (1) each sample trajectory in  $\mathcal{M}_{\text{subspace}}$  will mix sufficiently and nearly reach stationarity (when constrained to  $t \in \Omega_1 \cup \Omega_2$ ), (2) the samples in  $\Omega_1$  are nearly independent of those in  $\Omega_2$ . Therefore, defining

$$\mathbf{g}_{m,i,j} := \frac{1}{|\Omega_j|} \sum_{t \in \Omega_j} (\mathbf{x}_{m,t+1})_i \mathbf{x}_{m,t} \quad (7)$$

for  $1 \leq i \leq d$  and  $j \in \{1, 2\}$ , it is easy to check that  $\mathbb{E}[\mathbf{g}_{m,i,1} \mathbf{g}_{m,i,2}^\top] \approx \mathbb{E}[\mathbf{g}_{m,i,1}] \mathbb{E}[\mathbf{g}_{m,i,2}]^\top \approx (\mathbf{Y}^{(k_m)})_i (\mathbf{Y}^{(k_m)})_i^\top$ . Taking the average over  $m \in \mathcal{M}_{\text{subspace}}$ , we have constructed a matrix  $\widehat{\mathbf{G}}_i$  such that

$\mathbb{E}[\widehat{\mathbf{G}}_i] \approx \sum_{k=1}^K p^{(k)} (\mathbf{Y}^{(k)})_i (\mathbf{Y}^{(k)})_i^\top =: \mathbf{G}_i$ , where  $p^{(k)}$  denotes the fraction of sample trajectories generated by the  $k$ -th model. As a consequence, if  $T_{\text{subspace}}$  and  $|\mathcal{M}_{\text{subspace}}|$  are sufficiently large, then one might expect  $\widehat{\mathbf{G}}_i$  to be a good approximation of  $\mathbf{G}_i$ , the latter of which is symmetric, has rank at most  $K$ , and has  $\mathbf{U}_i^*$  as its eigenspace. All this motivates us to compute the rank- $K$  eigenspace of  $\widehat{\mathbf{G}}_i + \widehat{\mathbf{G}}_i^\top$ .

## 2.3. Clustering

In this step, we propose to construct a similarity matrix  $\mathbf{S}$  for the sample trajectories in  $\mathcal{M}_{\text{clustering}}$ , based on their pairwise comparisons; then we can apply any mainstream clustering algorithm (e.g. spectral clustering (Chen et al., 2021b)) to  $\mathbf{S}$ , dividing  $\mathcal{M}_{\text{clustering}}$  into  $K$  disjoint clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  such that the trajectories in each cluster are primarily generated by the same LDS model.

The remaining question is how to perform the pairwise comparisons; we intend to achieve this by comparing the autocovariance matrices associated with the sample trajectories. Note that, even though  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  for  $k \neq \ell$ , it is indeed possible that  $\mathbf{\Gamma}^{(k)} = \mathbf{\Gamma}^{(\ell)}$  or  $\mathbf{Y}^{(k)} = \mathbf{Y}^{(\ell)}$ . Fortunately, the following fact ensures the separation of  $(\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)})$  versus  $(\mathbf{\Gamma}^{(\ell)}, \mathbf{Y}^{(\ell)})$ .

**Fact 2.1.** *If  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$ , then we have either  $\mathbf{\Gamma}^{(k)} \neq \mathbf{\Gamma}^{(\ell)}$  or  $\mathbf{Y}^{(k)} \neq \mathbf{Y}^{(\ell)}$  (or both).*

Now, let us compare the  $m$ -th and  $n$ -th trajectories for some  $m, n \in \mathcal{M}_{\text{clustering}}$ . In order to determine whether they have the same label (namely  $k_m = k_n$ ), we propose to estimate the quantity  $\|\mathbf{\Gamma}^{(k_m)} - \mathbf{\Gamma}^{(k_n)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k_m)} - \mathbf{Y}^{(k_n)}\|_{\mathbb{F}}^2$  using the data samples  $\{\mathbf{x}_{m,t}\}$  and  $\{\mathbf{x}_{n,t}\}$ , which is expected to be small (resp. large) if  $k_m = k_n$  (resp.  $k_m \neq k_n$ ). To do so, let us divide  $\{0, 1, \dots, T_{\text{clustering}}\}$  evenly into four segments, and denote by  $\Omega_1$  (resp.  $\Omega_2$ ) the 2nd (resp. 4th) segment. Recall our earlier definition of the  $\mathbf{g}$  vectors in (7). Assuming sufficient mixing and utilizing near independence between samples from  $\Omega_1$  and those from  $\Omega_2$ , we might resort to the following statistic:

$$\text{stat}_Y := \sum_{i=1}^d \left\langle \mathbf{U}_i^\top (\mathbf{g}_{m,i,1} - \mathbf{g}_{n,i,1}), \mathbf{U}_i^\top (\mathbf{g}_{m,i,2} - \mathbf{g}_{n,i,2}) \right\rangle,$$

whose expectation is given by

$$\begin{aligned} \mathbb{E}[\text{stat}_Y] &\approx \sum_{i=1}^d \left\| \mathbf{U}_i^\top ((\mathbf{Y}^{(k_m)})_i - (\mathbf{Y}^{(k_n)})_i) \right\|_2^2 \\ &\approx \sum_{i=1}^d \left\| (\mathbf{Y}^{(k_m)})_i - (\mathbf{Y}^{(k_n)})_i \right\|_2^2 = \|\mathbf{Y}^{(k_m)} - \mathbf{Y}^{(k_n)}\|_{\mathbb{F}}^2; \end{aligned}$$

here, the second approximation holds if each subspace  $\mathbf{U}_i$  is sufficiently close to  $\mathbf{U}_i^* = \text{span}\{(\mathbf{Y}^{(j)})_i, 1 \leq j \leq K\}$ .

The purpose of utilizing  $\{U_i\}$  is to reduce the variance of  $\text{stat}_Y$ . Similarly, we can compute another statistic  $\text{stat}_\Gamma$  with  $\mathbb{E}[\text{stat}_\Gamma] \approx \|\Gamma^{(k_m)} - \Gamma^{(k_n)}\|_F^2$ . Consequently, we shall declare  $k_m \neq k_n$  if  $\text{stat}_\Gamma + \text{stat}_Y$  exceeds some appropriate threshold.

*Remark 2.2.* For subspace estimation and clustering, we split each trajectory evenly into four parts, and only utilize the 2nd and 4th. In theory, this only affects sample complexities by a constant factor. The key benefit is that no prior knowledge of the mixing time  $t_{\text{mix}}$  is needed. In cases where  $t_{\text{mix}}$  is known/estimated, one might improve data efficiency by letting the 1st and 3rd parts have only length  $t_{\text{mix}}$ .

## 2.4. Model estimation

Suppose that we have obtained a good clustering accuracy in the previous step. Then, for each  $k$ , we use the samples  $\{\{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}\}_{m \in \mathcal{C}_k}$  to obtain an estimate  $\hat{\mathbf{A}}^{(k)}$  of the state transition matrix by solving a least-squares problem (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019). Finally, we can estimate the noise vector  $\hat{\mathbf{w}}_{m,t} := \mathbf{x}_{m,t+1} - \hat{\mathbf{A}}^{(k)} \mathbf{x}_{m,t} \approx \mathbf{w}_{m,t}$ , and use the empirical covariance of  $\{\hat{\mathbf{w}}_{m,t}\}$  as our estimate of the noise covariance matrix.

## 2.5. Classification

In the previous steps, we have obtained initial clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  and coarse model estimates  $\{\hat{\mathbf{A}}^{(k)}, \hat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ . With the assistance of additional sample trajectories in  $\mathcal{M}_{\text{classification}}$ , we can augment the clusters in the following way: for each new trajectory  $\{\mathbf{x}_t\}_{0 \leq t \leq T}$ , we infer its label as  $\hat{k} = \arg \min_k L(\hat{\mathbf{A}}^{(k)}, \hat{\mathbf{W}}^{(k)})$ , and then assign this trajectory to the  $\hat{k}$ -th cluster; here, the loss function  $L(\mathbf{A}, \mathbf{W}) := T \cdot \log \det(\mathbf{W}) + \sum_{t=0}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)^\top \mathbf{W}^{-1} (\mathbf{x}_{t+1} - \mathbf{A}\mathbf{x}_t)$  coincides with the negative log-likelihood under a Gaussian assumption. Once this is done, we can run Algorithm 4 again with the updated clusters  $\{\mathcal{C}_k\}$  to refine our model estimates, which is exactly the last step of Algorithm 1.

*Remark 2.3.* While this work focuses on mixtures of LDSs, we emphasize that the general principles of Algorithm 1 are applicable under much weaker assumptions. For Algorithms 2 and 3 to work, we essentially only require that each sample trajectory satisfies a certain *mixing* property, and that the autocovariances of different models are *sufficiently separated* (and hence distinguishable). As for Algorithms 4 and 5, we essentially require a *well-specified parametric form* of the time-series models. This observation might inspire future extensions (in theory or applications) of Algorithm 1 to much broader settings.

## 3. Main results

### 3.1. Model assumptions

To streamline the theoretical analysis, we focus on the case where the trajectories are driven by *Gaussian noise*; that is, for each  $1 \leq m \leq M, 0 \leq t \leq T_m$ , the noise vector  $\mathbf{w}_{m,t}$  in (1) is independently generated from the Gaussian distribution  $\mathcal{N}(\mathbf{0}, \mathbf{W}^{(k_m)})$ . Next, we assume for simplicity that the labels  $\{k_m\}_{1 \leq m \leq M}$  of the trajectories are pre-determined and fixed, although one might equivalently regard  $\{k_m\}$  as being random and independent of the noise vectors  $\{\mathbf{w}_{m,t}\}$ . Moreover, while our algorithms and analysis are largely insensitive to the initial states  $\{\mathbf{x}_{m,0}\}_{1 \leq m \leq M}$ , we focus on two canonical cases for concreteness: (i) the trajectories start at *zero* state, or (ii) they are segments of *one* continuous long trajectory. This is formalized as follows:

$$\text{Case 0: } \mathbf{x}_{m,0} = \mathbf{0}, \quad 1 \leq m \leq M; \quad (8a)$$

$$\text{Case 1: } \mathbf{x}_{1,0} = \mathbf{0}, \quad \text{and}$$

$$\mathbf{x}_{m+1,0} = \mathbf{x}_{m,T_m}, \quad 1 \leq m \leq M-1. \quad (8b)$$

We further define the total trajectory length as  $T_{\text{total}} := \sum_{1 \leq m \leq M} T_m = \sum_{\circ} T_{\text{total},\circ}$ , where

$$T_{\text{total},\circ} := T_{\circ} \cdot |\mathcal{M}_{\circ}|,$$

$$\circ \in \{\text{subspace, clustering, classification}\}.$$

Additionally, we assume that each model occupies a non-degenerate fraction of the data; in other words, there exists some  $0 < p_{\min} \leq 1/K$  such that for all  $1 \leq k \leq K$  and  $\circ \in \{\text{subspace, clustering, classification}\}$ ,

$$p_{\min} \leq p_{\circ}^{(k)} := \frac{1}{|\mathcal{M}_{\circ}|} \sum_{m \in \mathcal{M}_{\circ}} \mathbb{1}(k_m = k).$$

Finally, we make the following assumptions about the ground-truth LDS models, where we recall that the autocovariance matrices  $\{\Gamma^{(k)}, \mathbf{Y}^{(k)}\}$  have been defined in (5).

**Assumption 3.1.** *The LDS models  $\{\mathbf{A}^{(k)}, \mathbf{W}^{(k)}\}_{1 \leq k \leq K}$  satisfy the following conditions:*

1. *There exist  $\kappa_A \geq 1$  and  $0 \leq \rho < 1$  such that for any  $1 \leq k \leq K$ ,  $\|(\mathbf{A}^{(k)})^t\| \leq \kappa_A \cdot \rho^t, t = 1, 2, \dots$ ;*
2. *There exist  $\Gamma_{\max} \geq W_{\max} \geq W_{\min} > 0$  and  $\kappa_{w,\text{cross}} \geq \kappa_w \geq 1$  such that for any  $1 \leq k \leq K$ ,*
  - (i)  $\lambda_{\max}(\Gamma^{(k)}) \leq \Gamma_{\max}$ ,
  - (ii)  $W_{\min} \leq \lambda_{\min}(\mathbf{W}^{(k)}) \leq \lambda_{\max}(\mathbf{W}^{(k)}) \leq W_{\max}$ ,
  - (iii)  $W_{\max}/W_{\min} =: \kappa_{w,\text{cross}}$ ,
  - (iv)  $\kappa(\mathbf{W}^{(k)}) = \lambda_{\max}(\mathbf{W}^{(k)})/\lambda_{\min}(\mathbf{W}^{(k)}) \leq \kappa_w$ ;
3. *There exist  $\Delta_{\Gamma,Y}, \Delta_{A,W} > 0$  such that for any  $1 \leq k < \ell \leq K$ ,*

$$\|\Gamma^{(k)} - \Gamma^{(\ell)}\|_F^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_F^2 \geq \Delta_{\Gamma,Y}^2,$$

$$\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_F^2 + \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_F^2}{W_{\max}^2} \geq \Delta_{A,W}^2.$$

The first assumption states that each state transition matrix  $\mathbf{A}^{(k)}$  is *exponentially stable*, which is a quantified version of stability and has appeared in various forms in the literature of LDS (Kailath et al., 2000; Cohen et al., 2018); here,  $\kappa_A$  can be regarded as a condition number, while  $\rho$  is a contraction rate. The second assumption ensures that the noise covariance matrices  $\{\mathbf{W}^{(k)}\}$  are well conditioned, and the autocovariance matrices  $\{\mathbf{\Gamma}^{(k)}\}$  are bounded. The third assumption quantifies the separation between different LDS models. It is important that we consider the separation of  $(\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)})$  versus  $(\mathbf{\Gamma}^{(\ell)}, \mathbf{Y}^{(\ell)})$  jointly, which guarantees that  $\Delta_{\Gamma, Y}$  is always strictly positive (thanks to Fact 2.1), despite the possibility of  $\mathbf{\Gamma}^{(k)} = \mathbf{\Gamma}^{(\ell)}$  or  $\mathbf{Y}^{(k)} = \mathbf{Y}^{(\ell)}$ ; the reasoning for our definition of  $\Delta_{A, W}$  is similar.

**Remark 3.2.** Since the separation parameters  $\Delta_{A, W}$ ,  $\Delta_{\Gamma, Y}$  are defined with respect to the Frobenius norm, we may naturally regard them as  $\Delta_{A, W} = \sqrt{d}\delta_{A, W}$ ,  $\Delta_{\Gamma, Y} = \sqrt{d}\delta_{\Gamma, Y}$ , where  $\delta_{A, W}$ ,  $\delta_{\Gamma, Y}$  are the *canonical separation parameters* (in terms of the spectral norm). For example, in the simple setting where  $K = 2$ ,  $\mathbf{W}^{(1)} = \mathbf{W}^{(2)}$ ,  $\mathbf{A}^{(1)} = 0.5\mathbf{I}_d$  and  $\mathbf{A}^{(2)} = (0.5 + \delta)\mathbf{I}_d$  for some canonical separation parameter  $\delta$ , we have  $\Delta_{A, W} = \|\mathbf{A}^{(1)} - \mathbf{A}^{(2)}\|_F = \|\delta\mathbf{I}_d\|_F = \sqrt{d}\delta$ . This observation will be crucial for obtaining the correct dependence on the dimension  $d$  in our subsequent analysis.

**Remark 3.3.** Most of the parameters in Assumption 3.1 come directly from the ground-truth LDS models  $\{\mathbf{A}^{(k)}, \mathbf{W}^{(k)}\}$ , except  $\Gamma_{\max}$  and  $\Delta_{\Gamma, Y}$ . In fact, they can be bounded by  $\Gamma_{\max} \lesssim W_{\max}$  and  $\Delta_{\Gamma, Y} \gtrsim \Delta_{A, W}$ ; see Fact D.4 for the formal statements. However, the pre-factors in these bounds can be pessimistic in general cases, therefore we choose to preserve  $\Gamma_{\max}$  and  $\Delta_{\Gamma, Y}$  in our analysis.

### 3.2. Theoretical guarantees

We are now ready to present our end-to-end performance guarantees for Algorithm 1. Our main results for Cases 0 and 1 (defined in (8)) are summarized in Theorems 3.4 and 3.5 below, with proofs deferred to Appendix B.

**Theorem 3.4** (Case 0). *There exist positive constants  $\{C_i\}_{1 \leq i \leq 8}$  such that the following holds for any fixed  $0 < \delta < 1/2$ . Consider the model (1) under the assumptions in Sections 2 and 3.1, with a focus on Case 0 (8a). Suppose that we run Algorithm 1 with parameters  $\tau, G$  that satisfy  $1/8 < \tau/\Delta_{\Gamma, Y}^2 < 3/8$ ,  $G \geq C_1\iota_1$ , and data  $\{\mathbf{X}_m\}_{1 \leq m \leq M}$  (where  $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}$ ) that satisfies*

$$T_{\text{subspace}} \geq C_2 \frac{\iota_1}{1 - \rho},$$

$$T_{\text{total,subspace}} \geq C_3 \frac{d}{1 - \rho} \left( \left( \frac{\Gamma_{\max} \sqrt{d}}{\Delta_{\Gamma, Y}} \right)^4 \frac{K^2}{p_{\min}^2} + 1 \right) \cdot \iota_1^4, \quad (9a)$$

$$T_{\text{clustering}} \geq C_4 \frac{G}{1 - \rho} \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\Delta_{\Gamma, Y}^2} + 1 \right) \iota_2,$$

$$T_{\text{total,clustering}} \geq C_5 \frac{d \kappa_{w, \text{cross}}^2}{p_{\min}} \left( \frac{d}{\Delta_{A, W}^2} \frac{\Gamma_{\max}}{W_{\min}} + 1 \right) \iota_3^2, \quad (9b)$$

$$T_{\text{classification}} \geq C_6 \left( \kappa_w^2 + \frac{\kappa_{w, \text{cross}}^6}{\Delta_{A, W}^2} \right) \iota_1^2, \quad (9c)$$

where we define the logarithmic terms  $\iota_1 := \log\left(\frac{d \kappa_A T_{\text{total}}}{\delta}\right)$ ,  $\iota_2 := \log\left(\left(\frac{\Gamma_{\max}}{\Delta_{\Gamma, Y}} + 2\right) \frac{d \kappa_A T_{\text{total}}}{\delta}\right)$ ,  $\iota_3 := \log\left(\frac{\Gamma_{\max}}{W_{\min}} \frac{d \kappa_A T_{\text{total}}}{\delta}\right)$ . Then, with probability at least  $1 - \delta$ , Algorithm 1 achieves exact clustering in Line 4 and exact classification in Line 7; moreover, there exists some permutation  $\pi : \{1, \dots, K\} \rightarrow \{1, \dots, K\}$  such that the final model estimation  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$  in Line 8 obeys

$$1 \leq k \leq K, \quad \|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(\pi(k))}\| \leq C_7 \sqrt{\frac{d \kappa_w \iota_3}{p_{\min} T}},$$

$$\frac{\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(\pi(k))}\|}{\|\mathbf{W}^{(\pi(k))}\|} \leq C_8 \sqrt{\frac{d \iota_3}{p_{\min} T}}, \quad (10)$$

where  $T := T_{\text{total,clustering}} + T_{\text{total,classification}}$ .

**Theorem 3.5** (Case 1). *Consider the same setting of Theorem 3.4, except that we focus on Case 1 (cf. (8b)) instead. Then the same performance guarantees continue to hold, if we replace the conditions on  $T_{\text{total,clustering}}$  and  $T_{\text{classification}}$  in (9) with*

$$T_{\text{total,clustering}} \geq C_5 \frac{d \kappa_{w, \text{cross}}^2}{p_{\min}} \left( \frac{d}{\Delta_{A, W}^2} \frac{\Gamma_{\max}}{W_{\min}} \kappa_A^2 + 1 \right) \iota_3^2, \quad (11a)$$

$$T_{\text{classification}} \geq C_6 \left( \kappa_w^2 + \frac{\kappa_{w, \text{cross}}^6}{\Delta_{A, W}^2} \right) \iota_1^2 + \frac{1}{1 - \rho} \log(2\kappa_A). \quad (11b)$$

**Remark 3.6.** It is worth noting that all the short trajectory lengths  $\{T_o\}$  in the above theorems have sublinear dependence on the dimension  $d$ . In particular, the  $\sqrt{dK}$  dependence in  $T_{\text{clustering}}$  (9) is due to variance reduction achieved by subspace projection (see Section 2.3), without which this dependence would have been linear in  $d$ .

While Theorems 3.4 and 3.5 guarantee that Algorithm 1 successfully learns a mixture of LDS models with a polynomial number of samples, these results involve many parameters and may be somewhat difficult to interpret. In the following corollary, we make some simplifications and focus on the most important parameters.

**Corollary 3.7.** *Consider the same setting of Theorems 3.4 and 3.5. For simplicity, suppose that the condition numbers  $\kappa_A, \kappa_w, \kappa_{w, \text{cross}} \asymp 1$ , and the fractions of data generated by different LDS models are balanced, namely  $p_{\min} \asymp 1/K$ . Moreover, define the canonical separation parameters  $\delta_{A, W} := \Delta_{A, W}/\sqrt{d}$  and  $\delta_{\Gamma, Y} := \Delta_{\Gamma, Y}/\sqrt{d}$ ,*

as suggested by Remark 3.2. Finally, define the mixing time  $t_{\text{mix}} := 1/(1 - \rho)$ . Then we can rewrite the sample complexities in Theorems 3.4 and 3.5 as follows (where we hide the logarithmic terms  $\{\iota_i\}_{1 \leq i \leq 3}$ ): if

$$\begin{aligned} T_{\text{subspace}} &\gtrsim t_{\text{mix}}, \\ T_{\text{total,subspace}} &\gtrsim t_{\text{mix}} d \left( \left( \frac{\Gamma_{\max} K}{\delta_{\Gamma, Y}} \right)^4 + 1 \right), \\ T_{\text{clustering}} &\gtrsim t_{\text{mix}} \left( \left( \frac{\Gamma_{\max}}{\delta_{\Gamma, Y}} \right)^2 \sqrt{\frac{K}{d}} + 1 \right), \\ T_{\text{total,clustering}} &\gtrsim K d \left( \frac{1}{\delta_{A, W}^2} \frac{\Gamma_{\max}}{W_{\min}} + 1 \right), \\ T_{\text{classification}} &\gtrsim \begin{cases} \frac{1}{d \delta_{A, W}^2} + 1 & \text{for Case 0,} \\ \frac{1}{d \delta_{A, W}^2} + t_{\text{mix}} & \text{for Case 1,} \end{cases} \\ T_{\text{total,clustering}} + T_{\text{total,classification}} &\gtrsim \frac{K d}{\epsilon^2}, \end{aligned}$$

then with high probability, Algorithm 1 achieves exact clustering in Line 4, exact classification in Line 7, and final model estimation errors

$$\|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(\pi(k))}\| \leq \epsilon, \quad \frac{\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(\pi(k))}\|}{\|\mathbf{W}^{(\pi(k))}\|} \leq \epsilon.$$

Below are some key implications of Corollary 3.7.

- *Dimension  $d$  and targeted error  $\epsilon$ :* (1) Our algorithms allow the  $T_o$ 's to be much smaller than (and even decrease with)  $d$ . (2) Each  $T_{\text{total},o}$  shall grow linearly with  $d$ , and it takes an order of  $Kd/\epsilon^2$  samples in total to learn  $K$  models in up to  $\epsilon \rightarrow 0$  errors (in the spectral norm), which is just  $K$  times the usual parametric rate ( $d/\epsilon^2$ ) of estimating a single model.
- *Mixing time  $t_{\text{mix}}$ :* (1)  $T_{\text{subspace}}$  and  $T_{\text{clustering}}$  are linear in  $t_{\text{mix}}$ , which ensures sufficient mixing and thus facilitates our algorithms for Stage 1. In contrast,  $T_{\text{classification}}$  depends on  $t_{\text{mix}}$  only for Case 1, with the sole and different purpose of ensuring that the states  $\{\mathbf{x}_{m,t}\}$  are bounded throughout (see Example D.5 for an explanation). (2) While  $T_{\text{total,subspace}}$  needs to grow linearly with  $t_{\text{mix}}$ , this is not required for  $T_{\text{total,clustering}}$  and  $T_{\text{total,classification}}$ , because our methods for model estimation (Algorithm 4) do not rely on mixing (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019).
- *Canonical separation parameters  $\delta_{A,W}, \delta_{\Gamma,Y}$ :* (1)  $T_{\text{clustering}} \gtrsim 1/\delta_{\Gamma,Y}^2$  guarantees exact clustering of the trajectories, while  $T_{\text{classification}} \gtrsim 1/\delta_{A,W}^2$  guarantees exact classification. (2)  $T_{\text{total,subspace}} \gtrsim 1/\delta_{\Gamma,Y}^4$  leads to sufficiently accurate subspaces, while

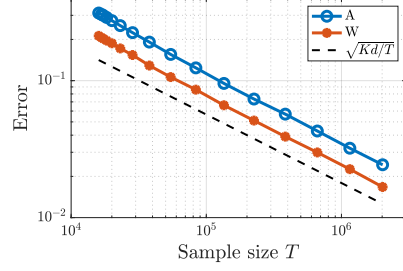


Figure 2. The model estimation errors of Algorithm 1 versus the total sample size (excluding  $\mathcal{M}_{\text{subspace}}$ ). Each curve is an average over 12 independent trials.

$T_{\text{total,clustering}} \gtrsim 1/\delta_{A,W}^2$  leads to accurate initial model estimation.<sup>1</sup>

*Remark 3.8.* It is worth noting that  $T_{\text{clustering}} \gg T_{\text{classification}}$  in Corollary 3.7, i.e. clustering requires a larger trajectory length than classification does. This justifies the benefit of the proposed two-stage procedure, compared with simply doing a large clustering of all trajectories and then one-shot model estimation.

### 3.3. Numerical experiments

We now validate our theoretical findings with a series of numerical experiments, confirming that Algorithm 1 successfully solves the mixed LDSs problem. In these experiments, we fix  $d = 80$ ,  $K = 4$ ; moreover, let  $T_{\text{subspace}} = 20$ ,  $T_{\text{clustering}} = 20$  and  $T_{\text{classification}} = 5$ , all of which are much smaller than  $d$ . We take  $|\mathcal{M}_{\text{subspace}}| = 30d$ ,  $|\mathcal{M}_{\text{clustering}}| = 10d$ , and vary  $|\mathcal{M}_{\text{classification}}|$  between  $[0, 5000d]$ . Our experiments focus on Case 1 as defined in (8b), and we generate the labels of the sample trajectories uniformly at random. The ground-truth LDS models are generated in the following manner:  $\mathbf{A}^{(k)} = \rho \mathbf{R}^{(k)}$ , where  $\rho = 0.5$  and  $\mathbf{R}^{(k)} \in \mathbb{R}^{d \times d}$  is a random orthogonal matrix;  $\mathbf{W}^{(k)}$  has eigendecomposition  $\mathbf{U}^{(k)} \mathbf{\Lambda}^{(k)} (\mathbf{U}^{(k)})^\top$ , where  $\mathbf{U}^{(k)}$  is a random orthogonal matrix, and the diagonal entries of  $\mathbf{\Lambda}^{(k)}$  are independently drawn from the uniform distribution on  $[1, 2]$ .

Our experimental results are illustrated in Figure 2. Here, the horizontal axis represents the sample size  $T = T_{\text{total,clustering}} + T_{\text{total,classification}}$  for model estimation, and the vertical axis represents the estimation errors, measured by  $\max_k \|\widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(\pi(k))}\|$  (plotted in blue) and  $\max_k \|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(\pi(k))}\| / \|\mathbf{W}^{(\pi(k))}\|$  (plotted in orange). The results confirm our main theoretical prediction: Algorithm 1 recovers the LDS models based on a mixture of short trajectories with length  $T_o \ll d$ , and achieves an error rate of  $1/\sqrt{T}$ . In addition, we observe in our experiments

<sup>1</sup>It is possible to improve the  $1/\delta_{\Gamma,Y}^4$  factor in  $T_{\text{total,subspace}}$  to  $1/\delta_{\Gamma,Y}^2$ , if one is willing to pay for some extra factors of eigen-gaps; see Appendix B for a detailed discussion.

that the final outputs of Algorithm 1 are robust to a small number of mis-clustered or mis-classified trajectories during the intermediate steps; this is clearly an appealing property to have, especially when the  $T_o$ 's and  $|\mathcal{M}_o|$ 's in reality are slightly smaller than what our theory requires. Additional experimental results can be found in Appendix F.

#### 4. Related works

**Meta-learning for mixed linear regression.** Our work is closely related to the recent papers (Kong et al., 2020b;a) that bridge mixed linear regression (Quandt & Ramsey, 1978; Yi et al., 2014; Chen et al., 2017; Li & Liang, 2018; Chen et al., 2020; Kwon et al., 2021b; Diakonikolas & Kane, 2020; Yin et al., 2018; Pal & Mazumdar, 2020; Chen et al., 2021c; Diamandis et al., 2021) and meta-learning (Finn et al., 2017; Harrison et al., 2020; Snell et al., 2017; Du et al., 2021; Pan & Yang, 2009; Tripuraneni et al., 2020; Baxter, 2000; Maurer et al., 2016). In this setting of *meta-learning for mixed linear regression* (Kong et al., 2020b;a), one has access to multiple tasks, each containing a few independent and identically distributed samples generated by an unknown model; the inductive bias (i.e. the common structure underlying these tasks) for meta-learning is that there is only a small discrete set of ground-truth linear regression models, akin to the case of mixed linear regression. While the high-level idea of our Algorithm 1 is largely inspired by the work of (Kong et al., 2020b;a), our *detailed implementations* are substantially different, due to the lack of ideal i.i.d. assumption (among others) in our case of time-series data; in addition, we improve upon some of the analyses in (Kong et al., 2020b;a).<sup>2</sup>

**Mixtures of time-series models and trajectories.** Mixture models for time series have achieved empirical success in the study of psychology (Bulteel et al., 2016; Takano et al., 2020), neuroscience (Albert, 1991; Mezer et al., 2009), biology (Wong & Li, 2000), air pollution (D'Urso et al., 2015), economics (McCulloch & Tsay, 1994; Maharaj, 2000; Kalliovirta et al., 2016), automobile sensors (Hallac et al., 2017), and many other domains. Some specific aspects of mixture models include hypothesis testing for a pair of trajectories (Maharaj, 2000), or clustering of multiple trajectories (Liao, 2005; Aghabozorgi et al., 2015; Pathak

et al., 2021; Huang et al., 2021). In addition to mixture models, other related yet different models include time-varying systems (Qu et al., 2021; Minasyan et al., 2021), systems with random parameters (Du et al., 2020), switching systems (Sun, 2006; Sarkar et al., 2019; Ansari et al., 2021), switching state-space models (Ghahramani & Hinton, 2000; Linderman et al., 2017), Markovian jump systems (Shi & Li, 2015; Zhao et al., 2019), and event-triggered systems (Sedghi et al., 2020; Schluter et al., 2020), to name just a few. There are even more related models in reinforcement learning (RL), such as latent bandit (Maillard & Mannor, 2014; Hong et al., 2020), multi-task learning (Wilson et al., 2007; Brunskill & Li, 2013; Liu et al., 2016; Sodhani et al., 2021) / meta-learning (Finn et al., 2017) / transfer-learning (Taylor & Stone, 2009; Tirinzoni et al., 2020) for RL, latent Markov decision processes (Kwon et al., 2021a; Brunskill et al., 2009; Steimle et al., 2021), and so on. What distinguishes our work from this extensive literature is that, we design algorithms and prove non-asymptotic sample complexities for model estimation, in the specific setting of mixture models that features (1) a finite set of underlying time-series models, and (2) unknown labels of the trajectories, with no probabilistic assumptions imposed on these latent variables.

**Linear dynamical systems.** LDS (also called vector autoregressive models) is one of the most fundamental models in system identification and optimal control (Ljung, 1998; Khalil et al., 1996). Recently, there has been a surge of studies about non-asymptotic theoretical analyses of various learning procedures for the basic LDS model (Faradonbeh et al., 2018; Simchowitz et al., 2018; Sarkar & Rakhlin, 2019), linear-quadratic regulators (Dean et al., 2020; Cohen et al., 2018; Jedra & Proutiere, 2019; Faradonbeh et al., 2020; Mania et al., 2019; Simchowitz & Foster, 2020; Fazel et al., 2018; Malik et al., 2019), and LDSs with partial observations (Oymak & Ozay, 2019; Simchowitz et al., 2019; Sarkar et al., 2021; Sun et al., 2020; Tsiamis et al., 2020; Lale et al., 2020; Zheng et al., 2021). In particular, it was only until recently that the authors of (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019) proved sharp error bounds of least squares for estimating the state transition matrix of a LDS model, using a single trajectory; our analysis of Algorithm 4 is largely inspired by their techniques.

#### 5. Discussion

This paper has developed a theoretical and algorithmic framework for learning multiple LDS models from a mixture of short, unlabeled sample trajectories. Our key contributions include a modular two-stage meta-algorithm, as well as theoretical analysis demonstrating its computational and statistical efficiency. We would like to invite the readers to contribute to this important topic by, say, further strengthening the theoretical analysis and algorithmic design. For exam-

<sup>2</sup>There is a concurrent preprint (Modi et al., 2021) that extends multi-task learning (Du et al., 2021; Tripuraneni et al., 2020) to time-series data, and the setting therein includes mixed LDSs as a special case. However, the authors of (Modi et al., 2021) assume oracle access to the global optimum of a non-convex optimization problem, without providing a practical algorithm that can provably find it; moreover, with the short trajectory length fixed, the estimation error bounds in that work will remain bounded away from zero, even if the number of trajectories grows to infinity. In comparison, we consider a simpler problem setting, and propose computationally efficient algorithms with better error bounds.



ple, in certain cases  $T_{\text{clustering}}$  can be a bottleneck compared with  $T_{\text{subspace}}$  and  $T_{\text{classification}}$ , and thus one might hope to achieve a better dependence on  $T_{\text{total,clustering}}$  (e.g. allowing  $T_{\text{total,clustering}} \ll d$ ), by replacing Line 5 in Algorithm 1 with a different method (e.g. adapting Algorithm 3 of (Kong et al., 2020b) to our setting). As another example, in some practical scenarios, the data is a single continuous trajectory, and the time steps when the underlying model changes are *unknown* (Hallac et al., 2017; Harrison et al., 2020); in order to accommodate such a case, one might need to incorporate change-point detection into the learning process.

Moving beyond the current setting of mixed LDSs, we remark that there are plenty of opportunities for future studies. For instance, while our methods in Stage 1 rely on the mixing property of the LDS models, it is worth exploring whether it is feasible to handle the non-mixing case (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019). Another potential direction is to consider the robustness against outliers and adversarial noise (Chen et al., 2021a; Kong et al., 2020a). One might even go further and extend the ideas to learning mixtures of other time-series models (see Remark 2.3), such as LDS with partial or nonlinear observations (Mhammedi et al., 2020), or nonlinear dynamical systems (Mania et al., 2020; Kakade et al., 2020; Foster et al., 2020). Ultimately, it would be of great importance to consider the case with controlled inputs, such as learning mixtures of linear-quadratic regulators, or latent Markov decision processes (Kwon et al., 2021a) that arises in reinforcement learning.

## Acknowledgements

Y. Chen is supported in part by the ARO grant W911NF-20-1-0097, the NSF grants CCF-1907661 and IIS-1900140, and the AFOSR grant FA9550-19-1-0030. H. V. Poor is supported in part by the NSF under Grant CCF-1908308. We would like to thank Yuxin Chen and Gen Li for numerous helpful discussions.

## References

- Aghabozorgi, S., Shirkhorshidi, A. S., and Wah, T. Y. Time-series clustering – a decade review. *Information Systems*, 53:16–38, 2015.
- Albert, P. S. A two-state markov mixture model for a time series of epileptic seizure counts. *Biometrics*, pp. 1371–1381, 1991.
- Ansari, A. F., Benidis, K., Kurle, R., Turkmen, A. C., Soh, H., Smola, A. J., Wang, B., and Januschowski, T. Deep explicit duration switching models for time series. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Baxter, J. A model of inductive bias learning. *Journal of Artificial Intelligence Research*, 12:149–198, 2000.
- Brunskill, E. and Li, L. Sample complexity of multi-task reinforcement learning. In *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, pp. 122–131, Arlington, Virginia, USA, 2013. AUAI Press.
- Brunskill, E., Leffler, B. R., Li, L., Littman, M. L., and Roy, N. Provably efficient learning with typed parametric models. *Journal of Machine Learning Research*, 10(68):1955–1988, 2009.
- Bulteel, K., Tuerlinckx, F., Brose, A., and Ceulemans, E. Clustering vector autoregressive models: Capturing qualitative differences in within-person dynamics. *Frontiers in Psychology*, 7:1540, 2016.
- Chen, S., Li, J., and Song, Z. Learning mixtures of linear regressions in subexponential time via fourier moments. In *Proceedings of the 52nd Annual ACM SIGACT Symposium on Theory of Computing*, pp. 587–600, 2020.
- Chen, S., Koehler, F., Moitra, A., and Yau, M. Kalman filtering with adversarial corruptions. *arXiv preprint arXiv:2111.06395*, 2021a.
- Chen, Y., Yi, X., and Caramanis, C. Convex and nonconvex formulations for mixed regression with two components: Minimax optimal rates. *IEEE Transactions on Information Theory*, 64(3):1738–1766, 2017.
- Chen, Y., Chi, Y., Fan, J., Ma, C., et al. Spectral methods for data science: A statistical perspective. *Foundations and Trends® in Machine Learning*, 14(5):566–806, 2021b.
- Chen, Y., Ma, C., Poor, H. V., and Chen, Y. Learning mixtures of low-rank models. *IEEE Transactions on Information Theory*, 67(7):4613–4636, 2021c. doi: 10.1109/TIT.2021.3065700.
- Cohen, A., Hasidim, A., Koren, T., Lazic, N., Mansour, Y., and Talwar, K. Online linear quadratic control. In *Proceedings of the International Conference on Machine Learning*, pp. 1029–1038. PMLR, 2018.
- Davis, C. and Kahan, W. M. The rotation of eigenvectors by a perturbation. iii. *SIAM Journal on Numerical Analysis*, 7(1):1–46, 1970.
- Dean, S., Mania, H., Matni, N., Recht, B., and Tu, S. On the sample complexity of the linear quadratic regulator. *Foundations of Computational Mathematics*, 20(4):633–679, 2020.
- Diakonikolas, I. and Kane, D. M. Small covers for near-zero sets of polynomials and learning latent variable models. In *Proceedings of the 2020 IEEE 61st Annual Symposium*

- on *Foundations of Computer Science (FOCS)*, pp. 184–195. IEEE, 2020.
- Diamandis, T., Eldar, Y., Fallah, A., Farnia, F., and Ozdaglar, A. A Wasserstein minimax framework for mixed linear regression. In Meila, M. and Zhang, T. (eds.), *Proceedings of the International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 2697–2706. PMLR, 18–24 Jul 2021.
- Du, K., Meng, Q., and Zhang, F. A q-learning algorithm for discrete-time linear-quadratic control with random parameters of unknown distribution: convergence and stabilization. *arXiv preprint arXiv:2011.04970*, 2020.
- Du, S. S., Hu, W., Kakade, S. M., Lee, J. D., and Lei, Q. Few-shot learning via learning the representation, provably. In *International Conference on Learning Representations*, 2021.
- D’Urso, P., De Giovanni, L., and Massari, R. Time series clustering by a robust autoregressive metric with application to air pollution. *Chemometrics and Intelligent Laboratory Systems*, 141:107–124, 2015.
- Faradonbeh, M. K. S., Tewari, A., and Michailidis, G. Finite time identification in unstable linear systems. *Automatica*, 96:342–353, 2018.
- Faradonbeh, M. K. S., Tewari, A., and Michailidis, G. On adaptive linear–quadratic regulators. *Automatica*, 117:108982, 2020.
- Fazel, M., Ge, R., Kakade, S., and Mesbahi, M. Global convergence of policy gradient methods for the linear quadratic regulator. In *Proceedings of the International Conference on Machine Learning*, pp. 1467–1476. PMLR, 2018.
- Finn, C., Abbeel, P., and Levine, S. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the International Conference on Machine Learning - Volume 70*, pp. 1126–1135. JMLR.org, 2017.
- Foster, D., Sarkar, T., and Rakhlin, A. Learning nonlinear dynamical systems from a single trajectory. In *Learning for Dynamics and Control*, pp. 851–861. PMLR, 2020.
- Ghahramani, Z. and Hinton, G. E. Variational learning for switching state-space models. *Neural computation*, 12(4):831–864, 2000.
- Hallac, D., Vare, S., Boyd, S., and Leskovec, J. Toeplitz inverse covariance-based clustering of multivariate time series data. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 215–223, 2017.
- Harrison, J., Sharma, A., Finn, C., and Pavone, M. Continuous meta-learning without tasks. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Hong, J., Kveton, B., Zaheer, M., Chow, Y., Ahmed, A., and Boutilier, C. Latent bandits revisited. In *Advances in Neural Information Processing Systems*, volume 33, pp. 13423–13433. Curran Associates, Inc., 2020.
- Huang, L., Sudhir, K., and Vishnoi, N. Coresets for time series clustering. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Jedra, Y. and Proutiere, A. Sample complexity lower bounds for linear system identification. In *Proceedings of the 2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 2676–2681. IEEE, 2019.
- Kailath, T., Sayed, A. H., and Hassibi, B. *Linear Estimation*. Prentice Hall, 2000.
- Kakade, S., Krishnamurthy, A., Lowrey, K., Ohnishi, M., and Sun, W. Information theoretic regret bounds for online nonlinear control. In *Advances in Neural Information Processing Systems*, volume 33, pp. 15312–15325. Curran Associates, Inc., 2020.
- Kalliovirta, L., Meitz, M., and Saikkonen, P. Gaussian mixture vector autoregression. *Journal of Econometrics*, 192(2):485–498, 2016.
- Khalil, I. S., Doyle, J., and Glover, K. *Robust and Optimal Control*. Prentice Hall, New Jersey, 1996.
- Kong, W., Somani, R., Kakade, S., and Oh, S. Robust meta-learning for mixed linear regression with small batches. In *Advances in Neural Information Processing Systems*, volume 33, pp. 4683–4696. Curran Associates, Inc., 2020a.
- Kong, W., Somani, R., Song, Z., Kakade, S., and Oh, S. Meta-learning for mixed linear regression. In *Proceedings of the International Conference on Machine Learning*, pp. 5394–5404. PMLR, 2020b.
- Kwon, J., Efroni, Y., Caramanis, C., and Mannor, S. RL for latent MDPs: Regret guarantees and a lower bound. *arXiv preprint arXiv:2102.04939*, 2021a.
- Kwon, J., Ho, N., and Caramanis, C. On the minimax optimality of the em algorithm for learning two-component mixed linear regression. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pp. 1405–1413. PMLR, 2021b.
- Lale, S., Azzadenesheli, K., Hassibi, B., and Anandkumar, A. Logarithmic regret bound in partially observable linear dynamical systems. In *Advances in Neural Information Processing Systems*, volume 33, pp. 20876–20888. Curran Associates, Inc., 2020.

- Li, Y. and Liang, Y. Learning mixtures of linear regressions with nearly optimal complexity. In *Proceedings of the Conference On Learning Theory*, pp. 1125–1144, 2018.
- Liao, T. W. Clustering of time series data – a survey. *Pattern Recognition*, 38(11):1857–1874, 2005.
- Liberzon, D. *Switching in Systems and Control*. Springer Science & Business Media, 2003.
- Linderman, S., Johnson, M., Miller, A., Adams, R., Blei, D., and Paninski, L. Bayesian learning and inference in recurrent switching linear dynamical systems. In *Artificial Intelligence and Statistics*, pp. 914–922. PMLR, 2017.
- Liu, Y., Guo, Z., and Brunskill, E. Pac continuous state online multitask reinforcement learning with identification. In *Proceedings of the 2016 International Conference on Autonomous Agents & Multiagent Systems*, pp. 438–446, 2016.
- Ljung, L. System identification. In *Signal analysis and prediction*, pp. 163–173. Springer, 1998.
- Lugosi, G. and Mendelson, S. Mean estimation and regression under heavy-tailed distributions: A survey. *Foundations of Computational Mathematics*, 19(5):1145–1190, 2019.
- Maharaj, E. A. Cluster of time series. *Journal of Classification*, 17(2):297–314, 2000.
- Maillard, O.-A. and Mannor, S. Latent bandits. In *Proceedings of the International Conference on Machine Learning*, pp. 136–144. PMLR, 2014.
- Malekzadeh, M., Clegg, R. G., Cavallaro, A., and Haddadi, H. Mobile sensor data anonymization. In *Proceedings of the International Conference on Internet of Things Design and Implementation*, IoTDI ’19, pp. 49–58. ACM, 2019.
- Malik, D., Pananjady, A., Bhatia, K., Khamaru, K., Bartlett, P., and Wainwright, M. Derivative-free methods for policy optimization: Guarantees for linear quadratic systems. In *Proceedings of the International Conference on Artificial Intelligence and Statistics*, pp. 2916–2925. PMLR, 2019.
- Mania, H., Tu, S., and Recht, B. Certainty equivalence is efficient for linear quadratic control. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pp. 10154–10164, 2019.
- Mania, H., Jordan, M. I., and Recht, B. Active learning for nonlinear system identification with guarantees. *arXiv preprint arXiv:2006.10277*, 2020.
- Maurer, A., Pontil, M., and Romera-Paredes, B. The benefit of multitask representation learning. *The Journal of Machine Learning Research*, 17(1):2853–2884, 2016.
- McCulloch, R. E. and Tsay, R. S. Statistical analysis of economic time series via markov switching models. *Journal of Time Series Analysis*, 15(5):523–539, 1994.
- Mezer, A., Yovel, Y., Pasternak, O., Gorfine, T., and Asfari, Y. Cluster analysis of resting-state fmri time series. *Neuroimage*, 45(4):1117–1125, 2009.
- Mhammedi, Z., Foster, D. J., Simchowitz, M., Misra, D., Sun, W., Krishnamurthy, A., Rakhlin, A., and Langford, J. Learning the linear quadratic regulator from nonlinear observations. *arXiv preprint arXiv:2010.03799*, 2020.
- Minasyan, E., Gradu, P., Simchowitz, M., and Hazan, E. Online control of unknown time-varying dynamical systems. In *Advances in Neural Information Processing Systems*, volume 34, 2021.
- Modi, A., Faradonbeh, M. K. S., Tewari, A., and Michailidis, G. Joint learning of linear time-invariant dynamical systems. *arXiv preprint arXiv:2112.10955*, 2021.
- Oymak, S. and Ozay, N. Non-asymptotic identification of LTI systems from a single trajectory. In *Proceedings of the 2019 American Control Conference (ACC)*, pp. 5655–5661. IEEE, 2019.
- Pal, S. and Mazumdar, A. Recovery of sparse signals from a mixture of linear samples. In *Proceedings of the International Conference on Machine Learning*, pp. 7466–7475. PMLR, 2020.
- Pan, S. J. and Yang, Q. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2009.
- Pathak, R., Sen, R., Rao, N., Erichson, N. B., Jordan, M. I., and Dhillon, I. S. Cluster-and-conquer: A framework for time-series forecasting. *arXiv preprint arXiv:2110.14011*, 2021.
- Qu, G., Shi, Y., Lale, S., Anandkumar, A., and Wierman, A. Stable online control of linear time-varying systems. In *Learning for Dynamics and Control*, pp. 742–753. PMLR, 2021.
- Quandt, R. E. and Ramsey, J. B. Estimating mixtures of normal distributions and switching regressions. *Journal of the American Statistical Association*, 73(364):730–738, 1978.
- Sarkar, T. and Rakhlin, A. Near optimal finite time identification of arbitrary linear dynamical systems. In *Proceedings of the International Conference on Machine Learning*, pp. 5610–5618. PMLR, 2019.

- Sarkar, T., Rakhlin, A., and Dahleh, M. Nonparametric system identification of stochastic switched linear systems. In *Proceedings of the 2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 3623–3628. IEEE, 2019.
- Sarkar, T., Rakhlin, A., and Dahleh, M. A. Finite time LTI system identification. *Journal of Machine Learning Research*, 22(26):1–61, 2021.
- Schluter, H., Solowjow, F., and Trimpe, S. Event-triggered learning for linear quadratic control. *IEEE Transactions on Automatic Control*, 2020.
- Sedghi, L., Ijaz, Z., Witheephanich, K., Pesch, D., et al. Machine learning in event-triggered control: Recent advances and open issues. *arXiv preprint arXiv:2009.12783*, 2020.
- Shi, P. and Li, F. A survey on markovian jump systems: modeling and design. *International Journal of Control, Automation and Systems*, 13(1):1–16, 2015.
- Simchowitz, M. and Foster, D. Naive exploration is optimal for online lqr. In *Proceedings of the International Conference on Machine Learning*, pp. 8937–8948. PMLR, 2020.
- Simchowitz, M., Mania, H., Tu, S., Jordan, M. I., and Recht, B. Learning without mixing: Towards a sharp analysis of linear system identification. In *Proceedings of the Conference On Learning Theory*, pp. 439–473. PMLR, 2018.
- Simchowitz, M., Boczar, R., and Recht, B. Learning linear dynamical systems with semi-parametric least squares. In *Proceedings of the Conference on Learning Theory*, pp. 2714–2802. PMLR, 2019.
- Snell, J., Swersky, K., and Zemel, R. Prototypical networks for few-shot learning. In *Advances in Neural Information Processing Systems*, pp. 4077–4087, 2017.
- Sodhani, S., Zhang, A., and Pineau, J. Multi-task reinforcement learning with context-based representations. In Meila, M. and Zhang, T. (eds.), *Proceedings of the International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pp. 9767–9779. PMLR, 18–24 Jul 2021.
- Steimle, L. N., Kaufman, D. L., and Denton, B. T. Multi-model markov decision processes. *IJSE Transactions*, pp. 1–16, 2021.
- Sun, Y., Oymak, S., and Fazel, M. Finite sample system identification: Optimal rates and the role of regularization. In *Proceedings of the 2nd Conference on Learning for Dynamics and Control*, volume 120 of *Proceedings of Machine Learning Research*, pp. 16–25. PMLR, 10–11 Jun 2020.
- Sun, Z. *Switched Linear Systems: Control and Design*. Springer Science & Business Media, 2006.
- Takano, K., Stefanovic, M., Rosenkranz, T., and Ehring, T. Clustering individuals on limited features of a vector autoregressive model. *Multivariate Behavioral Research*, pp. 1–19, 2020.
- Taylor, M. E. and Stone, P. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(7), 2009.
- Tirinzone, A., Poiani, R., and Restelli, M. Sequential transfer in reinforcement learning with a generative model. In *Proceedings of the International Conference on Machine Learning*, pp. 9481–9492. PMLR, 2020.
- Tripuraneni, N., Jordan, M. I., and Jin, C. On the theory of transfer learning: The importance of task diversity. In *Advances in Neural Information Processing Systems*, volume 33, 2020.
- Tsiamis, A., Matni, N., and Pappas, G. Sample complexity of kalman filtering for unknown systems. In *Learning for Dynamics and Control*, pp. 435–444. PMLR, 2020.
- Vershynin, R. *High-dimensional Probability: An Introduction with Applications in Data Science*, volume 47. Cambridge university press, 2018.
- Wilson, A., Fern, A., Ray, S., and Tadepalli, P. Multi-task reinforcement learning: a hierarchical Bayesian approach. In *Proceedings of the International Conference on Machine Learning*, pp. 1015–1022, 2007.
- Wong, C. S. and Li, W. K. On a mixture autoregressive model. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 62(1):95–115, 2000.
- Yi, X., Caramanis, C., and Sanghavi, S. Alternating minimization for mixed linear regression. In *Proceedings of the International Conference on Machine Learning*, pp. 613–621, 2014.
- Yin, D., Pedarsani, R., Chen, Y., and Ramchandran, K. Learning mixtures of sparse linear regressions using sparse graph codes. *IEEE Transactions on Information Theory*, 65(3):1430–1451, 2018.
- Zhao, P., Kang, Y., and Zhao, Y.-B. A brief tutorial and survey on markovian jump systems: Stability and control. *IEEE Systems, Man, and Cybernetics Magazine*, 5(2): 37–C3, 2019.

Zheng, Y., Frieri, L., Kamgarpour, M., and Li, N. Sample complexity of linear quadratic gaussian (lqg) control for output feedback systems. In *Learning for Dynamics and Control*, pp. 559–570. PMLR, 2021.

Table 1. A list of notation and parameters. In the subscripts of  $\mathcal{M}_o, T_o, T_{\text{total},o}$ , the symbol  $o$  takes value in {subspace, clustering, classification}.

Notation	Explanation
$d$	State dimension
$K$	Number of LDS models
$k_m$	The unknown label (latent variable) of the $m$ -th trajectory
$\mathcal{M}_o$	Subsets of $M$ trajectories, $\{1, \dots, M\} = \mathcal{M}_{\text{subspace}} \cup \mathcal{M}_{\text{clustering}} \cup \mathcal{M}_{\text{classification}}$
$p_{\min}$	A lower bound for the fraction of trajectories generated by each model
$T_o$	Short trajectory length for $\mathcal{M}_o$
$T_{\text{total},o}$	Total trajectory length, $T_{\text{total},o} = T_o \cdot  \mathcal{M}_o $
$\mathbf{A}^{(k)}, \mathbf{W}^{(k)}$	State transition matrix and noise covariance matrix of the $k$ -th LDS model
$\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)}$	Order-0 and order-1 stationary autocovariance matrices of the $k$ -th LDS model
$\kappa_A, \rho$	$\ (\mathbf{A}^{(k)})^t\  \leq \kappa_A \cdot \rho^t, \quad t = 1, 2, \dots$
$\Delta_{\Gamma, Y}, \Delta_{A, W}$	Model separation parameters (see Assumption 3.1)
$W_{\min}, W_{\max}$	$W_{\min} \leq \lambda_{\min}(\mathbf{W}^{(k)}) \leq \lambda_{\max}(\mathbf{W}^{(k)}) \leq W_{\max}, \quad 1 \leq k \leq K$
$\kappa_{w, \text{cross}}, \kappa_w$	$\kappa_{w, \text{cross}} = W_{\max}/W_{\min}; \quad \kappa(\mathbf{W}^{(k)}) \leq \kappa_w, \quad 1 \leq k \leq K$
$\Gamma_{\max}$	$\ \mathbf{\Gamma}^{(k)}\  \leq \Gamma_{\max}, \quad 1 \leq k \leq K$

This appendix is organized as follows. Appendix A complements Section 2, providing full details of the proposed Algorithms 2 – 5. Appendix B includes detailed theoretical results for the algorithms, with proofs postponed to Appendix C; it also provides a proof for our main results in Section 3. Appendix D collects some miscellaneous results. In Appendix E we discuss on some potential extensions of our methods. Finally, Appendix F includes additional experimental results. For the readers’ convenience, we include Table 1 for a quick reference to the key notation and parameters used in our analysis.

**Additional notation.** Let  $\text{vec}(\cdot)$  denote the vectorization of a matrix. For two matrices  $\mathbf{A}, \mathbf{B}$ , let  $\mathbf{A} \otimes \mathbf{B}$  denote their Kronecker product. Given a sequence of real numbers  $\{x_i\}_{1 \leq i \leq N}$ , we denote its median as  $\text{median}\{x_i, 1 \leq i \leq N\}$ . We shall also let  $\text{poly}(n)$  denote some polynomial in  $n$  of a constant degree. For a positive integer  $n$ , we denote  $[n] := \{1, 2, \dots, n\}$ .

## A. Detailed Algorithms

### A.1. Subspace estimation

**Procedure.** Recall that the notation  $(\mathbf{\Gamma}^{(k)})_i \in \mathbb{R}^d$  (resp.  $(\mathbf{Y}^{(k)})_i$ ) represents the transpose of the  $i$ -th row of  $\mathbf{\Gamma}^{(k)}$  (resp.  $\mathbf{Y}^{(k)}$ ) defined in (5). With this set of notation in place, let us define the following subspaces:

$$\mathbf{V}_i^* := \text{span}\left\{(\mathbf{\Gamma}^{(k)})_i, 1 \leq k \leq K\right\}, \quad \mathbf{U}_i^* := \text{span}\left\{(\mathbf{Y}^{(k)})_i, 1 \leq k \leq K\right\}, \quad 1 \leq i \leq d. \quad (12)$$

It is easily seen from the construction that each of these subspaces has rank at most  $K$ .

As it turns out, the collection of  $2d$  subspaces defined in (12) provides crucial low-dimensional information about the linear dynamical systems of interest. This motivates us to develop a data-driven method to estimate these subspaces, which will in turn allow for proper dimension reduction in subsequent steps. Towards this end, we propose to employ a spectral method for subspace estimation using sample trajectories in  $\mathcal{M}_{\text{subspace}}$ :

- (i) divide  $\{0, 1, \dots, T_{\text{subspace}}\}$  into four segments of the same size, and denote the 2nd (resp. 4th) segment as  $\Omega_1$  (resp.  $\Omega_2$ );

(ii) for each  $m \in \mathcal{M}_{\text{subspace}}$ ,  $1 \leq i \leq d, j \in \{1, 2\}$ , compute

$$\mathbf{h}_{m,i,j} := \frac{1}{|\Omega_j|} \sum_{t \in \Omega_j} (\mathbf{x}_{m,t})_i \mathbf{x}_{m,t}, \quad \mathbf{g}_{m,i,j} := \frac{1}{|\Omega_j|} \sum_{t \in \Omega_j} (\mathbf{x}_{m,t+1})_i \mathbf{x}_{m,t}; \quad (13)$$

(iii) for each  $1 \leq i \leq d$ , compute the following matrices

$$\widehat{\mathbf{H}}_i := \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \mathbf{h}_{m,i,1} \mathbf{h}_{m,i,2}^\top, \quad \widehat{\mathbf{G}}_i := \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \mathbf{g}_{m,i,1} \mathbf{g}_{m,i,2}^\top, \quad (14)$$

and let  $\mathbf{V}_i \in \mathbb{R}^{d \times K}$  (resp.  $\mathbf{U}_i$ ) be the top- $K$  eigenspace of  $\widehat{\mathbf{H}}_i + \widehat{\mathbf{H}}_i^\top$  (resp.  $\widehat{\mathbf{G}}_i + \widehat{\mathbf{G}}_i^\top$ ).

The output  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$  will then serve as our estimate of  $\{\mathbf{V}_i^*, \mathbf{U}_i^*\}_{1 \leq i \leq d}$ . This spectral approach is summarized in Algorithm 2.

**Rationale.** According to the mixing property of LDS, if  $T_{\text{subspace}}$  is larger than some appropriately defined mixing time, then each sample trajectory in  $\mathcal{M}_{\text{subspace}}$  will mix sufficiently and nearly reach stationarity (when constrained to  $t \in \Omega_1 \cup \Omega_2$ ). In this case, it is easy to check that

$$\mathbb{E}[\mathbf{h}_{m,i,j}] \approx (\mathbf{\Gamma}^{(k_m)})_i, \quad \mathbb{E}[\mathbf{g}_{m,i,j}] \approx (\mathbf{Y}^{(k_m)})_i, \quad 1 \leq i \leq d, \quad j \in \{1, 2\}.$$

Moreover, the samples in  $\Omega_1$  are nearly independent of those in  $\Omega_2$  as long as the spacing between them exceeds the mixing time, and therefore,

$$\mathbb{E}[\widehat{\mathbf{H}}_i] \approx \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} (\mathbf{\Gamma}^{(k_m)})_i (\mathbf{\Gamma}^{(k_m)})_i^\top = \sum_{k=1}^K p_{\text{subspace}}^{(k)} (\mathbf{\Gamma}^{(k)})_i (\mathbf{\Gamma}^{(k)})_i^\top =: \mathbf{H}_i, \quad (15a)$$

$$\mathbb{E}[\widehat{\mathbf{G}}_i] \approx \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} (\mathbf{Y}^{(k_m)})_i (\mathbf{Y}^{(k_m)})_i^\top = \sum_{k=1}^K p_{\text{subspace}}^{(k)} (\mathbf{Y}^{(k)})_i (\mathbf{Y}^{(k)})_i^\top =: \mathbf{G}_i, \quad (15b)$$

where  $p_{\text{subspace}}^{(k)}$  denotes the fraction of sample trajectories generated by the  $k$ -th model, namely,

$$p_{\text{subspace}}^{(k)} := \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \mathbb{1}(k_m = k), \quad 1 \leq k \leq K. \quad (16)$$

As a consequence, if  $T_{\text{subspace}}$  and  $|\mathcal{M}_{\text{subspace}}|$  are both sufficiently large, then one might expect  $\widehat{\mathbf{H}}_i$  (resp.  $\widehat{\mathbf{G}}_i$ ) to be a reasonably good approximation of  $\mathbf{H}_i$  (resp.  $\mathbf{G}_i$ ), the latter of which has rank at most  $K$  and has  $\mathbf{V}_i^*$  (resp.  $\mathbf{U}_i^*$ ) as its eigenspace. All this motivates us to compute the rank- $K$  eigenspaces of  $\widehat{\mathbf{H}}_i + \widehat{\mathbf{H}}_i^\top$  and  $\widehat{\mathbf{G}}_i + \widehat{\mathbf{G}}_i^\top$  in Algorithm 2.

## A.2. Clustering

This step seeks to divide the sample trajectories in  $\mathcal{M}_{\text{clustering}}$  into  $K$  clusters (albeit not perfectly), such that the trajectories in each cluster are primarily generated by the same LDS model. We intend to achieve this by performing pairwise comparisons of the autocovariance matrices associated with the sample trajectories.

**Key observation.** Even though  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  for  $k \neq \ell$ , it is indeed possible that  $\mathbf{\Gamma}^{(k)} = \mathbf{\Gamma}^{(\ell)}$  or  $\mathbf{Y}^{(k)} = \mathbf{Y}^{(\ell)}$ . Therefore, in order to differentiate sample trajectories generated by different systems based on  $\mathbf{\Gamma}(\mathbf{A}, \mathbf{W})$  and  $\mathbf{Y}(\mathbf{A}, \mathbf{W})$ , it is important to ensure separation of  $(\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)})$  and  $(\mathbf{\Gamma}^{(\ell)}, \mathbf{Y}^{(\ell)})$  when  $k \neq \ell$ , which can be guaranteed by the following fact.

**Fact A.1.** *If  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$ , then we have either  $\mathbf{\Gamma}^{(k)} \neq \mathbf{\Gamma}^{(\ell)}$  or  $\mathbf{Y}^{(k)} \neq \mathbf{Y}^{(\ell)}$  (or both).*

*Proof.* First, observe that the definitions of  $\{\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)}\}$  in (5) suggest that we can recover  $\mathbf{A}^{(k)}, \mathbf{W}^{(k)}$  from  $\mathbf{\Gamma}^{(k)}, \mathbf{Y}^{(k)}$  as follows:

$$\mathbf{A}^{(k)} = \mathbf{Y}^{(k)} \mathbf{\Gamma}^{(k)-1}, \quad \mathbf{W}^{(k)} = \mathbf{\Gamma}^{(k)} - \mathbf{A}^{(k)} \mathbf{\Gamma}^{(k)} \mathbf{A}^{(k)\top}. \quad (17)$$

**Algorithm 2** Subspace estimation

- 1: **Input:** short trajectories  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{subspace}}}$ , where  $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_{\text{subspace}}}$ .
- 2: Let  $N \leftarrow \lfloor T_{\text{subspace}}/4 \rfloor$ , and  $\Omega_1 \leftarrow \{N + j, 1 \leq j \leq N\}$ ,  $\Omega_2 \leftarrow \{3N + j, 1 \leq j \leq N\}$ .
- 3: **for**  $(m, i, j) \in \mathcal{M}_{\text{subspace}} \times [d] \times [2]$  **do**
- 4:   Compute

$$\mathbf{h}_{m,i,j} \leftarrow \frac{1}{|\Omega_j|} \sum_{t \in \Omega_j} (\mathbf{x}_{m,t})_i \mathbf{x}_{m,t}, \quad \mathbf{g}_{m,i,j} \leftarrow \frac{1}{|\Omega_j|} \sum_{t \in \Omega_j} (\mathbf{x}_{m,t+1})_i \mathbf{x}_{m,t}.$$

- 5: **end for**
- 6: **for**  $i = 1, \dots, d$  **do**
- 7:   Compute

$$\widehat{\mathbf{H}}_i \leftarrow \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \mathbf{h}_{m,i,1} \mathbf{h}_{m,i,2}^\top, \quad \widehat{\mathbf{G}}_i \leftarrow \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \mathbf{g}_{m,i,1} \mathbf{g}_{m,i,2}^\top,$$

- 8:   Let  $\mathbf{V}_i \in \mathbb{R}^{d \times K}$  (resp.  $\mathbf{U}_i$ ) be the top- $K$  eigenspace of  $\widehat{\mathbf{H}}_i + \widehat{\mathbf{H}}_i^\top$  (resp.  $\widehat{\mathbf{G}}_i + \widehat{\mathbf{G}}_i^\top$ ).
- 9: **end for**
- 10: **Output:** subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$ .

With this, we can prove the fact by contradiction. Suppose instead that  $\Gamma^{(k)} = \Gamma^{(\ell)}$  and  $\mathbf{Y}^{(k)} = \mathbf{Y}^{(\ell)}$ , then (17) yields

$$\begin{aligned} \mathbf{A}^{(k)} &= \mathbf{Y}^{(k)} \Gamma^{(k)-1} = \mathbf{Y}^{(\ell)} \Gamma^{(\ell)-1} = \mathbf{A}^{(\ell)}, \quad \text{and} \\ \mathbf{W}^{(k)} &= \Gamma^{(k)} - \mathbf{A}^{(k)} \Gamma^{(k)} \mathbf{A}^{(k)\top} = \Gamma^{(\ell)} - \mathbf{A}^{(\ell)} \Gamma^{(\ell)} \mathbf{A}^{(\ell)\top} = \mathbf{W}^{(\ell)}, \end{aligned}$$

which is contradictory to the assumption that  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$ .  $\square$

**Idea.** Let us compare a pair of sample trajectories  $\{\mathbf{x}_t\}_{0 \leq t \leq T_{\text{clustering}}}$  and  $\{\mathbf{z}_t\}_{0 \leq t \leq T_{\text{clustering}}}$ , where  $\{\mathbf{x}_t\}$  is generated by the system  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  and  $\{\mathbf{z}_t\}$  by the system  $(\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$ . In order to determine whether  $k = \ell$ , we propose to estimate the quantity  $\|\Gamma^{(k)} - \Gamma^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2$  using the data samples  $\{\mathbf{x}_t\}$  and  $\{\mathbf{z}_t\}$ , which is expected to be small (resp. large) if  $k = \ell$  (resp.  $k \neq \ell$ ). To do so, let us divide  $\{0, 1, \dots, T_{\text{clustering}}\}$  evenly into four segments, and denote by  $\Omega_1$  (resp.  $\Omega_2$ ) the 2nd (resp. 4th) segment, akin to what we have done in the previous step. Observe that

$$\mathbf{U}_i^\top \mathbb{E} \left[ ((\mathbf{x}_{t+1})_i \mathbf{x}_t - (\mathbf{z}_{t+1})_i \mathbf{z}_t) \right] \approx \mathbf{U}_i^\top ((\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i)$$

for all  $1 \leq i \leq d$  and  $t \in \Omega_1 \cup \Omega_2$ . Assuming sufficient mixing and utilizing (near) statistical independence due to sample splitting, we might resort to the following statistic

$$\text{stat}_Y := \sum_{i=1}^d \left\langle \mathbf{U}_i^\top \frac{1}{|\Omega_1|} \sum_{t \in \Omega_1} ((\mathbf{x}_{t+1})_i \mathbf{x}_t - (\mathbf{z}_{t+1})_i \mathbf{z}_t), \mathbf{U}_i^\top \frac{1}{|\Omega_2|} \sum_{t \in \Omega_2} ((\mathbf{x}_{t+1})_i \mathbf{x}_t - (\mathbf{z}_{t+1})_i \mathbf{z}_t) \right\rangle, \quad (19)$$

whose expectation is given by

$$\begin{aligned} \mathbb{E}[\text{stat}_Y] &\approx \sum_{i=1}^d \left\langle \mathbf{U}_i^\top ((\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i), \mathbf{U}_i^\top ((\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i) \right\rangle \\ &= \sum_{i=1}^d \left\| \mathbf{U}_i^\top ((\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i) \right\|_2^2 \approx \sum_{i=1}^d \left\| (\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i \right\|_2^2 = \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2; \end{aligned}$$

here, the first approximation is due to the near independence between samples from  $\Omega_1$  and those from  $\Omega_2$ , whereas the second approximation holds if each subspace  $\mathbf{U}_i$  is sufficiently close to  $\mathbf{U}_i^* = \text{span}\{(\mathbf{Y}^{(j)})_i, 1 \leq j \leq K\}$ . The purpose of



**Algorithm 3** Clustering

- 1: **Input:** short trajectories  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{clustering}}}$ , where  $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_{\text{clustering}}}$ ; subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$ ; testing threshold  $\tau$ ; number of copies  $G$ .
- 2: Let  $N \leftarrow \lfloor T_{\text{clustering}}/4G \rfloor$ , and

$$\Omega_{g,1} \leftarrow \{(4g-3)N + j, 1 \leq j \leq N\}, \quad \Omega_{g,2} \leftarrow \{(4g-1)N + j, 1 \leq j \leq N\}, \quad 1 \leq g \leq G.$$

- 3: **for**  $(m, i, g, j) \in \mathcal{M}_{\text{clustering}} \times [d] \times [G] \times [2]$  **do**
- 4:   Compute

$$\mathbf{h}_{m,i,g,j} \leftarrow \frac{1}{|\Omega_{g,j}|} \sum_{t \in \Omega_{g,j}} (\mathbf{x}_{m,t})_i \mathbf{x}_{m,t}, \quad \mathbf{g}_{m,i,g,j} \leftarrow \frac{1}{|\Omega_{g,j}|} \sum_{t \in \Omega_{g,j}} (\mathbf{x}_{m,t+1})_i \mathbf{x}_{m,t}.$$

- 5: **end for**
- 6: // Compute the similarity matrix  $\mathbf{S}$ :
- 7: **for**  $(m, n) \in \mathcal{M}_{\text{clustering}} \times \mathcal{M}_{\text{clustering}}$  **do**
- 8:   **for**  $g = 1, \dots, G$  **do**
- 9:     Compute

$$\text{stat}_{\Gamma,g} \leftarrow \sum_{i=1}^d \left\langle \mathbf{V}_i^\top (\mathbf{h}_{m,i,g,1} - \mathbf{h}_{n,i,g,1}), \mathbf{V}_i^\top (\mathbf{h}_{m,i,g,2} - \mathbf{h}_{n,i,g,2}) \right\rangle, \quad (18a)$$

$$\text{stat}_{\Upsilon,g} \leftarrow \sum_{i=1}^d \left\langle \mathbf{U}_i^\top (\mathbf{g}_{m,i,g,1} - \mathbf{g}_{n,i,g,1}), \mathbf{U}_i^\top (\mathbf{g}_{m,i,g,2} - \mathbf{g}_{n,i,g,2}) \right\rangle, \quad (18b)$$

- 10: **end for**
- 11:   Set  $S_{m,n} \leftarrow \mathbb{1} \left( \text{median} \{ \text{stat}_{\Gamma,g}, 1 \leq g \leq G \} + \text{median} \{ \text{stat}_{\Upsilon,g}, 1 \leq g \leq G \} \leq \tau \right)$ .
- 12: **end for**
- 13: Divide  $\mathcal{M}_{\text{clustering}}$  into  $K$  clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  according to  $\{S_{m,n}\}_{m,n \in \mathcal{M}_{\text{clustering}}}$ .
- 14: **Output:** clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .

utilizing  $\{\mathbf{U}_i\}$  is to reduce the variance of  $\text{stat}_{\Upsilon}$ . Similarly, we can compute another statistic (by replacing  $\{\mathbf{U}_i\}$  with  $\{\mathbf{V}_i\}$  and  $(\mathbf{x}_{t+1})_i, (\mathbf{z}_{t+1})_i$  with  $(\mathbf{x}_t)_i, (\mathbf{z}_t)_i$ ) as follows:

$$\text{stat}_{\Gamma} := \sum_{i=1}^d \left\langle \mathbf{V}_i^\top \frac{1}{|\Omega_1|} \sum_{t \in \Omega_1} ((\mathbf{x}_t)_i \mathbf{x}_t - (\mathbf{z}_t)_i \mathbf{z}_t), \mathbf{V}_i^\top \frac{1}{|\Omega_2|} \sum_{t \in \Omega_2} ((\mathbf{x}_t)_i \mathbf{x}_t - (\mathbf{z}_t)_i \mathbf{z}_t) \right\rangle, \quad (20)$$

which has expectation

$$\mathbb{E}[\text{stat}_{\Gamma}] \approx \sum_{i=1}^d \|\mathbf{V}_i^\top ((\mathbf{\Gamma}^{(k)})_i - (\mathbf{\Gamma}^{(\ell)})_i)\|_2^2 \approx \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}}^2.$$

Consequently, we shall declare  $k \neq \ell$  if  $\text{stat}_{\Gamma} + \text{stat}_{\Upsilon}$  exceeds some appropriate threshold  $\tau$ .

**Procedure.** We are now positioned to describe the proposed clustering procedure. We first compute the statistics  $\text{stat}_{\Gamma}$  and  $\text{stat}_{\Upsilon}$  for each pair of sample trajectories in  $\mathcal{M}_{\text{clustering}}$  by means of the method described above, and then construct a similarity matrix  $\mathbf{S}$ , in a way that  $S_{m,n}$  is set to 0 if  $\text{stat}_{\Gamma} + \text{stat}_{\Upsilon}$  (computed for the  $m$ -th and  $n$ -th trajectories) is larger than a threshold  $\tau$ , or set to 1 otherwise. In order to enhance the robustness of these statistics, we divide  $\{0, 1, \dots, T_{\text{clustering}}\}$  into  $4G$  (instead of 4) segments, compute  $G$  copies of  $\text{stat}_{\Gamma}$  and  $\text{stat}_{\Upsilon}$ , and take the medians of these values. Next, we apply a mainstream clustering algorithm (e.g. spectral clustering (Chen et al., 2021b)) to the similarity matrix  $\mathbf{S}$ , and divide  $\mathcal{M}_{\text{clustering}}$  into  $K$  disjoint clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ . The complete clustering procedure is provided in Algorithm 3. The threshold  $\tau$  shall be chosen to be on the order of  $\min_{1 \leq k < \ell \leq K} \{\|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2\}$  (which is strictly positive due to Fact 2.1), and it suffices to set the number of copies  $G$  to be on some logarithmic order. Our choice of these parameters are specified in Theorem 3.4.

**Algorithm 4** Least squares and covariance estimation

- 1: **Input:** clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .
- 2: **for**  $k = 1, \dots, K$  **do**
- 3:   Compute

$$\widehat{\mathbf{A}}^{(k)} \leftarrow \left( \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \mathbf{x}_{m,t+1} \mathbf{x}_{m,t}^\top \right) \left( \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \mathbf{x}_{m,t} \mathbf{x}_{m,t}^\top \right)^{-1}, \quad \text{and} \quad (22a)$$

$$\widehat{\mathbf{W}}^{(k)} \leftarrow \frac{1}{\sum_{m \in \mathcal{C}_k} T_m} \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \widehat{\mathbf{w}}_{m,t} \widehat{\mathbf{w}}_{m,t}^\top, \quad \text{where} \quad \widehat{\mathbf{w}}_{m,t} = \mathbf{x}_{m,t+1} - \widehat{\mathbf{A}}^{(k)} \mathbf{x}_{m,t}. \quad (22b)$$

- 4: **end for**
- 5: **Output:** estimated models  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ .

**A.3. Model estimation**

Suppose that we have obtained reasonably good clustering accuracy in the previous step, namely for each  $1 \leq k \leq K$ , the cluster  $\mathcal{C}_k$  output by Algorithm 3 contains mostly trajectories generated by the same model  $(\mathbf{A}^{(\pi(k))}, \mathbf{W}^{(\pi(k))})$ , with  $\pi$  representing some permutation function of  $\{1, \dots, K\}$ . We propose to obtain a coarse model estimation and covariance estimation, as summarized in Algorithm 4. More specifically, for each  $k$ , we use the samples  $\{\{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}\}_{m \in \mathcal{C}_k}$  to obtain an estimate of  $\mathbf{A}^{(\pi(k))}$  by solving the following least-squares problem (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019)

$$\widehat{\mathbf{A}}^{(k)} := \arg \min_{\mathbf{A}} \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \|\mathbf{x}_{m,t+1} - \mathbf{A} \mathbf{x}_{m,t}\|_2^2, \quad (21)$$

whose closed-form solution is given in (22a). Next, we can use  $\widehat{\mathbf{A}}^{(k)}$  to estimate the noise vector

$$\widehat{\mathbf{w}}_{m,t} := \mathbf{x}_{m,t+1} - \widehat{\mathbf{A}}^{(k)} \mathbf{x}_{m,t} \approx \mathbf{w}_{m,t},$$

and finally estimate the noise covariance  $\mathbf{W}^{(\pi(k))}$  with the empirical covariance of  $\{\widehat{\mathbf{w}}_{m,t}\}$ , as shown in (22b).

**A.4. Classification**

**Procedure.** In the previous steps, we have obtained initial clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  and coarse model estimates  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ . With the assistance of additional sample trajectories  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{classification}}}$ , we can infer their latent labels and assign them to the corresponding clusters; the procedure is stated in Algorithm 5 and will be explained shortly. Once this is done, we can run Algorithm 4 again with the updated clusters  $\{\mathcal{C}_k\}$  to refine our model estimates, which is exactly the last step of Algorithm 1.

**Rationale.** The strategy of inferring labels in Algorithm 5 can be derived from the maximum likelihood estimator, under the assumption that the noise vectors  $\{\mathbf{w}_{m,t}\}$  follow Gaussian distributions. Note, however, that even in the absence of Gaussian assumptions, Algorithm 5 remains effective in principle. To see this, consider a short trajectory  $\{\mathbf{x}_t\}_{0 \leq t \leq T}$  generated by model  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$ , i.e.  $\mathbf{x}_{t+1} = \mathbf{A}^{(k)} \mathbf{x}_t + \mathbf{w}_t$  where  $\mathbb{E}[\mathbf{w}_t] = \mathbf{0}$ ,  $\text{cov}(\mathbf{w}_t) = \mathbf{W}^{(k)}$ . Let us define the following loss function

$$L(\mathbf{A}, \mathbf{W}) := T \cdot \log \det(\mathbf{W}) + \sum_{t=0}^{T-1} (\mathbf{x}_{t+1} - \mathbf{A} \mathbf{x}_t)^\top \mathbf{W}^{-1} (\mathbf{x}_{t+1} - \mathbf{A} \mathbf{x}_t). \quad (23)$$

With some elementary calculation, we can easily check that for any incorrect label  $\ell \neq k$ , it holds that  $\mathbb{E}[L(\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)}) - L(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})] > 0$ , with the proviso that  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  and  $\{\mathbf{x}_t\}$  are non-degenerate in some sense; in other words, the correct model  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  achieves the minimal expected loss (which, due to the quadratic form of the loss function, depends solely on the first and second moments of the distributions of  $\{\mathbf{w}_t\}$ , as well as the initial state  $\mathbf{x}_0$ ). This justifies the proposed procedure for inferring unknown labels in Algorithm 5, provided that  $T_m$  is large enough and that the estimated LDS models are sufficiently reliable.

**Algorithm 5** Classification

- 1: **Input:** short trajectories  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{classification}}}$ , where  $\mathbf{X}_m = \{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}$ ; coarse models  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$ ; clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .
- 2: **for**  $m \in \mathcal{M}_{\text{classification}}$  **do**
- 3: Infer the label of the  $m$ -th trajectory by

$$\widehat{k}_m \leftarrow \arg \min_{\ell} \left\{ T_m \cdot \log \det(\widehat{\mathbf{W}}^{(\ell)}) + \sum_{t=0}^{T_m-1} (\mathbf{x}_{m,t+1} - \widehat{\mathbf{A}}^{(\ell)} \mathbf{x}_{m,t})^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{x}_{m,t+1} - \widehat{\mathbf{A}}^{(\ell)} \mathbf{x}_{m,t}) \right\},$$

then add  $m$  to cluster  $\mathcal{C}_{\widehat{k}_m}$ .

- 4: **end for**
- 5: **Output:** updated clusters  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$ .

**B. Detailed analysis**

This section provides detailed, modular theoretical results for the performance of Algorithms 2–5, and concludes with a proof for the main theorems in Section 3.2 (i.e. the performance guarantees of Algorithm 1).

**Subspace estimation.** The following theorem provides upper bounds on the errors of subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  output by Algorithm 2, assuming sufficient mixing of the short trajectories.

**Theorem B.1.** *Consider the model (1) under the assumptions in Sections 2 and 3.1. There exist some universal constants  $C_1, C_2, C_3 > 0$  such that the following holds. Suppose that we run Algorithm 2 with data  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{subspace}}}$  obeying*

$$T_{\text{subspace}} \geq C_1 \cdot t_{\text{mix}}, \quad T_{\text{total,subspace}} \geq C_2 \cdot t_{\text{mix}} d \cdot \log \frac{T_{\text{total}} d}{\delta}, \quad \text{where } t_{\text{mix}} := \frac{1}{1-\rho} \cdot \log \left( \frac{d \kappa_A T_{\text{total}}}{\delta} \right).$$

Then with probability at least  $1 - \delta$ , Algorithm 2 outputs  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$  satisfying the following: for all  $1 \leq k \leq K$  and  $1 \leq i \leq d$ ,

$$\begin{aligned} & \max \left\{ \|\mathbf{(\Gamma}^{(k)})_i - \mathbf{V}_i \mathbf{V}_i^\top (\mathbf{\Gamma}^{(k)})_i\|_2, \|\mathbf{(Y}^{(k)})_i - \mathbf{U}_i \mathbf{U}_i^\top (\mathbf{Y}^{(k)})_i\|_2 \right\} \\ & \leq C_3 \cdot \Gamma_{\max} \left( \frac{K}{p_{\min}} \right)^{1/2} \left( \frac{t_{\text{mix}} d}{T_{\text{total,subspace}}} \log^3 \frac{T_{\text{total}} d}{\delta} \right)^{1/4}. \end{aligned} \quad (24)$$

Our proof (deferred to Appendix C.2) includes a novel perturbation analysis; the resulted error bound (24) has a  $1/T_{\text{total,subspace}}^{1/4}$  dependence and is gap-free (i.e. independent of the eigenvalue gaps of the ground-truth low-rank matrices, which can be arbitrarily close to zero in the worst case). It is possible to adapt the existing perturbation results in (Kong et al., 2020b;a) to our setting (which we include in Lemma C.3 in the appendix for completeness); however, one of them is dependent on the eigenvalue gaps, while the other one incurs a worse  $1/T_{\text{total,subspace}}^{1/6}$  dependence. It would be interesting future work to investigate whether a gap-free bound with a  $1/T_{\text{total,subspace}}^{1/2}$  dependence is possible.

**Clustering.** Our next theorem shows that Algorithm 3 achieves exact clustering of  $\mathcal{M}_{\text{clustering}}$ , if  $T_{\text{clustering}}$  is sufficiently large and subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  are accurate. The proof is deferred to Appendix C.3.

**Theorem B.2.** *Consider the model (1) under the assumptions in Sections 2 and 3.1. There exist universal constants  $C_1, C_2, c_3 > 0$  such that the following holds. Suppose that we run Algorithm 3 with data  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{clustering}}}$ , independent subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$  and parameters  $\tau, G$  that satisfy the following:*

- The threshold  $\tau$  obeys  $1/8 < \tau / \Delta_{\Gamma, Y}^2 < 3/8$ ;
- The short trajectory length  $T_{\text{clustering}} = NG$ , where

$$G \geq C_1 \cdot \log \frac{|\mathcal{M}_{\text{clustering}}|}{\delta}, \quad N \geq C_2 \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\Delta_{\Gamma, Y}^2} + 1 \right) \frac{1}{1-\rho} \log \left( \left( \frac{\Gamma_{\max}}{\Delta_{\Gamma, Y}} + 2 \right) \frac{d \kappa_A T_{\text{total}}}{\delta} \right);$$

- The subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}_{1 \leq i \leq d}$  satisfy that, for all  $1 \leq i \leq d$  and  $1 \leq k \leq K$ ,

$$\max \left\{ \left\| (\mathbf{\Gamma}^{(k)})_i - \mathbf{V}_i \mathbf{V}_i^\top (\mathbf{\Gamma}^{(k)})_i \right\|_2, \left\| (\mathbf{Y}^{(k)})_i - \mathbf{U}_i \mathbf{U}_i^\top (\mathbf{Y}^{(k)})_i \right\|_2 \right\} \leq c_3 \frac{\Delta_{\Gamma, Y}}{\sqrt{d}}. \quad (25)$$

Then with probability at least  $1 - \delta$ , Algorithm 3 achieves exact clustering: for all  $m_1, m_2 \in \mathcal{M}_{\text{clustering}}$ ,  $S_{m_1, m_2} = 1$  if and only if the  $m_1$ -th and  $m_2$ -th trajectories are generated by the same model, i.e. they have the same label  $k_{m_1} = k_{m_2}$ .

**Least squares and covariance estimation.** The next result controls the model estimation errors of Algorithm 4, under the assumption that every cluster is pure.

**Theorem B.3.** Consider the model (1) under the assumptions in Section 3.1. There exist universal constants  $C_1, C_2, C_3 > 0$  such that the following holds. Let  $\{\mathcal{C}_k\}_{1 \leq k \leq K}$  be subsets of  $\mathcal{M}_{\text{clustering}} \cup \mathcal{M}_{\text{classification}}$  such that for all  $1 \leq k \leq K$ ,  $\mathcal{C}_k$  contains only short trajectories generated by model  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$ , namely  $k_m = k$  for all  $m \in \mathcal{C}_k$ . Suppose that for all  $m \in \mathcal{M}_{\text{clustering}} \cup \mathcal{M}_{\text{classification}}$ ,  $T_m \geq 4$ , and for all  $1 \leq k \leq K$ ,

$$T_{\text{total}}^{(k)} := \sum_{m \in \mathcal{C}_k} T_m \geq C_1 d \kappa_w^2 \iota, \quad \text{where } \iota := \log \left( \frac{\Gamma_{\max} d \kappa_A T_{\text{total}}}{W_{\min} \delta} \right).$$

Let  $\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}$  be computed by (22) in Algorithm 4. Then with probability at least  $1 - \delta$ , one has

$$\left\| \widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(k)} \right\| \leq C_2 \sqrt{\frac{d \kappa_w \iota}{T_{\text{total}}^{(k)}}}, \quad \frac{\left\| \widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)} \right\|}{\left\| \mathbf{W}^{(k)} \right\|} \leq C_3 \sqrt{\frac{d \iota}{T_{\text{total}}^{(k)}}}, \quad 1 \leq k \leq K.$$

Our proof (postponed to Appendix C.4) is based on the techniques of (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019), but with two major differences. First, the authors of (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019) consider the setting where ordinary least squares is applied to a single continuous trajectory generated by a single LDS model; this is not the case for our setting, and thus our proof and results are different from theirs. Second, the noise covariance matrix  $\mathbf{W}$  is assumed to be  $\sigma^2 \mathbf{I}_d$  in (Simchowitz et al., 2018; Sarkar & Rakhlin, 2019), while in our case,  $\{\mathbf{W}^{(k)}\}_{1 \leq k \leq K}$  are unknown and need to be estimated.

**Classification.** Our last theorem shows that Algorithm 5 correctly classifies all trajectories in  $\mathcal{M}_{\text{classification}}$ , as long as the coarse models are sufficiently accurate and the short trajectory lengths are large enough; these conditions are slightly different for Cases 0 and 1 defined in (8). See Appendix C.5 for the proof.

**Theorem B.4.** Consider the model (1) under the assumptions in Section 3.1. There exist universal constants  $c_1, c_2, C_3 > 0$  such that the following holds. Suppose that we run Algorithm 5 with data  $\{\mathbf{X}_m\}_{m \in \mathcal{M}_{\text{classification}}}$  and independent coarse models  $\{\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}\}_{1 \leq k \leq K}$  satisfying  $\left\| \widehat{\mathbf{A}}^{(k)} - \mathbf{A}^{(k)} \right\| \leq \epsilon_A, \left\| \widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)} \right\| \leq \epsilon_W$  for all  $k$ . Then with probability at least  $1 - \delta$ , Algorithm 5 correctly classifies all trajectories in  $\mathcal{M}_{\text{classification}}$ , provided that

$$\text{For Case 0: } \epsilon_A \leq c_1 \Delta_{A, W} \sqrt{\frac{W_{\min}}{\Gamma_{\max} \kappa_{w, \text{cross}} (d + \iota)}}, \quad \frac{\epsilon_W}{W_{\min}} \leq c_2 \cdot \min \left\{ 1, \frac{\Delta_{A, W}}{\sqrt{\kappa_{w, \text{cross}} d}} \right\}, \quad (26a)$$

$$T_m \geq C_3 \left( \kappa_w^2 + \frac{\kappa_{w, \text{cross}}^6}{\Delta_{A, W}^2} \right) \iota^2, \quad m \in \mathcal{M}_{\text{classification}}; \quad (26b)$$

$$\text{For Case 1: } \epsilon_A \leq c_1 \Delta_{A, W} \sqrt{\frac{W_{\min}}{\Gamma_{\max} \kappa_{w, \text{cross}} \kappa_A^2 (d + \iota)}}, \quad \frac{\epsilon_W}{W_{\min}} \leq c_2 \cdot \min \left\{ 1, \frac{\Delta_{A, W}}{\sqrt{\kappa_{w, \text{cross}} d}} \right\}, \quad (26c)$$

$$T_m \geq C_3 \left( \kappa_w^2 + \frac{\kappa_{w, \text{cross}}^6}{\Delta_{A, W}^2} \right) \iota^2 + \frac{1}{1 - \rho} \log(2 \kappa_A), \quad m \in \mathcal{M}_{\text{classification}}, \quad (26d)$$

where  $\iota := \log \frac{T_{\text{total}}}{\delta}$  is a logarithmic term.

**Proof of Theorems 3.4 and 3.5.** Our main theorems are direct implications of the above guarantees for the individual steps.

- *Stage 1:* To begin with, according to Theorem B.1, if the condition (9a) on  $T_{\text{subspace}}$  and  $T_{\text{total,subspace}}$  hold, then the subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  output by Line 3 of Algorithm 1 are guaranteed to satisfy the error bounds (25) required by Theorem B.2. This together with the condition (9b) on  $T_{\text{clustering}}$  ensures exact clustering in Line 4.
- *Stage 2:* Based on this, if we further know that  $T_{\text{total,clustering}}$  obeys condition (9b) for Case 0 or (11a) for Case 1, then Theorem B.3 tells us that the coarse model estimation in Line 5 satisfies the error bounds (26a) or (26c) required by Theorem B.4. Together with the assumption (9c) or (11b) on  $T_{\text{classification}}$ , this guarantees exact classification in Line 7 of Algorithm 1, for either Case 0 or 1. At the end, the final model estimation errors (10) follow immediately from Theorem B.3.

Note that all the statements above are high-probability guarantees; it suffices to take the union bound over these steps, so that the performance guarantees of Theorems 3.4 and 3.5 hold with probability at least  $1 - \delta$ . This finishes the proof of our main theorems.

*Remark B.5.* The step of subspace estimation in Algorithm 1 is non-essential and optional; it allows for a smaller  $T_{\text{clustering}}$ , but comes at the price of complicating the overall algorithm. For practitioners who prefer a simpler algorithm, they might simply remove this step (i.e. Line 3 of Algorithm 1), and replace the rank- $K$  subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  with  $\mathbf{I}_d$  (i.e. no dimensionality reduction for the clustering step). The theoretical guarantees continue to hold with minor modification: in Corollary 3.7, one simply needs to remove the conditions on  $T_{\text{subspace}}$ ,  $T_{\text{total,subspace}}$ , and in the condition for  $T_{\text{clustering}}$ , replace the factor  $\sqrt{K/d}$  (where  $K$  is due to dimensionality reduction) with  $\sqrt{d/d} = 1$ . This is one example regarding how our modular algorithms and theoretical analysis can be easily modified and adapted to accommodate different situations.

## C. Proofs for Appendix B

This section starts with some preliminaries about linear dynamical systems that will be helpful later. Then, it provides the main proofs for the theorems in Appendix B.

### C.1. Preliminaries

**Truncation of autocovariance.** Recall the notation  $\mathbf{\Gamma}^{(k)} = \sum_{i=0}^{\infty} \mathbf{A}^{(k)i} \mathbf{W}^{(k)} (\mathbf{A}^{(k)i})^\top$  from (3) and (5). We add a subscript  $t$  to represent its  $t$ -step truncation:

$$\mathbf{\Gamma}_t^{(k)} := \sum_{i=0}^{t-1} \mathbf{A}^{(k)i} \mathbf{W}^{(k)} (\mathbf{A}^{(k)i})^\top. \quad (27)$$

Also recall the assumption of exponential stability in Assumption 3.1, namely  $\|(\mathbf{A}^{(k)})^t\| \leq \kappa_A \rho^t$ . As a result,  $\mathbf{\Gamma}_t^{(k)}$  is close to  $\mathbf{\Gamma}^{(k)}$ :

$$\begin{aligned} \mathbf{0} &\preceq \mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}_t^{(k)} = \sum_{i=t}^{\infty} \mathbf{A}^{(k)i} \mathbf{W}^{(k)} (\mathbf{A}^{(k)i})^\top = \mathbf{A}^{(k)t} \mathbf{\Gamma}^{(k)} (\mathbf{A}^{(k)t})^\top, \\ \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}_t^{(k)}\| &\leq \|\mathbf{\Gamma}^{(k)}\| \|\mathbf{A}^{(k)t}\|^2 \leq \|\mathbf{\Gamma}^{(k)}\| \kappa_A^2 \rho^{2t} \leq \Gamma_{\max} \kappa_A^2 \rho^{2t}. \end{aligned} \quad (28)$$

Moreover, let  $\mathbf{Y}_t^{(k)} := \mathbf{A}^{(k)} \mathbf{\Gamma}_t^{(k)}$ , then  $\mathbf{Y}_t^{(k)}$  is also close to  $\mathbf{Y}^{(k)}$ :

$$\|\mathbf{Y}^{(k)} - \mathbf{Y}_t^{(k)}\| \leq \|\mathbf{A}^{(k)}\| \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}_t^{(k)}\| \leq \Gamma_{\max} \kappa_A^3 \rho^{2t}. \quad (29)$$

**“Independent version” of states.** Given some mixing time  $t_{\text{mix}}$ , we define  $\tilde{\mathbf{x}}_{m,t}(t_{\text{mix}})$  as the  $(t_{\text{mix}} - 1)$ -step approximation of  $\mathbf{x}_{m,t}$ :

$$\tilde{\mathbf{x}}_{m,t} = \tilde{\mathbf{x}}_{m,t}(t_{\text{mix}}) := \sum_{i=0}^{t_{\text{mix}}-2} (\mathbf{A}^{(k_m)})^i \mathbf{w}_{m,t-i-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k_m)}), \quad t_{\text{mix}} \leq t \leq T_m, \quad 1 \leq m \leq M. \quad (30)$$

Since  $\tilde{\mathbf{x}}_{m,t}$  consists of only the most recent noise vectors, it is independent of the history up to  $\mathbf{x}_{m,t-t_{\text{mix}}+1}$ . Our proofs for Stage 1 will rely on this “independent version”  $\{\tilde{\mathbf{x}}_{m,t}\}$  of states  $\{\mathbf{x}_{m,t}\}$ ; the basic idea is that, for an appropriately chosen  $t_{\text{mix}}$ , a trajectory of length  $T$  can be regarded as a collection of  $T/t_{\text{mix}}$  independent samples. We will often use the notation  $\tilde{\cdot}$  to represent the “independent version” of other variables as well.

**Boundedness of states.** The following lemma provides upper bounds for  $\{\|\tilde{\mathbf{x}}_{m,t}\|_2\}$  and  $\{\|\mathbf{x}_{m,t}\|_2\}$ . This will help to control the effects of mixing errors and model estimation errors in the analyses later.

**Lemma C.1** (Bounded states). *Consider the model (1) under the assumptions in Sections 2 and 3.1. Fix any  $t_{\text{mix}} \geq 3$ . Then with probability at least  $1 - \delta$ , we have  $\|\tilde{\mathbf{x}}_{m,t}\|_2 \leq C_0 \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $1 \leq m \leq M, t_{\text{mix}} \leq t \leq T_m$ , where  $C_0 > 0$  is some universal constant; moreover, for both cases of initial states defined in (8), all states  $\{\{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}\}_{1 \leq m \leq M}$  are bounded throughout:*

- *Case 0:* with probability at least  $1 - \delta$ , we have  $\|\mathbf{x}_{m,t}\|_2 \leq C_0 \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $m, t$ .
- *Case 1:* suppose that  $t_{\text{mix}} \geq \frac{1}{1-\rho} \log(\sqrt{2}\kappa_A)$ , and  $T_m \geq t_{\text{mix}}$  for all  $m$ , then with probability at least  $1 - \delta$ , we have  $\|\mathbf{x}_{m,t}\|_2 \leq 3C_0 \kappa_A \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $m, t$ , and  $\|\mathbf{x}_{m,t}\|_2 \leq 2C_0 \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $t \geq t_{\text{mix}}$  or  $t = 0$ .

*Proof.* First, recall from Corollary 7.3.3 of (Vershynin, 2018) that, if random vector  $\mathbf{a} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , then for all  $u \geq 0$ , we have  $\mathbb{P}(\|\mathbf{a}\|_2 \geq 2\sqrt{d} + u) \leq 2 \exp(-cu^2)$ . Since  $\tilde{\mathbf{x}}_{m,t} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k_m)})$ , where  $\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k_m)} \preceq \mathbf{\Gamma}^{(k_m)} \preceq \Gamma_{\max} \mathbf{I}_d$ , we have  $\mathbb{P}(\|\tilde{\mathbf{x}}_{m,t}\|_2 \geq \sqrt{\Gamma_{\max}}(2\sqrt{d} + u)) \leq 2 \exp(-cu^2)$ . Taking the union bound, we have

$$\begin{aligned} & \mathbb{P}\left(\text{there exists } m, t \text{ such that } \|\tilde{\mathbf{x}}_{m,t}\|_2 \geq \sqrt{\Gamma_{\max}}(2\sqrt{d} + u)\right) \\ & \leq \sum_{m=1}^M \sum_{t=t_{\text{mix}}}^{T_m} \mathbb{P}\left(\|\tilde{\mathbf{x}}_{m,t}\|_2 \geq \sqrt{\Gamma_{\max}}(2\sqrt{d} + u)\right) \leq 2T_{\text{total}} \exp(-cu^2) \leq \delta, \end{aligned}$$

where the last inequality holds if we pick  $u \geq \sqrt{\frac{1}{c} \log \frac{2T_{\text{total}}}{\delta}}$ . This finishes the proof of the first claim in the lemma. Next, we prove the boundedness of  $\{\|\mathbf{x}_{m,t}\|_2\}$ .

- *Case 0:* It is easy to check that  $\mathbf{x}_{m,t} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_t^{(k_m)})$ , where  $\mathbf{\Gamma}_t^{(k_m)} \preceq \mathbf{\Gamma}^{(k_m)} \preceq \Gamma_{\max} \mathbf{I}_d$ . The boundedness of  $\{\|\mathbf{x}_{m,t}\|_2\}$  can be proved by a similar argument as before, which we omit for brevity.
- *Case 1:* Define  $\xi_{m,t} := \mathbf{x}_{m,t} - (\mathbf{A}^{(k_m)})^t \mathbf{x}_{m,0} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_t^{(k_m)})$ . By a similar argument as before, we have with probability at least  $1 - \delta$ ,  $\|\xi_{m,t}\|_2 \leq C_0 \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $m, t$ . Moreover, for any  $t \geq t_{\text{mix}} \geq \frac{1}{1-\rho} \log(\sqrt{2}\kappa_A)$ , we have  $\|(\mathbf{A}^{(k_m)})^t\| \leq \kappa_A \rho^t \leq 1/2$ . With this in place, we have

$$\mathbf{x}_{m+1,0} = \mathbf{x}_{m,T_m} = (\mathbf{A}^{(k_m)})^{T_m} \mathbf{x}_{m,0} + \xi_{m,T_m}, \quad \text{and thus}$$

$$\|\mathbf{x}_{m+1,0}\|_2 \leq \|(\mathbf{A}^{(k_m)})^{T_m}\| \|\mathbf{x}_{m,0}\|_2 + \|\xi_{m,T_m}\|_2 \leq \frac{1}{2} \|\mathbf{x}_{m,0}\|_2 + C_0 \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}.$$

Recall the assumption that  $\mathbf{x}_{1,0} = \mathbf{0}$ ; by induction, we have  $\|\mathbf{x}_{m,0}\|_2 \leq 2C_0 \sqrt{\Gamma_{\max}(d + \log(T_{\text{total}}/\delta))}$  for all  $1 \leq m \leq M$ . Now, we have for all  $m, t$ ,

$$\begin{aligned} \|\mathbf{x}_{m,t}\|_2 & \leq \|(\mathbf{A}^{(k_m)})^t \mathbf{x}_{m,0}\|_2 + \|\xi_{m,t}\|_2 \\ & \leq \kappa_A \|\mathbf{x}_{m,0}\|_2 + C_0 \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})} \leq 3C_0 \kappa_A \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}; \end{aligned}$$

moreover, for  $t \geq t_{\text{mix}}$ , since  $\|(\mathbf{A}^{(k_m)})^t\| \leq 1/2$ , we obtain a better bound

$$\begin{aligned} \|\mathbf{x}_{m,t}\|_2 & \leq \|(\mathbf{A}^{(k_m)})^t \mathbf{x}_{m,0}\|_2 + \|\xi_{m,t}\|_2 \\ & \leq \frac{1}{2} \|\mathbf{x}_{m,0}\|_2 + C_0 \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})} \leq 2C_0 \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}. \end{aligned}$$

This shows the boundedness of  $\{\|\mathbf{x}_{m,t}\|_2\}$  and completes our proof of the lemma.  $\square$

## C.2. Proof of Theorem B.1

Theorem B.1 is an immediate consequence of Lemmas C.2 and C.3 below. The former shows the concentration of  $\widehat{\mathbf{H}}_i, \widehat{\mathbf{G}}_i$  around the targeted low-rank matrices  $\mathbf{H}_i, \mathbf{G}_i$ , while the latter is a result of perturbation analysis.

**Lemma C.2.** *Under the setting of Theorem B.1, with probability at least  $1 - \delta$ , we have for all  $1 \leq i \leq d$ ,*

$$\max \left\{ \|\widehat{\mathbf{H}}_i - \mathbf{H}_i\|, \|\widehat{\mathbf{G}}_i - \mathbf{G}_i\| \right\} \lesssim \Gamma_{\max}^2 \sqrt{\frac{t_{\text{mix}} d}{T_{\text{total, subspace}}} \log^3 \frac{T_{\text{total}} d}{\delta}}.$$

**Lemma C.3.** *Consider the matrix  $\mathbf{M}_\star = \sum_{k=1}^K p^{(k)} \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \in \mathbb{R}^{d \times d}$ , where  $0 < p^{(k)} < 1$ ,  $\sum_{k=1}^K p^{(k)} = 1$ , and  $\mathbf{y}^{(k)} \in \mathbb{R}^d$ . Let  $\mathbf{M}$  be symmetric and satisfy  $\|\mathbf{M} - \mathbf{M}_\star\| \leq \epsilon$ , and  $\mathbf{U} \in \mathbb{R}^{d \times K}$  be the top- $K$  eigenspace of  $\mathbf{M}$ . Then we have*

$$\sum_{k=1}^K p^{(k)} \|\mathbf{y}^{(k)} - \mathbf{U} \mathbf{U}^\top \mathbf{y}^{(k)}\|_2^2 \leq 2K\epsilon, \quad (31)$$

and for all  $1 \leq k \leq K$ , it holds that

$$\|\mathbf{y}^{(k)} - \mathbf{U} \mathbf{U}^\top \mathbf{y}^{(k)}\|_2 \leq \min \left\{ \left( \frac{2K\epsilon}{p^{(k)}} \right)^{1/2}, \frac{2\epsilon}{\lambda_{\min}(\mathbf{M})} \|\mathbf{y}^{(k)}\|_2, \sqrt{2} \left( \frac{\epsilon}{p^{(k)}} \|\mathbf{y}^{(k)}\|_2 \right)^{1/3} \right\}. \quad (32)$$

For our main analyses in Sections 3.2 and B, we choose the first term on the right-hand side of (32).

### C.2.1. PROOF OF LEMMA C.2

We first analyze the idealized case with i.i.d. samples; then we make a connection between this i.i.d. case and the actual case of mixed LDSs, by utilizing the mixing property of linear dynamical systems. We prove the result of Lemma C.2 for  $\|\widehat{\mathbf{G}}_i - \mathbf{G}_i\|$  only, since the analysis for  $\|\widehat{\mathbf{H}}_i - \mathbf{H}_i\|$  is mostly the same (and simpler in fact).

**Step 1: the idealized i.i.d. case.** With some abuse of notation, suppose that for all  $1 \leq m \leq M$ , we have for some  $k_m \in \{1, \dots, K\}$ ,

$$\begin{aligned} \mathbf{x}_{m,t}, \mathbf{z}_{m,t} &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \widetilde{\mathbf{\Gamma}}^{(k_m)}), & \mathbf{w}_{m,t}, \mathbf{v}_{m,t} &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{W}^{(k_m)}), \\ \mathbf{x}_{m,t}' &= \mathbf{A}^{(k_m)} \mathbf{x}_{m,t} + \mathbf{w}_{m,t}, & \mathbf{z}_{m,t}' &= \mathbf{A}^{(k_m)} \mathbf{z}_{m,t} + \mathbf{v}_{m,t}, \quad 1 \leq t \leq N, \end{aligned}$$

where for all  $1 \leq k \leq K$ , it holds that  $\mathbf{W}^{(k)}, \widetilde{\mathbf{\Gamma}}^{(k)} \preceq \mathbf{\Gamma}^{(k)} \preceq \Gamma_{\max} \mathbf{I}_d$ . Notice that  $\text{cov}(\mathbf{x}_{m,t}') = \mathbf{A}^{(k_m)} \widetilde{\mathbf{\Gamma}}^{(k_m)} (\mathbf{A}^{(k_m)})^\top + \mathbf{W}^{(k_m)} \preceq \mathbf{A}^{(k_m)} \mathbf{\Gamma}^{(k_m)} (\mathbf{A}^{(k_m)})^\top + \mathbf{W}^{(k_m)} = \mathbf{\Gamma}^{(k_m)} \preceq \Gamma_{\max} \mathbf{I}_d$ . Consider the i.i.d. version of matrix  $\widehat{\mathbf{G}}_i$  and its expectation  $\mathbf{G}_i$  defined as follows:

$$\widehat{\mathbf{G}}_i := \frac{1}{MN} \sum_{m=1}^M \sum_{t=1}^N \left( (\mathbf{x}_{m,t}')_i \mathbf{x}_{m,t} \right) \left( (\mathbf{z}_{m,t}')_i \mathbf{z}_{m,t} \right)^\top, \quad \mathbf{G}_i = \sum_{k=1}^K p^{(k)} (\widetilde{\mathbf{Y}}^{(k)})_i (\widetilde{\mathbf{Y}}^{(k)})_i^\top,$$

where  $(\widetilde{\mathbf{Y}}^{(k)})_i$  is the transpose of the  $i$ -th row of  $\widetilde{\mathbf{Y}}^{(k)} := \mathbf{A}^{(k)} \widetilde{\mathbf{\Gamma}}^{(k)}$ . For this i.i.d. setting, we claim that, if the i.i.d. sample size  $MN$  satisfies  $MN \gtrsim d \cdot \log(MNd/\delta)$ , then with probability at least  $1 - \delta$ ,

$$\|\widehat{\mathbf{G}}_i - \mathbf{G}_i\| \lesssim \Gamma_{\max}^2 \sqrt{\frac{d}{MN} \log^3 \frac{MNd}{\delta}}, \quad 1 \leq i \leq d. \quad (33)$$

This claim can be proved by a standard covering argument with truncation; we will provide a proof later for completeness.

**Step 2: back to the actual case of mixed LDSs.** Now we turn to the analysis of  $\widehat{\mathbf{G}}_i$  defined in (14) versus its expectation  $\mathbf{G}_i$  defined in (15b), for the mixed LDSs setting. We first show that  $\widehat{\mathbf{G}}_i$  can be written as a *weighted average* of some matrices, each of which can be further decomposed into an i.i.d. *part* (as in Step 1) plus a *negligible mixing error term*. Then we analyze each term in the decomposition, and finally put pieces together to show that  $\widehat{\mathbf{G}}_i \approx \mathbf{G}_i$ .

**Step 2.1: decomposition of the index set  $\Omega_1 \times \Omega_2$ .** Recall the definition of index sets  $\Omega_1, \Omega_2$  in Algorithm 2. Denote the first index of  $\Omega_1$  (resp.  $\Omega_2$ ) as  $\tau_1 + 1$  (resp.  $\tau_2 + 1$ ), and let  $\Delta := \tau_2 - \tau_1$  be their distance. Also denote  $N := |\Omega_1| = |\Omega_2| \asymp T_{\text{subspace}}$ . For any  $t \in \Omega_1$  and  $1 \leq j \leq N$ , define

$$s_j(t) := \text{Cycle}(t + \Delta + j; \Omega_2) = \begin{cases} t + \Delta + j & \text{if } t + \Delta + j \leq \tau_2 + N, \\ t + \Delta + j - N & \text{otherwise,} \end{cases}$$

where  $\text{Cycle}(i; \Omega)$  represents the cyclic indexing of value  $i$  on set  $\Omega$ . Then we have

$$\Omega_1 \times \Omega_2 = \{(t_1, t_2), t_1 \in \Omega_1, t_2 \in \Omega_2\} = \cup_{j=1}^N \{(t, s_j(t)), t \in \Omega_1\}.$$

We further define

$$\mathcal{S}_\tau := \{\tau_1 + \tau + f \cdot t_{\text{mix}} : f \geq 0, \tau + f \cdot t_{\text{mix}} \leq N\}, \quad 1 \leq \tau \leq t_{\text{mix}},$$

so that  $\Omega_1 = \cup_{\tau=1}^{t_{\text{mix}}} \mathcal{S}_\tau$ . Notice that for each  $\tau$ , the elements of  $\mathcal{S}_\tau$  are at least  $t_{\text{mix}}$  far apart. Putting together, we have

$$\Omega_1 \times \Omega_2 = \cup_{j=1}^N \cup_{\tau=1}^{t_{\text{mix}}} \{(t, s_j(t)), t \in \mathcal{S}_\tau\}. \quad (34)$$

**Step 2.2: decomposition of  $\widehat{\mathbf{G}}_i$ .** In the remaining proof, we denote  $\mathbf{x}_{m,t}' := \mathbf{x}_{m,t+1}$  for notational consistency. Using the decomposition (34) of  $\Omega_1 \times \Omega_2$ , we can rewrite  $\widehat{\mathbf{G}}_i$  defined in (14) as a weighted average of  $N t_{\text{mix}}$  matrices:

$$\widehat{\mathbf{G}}_i = \frac{1}{|\mathcal{M}_{\text{subspace}}|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \frac{1}{N^2} \sum_{(t_1, t_2) \in \Omega_1 \times \Omega_2} (\mathbf{x}'_{m,t_1})_i \mathbf{x}_{m,t_1} \cdot (\mathbf{x}'_{m,t_2})_i \mathbf{x}_{m,t_2}^\top = \sum_{j=1}^N \sum_{\tau=1}^{t_{\text{mix}}} \frac{|\mathcal{S}_\tau|}{N^2} \cdot \mathbf{F}_{i,j,\tau},$$

where

$$\mathbf{F}_{i,j,\tau} := \frac{1}{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \sum_{t \in \mathcal{S}_\tau} (\mathbf{x}_{m,t}')_i \mathbf{x}_{m,t} \cdot (\mathbf{x}'_{m,s_j(t)})_i \mathbf{x}_{m,s_j(t)}^\top. \quad (35)$$

Recalling the definition of  $\tilde{\mathbf{x}}_{m,t}$  in (30) and  $\mathbf{\Gamma}_t^{(k)}$  in (27), we have

$$\mathbf{x}_{m,t} = \tilde{\mathbf{x}}_{m,t} + \underbrace{(\mathbf{A}^{(k_m)})^{t_{\text{mix}}-1} \mathbf{x}_{m,t-t_{\text{mix}}+1}}_{=:\boldsymbol{\delta}_{m,t}} = \tilde{\mathbf{x}}_{m,t} + \boldsymbol{\delta}_{m,t},$$

where

$$\|\boldsymbol{\delta}_{m,t}\|_2 \leq \|(\mathbf{A}^{(k_m)})^{t_{\text{mix}}-1}\| \cdot \|\mathbf{x}_{m,t-t_{\text{mix}}+1}\|_2 \leq \kappa_A \rho^{t_{\text{mix}}-1} \|\mathbf{x}_{m,t-t_{\text{mix}}+1}\|_2, \quad (36)$$

and  $\tilde{\mathbf{x}}_{m,t} \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)})$  is independent of  $\boldsymbol{\delta}_{m,t}$ . Moreover,

$$\mathbf{x}_{m,t}' = \mathbf{A}^{(k_m)} \mathbf{x}_{m,t} + \mathbf{w}_{m,t} = \tilde{\mathbf{x}}_{m,t}' + \mathbf{A}^{(k_m)} \boldsymbol{\delta}_{m,t}, \quad \text{where } \tilde{\mathbf{x}}_{m,t}' := \mathbf{A}^{(k_m)} \tilde{\mathbf{x}}_{m,t} + \mathbf{w}_{m,t}.$$

We can rewrite  $\mathbf{x}_{m,s_j(t)} = \tilde{\mathbf{x}}_{m,s_j(t)} + \boldsymbol{\delta}_{m,s_j(t)}$  and  $\mathbf{x}'_{m,s_j(t)} = \tilde{\mathbf{x}}_{m,s_j(t)}' + \mathbf{A}^{(k_m)} \boldsymbol{\delta}_{m,s_j(t)}$  in a similar manner. Putting this back to (35), one has

$$\begin{aligned} \mathbf{F}_{i,j,\tau} &= \frac{1}{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \sum_{t \in \mathcal{S}_\tau} (\mathbf{x}_{m,t}')_i \mathbf{x}_{m,t} \cdot (\mathbf{x}'_{m,s_j(t)})_i \mathbf{x}_{m,s_j(t)}^\top \\ &= \frac{1}{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \sum_{t \in \mathcal{S}_\tau} \\ &\quad (\tilde{\mathbf{x}}_{m,t}' + \mathbf{A}^{(k_m)} \boldsymbol{\delta}_{m,t})_i (\tilde{\mathbf{x}}_{m,t} + \boldsymbol{\delta}_{m,t}) \cdot (\tilde{\mathbf{x}}_{m,s_j(t)}' + \mathbf{A}^{(k_m)} \boldsymbol{\delta}_{m,s_j(t)})_i (\tilde{\mathbf{x}}_{m,s_j(t)} + \boldsymbol{\delta}_{m,s_j(t)})^\top \\ &= \frac{1}{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau|} \sum_{m \in \mathcal{M}_{\text{subspace}}} \sum_{t \in \mathcal{S}_\tau} \underbrace{(\tilde{\mathbf{x}}_{m,t}'_i \tilde{\mathbf{x}}_{m,t} \cdot (\tilde{\mathbf{x}}_{m,s_j(t)}')_i \tilde{\mathbf{x}}_{m,s_j(t)})^\top}_{=:\tilde{\mathbf{F}}_{i,j,\tau}} + \boldsymbol{\Delta}_{i,j,\tau}, \end{aligned} \quad (37)$$

where  $\boldsymbol{\Delta}_{i,j,\tau}$  contains all the  $\{\boldsymbol{\delta}_{m,t}\}$  terms in the expansion. The key observation here is that, by our definition of index set  $\mathcal{S}_\tau$ , the  $\{\tilde{\mathbf{x}}_{m,t}\}$  terms in  $\tilde{\mathbf{F}}_{i,j,\tau}$  are independent, and thus we can utilize our earlier analysis of the i.i.d. case in Step 1 to study  $\tilde{\mathbf{F}}_{i,j,\tau}$ .



**Step 2.3: analysis for each term of the decomposition.** Towards showing  $\widehat{\mathbf{G}}_i \approx \mathbf{G}_i$ , we prove in the following that  $\widetilde{\mathbf{F}}_{i,j,\tau}$  concentrates around its expectation  $\widetilde{\mathbf{G}}_i$ , which is close to  $\mathbf{G}_i$ ; moreover, the error term  $\Delta_{i,j,\tau}$  becomes exponentially small as  $t_{\text{mix}}$  grows. More specifically, we have the following:

- Recall the notation  $\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} = \mathbf{A}^{(k)} \Gamma_{t_{\text{mix}}-1}^{(k)}$ . It holds that

$$\mathbb{E}[\widetilde{\mathbf{F}}_{i,j,\tau}] = \widetilde{\mathbf{G}}_i := \sum_{k=1}^K p^{(k)} (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i^\top.$$

According to our result (33) for the i.i.d. case, for fixed  $j, \tau$ , we have with probability at least  $1 - \delta$ ,

$$\begin{aligned} 1 \leq i \leq d, \quad \|\widetilde{\mathbf{F}}_{i,j,\tau} - \widetilde{\mathbf{G}}_i\| &\lesssim \Gamma_{\max}^2 \sqrt{\frac{d}{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau|} \log^3 \frac{|\mathcal{M}_{\text{subspace}}| \cdot |\mathcal{S}_\tau| d}{\delta}} \\ &\lesssim \Gamma_{\max}^2 \sqrt{\frac{t_{\text{mix}} d}{T_{\text{total,subspace}}} \log^3 \frac{T_{\text{total,subspace}} d}{\delta}}. \end{aligned} \quad (38)$$

- Recall from (29) that  $\|\mathbf{Y}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(k)}\| \leq \Gamma_{\max} \kappa_A^3 \rho^{2(t_{\text{mix}}-1)}$ . Therefore,

$$\begin{aligned} &\|(\mathbf{Y}^{(k)})_i (\mathbf{Y}^{(k)})_i^\top - (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i^\top\| \\ &\leq (\|(\mathbf{Y}^{(k)})_i\|_2 + \|(\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i\|_2) \|(\mathbf{Y}^{(k)})_i - (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i\|_2 \\ &\leq (2\|\mathbf{Y}^{(k)}\| + \|\mathbf{Y}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(k)}\|) \|\mathbf{Y}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(k)}\| \\ &\leq (2\Gamma_{\max} \kappa_A + \Gamma_{\max} \kappa_A^3 \rho^{2(t_{\text{mix}}-1)}) \cdot \Gamma_{\max} \kappa_A^3 \rho^{2(t_{\text{mix}}-1)} \\ &= (2 + \kappa_A^2 \rho^{2(t_{\text{mix}}-1)}) \cdot \Gamma_{\max}^2 \kappa_A^4 \rho^{2(t_{\text{mix}}-1)} \leq 3\Gamma_{\max}^2 \kappa_A^4 \rho^{2(t_{\text{mix}}-1)}, \end{aligned}$$

where we use the mild assumption that  $t_{\text{mix}} \geq 1 + \frac{\log \kappa_A}{1-\rho}$ , and the fact that  $\|\mathbf{Y}^{(k)}\|, \|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)}\| \leq \Gamma_{\max} \kappa_A$ . Consequently,

$$\|\mathbf{G}_i - \widetilde{\mathbf{G}}_i\| \leq \sum_{k=1}^K p^{(k)} \|(\mathbf{Y}^{(k)})_i (\mathbf{Y}^{(k)})_i^\top - (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i (\mathbf{Y}_{t_{\text{mix}}-1}^{(k)})_i^\top\| \leq 3\Gamma_{\max}^2 \kappa_A^4 \rho^{2(t_{\text{mix}}-1)}. \quad (39)$$

- By Lemma C.1, if  $t_{\text{mix}} \gtrsim \frac{1}{1-\rho} \log(2\kappa_A)$ , then we have with probability at least  $1 - \delta$ , all the  $\mathbf{x}_{m,t}$ 's and  $\widetilde{\mathbf{x}}_{m,t}$ 's involved in the definition of  $\Delta_{i,j,\tau}$  in (37) have  $\ell_2$  norm bounded by  $\sqrt{\Gamma_{\max}} \text{poly}(d, \kappa_A, \log(T_{\text{total}}/\delta))$ . This together with the upper bound on  $\|\delta_{m,t}\|_2$  in (36) implies that for all  $i, j, \tau$ , it holds that

$$\|\Delta_{i,j,\tau}\| \leq \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}. \quad (40)$$

**Step 2.4: putting pieces together.** With (38), (39) and (40) in place and taking the union bound, we have with probability at least  $1 - \delta$ , for all  $1 \leq i \leq d$ ,

$$\begin{aligned} \|\widehat{\mathbf{G}}_i - \mathbf{G}_i\| &= \left\| \sum_{j=1}^N \sum_{\tau=1}^{t_{\text{mix}}} \frac{|\mathcal{S}_\tau|}{N^2} \cdot \mathbf{F}_{i,j,\tau} - \mathbf{G}_i \right\| \leq \max_{j,\tau} \|\widetilde{\mathbf{F}}_{i,j,\tau} - \widetilde{\mathbf{G}}_i\| + \|\widetilde{\mathbf{G}}_i - \mathbf{G}_i\| + \max_{j,\tau} \|\Delta_{i,j,\tau}\| \\ &\lesssim \Gamma_{\max}^2 \sqrt{\frac{t_{\text{mix}} d}{T_{\text{total,subspace}}} \log^3 \frac{T_{\text{total,subspace}} d}{\delta}} + \Gamma_{\max}^2 \kappa_A^4 \rho^{2(t_{\text{mix}}-1)} + \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{dT_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1} \\ &\lesssim \Gamma_{\max}^2 \sqrt{\frac{t_{\text{mix}} d}{T_{\text{total,subspace}}} \log^3 \frac{T_{\text{total}} d}{\delta}}, \end{aligned}$$

where the last inequality holds if  $t_{\text{mix}} \gtrsim \frac{1}{1-\rho} \log\left(\frac{d\kappa_A T_{\text{total}}}{\delta}\right)$ . This finishes the proof of Lemma C.2.

*Proof of (33).* Define the truncating operator

$$\text{Trunc}(x; D) := x \cdot \mathbb{1}(|x| \leq D), \quad x \in \mathbb{R}, \quad D \geq 0.$$

Consider the following truncated version of  $\widehat{\mathbf{G}}_i$ :

$$\widehat{\mathbf{G}}_i^{\text{Trunc}} := \frac{1}{MN} \sum_{m=1}^M \sum_{t=1}^N \left( \text{Trunc}\left((\mathbf{x}_{m,t}')_i; D_0\right) \mathbf{x}_{m,t} \right) \left( \text{Trunc}\left((\mathbf{z}_{m,t}')_i; D_0\right) \mathbf{z}_{m,t} \right)^\top$$

(the truncating level  $D_0$  will be specified later), and let  $\mathbf{E}_i^{\text{Trunc}} := \mathbb{E}[\widehat{\mathbf{G}}_i^{\text{Trunc}}]$  be its expectation. In the following, we first show that  $\widehat{\mathbf{G}}_i^{\text{Trunc}}$  concentrates around  $\mathbf{E}_i^{\text{Trunc}}$ , and then prove that  $\mathbf{E}_i^{\text{Trunc}} \approx \mathbf{G}_i$ .

- By a standard covering argument, we have

$$\|\widehat{\mathbf{G}}_i^{\text{Trunc}} - \mathbf{E}_i^{\text{Trunc}}\| = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{S}^{d-1}} \mathbf{u}^\top \left( \widehat{\mathbf{G}}_i^{\text{Trunc}} - \mathbf{E}_i^{\text{Trunc}} \right) \mathbf{v} \leq 4 \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{N}_{1/8}} \mathbf{u}^\top \left( \widehat{\mathbf{G}}_i^{\text{Trunc}} - \mathbf{E}_i^{\text{Trunc}} \right) \mathbf{v},$$

where  $\mathcal{N}_{1/8}$  denotes the  $1/8$ -covering of the unit sphere  $\mathcal{S}^{d-1}$  and has cardinality  $|\mathcal{N}_{1/8}| \leq 32^d$ . For fixed  $\mathbf{u}, \mathbf{v} \in \mathcal{N}_{1/8}$ , one has

$$\mathbf{u}^\top \widehat{\mathbf{G}}_i^{\text{Trunc}} \mathbf{v} = \frac{1}{MN} \sum_{m=1}^M \sum_{t=1}^N \text{Trunc}\left((\mathbf{x}_{m,t}')_i; D_0\right) \mathbf{u}^\top \mathbf{x}_{m,t} \cdot \text{Trunc}\left((\mathbf{z}_{m,t}')_i; D_0\right) \mathbf{v}^\top \mathbf{z}_{m,t},$$

where (cf. Chapter 2 of (Vershynin, 2018) for the definitions of subgaussian norm  $\|\cdot\|_{\psi_2}$  and subexponential norm  $\|\cdot\|_{\psi_1}$ )

$$\left| \text{Trunc}\left((\mathbf{x}_{m,t}')_i; D_0\right) \right|, \left| \text{Trunc}\left((\mathbf{z}_{m,t}')_i; D_0\right) \right| \leq D_0, \quad \|\mathbf{u}^\top \mathbf{x}_{m,t}\|_{\psi_2}, \|\mathbf{v}^\top \mathbf{z}_{m,t}\|_{\psi_2} \lesssim \sqrt{\Gamma_{\max}}.$$

Hence

$$\left\| \text{Trunc}\left((\mathbf{x}_{m,t}')_i; D_0\right) \mathbf{u}^\top \mathbf{x}_{m,t} \cdot \text{Trunc}\left((\mathbf{z}_{m,t}')_i; D_0\right) \mathbf{v}^\top \mathbf{z}_{m,t} \right\|_{\psi_1} \lesssim D_0^2 \Gamma_{\max},$$

and by Bernstein's inequality (Corollary 2.8.3 of (Vershynin, 2018)), we have

$$\mathbb{P}\left(\left| \mathbf{u}^\top \left( \widehat{\mathbf{G}}_i^{\text{Trunc}} - \mathbf{E}_i^{\text{Trunc}} \right) \mathbf{v} \right| \geq \tau\right) \leq 2 \exp\left(-c_0 MN \left(\frac{\tau}{D_0^2 \Gamma_{\max}}\right)^2\right)$$

for all  $0 \leq \tau \leq D_0^2 \Gamma_{\max}$ . Taking the union bound over  $\mathbf{u}, \mathbf{v} \in \mathcal{N}_{1/8}$ , we have with probability at least  $1 - \delta/2$ ,

$$\|\widehat{\mathbf{G}}_i^{\text{Trunc}} - \mathbf{E}_i^{\text{Trunc}}\| \lesssim D_0^2 \Gamma_{\max} \sqrt{\frac{d + \log \frac{1}{\delta}}{MN}}, \quad \text{provided that } MN \gtrsim d + \log \frac{1}{\delta}.$$

- Note that

$$\|\mathbf{G}_i - \mathbf{E}_i^{\text{Trunc}}\| = \left\| \sum_{k=1}^K p^{(k)} \mathbb{E}_{\mathbf{x}_t, \mathbf{z}_t \sim \mathcal{N}(0, \widetilde{\Gamma}^{(k)})} \left[ \left( (\mathbf{x}_t')_i (\mathbf{z}_t')_i - \text{Trunc}((\mathbf{x}_t')_i) \text{Trunc}((\mathbf{z}_t')_i) \right) \mathbf{x}_t \mathbf{z}_t^\top \right] \right\|,$$

where for each  $k$ ,

$$\begin{aligned} & \left\| \mathbb{E}_{\mathbf{x}_t, \mathbf{z}_t \sim \mathcal{N}(0, \widetilde{\Gamma}^{(k)})} \left[ \left( (\mathbf{x}_t')_i (\mathbf{z}_t')_i - \text{Trunc}((\mathbf{x}_t')_i) \text{Trunc}((\mathbf{z}_t')_i) \right) \mathbf{x}_t \mathbf{z}_t^\top \right] \right\| \\ &= \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{S}^{d-1}} \left| \mathbb{E}_{\mathbf{x}_t, \mathbf{z}_t \sim \mathcal{N}(0, \widetilde{\Gamma}^{(k)})} \left[ (\mathbf{x}_t')_i (\mathbf{z}_t')_i \mathbf{u}^\top \mathbf{x}_t \mathbf{v}^\top \mathbf{z}_t \cdot \left( 1 - \mathbb{1}(|(\mathbf{x}_t')_i| \leq D_0, |(\mathbf{z}_t')_i| \leq D_0) \right) \right] \right| \end{aligned}$$

$$\begin{aligned}
 &\leq \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{S}^{d-1}} \sqrt{\mathbb{E} \left( (\mathbf{x}_{t'})_i (\mathbf{z}_{t'})_i \mathbf{u}^\top \mathbf{x}_t \mathbf{v}^\top \mathbf{z}_t \right)^2} \sqrt{\mathbb{E} \left( 1 - \mathbb{1}(|(\mathbf{x}_{t'})_i| \leq D_0, |(\mathbf{z}_{t'})_i| \leq D_0) \right)^2} \\
 &\lesssim \Gamma_{\max}^2 \sqrt{\mathbb{P}(|\mathcal{N}(0, \Gamma_{\max})| > D_0)} \\
 &\lesssim \Gamma_{\max}^2 \exp \left( -c_1 \frac{D_0^2}{\Gamma_{\max}} \right).
 \end{aligned}$$

Finally, notice that if the truncating level  $D_0$  is sufficiently large, then we have  $\widehat{\mathbf{G}}_i = \widehat{\mathbf{G}}_i^{\text{Trunc}}$  with high probability. More formally, we have shown that for fixed  $1 \leq i \leq d$ , if  $MN \gtrsim d + \log(1/\delta)$ , then

$$\begin{aligned}
 \mathbb{P} \left( \|\widehat{\mathbf{G}}_i - \mathbf{G}_i\| \lesssim D_0^2 \Gamma_{\max} \sqrt{\frac{d + \log \frac{1}{\delta}}{MN}} + \Gamma_{\max}^2 \exp \left( -c_1 \frac{D_0^2}{\Gamma_{\max}} \right) \right) \\
 \geq 1 - \frac{\delta}{2} - \sum_{m,t} \mathbb{P} \left( |(\mathbf{x}_{m,t'})_i| > D_0 \text{ or } |(\mathbf{z}_{m,t'})_i| > D_0 \right).
 \end{aligned}$$

If we pick the truncating level  $D_0 \asymp \sqrt{\Gamma_{\max} \log(MN/\delta)}$ , then it is easy to check that with probability at least  $1 - \delta$ ,

$$\|\widehat{\mathbf{G}}_i - \mathbf{G}_i\| \lesssim D_0^2 \Gamma_{\max} \sqrt{\frac{d + \log \frac{1}{\delta}}{MN}} \lesssim \Gamma_{\max}^2 \sqrt{\frac{d}{MN} \log^3 \frac{MN}{\delta}}.$$

Taking the union bound over  $1 \leq i \leq d$  finishes our proof of (33).  $\square$

### C.2.2. PROOF OF LEMMA C.3

Define  $\Delta := M - M_*$ , and denote the eigendecomposition of  $M_*$  and  $M$  as  $M_* = \mathbf{U}_* \Lambda_* \mathbf{U}_*^\top$  and  $M = \mathbf{U} \Lambda \mathbf{U}^\top + \mathbf{U}_\perp \Lambda_\perp \mathbf{U}_\perp^\top$ , where diagonal matrix  $\Lambda$  (resp.  $\Lambda_*$ ) contains the top- $K$  eigenvalues of  $M$  (resp.  $M_*$ ). Then we have

$$\begin{aligned}
 \Lambda &= \mathbf{U}^\top M \mathbf{U} = \mathbf{U}^\top M_* \mathbf{U} + \mathbf{U}^\top \Delta \mathbf{U} = \sum_{k=1}^K p^{(k)} \mathbf{U}^\top \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{U} + \mathbf{U}^\top \Delta \mathbf{U}, \\
 \Lambda_* &= \mathbf{U}_*^\top M_* \mathbf{U}_* = \sum_{k=1}^K p^{(k)} \mathbf{U}_*^\top \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{U}_*.
 \end{aligned}$$

Subtracting these two equations gives

$$\sum_{k=1}^K p^{(k)} \left( \mathbf{U}_*^\top \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{U}_* - \mathbf{U}^\top \mathbf{y}^{(k)} \mathbf{y}^{(k)\top} \mathbf{U} \right) = \Lambda_* - \Lambda + \mathbf{U}^\top \Delta \mathbf{U}.$$

Taking the trace of both sides, we get

$$\sum_{k=1}^K p^{(k)} \left( \|\mathbf{U}_*^\top \mathbf{y}^{(k)}\|_2^2 - \|\mathbf{U}^\top \mathbf{y}^{(k)}\|_2^2 \right) = \text{Tr}(\Lambda_* - \Lambda) + \text{Tr}(\mathbf{U}^\top \Delta \mathbf{U}).$$

On the left-hand side,

$$\|\mathbf{U}_*^\top \mathbf{y}^{(k)}\|_2^2 - \|\mathbf{U}^\top \mathbf{y}^{(k)}\|_2^2 = \|\mathbf{y}^{(k)}\|_2^2 - \|\mathbf{U}^\top \mathbf{y}^{(k)}\|_2^2 = \|\mathbf{y}^{(k)} - \mathbf{U} \mathbf{U}^\top \mathbf{y}^{(k)}\|_2^2 \geq 0,$$

while on the right-hand side,

$$\text{Tr}(\Lambda_* - \Lambda) = \sum_{k=1}^K \left( \lambda_k(\Lambda_*) - \lambda_k(\Lambda) \right) \stackrel{(i)}{\leq} K \|\Delta\| \leq K\epsilon, \quad \text{Tr}(\mathbf{U}^\top \Delta \mathbf{U}) \leq \|\Delta\| \cdot \text{Tr}(\mathbf{U}^\top \mathbf{U}) = \|\Delta\| \cdot \text{Tr}(\mathbf{I}_K) \leq K\epsilon,$$

where (i) follows from Weyl's inequality. Putting things together, we have proved (31), which immediately leads to the first upper bound in (32). The second upper bound in (32) follows from a simple application of Davis-Kahan's sin  $\Theta$  theorem (Davis & Kahan, 1970), and the third term is due to Lemma A.11 of (Kong et al., 2020b); we skip the details for brevity.

### C.3. Proof of Theorem B.2

Our proof follows the three steps below:

1. Consider the idealized i.i.d. case, and characterize the expectations and variances of the testing statistics computed by Algorithm 3;
2. Go back to the actual case of mixed LDSs, and analyze one copy of  $\text{stat}_{\Gamma,g}$  or  $\text{stat}_{Y,g}$  defined in (18) for some fixed  $1 \leq g \leq G$ , by decomposing it into an i.i.d. part plus a negligible mixing error term;
3. Analyze  $\text{median}\{\text{stat}_{\Gamma,g}, 1 \leq g \leq G\}$  and  $\text{median}\{\text{stat}_{Y,g}, 1 \leq g \leq G\}$ , and prove the correct testing for each pair of trajectories, which implies that Algorithm 3 achieves exact clustering.

**Step 1: the idealized i.i.d. case.** Recall the definition of  $\text{stat}_Y$  in (19) when we first introduce our method for clustering. For notational convenience, we drop the subscript in  $\text{stat}_Y$ , and replace  $\mathbf{x}_{t+1}, \mathbf{z}_{t+1}$  with  $\mathbf{x}_t', \mathbf{z}_t'$ ; then, with some elementary linear algebra, (19) can be rewritten as

$$\begin{aligned} \text{stat} &= \sum_{i=1}^d \left\langle \mathbf{U}_i^\top \frac{1}{|\Omega_1|} \sum_{t \in \Omega_1} ((\mathbf{x}_t')_i \mathbf{x}_t - (\mathbf{z}_t')_i \mathbf{z}_t), \mathbf{U}_i^\top \frac{1}{|\Omega_2|} \sum_{t \in \Omega_2} ((\mathbf{x}_t')_i \mathbf{x}_t - (\mathbf{z}_t')_i \mathbf{z}_t) \right\rangle \\ &= \left\langle \mathbf{U}^\top \frac{1}{|\Omega_1|} \sum_{t \in \Omega_1} \text{vec}(\mathbf{x}_t (\mathbf{x}_t')^\top - \mathbf{z}_t (\mathbf{z}_t')^\top), \mathbf{U}^\top \frac{1}{|\Omega_2|} \sum_{t \in \Omega_2} \text{vec}(\mathbf{x}_t (\mathbf{x}_t')^\top - \mathbf{z}_t (\mathbf{z}_t')^\top) \right\rangle, \end{aligned}$$

where we define a large orthonormal matrix

$$\mathbf{U} := \begin{bmatrix} \mathbf{U}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{U}_2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0} \\ \mathbf{0} & \dots & \mathbf{0} & \mathbf{U}_d \end{bmatrix} \in \mathbb{R}^{d^2 \times dK}, \quad (41)$$

In this step, we consider the idealized i.i.d. case:

$$\begin{aligned} t \in \Omega_1 \cup \Omega_2, \quad \mathbf{x}_t &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{\Gamma}}^{(k)}), \quad \mathbf{w}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{W}}^{(k)}), \quad \mathbf{x}_t' = \tilde{\mathbf{A}}^{(k)} \mathbf{x}_t + \mathbf{w}_t, \\ \mathbf{z}_t &\stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{\Gamma}}^{(l)}), \quad \mathbf{v}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \tilde{\mathbf{W}}^{(l)}), \quad \mathbf{z}_t' = \tilde{\mathbf{A}}^{(l)} \mathbf{z}_t + \mathbf{v}_t, \end{aligned}$$

where  $\tilde{\mathbf{\Gamma}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(l)}, \tilde{\mathbf{W}}^{(k)}, \tilde{\mathbf{W}}^{(l)}$  are  $d \times d$  covariance matrices, and  $\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{A}}^{(l)}$  are  $d \times d$  state transition matrix. Our goal is to characterize the expectation and variance of  $\text{stat}$  in this i.i.d. case.

Before we present the results, we need some additional notation. First, let  $\{\mathbf{e}_i\}_{1 \leq i \leq d}$  be the canonical basis of  $\mathbb{R}^d$ , and define

$$\begin{aligned} \tilde{\mathbf{Y}}^{(k)} &:= \tilde{\mathbf{A}}^{(k)} \tilde{\mathbf{\Gamma}}^{(k)}, \\ \Sigma^{(k)} &:= \left( \tilde{\mathbf{A}}^{(k)} \otimes \mathbf{I}_d \right) \left( (\tilde{\mathbf{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\mathbf{\Gamma}}^{(k)})^{1/2} \right) \left( \mathbf{I}_{d^2} + \mathbf{P} \right) \left( (\tilde{\mathbf{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\mathbf{\Gamma}}^{(k)})^{1/2} \right) \left( (\tilde{\mathbf{A}}^{(k)})^\top \otimes \mathbf{I}_d \right), \end{aligned}$$

where  $\mathbf{P} \in \mathbb{R}^{d^2 \times d^2}$  is a symmetric permutation matrix, whose  $(i, j)$ -th block is  $\mathbf{e}_j \mathbf{e}_i^\top \in \mathbb{R}^{d \times d}$ ,  $1 \leq i, j \leq d$ . Let  $\tilde{\mathbf{Y}}^{(l)}, \Sigma^{(l)}$  be defined similarly, with  $\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(k)}$  replaced by  $\tilde{\mathbf{A}}^{(l)}, \tilde{\mathbf{\Gamma}}^{(l)}$ . Moreover, define

$$\boldsymbol{\mu}_{k,l} := \mathbf{U}^\top \text{vec} \left( (\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)})^\top \right) \in \mathbb{R}^{dK}, \quad (42a)$$

$$\Sigma_{k,l} := \mathbf{U}^\top \left( \Sigma^{(k)} + \tilde{\mathbf{W}}^{(k)} \otimes \tilde{\mathbf{\Gamma}}^{(k)} + \Sigma^{(l)} + \tilde{\mathbf{W}}^{(l)} \otimes \tilde{\mathbf{\Gamma}}^{(l)} \right) \mathbf{U} \in \mathbb{R}^{dK \times dK}. \quad (42b)$$

Now we are ready to present our results for the i.i.d. case. The first lemma below gives a precise characterization of  $\mathbb{E}[\text{stat}]$  and  $\text{var}(\text{stat})$ , in terms of  $\boldsymbol{\mu}_{k,l}$  and  $\Sigma_{k,l}$ ; the second lemma provides some upper and lower bounds, which will be handy for our later analyses.

**Lemma C.4.** Denote  $N = \min\{|\Omega_1|, |\Omega_2|\}$ . For the i.i.d. case just described, one has

$$\mathbb{E}[\text{stat}] = \|\boldsymbol{\mu}_{k,l}\|_2^2, \quad \text{var}(\text{stat}) \leq \frac{1}{N^2} \text{Tr}(\boldsymbol{\Sigma}_{k,l}^2) + \frac{2}{N} \boldsymbol{\mu}_{k,l}^\top \boldsymbol{\Sigma}_{k,l} \boldsymbol{\mu}_{k,l},$$

and the inequality becomes an equality if  $|\Omega_1| = |\Omega_2| = N$ .

**Lemma C.5.** Consider the same setting of Lemma C.4. Furthermore, suppose that

$$\tilde{\boldsymbol{\Gamma}}^{(k)}, \tilde{\boldsymbol{\Gamma}}^{(l)} \preceq \Gamma_{\max} \mathbf{I}_d, \quad \tilde{\mathbf{W}}^{(k)}, \tilde{\mathbf{W}}^{(l)} \preceq W_{\max} \mathbf{I}_d, \quad \|\tilde{\mathbf{A}}^{(k)}\|, \|\tilde{\mathbf{A}}^{(l)}\| \leq \kappa_A$$

for some  $0 < W_{\max} \leq \Gamma_{\max}$  and  $\kappa_A \geq 1$ . Then the following holds true.

- (Upper bound on expectation) It holds that

$$\mathbb{E}[\text{stat}] = \|\boldsymbol{\mu}_{k,l}\|_2^2 \leq \|\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)}\|_{\mathbb{F}}^2. \quad (43)$$

- (Lower bound on expectation) If  $\tilde{\mathbf{Y}}^{(k)} \neq \tilde{\mathbf{Y}}^{(l)}$  and subspaces  $\{\mathbf{U}_i\}_{1 \leq i \leq d}$  satisfy

$$1 \leq i \leq d, \quad \max \left\{ \|\tilde{\mathbf{Y}}^{(k)}_i - \mathbf{U}_i \mathbf{U}_i^\top \tilde{\mathbf{Y}}^{(k)}_i\|_2, \|\tilde{\mathbf{Y}}^{(l)}_i - \mathbf{U}_i \mathbf{U}_i^\top \tilde{\mathbf{Y}}^{(l)}_i\|_2 \right\} \leq \epsilon, \quad (44)$$

for some  $\epsilon \geq 0$ , then we have

$$\mathbb{E}[\text{stat}] = \|\boldsymbol{\mu}_{k,l}\|_2^2 \geq \|\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)}\|_{\mathbb{F}}^2 - 4\epsilon \sum_{i=1}^d \|\tilde{\mathbf{Y}}^{(k)}_i + \tilde{\mathbf{Y}}^{(l)}_i\|_2.$$

- (Upper bound on variance) The matrix  $\boldsymbol{\Sigma}_{k,l}$  is symmetric and satisfies

$$\mathbf{0} \preceq \boldsymbol{\Sigma}_{k,l} \preceq 6\Gamma_{\max}^2 \kappa_A^2 \mathbf{I}_{dK};$$

this, together with the earlier upper bound (43) on  $\|\boldsymbol{\mu}_{k,l}\|_2^2$ , implies that

$$\text{var}(\text{stat}) \leq \frac{1}{N^2} \text{Tr}(\boldsymbol{\Sigma}_{k,l}^2) + \frac{2}{N} \boldsymbol{\mu}_{k,l}^\top \boldsymbol{\Sigma}_{k,l} \boldsymbol{\mu}_{k,l} \lesssim \left( \frac{\Gamma_{\max}^2 \kappa_A^2}{N} \right)^2 dK + \frac{\Gamma_{\max}^2 \kappa_A^2}{N} \|\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)}\|_{\mathbb{F}}^2.$$

**Step 2: one copy of  $\text{stat}_{Y,g}, \text{stat}_{\Gamma,g}$  for a fixed  $g$ .** Now we turn back to the mixed LDSs setting and prove Theorem B.2. Let us focus on the testing of one pair of short trajectories  $\{\mathbf{x}_{m_1,t}\}, \{\mathbf{x}_{m_2,t}\}$  for some  $m_1, m_2 \in \mathcal{M}_{\text{clustering}}, m_1 \neq m_2$ . For notational consistency, in this proof we rewrite these two trajectories as  $\{\mathbf{x}_t\}$  and  $\{\mathbf{z}_t\}$ , their labels  $k_{m_1}, k_{m_2}$  as  $k, \ell$ , and the trajectory length  $T_{\text{clustering}}$  as  $T$ , respectively. Also denote  $\mathbf{x}'_t := \mathbf{x}_{t+1}$  and  $\mathbf{z}'_t := \mathbf{z}_{t+1}$ . Recall the definition of  $\{\text{stat}_{Y,g}\}_{1 \leq g \leq G}$  in (18b); for now, we consider one specific element and ignore the subscript  $g$ . Recalling the definition of  $\mathbf{U} \in \mathbb{R}^{d^2 \times dK}$  in (41), we have

$$\begin{aligned} \text{stat}_Y &= \sum_{i=1}^d \left\langle \frac{1}{N} \sum_{t \in \Omega_1} \mathbf{U}_i^\top ((\mathbf{x}'_t)_i \mathbf{x}_t - (\mathbf{z}'_t)_i \mathbf{z}_t), \frac{1}{N} \sum_{t \in \Omega_2} \mathbf{U}_i^\top ((\mathbf{x}'_t)_i \mathbf{x}_t - (\mathbf{z}'_t)_i \mathbf{z}_t) \right\rangle \\ &= \left\langle \mathbf{U}^\top \frac{1}{N} \sum_{t \in \Omega_1} \text{vec}(\mathbf{x}_t (\mathbf{x}'_t)^\top - \mathbf{z}_t (\mathbf{z}'_t)^\top), \mathbf{U}^\top \frac{1}{N} \sum_{t \in \Omega_2} \text{vec}(\mathbf{x}_t (\mathbf{x}'_t)^\top - \mathbf{z}_t (\mathbf{z}'_t)^\top) \right\rangle, \end{aligned}$$

where  $N = \lfloor T/4G \rfloor = |\Omega_1| = |\Omega_2|$ .

In the following, we show how to decompose  $\text{stat}_Y$  into an i.i.d. term plus a negligible mixing error term, and then analyze each component of this decomposition; finally, we put pieces together to give a characterization of  $\text{stat}_Y$ , or  $\{\text{stat}_{Y,g}\}_{1 \leq g \leq G}$  when we put the subscript  $g$  back in at the end of this step.

**Step 2.1: decomposition of  $\text{stat}_Y$ .** Define  $\mathcal{S}^{1,\tau_1} := \{t_1 + \tau_1 + f \cdot t_{\text{mix}} : f \geq 0, \tau_1 + f \cdot t_{\text{mix}} \leq |\Omega_1|\}$ , where  $t_1 + 1$  is the first index of  $\Omega_1$ , and the mixing time  $t_{\text{mix}}$  will be specified later; define  $\mathcal{S}^{2,\tau_2}$  similarly. Note that for each  $\tau_1$ , the elements of  $\mathcal{S}^{1,\tau_1}$  are at least  $t_{\text{mix}}$  far apart; moreover, we have  $\Omega_1 = \cup_{\tau_1=1}^{t_{\text{mix}}} \mathcal{S}^{1,\tau_1}$ ,  $\Omega_2 = \cup_{\tau_2=1}^{t_{\text{mix}}} \mathcal{S}^{2,\tau_2}$ , and thus

$$\begin{aligned} \frac{1}{N} \sum_{t \in \Omega_1} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top) &= \sum_{\tau_1=1}^{t_{\text{mix}}} \frac{|\mathcal{S}^{1,\tau_1}|}{N} \cdot \frac{1}{|\mathcal{S}^{1,\tau_1}|} \sum_{t \in \mathcal{S}^{1,\tau_1}} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top), \\ \frac{1}{N} \sum_{t \in \Omega_2} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top) &= \sum_{\tau_2=1}^{t_{\text{mix}}} \frac{|\mathcal{S}^{2,\tau_2}|}{N} \cdot \frac{1}{|\mathcal{S}^{2,\tau_2}|} \sum_{t \in \mathcal{S}^{2,\tau_2}} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top). \end{aligned}$$

Therefore, we can rewrite  $\text{stat}_Y$  as a weighted average

$$\begin{aligned} \text{stat}_Y &= \sum_{\tau_1=1}^{t_{\text{mix}}} \sum_{\tau_2=1}^{t_{\text{mix}}} w^{\tau_1,\tau_2} \cdot \text{stat}_Y^{\tau_1,\tau_2}, \quad \text{where} \\ w^{\tau_1,\tau_2} &:= \frac{|\mathcal{S}^{1,\tau_1}| |\mathcal{S}^{2,\tau_2}|}{N^2}, \quad \sum_{\tau_1=1}^{t_{\text{mix}}} \sum_{\tau_2=1}^{t_{\text{mix}}} w^{\tau_1,\tau_2} = 1, \quad \text{and} \\ \text{stat}_Y^{\tau_1,\tau_2} &:= \left\langle \mathbf{U}^\top \frac{1}{|\mathcal{S}^{1,\tau_1}|} \sum_{t \in \mathcal{S}^{1,\tau_1}} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top), \mathbf{U}^\top \frac{1}{|\mathcal{S}^{2,\tau_2}|} \sum_{t \in \mathcal{S}^{2,\tau_2}} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top) \right\rangle. \quad (45) \end{aligned}$$

We can further decompose  $\text{stat}_Y^{\tau_1,\tau_2}$  into an i.i.d. term plus a small error term. To do this, recalling the definition of  $\tilde{\mathbf{x}}_{m,t}$  in (30) and dropping the subscript  $m$ , we have

$$\begin{aligned} \mathbf{x}_t &= \underbrace{\mathbf{A}^{(k)t_{\text{mix}}-1} \mathbf{x}_{t-t_{\text{mix}}+1}}_{=: \boldsymbol{\delta}_{x,t}} + \tilde{\mathbf{x}}_t = \boldsymbol{\delta}_{x,t} + \tilde{\mathbf{x}}_t, \\ \mathbf{x}_t' &= \mathbf{A}^{(k)} \mathbf{x}_t + \mathbf{w}_t = \mathbf{A}^{(k)} \boldsymbol{\delta}_{x,t} + \underbrace{(\mathbf{A}^{(k)} \tilde{\mathbf{x}}_t + \mathbf{w}_t)}_{=: \tilde{\mathbf{x}}_t'}, \end{aligned}$$

where  $\tilde{\mathbf{x}}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)})$ . Similarly, we rewrite  $\mathbf{z}_t = \boldsymbol{\delta}_{z,t} + \tilde{\mathbf{z}}_t$ ,  $\mathbf{z}_t' = \mathbf{A}^{(\ell)} \boldsymbol{\delta}_{z,t} + \tilde{\mathbf{z}}_t'$ . Plugging these into the right-hand side of (45) and expanding it, one has

$$\begin{aligned} \text{stat}_Y^{\tau_1,\tau_2} &= \widetilde{\text{stat}}_Y^{\tau_1,\tau_2} + \Delta_Y^{\tau_1,\tau_2}, \quad \text{where} \\ \widetilde{\text{stat}}_Y^{\tau_1,\tau_2} &:= \left\langle \mathbf{U}^\top \frac{1}{|\mathcal{S}^{1,\tau_1}|} \sum_{t \in \mathcal{S}^{1,\tau_1}} \text{vec}(\tilde{\mathbf{x}}_t(\tilde{\mathbf{x}}_t')^\top - \tilde{\mathbf{z}}_t(\tilde{\mathbf{z}}_t')^\top), \mathbf{U}^\top \frac{1}{|\mathcal{S}^{2,\tau_2}|} \sum_{t \in \mathcal{S}^{2,\tau_2}} \text{vec}(\tilde{\mathbf{x}}_t(\tilde{\mathbf{x}}_t')^\top - \tilde{\mathbf{z}}_t(\tilde{\mathbf{z}}_t')^\top) \right\rangle, \end{aligned}$$

and  $\Delta_Y^{\tau_1,\tau_2}$  involves  $\{\boldsymbol{\delta}_{x,t}, \boldsymbol{\delta}_{z,t}\}$  terms.

**Step 2.2: analysis of each component.** First, notice that the  $\{\tilde{\mathbf{x}}_t, \tilde{\mathbf{z}}_t\}$  terms in the definition of  $\widetilde{\text{stat}}_Y^{\tau_1,\tau_2}$  are independent; this suggests that we can characterize  $\widetilde{\text{stat}}_Y^{\tau_1,\tau_2}$  by applying our earlier analysis for the i.i.d. case in Step 1. Second,  $\Delta_Y^{\tau_1,\tau_2}$  involves  $\{\boldsymbol{\delta}_{x,t}, \boldsymbol{\delta}_{z,t}\}$  terms, which in turn involve  $\mathbf{A}^{(k)t_{\text{mix}}-1}$ ,  $\mathbf{A}^{(\ell)t_{\text{mix}}-1}$  and thus will be exponentially small as  $t_{\text{mix}}$  increases, thanks to Assumption 3.1. More formally, we have the following:

- Applying Lemmas C.4 and C.5 with  $(\tilde{\mathbf{\Gamma}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(\ell)}) = (\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)})$ ,  $(\widetilde{\mathbf{W}}^{(k)}, \widetilde{\mathbf{W}}^{(\ell)}) = (\mathbf{W}^{(k)}, \mathbf{W}^{(\ell)})$  and  $(\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{A}}^{(\ell)}) = (\mathbf{A}^{(k)}, \mathbf{A}^{(\ell)})$ , we have

$$\begin{aligned} \mathbb{E}[\widetilde{\text{stat}}_Y^{\tau_1,\tau_2}] &= \|\mathbf{U}^\top \text{vec}((\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)})^\top)\|_2^2, \\ \text{var}(\widetilde{\text{stat}}_Y^{\tau_1,\tau_2}) &\lesssim \left( \frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}} \right)^2 dK + \frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}} \|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)}\|_{\mathbf{F}}^2, \end{aligned}$$

where  $\tilde{N} := \min\{|\mathcal{S}^{1,\tau_1}|, |\mathcal{S}^{2,\tau_2}|\} \asymp N/t_{\text{mix}}$ .

- By Lemma C.1, with probability at least  $1 - \delta$ , all  $\{\mathbf{x}_t, \tilde{\mathbf{x}}_t, \mathbf{z}_t, \tilde{\mathbf{z}}_t\}$  terms involved in the definition of  $\Delta_Y^{\tau_1, \tau_2}$  has  $\ell_2$  norm bounded by  $\sqrt{\Gamma_{\max}} \text{poly}(d, \kappa_A, \log(T_{\text{total}}/\delta))$ . This implies that

$$|\Delta_Y^{\tau_1, \tau_2}| \leq \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}.$$

**Step 2.3: putting pieces together.** Putting the subscript  $g$  back in, we have already shown that

$$\begin{aligned} \text{stat}_{Y,g} &= \sum_{\tau_1=1}^{t_{\text{mix}}} \sum_{\tau_2=1}^{t_{\text{mix}}} w^{\tau_1, \tau_2} \cdot \text{stat}_{Y,g}^{\tau_1, \tau_2} \\ &= \underbrace{\sum_{\tau_1=1}^{t_{\text{mix}}} \sum_{\tau_2=1}^{t_{\text{mix}}} w^{\tau_1, \tau_2} \cdot \widetilde{\text{stat}}_{Y,g}^{\tau_1, \tau_2}}_{=:\widetilde{\text{stat}}_{Y,g}} + \underbrace{\sum_{\tau_1=1}^{t_{\text{mix}}} \sum_{\tau_2=1}^{t_{\text{mix}}} w^{\tau_1, \tau_2} \cdot \Delta_{Y,g}^{\tau_1, \tau_2}}_{=:\Delta_{Y,g}} = \widetilde{\text{stat}}_{Y,g} + \Delta_{Y,g}, \end{aligned}$$

where

$$\begin{aligned} \mathbb{E}[\widetilde{\text{stat}}_{Y,g}] &= \|\mathbf{U}^\top \text{vec}((\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)})^\top)\|_2^2, \\ \text{var}(\widetilde{\text{stat}}_{Y,g}) &\leq \max_{\tau_1, \tau_2} \text{var}(\widetilde{\text{stat}}_{Y,g}^{\tau_1, \tau_2}) \lesssim \left(\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}\right)^2 dK + \frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}} \|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)}\|_{\mathbb{F}}^2, \\ |\Delta_{Y,g}| &\leq \max_{\tau_1, \tau_2} |\Delta_{Y,g}^{\tau_1, \tau_2}| \leq \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}. \end{aligned}$$

So far in Step 2, we have focused on the analysis of  $\text{stat}_{Y,g}$ . We can easily adapt the argument to study  $\text{stat}_{\Gamma,g}$  as well; the major difference is that, in Step 2.2, we should apply Lemmas C.4 and C.5 with  $(\tilde{\mathbf{\Gamma}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(\ell)}) = (\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)}, \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)})$ ,  $(\tilde{\mathbf{W}}^{(k)}, \tilde{\mathbf{W}}^{(\ell)}) = (\mathbf{0}, \mathbf{0})$ ,  $(\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{A}}^{(\ell)}) = (\mathbf{I}_d, \mathbf{I}_d)$ , and subspaces  $\{\mathbf{U}_i\}$  replaced by  $\{\mathbf{V}_i\}$  instead. The final result is that, for all  $1 \leq g \leq G$ ,

$$\text{stat}_{\Gamma,g} = \widetilde{\text{stat}}_{\Gamma,g} + \Delta_{\Gamma,g},$$

where

$$\begin{aligned} \mathbb{E}[\widetilde{\text{stat}}_{\Gamma,g}] &= \|\mathbf{V}^\top \text{vec}((\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)})^\top)\|_2^2, \\ \text{var}(\widetilde{\text{stat}}_{\Gamma,g}) &\lesssim \left(\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}\right)^2 dK + \frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}} \|\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)}\|_{\mathbb{F}}^2, \\ |\Delta_{\Gamma,g}| &\leq \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}. \end{aligned}$$

**Step 3: analysis of the medians, and final results.** Recall the following standard result on the concentration of medians (or median-of-means in general; see Theorem 2 of (Lugosi & Mendelson, 2019)).

**Proposition C.6** (Concentration of medians). *Let  $X_1, \dots, X_G$  be i.i.d. random variables with mean  $\mu$  and variance  $\sigma^2$ . Then we have  $|\text{median}\{X_g, 1 \leq g \leq G\} - \mu| \leq 2\sigma$  with probability at least  $1 - e^{-c_0 G}$  for some constant  $c_0$ .*

Notice that by construction,  $\{\widetilde{\text{stat}}_{Y,g}\}_{1 \leq g \leq G}$  are i.i.d. (and so are  $\{\widetilde{\text{stat}}_{\Gamma,g}\}_{1 \leq g \leq G}$ ). Applying Proposition C.6 to our case, we know that if  $G \gtrsim \log(1/\delta)$ , then with probability at least  $1 - \delta$ , the following holds:

- If  $k = \ell$ , i.e. the two trajectories are generated by the same LDS model, then

$$\begin{aligned} &\text{median}\{\text{stat}_{\Gamma,g}, 1 \leq g \leq G\} + \text{median}\{\text{stat}_{Y,g}, 1 \leq g \leq G\} \\ &\leq \text{median}\{\widetilde{\text{stat}}_{\Gamma,g}\} + \text{median}\{\widetilde{\text{stat}}_{Y,g}\} + \max_g |\Delta_{\Gamma,g}| + \max_g |\Delta_{Y,g}| \\ &\leq 2\sqrt{\text{var}(\widetilde{\text{stat}}_{\Gamma,g})} + 2\sqrt{\text{var}(\widetilde{\text{stat}}_{Y,g})} + \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1} \\ &\lesssim \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\tilde{N}} + \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}, \end{aligned} \tag{46}$$

- On the other hand, if  $k \neq \ell$ , then

$$\begin{aligned}
 & \text{median}\{\text{stat}_{\Gamma,g}, 1 \leq g \leq G\} + \text{median}\{\text{stat}_{Y,g}, 1 \leq g \leq G\} \\
 & \geq \text{median}\{\widetilde{\text{stat}}_{\Gamma,g}\} + \text{median}\{\widetilde{\text{stat}}_{Y,g}\} - \left( \max_g |\Delta_{\Gamma,g}| + \max_g |\Delta_{Y,g}| \right) \\
 & \geq \mathbb{E}[\widetilde{\text{stat}}_{\Gamma,g}] + \mathbb{E}[\widetilde{\text{stat}}_{Y,g}] - 2 \left( \sqrt{\text{var}(\widetilde{\text{stat}}_{\Gamma,g})} + \sqrt{\text{var}(\widetilde{\text{stat}}_{Y,g})} \right) - \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1} \\
 & \geq \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)})^\top) \right\|_2^2 + \left\| \mathbf{U}^\top \text{vec}((\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)})^\top) \right\|_2^2 \\
 & \quad - C_0 \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\tilde{N}} + \sqrt{\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}} \left( \|\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}} + \|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}} \right) \right) \\
 & \quad - \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}. \tag{47}
 \end{aligned}$$

We need to further simplify the result (47) for the  $k \neq \ell$  case. According to (28) and (29), we have

$$\begin{aligned}
 \|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}} & \leq \Gamma_{\max} \sqrt{d} \kappa_A^3 \rho^{2(t_{\text{mix}}-1)} =: \epsilon_{\text{mix}}, \\
 \|\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}} & \leq \Gamma_{\max} \sqrt{d} \kappa_A^2 \rho^{2(t_{\text{mix}}-1)} \leq \epsilon_{\text{mix}},
 \end{aligned}$$

which implies that

$$\begin{aligned}
 \|\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}} & \leq \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}} + 2\epsilon_{\text{mix}}, \\
 \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}_{t_{\text{mix}}-1}^{(k)} - \mathbf{\Gamma}_{t_{\text{mix}}-1}^{(\ell)})^\top) \right\|_2^2 & \geq \max \left\{ \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2 - 2\epsilon_{\text{mix}}, 0 \right\}^2 \\
 & \geq \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2^2 - 4\epsilon_{\text{mix}} \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2 \\
 & \geq \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2^2 - 4\epsilon_{\text{mix}} \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}}.
 \end{aligned}$$

We can do a similar analysis for  $\|\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)}\|_{\text{F}}$  and  $\|\mathbf{U}^\top \text{vec}((\mathbf{Y}_{t_{\text{mix}}-1}^{(k)} - \mathbf{Y}_{t_{\text{mix}}-1}^{(\ell)})^\top)\|_2^2$ . Moreover, we claim (and prove later) that if the subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  satisfy the condition (25) in the theorem, then

$$\left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2^2 + \left\| \mathbf{U}^\top \text{vec}((\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^\top) \right\|_2^2 \geq \frac{1}{2} \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}}^2 \right). \tag{48}$$

Putting these back to (47), we have for the  $k \neq \ell$  case,

$$\begin{aligned}
 & \text{median}\{\text{stat}_{\Gamma,g}, 1 \leq g \leq G\} + \text{median}\{\text{stat}_{Y,g}, 1 \leq g \leq G\} \\
 & \geq \left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2^2 + \left\| \mathbf{U}^\top \text{vec}((\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^\top) \right\|_2^2 \\
 & \quad - 4\epsilon_{\text{mix}} \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}} + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} \right) \\
 & \quad - C_0 \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\tilde{N}} + \sqrt{\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}} \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}} + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} + 4\epsilon_{\text{mix}} \right) \right) \\
 & \quad - \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1} \\
 & \stackrel{(i)}{\geq} \frac{1}{2} \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}}^2 \right) - 0.01 \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}}^2 \right) \\
 & \quad - C_1 \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\tilde{N}} + \sqrt{\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}} \sqrt{\|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\text{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}}^2} \right) \\
 & \quad - \Gamma_{\max}^2 \cdot \text{poly}\left(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}\right) \cdot \rho^{t_{\text{mix}}-1}
 \end{aligned}$$



$$\begin{aligned}
 &\stackrel{(ii)}{\geq} 0.48 \left( \|\Gamma^{(k)} - \Gamma^{(\ell)}\|_F^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_F^2 \right) \\
 &\quad - C_1 \left( \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\tilde{N}} + \sqrt{\frac{\Gamma_{\max}^2 \kappa_A^2}{\tilde{N}}} \sqrt{\|\Gamma^{(k)} - \Gamma^{(\ell)}\|_F^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_F^2} \right), \tag{49}
 \end{aligned}$$

where (i) holds if  $t_{\text{mix}} \gtrsim \frac{1}{1-\rho} \log\left(\left(\frac{\Gamma_{\max}}{\Delta_{\Gamma,Y}} + 2\right)d\kappa_A\right)$  so that  $\epsilon_{\text{mix}} \leq 10^{-3}\Delta_{\Gamma,Y}$ , and (ii) holds if  $t_{\text{mix}} \gtrsim \frac{1}{1-\rho} \log\left(\left(\frac{\Gamma_{\max}}{\Delta_{\Gamma,Y}} + 2\right)\frac{d\kappa_A T_{\text{total}}}{\delta}\right)$  so that  $\Gamma_{\max}^2 \cdot \text{poly}(d, \kappa_A, \log \frac{T_{\text{total}}}{\delta}) \cdot \rho^{t_{\text{mix}}-1} \leq 10^{-3}\Delta_{\Gamma,Y}^2$ .

Putting (46) and (49) together, we can finally check that, if it further holds that  $\tilde{N} \asymp N/t_{\text{mix}} \gtrsim \frac{\Gamma_{\max}^2 \kappa_A^2 \sqrt{dK}}{\Delta_{\Gamma,Y}^2} + 1$ , then we have with probability at least  $1 - \delta$ ,

$$\text{median}\{\text{stat}_{\Gamma,g}\} + \text{median}\{\text{stat}_{\mathbf{Y},g}\} \begin{cases} \leq \frac{1}{8}\Delta_{\Gamma,Y}^2 & \text{if } k = \ell, \\ \geq \frac{3}{8}(\|\Gamma^{(k)} - \Gamma^{(\ell)}\|_F^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_F^2) \geq \frac{3}{8}\Delta_{\Gamma,Y}^2 & \text{if } k \neq \ell. \end{cases}$$

This together with our choice of testing threshold  $\tau \in [\Delta_{\Gamma,Y}^2/8, 3\Delta_{\Gamma,Y}^2/8]$  in Algorithm 3 implies correct testing of the two trajectories  $\{\mathbf{x}_t\}, \{\mathbf{z}_t\}$ . Finally, taking the union bound over all pairs of trajectories in  $\mathcal{M}_{\text{clustering}}$  leads to correct pairwise testing, which in turn implies exact clustering of  $\mathcal{M}_{\text{clustering}}$ ; this completes our proof of Theorem (B.2).

### C.3.1. PROOF OF LEMMA C.4

We first assume  $|\Omega_1| = |\Omega_2| = N$  for simplicity. Recall that

$$\text{stat} = \left\langle \underbrace{\mathbf{U}^\top \frac{1}{|\Omega_1|} \sum_{t \in \Omega_1} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top)}_{=: \mathbf{a}}, \underbrace{\mathbf{U}^\top \frac{1}{|\Omega_2|} \sum_{t \in \Omega_2} \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top)}_{=: \mathbf{b}} \right\rangle = \langle \mathbf{a}, \mathbf{b} \rangle,$$

where  $\mathbf{a}, \mathbf{b} \in \mathbb{R}^{dK}$  are i.i.d., and

$$\mathbb{E}[\mathbf{a}] = \mathbb{E}\left[\mathbf{U}^\top \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top)\right] = \mathbf{U}^\top \text{vec}\left((\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)})^\top\right) = \boldsymbol{\mu}_{k,l}.$$

Therefore, we have the expectation

$$\mathbb{E}[\text{stat}] = \langle \mathbb{E}[\mathbf{a}], \mathbb{E}[\mathbf{b}] \rangle = \|\boldsymbol{\mu}_{k,l}\|_2^2.$$

It remains to compute the variance  $\text{var}(\text{stat}) = \mathbb{E}[\text{stat}^2] - \mathbb{E}[\text{stat}]^2$ , where

$$\mathbb{E}[\text{stat}^2] = \mathbb{E}[(\mathbf{a}^\top \mathbf{b})^2] = \text{Tr}(\mathbb{E}[\mathbf{b}\mathbf{b}^\top]\mathbb{E}[\mathbf{a}\mathbf{a}^\top]) = \text{Tr}(\mathbb{E}[\mathbf{a}\mathbf{a}^\top]^2). \tag{50}$$

Here  $\mathbb{E}[\mathbf{a}\mathbf{a}^\top] = \mathbb{E}[\mathbf{a}]\mathbb{E}[\mathbf{a}]^\top + \text{cov}(\mathbf{a})$ , and since  $\mathbf{a}$  is an empirical average of  $N$  i.i.d. random vectors, we have

$$\text{cov}(\mathbf{a}) = \frac{1}{N} \text{cov}(\mathbf{f}), \quad \text{where } \mathbf{f} := \mathbf{U}^\top \text{vec}(\mathbf{x}_t(\mathbf{x}_t')^\top - \mathbf{z}_t(\mathbf{z}_t')^\top) \in \mathbb{R}^{dK}.$$

For now, we claim that

$$\text{cov}(\mathbf{f}) = \mathbf{U}^\top (\boldsymbol{\Sigma}^{(k)} + \widetilde{\mathbf{W}}^{(k)} \otimes \tilde{\Gamma}^{(k)} + \boldsymbol{\Sigma}^{(\ell)} + \widetilde{\mathbf{W}}^{(l)} \otimes \tilde{\Gamma}^{(l)}) \mathbf{U} = \boldsymbol{\Sigma}_{k,l}, \tag{51}$$

which will be proved soon later. Putting these back to (50), one has

$$\begin{aligned}
 \mathbb{E}[\mathbf{a}\mathbf{a}^\top] &= \text{cov}(\mathbf{a}) + \mathbb{E}[\mathbf{a}]\mathbb{E}[\mathbf{a}]^\top = \frac{1}{N} \boldsymbol{\Sigma}_{k,l} + \boldsymbol{\mu}_{k,l} \boldsymbol{\mu}_{k,l}^\top, \\
 \mathbb{E}[\mathbf{a}\mathbf{a}^\top]^2 &= \frac{1}{N^2} \boldsymbol{\Sigma}_{k,l}^2 + \frac{1}{N} (\boldsymbol{\Sigma}_{k,l} \boldsymbol{\mu}_{k,l} \boldsymbol{\mu}_{k,l}^\top + \boldsymbol{\mu}_{k,l} \boldsymbol{\mu}_{k,l}^\top \boldsymbol{\Sigma}_{k,l}) + \|\boldsymbol{\mu}_{k,l}\|_2^2 \boldsymbol{\mu}_{k,l} \boldsymbol{\mu}_{k,l}^\top,
 \end{aligned}$$

and finally

$$\text{var}(\text{stat}) = \mathbb{E}[\text{stat}^2] - \mathbb{E}[\text{stat}]^2 = \text{Tr}(\mathbb{E}[\mathbf{a}\mathbf{a}^\top]^2) - \|\boldsymbol{\mu}_{k,l}\|_2^4 = \frac{1}{N^2} \text{Tr}(\boldsymbol{\Sigma}_{k,l}^2) + \frac{2}{N} \boldsymbol{\mu}_{k,l}^\top \boldsymbol{\Sigma}_{k,l} \boldsymbol{\mu}_{k,l},$$

which completes our calculation of the variance for the case of  $|\Omega_1| = |\Omega_2| = N$ . For the more general case where (without loss of generality)  $|\Omega_1| = N \leq |\Omega_2|$ , we simply need to modify the equation (50) to an inequality  $\mathbb{E}[\text{stat}^2] = \text{Tr}(\mathbb{E}[\mathbf{b}\mathbf{b}^\top]\mathbb{E}[\mathbf{a}\mathbf{a}^\top]) \leq \text{Tr}(\mathbb{E}[\mathbf{a}\mathbf{a}^\top]^2)$ , and the remaining analysis is the same. This finishes the proof of Lemma C.4.

*Proof of (51).* For notational simplicity, we drop the subscripts  $t$  in the definition of  $\mathbf{f}$ . Then we have  $\mathbf{f} = \mathbf{U}^\top \text{vec}(\mathbf{x}(\mathbf{x}')^\top - \mathbf{z}(\mathbf{z}')^\top)$ , and hence

$$\text{cov}(\mathbf{f}) = \mathbf{U}^\top \text{cov}\left(\text{vec}(\mathbf{x}(\mathbf{x}')^\top - \mathbf{z}(\mathbf{z}')^\top)\right) \mathbf{U} = \mathbf{U}^\top \left( \text{cov}(\text{vec}(\mathbf{x}(\mathbf{x}')^\top)) + \text{cov}(\text{vec}(\mathbf{z}(\mathbf{z}')^\top)) \right) \mathbf{U}, \quad (52)$$

where the second equality uses the independence between  $(\mathbf{x}, \mathbf{x}')$  and  $(\mathbf{z}, \mathbf{z}')$ .

Let us focus on  $\text{cov}(\text{vec}(\mathbf{x}(\mathbf{x}')^\top))$ . For notational simplicity, for now we rewrite  $\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{W}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(k)}$  as  $\mathbf{A}, \mathbf{W}, \mathbf{\Gamma}$ . Define  $\mathbf{y} := \mathbf{\Gamma}^{-1/2} \mathbf{x} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , and recall that  $\mathbf{x}' = \mathbf{A}\mathbf{x} + \mathbf{w}$  where  $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{W})$ . Using the fact that  $\text{vec}(\mathbf{ABC}) = (\mathbf{C}^\top \otimes \mathbf{A})\text{vec}(\mathbf{B})$  for any matrices of compatible shape, we have

$$\begin{aligned} \mathbf{g} &:= \text{vec}\left(\mathbf{x}(\mathbf{x}')^\top\right) = \text{vec}(\mathbf{x}\mathbf{x}^\top \mathbf{A}^\top) + \text{vec}(\mathbf{x}\mathbf{w}^\top) \\ &= (\mathbf{A} \otimes \mathbf{I}_d) \text{vec}(\mathbf{x}\mathbf{x}^\top) + \text{vec}(\mathbf{x}\mathbf{w}^\top) = (\mathbf{A} \otimes \mathbf{I}_d) \text{vec}(\mathbf{\Gamma}^{1/2} \mathbf{y}\mathbf{y}^\top \mathbf{\Gamma}^{1/2}) + \text{vec}(\mathbf{x}\mathbf{w}^\top) \\ &= \underbrace{(\mathbf{A} \otimes \mathbf{I}_d)(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2}) \text{vec}(\mathbf{y}\mathbf{y}^\top)}_{=: \mathbf{g}_1} + \underbrace{\text{vec}(\mathbf{x}\mathbf{w}^\top)}_{=: \mathbf{g}_2}. \end{aligned}$$

Note that  $\mathbb{E}[\mathbf{g}_2] = \mathbf{0}$ ,  $\mathbb{E}[\mathbf{g}] = \mathbb{E}[\mathbf{g}_1]$ , and  $\mathbb{E}[\mathbf{g}_1 \mathbf{g}_2^\top] = \mathbf{0}$ . Hence

$$\text{cov}(\mathbf{g}) = \mathbb{E}[\mathbf{g}\mathbf{g}^\top] - \mathbb{E}[\mathbf{g}]\mathbb{E}[\mathbf{g}]^\top = \text{cov}(\mathbf{g}_1) + \mathbb{E}[\mathbf{g}_2 \mathbf{g}_2^\top]. \quad (53)$$

For the second term on the right-hand side, we have

$$\mathbf{g}_2 = \text{vec}(\mathbf{x}\mathbf{w}^\top) = \begin{bmatrix} w_1 \mathbf{x} \\ \vdots \\ w_d \mathbf{x} \end{bmatrix}, \quad \mathbb{E}[\mathbf{g}_2 \mathbf{g}_2^\top] = \mathbb{E}\left[[w_i w_j \mathbf{x}\mathbf{x}^\top]_{1 \leq i, j \leq d}\right] = \mathbf{W} \otimes \mathbf{\Gamma};$$

as for the first term, we claim (and prove soon later) that

$$\text{cov}\left(\text{vec}(\mathbf{y}\mathbf{y}^\top)\right) = \mathbf{I}_{d^2} + \mathbf{P}, \quad (54)$$

which implies that

$$\begin{aligned} \text{cov}(\mathbf{g}_1) &= \text{cov}\left((\mathbf{A} \otimes \mathbf{I}_d)(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2}) \text{vec}(\mathbf{y}\mathbf{y}^\top)\right) \\ &= (\mathbf{A} \otimes \mathbf{I}_d)(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2}) \text{cov}\left(\text{vec}(\mathbf{y}\mathbf{y}^\top)\right) (\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2})(\mathbf{A}^\top \otimes \mathbf{I}_d) \\ &= (\mathbf{A} \otimes \mathbf{I}_d)(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2})(\mathbf{I}_{d^2} + \mathbf{P})(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2})(\mathbf{A}^\top \otimes \mathbf{I}_d). \end{aligned}$$

Putting these back to (53), one has

$$\text{cov}(\text{vec}(\mathbf{x}(\mathbf{x}')^\top)) = \text{cov}(\mathbf{g}) = (\mathbf{A} \otimes \mathbf{I}_d)(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2})(\mathbf{I}_{d^2} + \mathbf{P})(\mathbf{\Gamma}^{1/2} \otimes \mathbf{\Gamma}^{1/2})(\mathbf{A}^\top \otimes \mathbf{I}_d) + \mathbf{W} \otimes \mathbf{\Gamma},$$

which is equal to  $\mathbf{\Sigma}^{(k)} + \tilde{\mathbf{W}}^{(k)} \otimes \tilde{\mathbf{\Gamma}}^{(k)}$  if we return to the original notation of  $\tilde{\mathbf{A}}^{(k)}, \tilde{\mathbf{W}}^{(k)}, \tilde{\mathbf{\Gamma}}^{(k)}$ . By a similar analysis, we can show that  $\text{cov}(\text{vec}(\mathbf{z}(\mathbf{z}')^\top)) = \mathbf{\Sigma}^{(l)} + \tilde{\mathbf{W}}^{(l)} \otimes \tilde{\mathbf{\Gamma}}^{(l)}$ . Putting these back to (52) finishes our calculation of  $\text{cov}(\mathbf{f})$ .

Finally, it remains to prove (54). Denote

$$\mathbf{u} := \text{vec}(\mathbf{y}\mathbf{y}^\top) = \begin{bmatrix} y_1 \mathbf{y} \\ \vdots \\ y_d \mathbf{y} \end{bmatrix}.$$

Then  $\mathbb{E}[\mathbf{u}] = \text{vec}(\mathbf{I}_d)$ , and thus  $\mathbb{E}[\mathbf{u}]\mathbb{E}[\mathbf{u}]^\top = [\mathbf{e}_i \mathbf{e}_j^\top]_{1 \leq i, j \leq d}$ . Next, consider

$$\mathbb{E}[\mathbf{u}\mathbf{u}^\top] = \mathbb{E}\left[[y_i y_j \mathbf{y}\mathbf{y}^\top]_{1 \leq i, j \leq d}\right].$$

- For  $i = j$ ,

$$\mathbb{E}[y_i^2 y_k y_\ell] = \begin{cases} 3 & \text{if } k = \ell = i, \\ 1 & \text{if } k = \ell \neq i, \\ 0 & \text{if } k \neq \ell, \end{cases}$$

and hence  $\mathbb{E}[y_i^2 \mathbf{y} \mathbf{y}^\top] = \mathbf{I}_d + 2\mathbf{e}_i \mathbf{e}_i^\top$ .

- For  $i \neq j$ ,

$$\mathbb{E}[y_i y_j y_k y_\ell] = \begin{cases} 1 & \text{if } k = i, \ell = j \text{ or } k = j, \ell = i, \\ 0 & \text{otherwise,} \end{cases}$$

and hence  $\mathbb{E}[y_i y_j \mathbf{y} \mathbf{y}^\top] = \mathbf{e}_i \mathbf{e}_j^\top + \mathbf{e}_j \mathbf{e}_i^\top$ .

Putting together, the  $(i, j)$ -th  $d \times d$  block of  $\text{cov}(\mathbf{u}) = \mathbb{E}[\mathbf{u} \mathbf{u}^\top] - \mathbb{E}[\mathbf{u}] \mathbb{E}[\mathbf{u}]^\top$  is equal to  $\mathbf{I}_d + \mathbf{e}_i \mathbf{e}_i^\top$  if  $i = j$ , and  $\mathbf{e}_j \mathbf{e}_i^\top$  if  $i \neq j$ . In other words,  $\text{cov}(\text{vec}(\mathbf{y} \mathbf{y}^\top)) = \text{cov}(\mathbf{u}) = \mathbf{I}_{d^2} + \mathbf{P}$ , where  $\mathbf{P} = [\mathbf{e}_j \mathbf{e}_i^\top]_{1 \leq i, j \leq d}$  is a symmetric permutation matrix; this completes our proof of (54).  $\square$

### C.3.2. PROOF OF LEMMA C.5

First, it holds that

$$\|\boldsymbol{\mu}_{k,l}\|_2^2 = \sum_{i=1}^d \left\| \mathbf{U}_i^\top ((\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i) \right\|_2^2 \leq \sum_{i=1}^d \left\| (\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i \right\|_2^2 = \|\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)}\|_{\mathbb{F}}^2,$$

which gives the upper bound on  $\mathbb{E}[\text{stat}] = \|\boldsymbol{\mu}_{k,l}\|_2^2$ . For the lower bound, the triangle inequality tells us that

$$\left\| \mathbf{U}_i^\top ((\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i) \right\|_2 = \left\| \mathbf{U}_i \mathbf{U}_i^\top (\tilde{\mathbf{Y}}^{(k)})_i - \mathbf{U}_i \mathbf{U}_i^\top (\tilde{\mathbf{Y}}^{(l)})_i \right\|_2 \geq \max \left\{ \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2 - 2\epsilon, 0 \right\},$$

which implies that

$$\left\| \mathbf{U}_i^\top ((\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i) \right\|_2^2 \geq \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2^2 - 4\epsilon \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2,$$

and hence

$$\begin{aligned} \|\boldsymbol{\mu}_{k,l}\|_2^2 &= \sum_{i=1}^d \left\| \mathbf{U}_i^\top ((\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i) \right\|_2^2 \geq \sum_{i=1}^d \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2^2 - 4\epsilon \sum_{i=1}^d \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2 \\ &= \|\tilde{\mathbf{Y}}^{(k)} - \tilde{\mathbf{Y}}^{(l)}\|_{\mathbb{F}}^2 - 4\epsilon \sum_{i=1}^d \|(\tilde{\mathbf{Y}}^{(k)})_i - (\tilde{\mathbf{Y}}^{(l)})_i\|_2. \end{aligned}$$

It remains to upper bound  $\boldsymbol{\Sigma}_{k,l}$ . Recall the definition

$$\boldsymbol{\Sigma}_{k,l} = \mathbf{U}^\top (\boldsymbol{\Sigma}^{(k)} + \tilde{\mathbf{W}}^{(k)} \otimes \tilde{\boldsymbol{\Gamma}}^{(k)} + \boldsymbol{\Sigma}^{(l)} + \tilde{\mathbf{W}}^{(l)} \otimes \tilde{\boldsymbol{\Gamma}}^{(l)}) \mathbf{U}.$$

We will utilize the following basic facts: (1) for square matrices  $\mathbf{A}$  and  $\mathbf{B}$  with eigenvalues  $\{\lambda_i\}$  and  $\{\mu_j\}$  respectively, their Kronecker product  $\mathbf{A} \otimes \mathbf{B}$  has eigenvalues  $\{\lambda_i \mu_j\}$ ; (2) For matrices  $\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}$  of compatible shapes, it holds that  $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{A}\mathbf{C}) \otimes (\mathbf{B}\mathbf{D})$ . These imply that

$$\mathbf{0} \preceq \tilde{\mathbf{W}}^{(k)} \otimes \tilde{\boldsymbol{\Gamma}}^{(k)} \preceq \|\tilde{\mathbf{W}}^{(k)}\| \|\tilde{\boldsymbol{\Gamma}}^{(k)}\| \mathbf{I}_{d^2} \preceq W_{\max} \Gamma_{\max} \mathbf{I}_{d^2},$$

and

$$\begin{aligned} \mathbf{0} \preceq \boldsymbol{\Sigma}^{(k)} &= \left( \tilde{\mathbf{A}}^{(k)} \otimes \mathbf{I}_d \right) \left( (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \right) \left( \mathbf{I}_{d^2} + \mathbf{P} \right) \left( (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \right) \left( (\tilde{\mathbf{A}}^{(k)})^\top \otimes \mathbf{I}_d \right) \\ &\preceq 2 \left( \tilde{\mathbf{A}}^{(k)} \otimes \mathbf{I}_d \right) \left( (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \right) \left( (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \otimes (\tilde{\boldsymbol{\Gamma}}^{(k)})^{1/2} \right) \left( (\tilde{\mathbf{A}}^{(k)})^\top \otimes \mathbf{I}_d \right) \end{aligned}$$

$$\begin{aligned}
 &= 2 \left( \tilde{\mathbf{A}}^{(k)} \otimes \mathbf{I}_d \right) \left( \tilde{\mathbf{\Gamma}}^{(k)} \otimes \tilde{\mathbf{\Gamma}}^{(k)} \right) \left( (\tilde{\mathbf{A}}^{(k)})^\top \otimes \mathbf{I}_d \right) \\
 &\preceq 2\Gamma_{\max}^2 \left( \tilde{\mathbf{A}}^{(k)} \otimes \mathbf{I}_d \right) \left( (\tilde{\mathbf{A}}^{(k)})^\top \otimes \mathbf{I}_d \right) \\
 &= 2\Gamma_{\max}^2 \left( \tilde{\mathbf{A}}^{(k)} (\tilde{\mathbf{A}}^{(k)})^\top \right) \otimes \mathbf{I}_d \preceq 2\Gamma_{\max}^2 \kappa_A^2 \mathbf{I}_{d^2}.
 \end{aligned}$$

Using the conditions  $W_{\max} \leq \Gamma_{\max}$  and  $\kappa_A \geq 1$ , we have

$$\mathbf{\Sigma}^{(k)} + \tilde{\mathbf{W}}^{(k)} \otimes \tilde{\mathbf{\Gamma}}^{(k)} \preceq (W_{\max} \Gamma_{\max} + 2\Gamma_{\max}^2 \kappa_A^2) \mathbf{I}_{d^2} \preceq 3\Gamma_{\max}^2 \kappa_A^2 \mathbf{I}_{d^2}.$$

We can upper bound  $\mathbf{\Sigma}^{(\ell)} + \tilde{\mathbf{W}}^{(\ell)} \otimes \tilde{\mathbf{\Gamma}}^{(\ell)}$  by the same analysis. As a result,

$$\mathbf{\Sigma}_{k,\ell} \preceq \mathbf{U}^\top (6\Gamma_{\max}^2 \kappa_A^2 \mathbf{I}_{d^2}) \mathbf{U} = 6\Gamma_{\max}^2 \kappa_A^2 \mathbf{U}^\top \mathbf{U} = 6\Gamma_{\max}^2 \kappa_A^2 \mathbf{I}_{dK},$$

which finishes the proof of Lemma C.5.

### C.3.3. PROOF OF (48)

Let  $\epsilon$  be the right-hand side of the condition (25) on the subspaces. Then, applying the second point of Lemma C.5 (with  $\tilde{\mathbf{Y}}^{(k)} = \mathbf{Y}^{(k)}$ ,  $\tilde{\mathbf{Y}}^{(\ell)} = \mathbf{Y}^{(\ell)}$ ) tells us that

$$\begin{aligned}
 \left\| \mathbf{U}^\top \text{vec}((\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^\top) \right\|_2^2 &\geq \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2 - 4\epsilon \sum_{i=1}^d \|(\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i\|_2 \\
 &\geq \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2 - 4\epsilon\sqrt{d} \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}},
 \end{aligned}$$

where the last line follows from the Cauchy-Schwarz inequality:

$$\sum_{i=1}^d \|(\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i\|_2 \leq \sqrt{d \sum_{i=1}^d \|(\mathbf{Y}^{(k)})_i - (\mathbf{Y}^{(\ell)})_i\|_2^2} = \sqrt{d} \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}.$$

We can lower bound  $\|\mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top)\|_2^2$  similarly. Putting pieces together, we have

$$\begin{aligned}
 &\left\| \mathbf{V}^\top \text{vec}((\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)})^\top) \right\|_2^2 + \left\| \mathbf{U}^\top \text{vec}((\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^\top) \right\|_2^2 \\
 &\geq \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2 - 4\epsilon\sqrt{d} (\|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}} + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}) \\
 &\geq \frac{1}{2} \left( \|\mathbf{\Gamma}^{(k)} - \mathbf{\Gamma}^{(\ell)}\|_{\mathbb{F}}^2 + \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\mathbb{F}}^2 \right),
 \end{aligned}$$

where the last inequality is due to the assumption  $\epsilon \lesssim \Delta_{\Gamma, Y} / \sqrt{d}$ . This completes our proof of (48).

### C.4. Proof of Theorem B.3

It suffices to prove the error bounds for one specific value of  $k$ , and then take the union bound over  $1 \leq k \leq K$ . For notational convenience, in this proof we rewrite  $T_{\text{total}}^{(k)}$ ,  $\mathbf{A}^{(k)}$ ,  $\hat{\mathbf{A}}^{(k)}$ ,  $\mathbf{W}^{(k)}$ ,  $\hat{\mathbf{W}}^{(k)}$  as  $T$ ,  $\mathbf{A}$ ,  $\hat{\mathbf{A}}$ ,  $\mathbf{W}$ ,  $\hat{\mathbf{W}}$ , respectively. We will investigate the close-form solution  $\hat{\mathbf{A}}$ , and prepare ourselves with a self-normalized concentration bound; this will be helpful in finally proving the error bounds for  $\|\hat{\mathbf{A}} - \mathbf{A}\|$  and  $\|\hat{\mathbf{W}} - \mathbf{W}\|$ .

**Step 1: preparation.** Recall the least-squares solution

$$\hat{\mathbf{A}} = \left( \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \mathbf{x}_{m,t+1} \mathbf{x}_{m,t}^\top \right) \left( \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \mathbf{x}_{m,t} \mathbf{x}_{m,t}^\top \right)^{-1}.$$

Using the notation

$$\mathbf{X} := \begin{bmatrix} \vdots \\ \mathbf{x}_{m,t}^\top \\ \vdots \end{bmatrix}_{0 \leq t \leq T_m - 1, m \in \mathcal{C}_k} \in \mathbb{R}^{T \times d}, \quad \mathbf{X}_+ := \begin{bmatrix} \vdots \\ \mathbf{x}_{m,t+1}^\top \\ \vdots \end{bmatrix}, \quad \mathbf{N} := \begin{bmatrix} \vdots \\ \mathbf{w}_{m,t}^\top \\ \vdots \end{bmatrix},$$

we have

$$\mathbf{X}_+^\top = \mathbf{A}\mathbf{X}^\top + \mathbf{N}^\top, \quad \widehat{\mathbf{A}} = \mathbf{X}_+^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} = \mathbf{A} + \mathbf{N}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1},$$

namely

$$\Delta_A := \widehat{\mathbf{A}} - \mathbf{A} = \mathbf{N}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1}. \quad (55)$$

We will utilize the following matrix form of self-normalized concentration; see Lemma C.4 of (Kakade et al., 2020).

**Lemma C.7** (Self-normalized concentration, matrix form). *Consider filtrations  $\{\mathcal{F}_t\}$ , and random vectors  $\{\mathbf{x}_t, \mathbf{z}_t\}$  satisfying  $\mathbf{x}_t \in \mathcal{F}_{t-1}$  and  $\mathbf{z}_t | \mathcal{F}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Let  $\mathbf{V} \in \mathbb{R}^{d \times d}$  be a fixed, symmetric, positive definite matrix, and denote  $\overline{\mathbf{V}}_T := \sum_{t=1}^T \mathbf{x}_t \mathbf{x}_t^\top + \mathbf{V}$ . Then with probability at least  $1 - \delta$ ,*

$$\left\| (\overline{\mathbf{V}}_T)^{-1/2} \left( \sum_{t=1}^T \mathbf{x}_t \mathbf{z}_t^\top \right) \right\| \lesssim \sqrt{d + \log \frac{1}{\delta} + \log \frac{\det(\overline{\mathbf{V}}_T)}{\det(\mathbf{V})}}.$$

**Step 2: estimation error of  $\widehat{\mathbf{A}}$ .** Let us rewrite (55) as

$$\Delta_A^\top = (\mathbf{X}^\top \mathbf{X})^{-1/2} \cdot (\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{X}^\top \mathbf{N}. \quad (56)$$

It is obvious that  $\mathbf{X}^\top \mathbf{X}$  plays a crucial role. Recall from Lemma C.1 that with probability at least  $1 - \delta$ ,  $\|\mathbf{x}_{m,t}\|_2 \leq D_{\text{vec}}$  for some  $D_{\text{vec}} \lesssim \sqrt{\Gamma_{\max}} \cdot \text{poly}(d, \kappa_A, \log(T_{\text{total}}/\delta))$ . Then trivially we have the upper bound

$$\mathbf{X}^\top \mathbf{X} \preceq D_{\text{vec}}^2 T \cdot \mathbf{I}_d =: \mathbf{V}_{\text{up}}.$$

For a lower bound, we claim (and prove later) that with probability at least  $1 - \delta$ ,

$$\mathbf{X}^\top \mathbf{X} \succeq \frac{1}{5} T \cdot \mathbf{W} =: \mathbf{V}_{\text{lb}}, \quad \text{provided that } T \gtrsim \kappa_w^2 d \cdot \log \left( \frac{\Gamma_{\max} \kappa_A d T_{\text{total}}}{W_{\min} \delta} \right). \quad (57)$$

Now we are ready to control  $\|\widehat{\mathbf{A}} - \mathbf{A}\| = \|\Delta_A\|$ . From (56), we have

$$\|\Delta_A\| \leq \|(\mathbf{X}^\top \mathbf{X})^{-1/2}\| \cdot \|(\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{X}^\top \mathbf{N}\|.$$

First, the lower bound (57) on  $\mathbf{X}^\top \mathbf{X}$  tells us that  $\|(\mathbf{X}^\top \mathbf{X})^{-1/2}\| \lesssim 1/\sqrt{T \cdot \lambda_{\min}(\overline{\mathbf{W}})}$ . Moreover, applying Lemma C.7 with  $\mathbf{V} = \mathbf{V}_{\text{lb}}$ , one has with probability at least  $1 - \delta$ ,

$$\begin{aligned} \|(\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{X}^\top \mathbf{N}\| &\lesssim \|(\mathbf{X}^\top \mathbf{X} + \mathbf{V})^{-1/2} \mathbf{X}^\top \mathbf{N}\| \lesssim \|\mathbf{W}^{1/2}\| \sqrt{d + \log \frac{1}{\delta} + \log \frac{\det(\mathbf{V}_{\text{up}} + \mathbf{V}_{\text{lb}})}{\det(\mathbf{V}_{\text{lb}})}} \\ &\lesssim \sqrt{\|\overline{\mathbf{W}}\|} \sqrt{d \cdot \log \left( \frac{\Gamma_{\max} d \kappa_A T_{\text{total}}}{W_{\min} \delta} \right)}. \end{aligned} \quad (58)$$

Putting these together, we have

$$\|\Delta_A\| \lesssim \frac{1}{\sqrt{T \cdot \lambda_{\min}(\overline{\mathbf{W}})}} \cdot \sqrt{\|\overline{\mathbf{W}}\|} \sqrt{d \cdot \log \left( \frac{\Gamma_{\max} d \kappa_A T_{\text{total}}}{W_{\min} \delta} \right)} \lesssim \sqrt{\frac{d \cdot \kappa_w}{T} \log \left( \frac{\Gamma_{\max} d \kappa_A T_{\text{total}}}{W_{\min} \delta} \right)},$$

which proves our upper bound for  $\|\widehat{\mathbf{A}} - \mathbf{A}\|$  in the theorem.

**Step 3: estimation error of  $\widehat{\mathbf{W}}$ .** By the definition of  $\widehat{\mathbf{w}}_{m,t} = \mathbf{x}_{m,t+1} - \widehat{\mathbf{A}} \mathbf{x}_{m,t} = \mathbf{w}_{m,t} - \Delta_A \mathbf{x}_{m,t}$ , we have

$$\begin{aligned} \widehat{\mathbf{N}} &:= \begin{bmatrix} \vdots \\ \widehat{\mathbf{w}}_{m,t}^\top \\ \vdots \end{bmatrix} \in \mathbb{R}^{T \times d}, \\ \widehat{\mathbf{N}}^\top &= \mathbf{N}^\top - \Delta_A \mathbf{X}^\top = \mathbf{N}^\top - \mathbf{N}^\top \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top = \mathbf{N}^\top (\mathbf{I}_T - \mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top). \end{aligned}$$

Notice that  $\mathbf{I}_T - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$  is a symmetric projection matrix. Therefore,

$$\widehat{\mathbf{W}} = \frac{1}{T} \widehat{\mathbf{N}}^\top \widehat{\mathbf{N}} = \frac{1}{T} \mathbf{N}^\top (\mathbf{I}_T - \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top) \mathbf{N},$$

and thus

$$\widehat{\mathbf{W}} - \mathbf{W} = \left( \frac{1}{T} \mathbf{N}^\top \mathbf{N} - \mathbf{W} \right) - \frac{1}{T} \mathbf{N}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{N}. \quad (59)$$

The first term on the right-hand side can be controlled by a standard result for covariance estimation (see Proposition D.1): with probability at least  $1 - \delta$ ,

$$\left\| \frac{1}{T} \mathbf{N}^\top \mathbf{N} - \mathbf{W} \right\| \lesssim \|\mathbf{W}\| \sqrt{\frac{d + \log \frac{1}{\delta}}{T}}.$$

As for the second term, we rewrite it as

$$\frac{1}{T} \mathbf{N}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{N} = \frac{1}{T} \left( (\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{X}^\top \mathbf{N} \right)^\top \left( (\mathbf{X}^\top \mathbf{X})^{-1/2} \mathbf{X}^\top \mathbf{N} \right).$$

Applying our earlier self-normalized concentration bound (58), we have

$$\left\| \frac{1}{T} \mathbf{N}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{N} \right\| \lesssim \frac{\|\mathbf{W}\| d \cdot \log\left(\frac{\Gamma_{\max}}{W_{\min}} \frac{d\kappa_A T_{\text{total}}}{\delta}\right)}{T}.$$

Putting these back to (59), we have

$$\begin{aligned} \|\widehat{\mathbf{W}} - \mathbf{W}\| &\leq \left\| \frac{1}{T} \mathbf{N}^\top \mathbf{N} - \mathbf{W} \right\| + \left\| \frac{1}{T} \mathbf{N}^\top \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{N} \right\| \\ &\lesssim \|\mathbf{W}\| \sqrt{\frac{d + \log \frac{1}{\delta}}{T}} + \|\mathbf{W}\| \frac{d \cdot \log\left(\frac{\Gamma_{\max}}{W_{\min}} \frac{d\kappa_A T_{\text{total}}}{\delta}\right)}{T} \lesssim \|\mathbf{W}\| \sqrt{\frac{d \cdot \log\left(\frac{\Gamma_{\max}}{W_{\min}} \frac{d\kappa_A T_{\text{total}}}{\delta}\right)}{T}}, \end{aligned}$$

where the last inequality uses  $T \gtrsim d \cdot \log\left(\frac{\Gamma_{\max}}{W_{\min}} \frac{d\kappa_A T_{\text{total}}}{\delta}\right)$ . This finishes our proof of Theorem B.3.

#### C.4.1. PROOF OF (57)

We start with the following decomposition:

$$\begin{aligned} \mathbf{X}^\top \mathbf{X} &= \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 1} \mathbf{x}_{m,t} \mathbf{x}_{m,t}^\top \gtrsim \sum_{m \in \mathcal{C}_k} \sum_{1 \leq t \leq T_m - 1} \mathbf{x}_{m,t} \mathbf{x}_{m,t}^\top \\ &= \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \mathbf{x}_{m,t+1} \mathbf{x}_{m,t+1}^\top = \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{A} \mathbf{x}_{m,t} + \mathbf{w}_{m,t}) (\mathbf{A} \mathbf{x}_{m,t} + \mathbf{w}_{m,t})^\top \\ &= \underbrace{\sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \mathbf{w}_{m,t} \mathbf{w}_{m,t}^\top}_{:=P} \\ &\quad + \underbrace{\sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \left( \mathbf{A} \mathbf{x}_{m,t} \mathbf{x}_{m,t}^\top \mathbf{A}^\top + \mathbf{A} \mathbf{x}_{m,t} \mathbf{w}_{m,t}^\top + \mathbf{w}_{m,t} \mathbf{x}_{m,t}^\top \mathbf{A}^\top \right)}_{:=Q} \\ &= P + Q. \end{aligned} \quad (60)$$

**Lower bound for  $P$ .** By Proposition D.1, we have with probability at least  $1 - \delta$ ,

$$\left\| \frac{1}{T - |\mathcal{C}_k|} P - \mathbf{W} \right\| \lesssim \|\mathbf{W}\| \sqrt{\frac{d + \log \frac{1}{\delta}}{T}} \lesssim \lambda_{\min}(\mathbf{W}) \cdot \mathbf{I}_d, \quad \text{provided that } T \gtrsim \kappa_w^2 (d + \log \frac{1}{\delta}).$$

As a result, we have  $\frac{1}{T - |\mathcal{C}_k|} P \gtrsim \frac{1}{2} \mathbf{W}$ , which implies  $P \gtrsim \frac{1}{4} T \cdot \mathbf{W}$ .

**Lower bound for  $\mathbf{Q}$ .** Let  $\mathcal{N}_\epsilon$  be an  $\epsilon$ -net of the unit sphere  $\mathcal{S}^{d-1}$  (where the value of  $0 < \epsilon < 1$  will be specified later), and let  $\pi$  be the projection onto  $\mathcal{N}_\epsilon$ . Recall the standard result that  $|\mathcal{N}_\epsilon| \leq (9/\epsilon)^d$ . Moreover, for any  $\mathbf{v} \in \mathcal{S}^{d-1}$ , denote  $\Delta_{\mathbf{v}} := \mathbf{v} - \pi(\mathbf{v})$ , which satisfies  $\|\Delta_{\mathbf{v}}\|_2 \leq \epsilon$ . Then we have

$$\begin{aligned} \lambda_{\min}(\mathbf{Q}) &= \inf_{\mathbf{v} \in \mathcal{S}^{d-1}} \mathbf{v}^\top \mathbf{Q} \mathbf{v} = \inf_{\mathbf{v} \in \mathcal{S}^{d-1}} (\pi(\mathbf{v}) + \Delta_{\mathbf{v}})^\top \mathbf{Q} (\pi(\mathbf{v}) + \Delta_{\mathbf{v}}) \\ &\geq \inf_{\mathbf{v} \in \mathcal{S}^{d-1}} \pi(\mathbf{v})^\top \mathbf{Q} \pi(\mathbf{v}) - (2\epsilon + \epsilon^2) \|\mathbf{Q}\| \geq \inf_{\mathbf{v} \in \mathcal{N}_\epsilon} \mathbf{v}^\top \mathbf{Q} \mathbf{v} - 3\epsilon \|\mathbf{Q}\|. \end{aligned} \quad (61)$$

For  $\|\mathbf{Q}\|$ , we simply use a crude upper bound, based on the boundedness of  $\{\|\mathbf{x}_{m,t}\|_2\}$  (Lemma C.1): with probability at least  $1 - \delta$ , one has

$$\|\mathbf{Q}\| \leq \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \left( \|\mathbf{A} \mathbf{x}_{m,t}\|_2^2 + 2 \|\mathbf{A} \mathbf{x}_{m,t}\|_2 \|\mathbf{w}_{m,t}\|_2 \right) \lesssim \Gamma_{\max} \cdot \text{poly}\left(\kappa_A, d, T_{\text{total}}, \frac{1}{\delta}\right). \quad (62)$$

Next, we lower bound  $\inf_{\mathbf{v} \in \mathcal{N}_\epsilon} \mathbf{v}^\top \mathbf{Q} \mathbf{v}$ . First, consider a fixed  $\mathbf{v} \in \mathcal{N}_\epsilon$ ; denoting  $\mathbf{y}_{m,t} := \mathbf{A} \mathbf{x}_{m,t}$  and  $\mathbf{u}_{m,t} := \mathbf{W}^{-1/2} \mathbf{w}_{m,t} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ , we have

$$\mathbf{v}^\top \mathbf{Q} \mathbf{v} = \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 + 2 \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \mathbf{v}^\top \mathbf{y}_{m,t} \cdot \mathbf{u}_{m,t}^\top \mathbf{W}^{1/2} \mathbf{v},$$

where  $\mathbf{u}_{m,t}^\top \mathbf{W}^{1/2} \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{v}^\top \mathbf{W} \mathbf{v})$ . Lemma D.2 (the scalar version of self-normalized concentration) tells us that, for any fixed  $\lambda > 0$ , with probability at least  $1 - \delta$ ,

$$\begin{aligned} \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} \mathbf{v}^\top \mathbf{y}_{m,t} \cdot \mathbf{u}_{m,t}^\top \mathbf{W}^{1/2} \mathbf{v} &\geq -\sqrt{\mathbf{v}^\top \mathbf{W} \mathbf{v}} \left( \frac{\lambda}{2} \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 + \frac{1}{\lambda} \log \frac{1}{\delta} \right) \\ &\geq -\sqrt{\|\mathbf{W}\|} \left( \frac{\lambda}{2} \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 + \frac{1}{\lambda} \log \frac{1}{\delta} \right). \end{aligned}$$

Replacing  $\delta$  with  $\delta/(9/\epsilon)^d$  and taking the union bound, we have with probability at least  $1 - \delta$ , for any  $\mathbf{v} \in \mathcal{N}_\epsilon$ ,

$$\begin{aligned} \mathbf{v}^\top \mathbf{Q} \mathbf{v} &\geq \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 - \sqrt{\|\mathbf{W}\|} \left( \lambda \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 + \frac{2}{\lambda} \left( d \cdot \log \frac{9}{\epsilon} + \log \frac{1}{\delta} \right) \right) \\ &= (1 - \sqrt{\|\mathbf{W}\|} \lambda) \sum_{m \in \mathcal{C}_k} \sum_{0 \leq t \leq T_m - 2} (\mathbf{v}^\top \mathbf{y}_{m,t})^2 - \sqrt{\|\mathbf{W}\|} \frac{2}{\lambda} \left( d \cdot \log \frac{9}{\epsilon} + \log \frac{1}{\delta} \right). \end{aligned}$$

With the choice of  $\lambda = 1/\sqrt{\|\mathbf{W}\|}$ , this implies

$$\mathbf{v}^\top \mathbf{Q} \mathbf{v} \geq -2\|\mathbf{W}\| \left( d \cdot \log \frac{9}{\epsilon} + \log \frac{1}{\delta} \right), \quad \text{for all } \mathbf{v} \in \mathcal{N}_\epsilon. \quad (63)$$

Putting (62) and (63) back to (61), we have with probability at least  $1 - \delta$ ,

$$\lambda_{\min}(\mathbf{Q}) \geq \inf_{\mathbf{v} \in \mathcal{N}_\epsilon} \mathbf{v}^\top \mathbf{Q} \mathbf{v} - 3\epsilon \|\mathbf{Q}\| \geq -C_0 \left( \|\mathbf{W}\| \left( d \cdot \log \frac{9}{\epsilon} + \log \frac{1}{\delta} \right) + \epsilon \Gamma_{\max} \cdot \text{poly}\left(\kappa_A, d, T_{\text{total}}, \frac{1}{\delta}\right) \right)$$

for some universal constant  $C_0 > 0$ .

**Putting things together.** Recall the decomposition  $\mathbf{X}^\top \mathbf{X} \succcurlyeq \mathbf{P} + \mathbf{Q}$  in (60). We have already shown that if  $T \gtrsim \kappa_w^2 (d + \log \frac{1}{\delta})$ , then with probability at least  $1 - \delta$ ,

$$\mathbf{X}^\top \mathbf{X} \succcurlyeq \mathbf{P} + \mathbf{Q} \succcurlyeq \frac{1}{4} T \cdot \mathbf{W} - C_0 \left( \|\mathbf{W}\| \left( d \cdot \log \frac{9}{\epsilon} + \log \frac{1}{\delta} \right) + \epsilon \Gamma_{\max} \cdot \text{poly}\left(\kappa_A, d, T_{\text{total}}, \frac{1}{\delta}\right) \right) \mathbf{I}_d.$$

It is easy to check that, if we further choose

$$\epsilon \asymp \frac{1}{\text{poly}\left(\kappa_A, d, T_{\text{total}}, \frac{1}{\delta}, \frac{\Gamma_{\max}}{W_{\min}}\right)}, \quad T \gtrsim \kappa_w^2 d \cdot \log \left( \frac{\Gamma_{\max} \kappa_A d T_{\text{total}}}{W_{\min} \delta} \right),$$

then we have  $\mathbf{X}^\top \mathbf{X} \succcurlyeq \frac{1}{5} T \cdot \mathbf{W}$ , which finishes the proof of (57).

### C.5. Proof of Theorem B.4

In this proof, we show the correct classification of one short trajectory  $\{\mathbf{x}_{m,t}\}_{0 \leq t \leq T_m}$  (with true label  $k_m$ ) for some  $m \in \mathcal{M}_{\text{classification}}$ ; then it suffices to take the union bound to prove the correct classification of all trajectories in  $\mathcal{M}_{\text{classification}}$ . For notational simplicity, we drop the subscript  $m$  and rewrite  $\{\mathbf{x}_{m,t}\}, T_m, k_m$  as  $\{\mathbf{x}_t\}, T, k$ , respectively. The basic idea of this proof is to show that

$$L(\widehat{\mathbf{A}}^{(\ell)}, \widehat{\mathbf{W}}^{(\ell)}) > L(\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}) \quad (64)$$

for any incorrect label  $\ell \neq k$ , where  $L$  is the loss function defined in (23) and used by Algorithm 5 for classification. Our proof below is simply a sequence of arguments that finally transform (64) into a sufficient condition in terms of the coarse model errors  $\epsilon_A, \epsilon_W$  and short trajectory length  $T$ .

Before we proceed, we record a few basic facts that will be useful later. First, the assumption  $\|\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)}\| \leq \epsilon_W \leq 0.1W_{\min}$  implies that  $\lambda_{\min}(\widehat{\mathbf{W}}^{(k)}) \geq 0.9\lambda_{\min}(\mathbf{W}^{(k)})$ , and  $\widehat{\mathbf{W}}^{(k)}$  is well conditioned with  $\kappa(\widehat{\mathbf{W}}^{(k)}) \lesssim \kappa(\mathbf{W}^{(k)}) \leq \kappa_w$ . Moreover, by Lemma C.1, with probability at least  $1 - \delta$ , we have for all  $0 \leq t \leq T$ ,

- In Case 0,  $\|\mathbf{x}_t\|_2 \leq D_x \lesssim \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}$ , provided that  $T \gtrsim 1$ ;
- In Case 1,  $\|\mathbf{x}_t\|_2 \leq D_x \lesssim \kappa_A \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}$ , provided that  $T \gtrsim \frac{1}{1-\rho} \log(2\kappa_A)$ .

Now we are ready to state our proof. Throughout our analyses, we will make some intermediate claims, whose proofs will be deferred to the end of this subsection.

**Step 1: a sufficient condition for correct classification.** In the following, we prove that for a fixed  $\ell \neq k$ , the condition (64) holds with high probability; at the end of the proof, we simply take the union bound over  $\ell \neq k$ . Using  $\mathbf{x}_{t+1} = \mathbf{A}^{(k)}\mathbf{x}_t + \mathbf{w}_t$  where  $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{W}^{(k)})$ , we can rewrite the loss function  $L$  as

$$\begin{aligned} L(\mathbf{A}, \mathbf{W}) &= T \cdot \log \det(\mathbf{W}) + \sum_{t=0}^{T-1} \mathbf{w}_t^\top \mathbf{W}^{-1} \mathbf{w}_t \\ &\quad + \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \mathbf{A})^\top \mathbf{W}^{-1} (\mathbf{A}^{(k)} - \mathbf{A}) \mathbf{x}_t + 2 \sum_{t=0}^{T-1} \mathbf{w}_t^\top \mathbf{W}^{-1} (\mathbf{A}^{(k)} - \mathbf{A}) \mathbf{x}_t. \end{aligned}$$

After some basic calculation, (64) can be equivalently written as

$$L(\widehat{\mathbf{A}}^{(\ell)}, \widehat{\mathbf{W}}^{(\ell)}) - L(\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}) = (\text{A}) + (\text{B}) - (\text{C}) > 0,$$

$$\text{where } (\text{A}) := T \cdot \left( \log \det(\widehat{\mathbf{W}}^{(\ell)}) - \log \det(\widehat{\mathbf{W}}^{(k)}) \right) + \sum_{t=0}^{T-1} \mathbf{w}_t^\top \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) \mathbf{w}_t$$

$$(\text{B}) := \left( \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t + 2 \sum_{t=0}^{T-1} \mathbf{w}_t^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t \right)$$

$$(\text{C}) := \left( \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)})^\top (\widehat{\mathbf{W}}^{(k)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}) \mathbf{x}_t + 2 \sum_{t=0}^{T-1} \mathbf{w}_t^\top (\widehat{\mathbf{W}}^{(k)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}) \mathbf{x}_t \right)$$

**Step 2: a lower bound for (A) + (B) - (C).** Intuitively, we expect that (A) + (B) should be large because the LDS models  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  and  $(\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  are well separated, while (C) should be small if  $(\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)}) \approx (\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$ . More formally, we claim that the following holds for some universal constants  $C_1, C_2, C_3 > 0$ :

- With probability at least  $1 - \delta$ ,

$$\begin{aligned} (\text{A}) &\geq T \cdot \left[ \log \det(\widehat{\mathbf{W}}^{(\ell)}) - \log \det(\widehat{\mathbf{W}}^{(k)}) + \text{Tr} \left( (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right) \right] \\ &\quad - C_1 \sqrt{T} \left\| \left( (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right) \right\|_{\text{F}} \log \frac{1}{\delta}; \end{aligned} \quad (65)$$



- With probability at least  $1 - \delta$ ,

$$(B) \geq C_2 \left( T \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}} - \kappa_w \kappa_{w,\text{cross}} \log \frac{1}{\delta} \right), \quad (66)$$

provided that  $T \gtrsim \kappa_w^2 \log^2(1/\delta)$ ;

- With probability at least  $1 - \delta$ ,

$$(C) \leq C_3 \left( T \frac{D_x^2 \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}\|^2}{W_{\min}} + \kappa_w \log \frac{1}{\delta} \right), \quad (67)$$

provided that  $T \gtrsim 1$  (under Case 0) or  $T \gtrsim \frac{1}{1-\rho} \log(2\kappa_A)$  (under Case 1).

Putting these together, we have with probability at least  $1 - \delta$ ,

$$\begin{aligned} (A) + (B) - (C) &\geq C_4 \cdot T \cdot \left[ \log \frac{\det(\widehat{\mathbf{W}}^{(\ell)})}{\det(\widehat{\mathbf{W}}^{(k)})} + \text{Tr} \left( \mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) \right) \right. \\ &\quad \left. + \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}} - \frac{D_x^2 \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}\|^2}{W_{\min}} \right] \\ &\quad - C_5 \left[ \sqrt{T} \left\| (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right\|_{\text{F}} + \kappa_w \kappa_{w,\text{cross}} \right] \log \frac{1}{\delta} \end{aligned}$$

for some universal constants  $C_4, C_5 > 0$ , provided that  $T \gtrsim \kappa_w^2 \log^2 \frac{1}{\delta}$  (under Case 0), or  $T \gtrsim \kappa_w^2 \log^2 \frac{1}{\delta} + \frac{1}{1-\rho} \log(2\kappa_A)$  (under Case 1). Now we have a lower bound of (A) + (B) - (C) as an order- $T$  term minus a low-order term. Therefore, to show (A) + (B) - (C)  $> 0$ , it suffices to prove that (a) the leading factor of order- $T$  term is positive and large, and (b) the low-order term is negligible compared with the order- $T$  term. More specifically, under the assumption that the coarse models  $\{\widehat{\mathbf{A}}^{(j)}, \widehat{\mathbf{W}}^{(j)}\}$  satisfy  $\|\widehat{\mathbf{A}}^{(j)} - \mathbf{A}^{(j)}\| \leq \epsilon_A$ ,  $\|\widehat{\mathbf{W}}^{(j)} - \mathbf{W}^{(j)}\| \leq \epsilon_W \leq 0.1W_{\min}$  for all  $1 \leq j \leq K$ , we make the following claims:

- (Order- $T$  term is large.) Define the leading factor  $\widehat{D}_{k,\ell}$  of the order- $T$  term and a related parameter  $D_{k,\ell}$  as follows:

$$\begin{aligned} \widehat{D}_{k,\ell} &:= \log \frac{\det(\widehat{\mathbf{W}}^{(\ell)})}{\det(\widehat{\mathbf{W}}^{(k)})} + \text{Tr} \left( \mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) \right) + \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}} - \frac{D_x^2 \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}\|^2}{W_{\min}}, \\ D_{k,\ell} &:= \frac{1}{\kappa_{w,\text{cross}}} \left( \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}}^2}{W_{\max}^2} + \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}^2 \right) \gtrsim \frac{\Delta_{\mathbf{A},\mathbf{W}}^2}{\kappa_{w,\text{cross}}}, \end{aligned} \quad (68)$$

where the last inequality follows from Assumption 3.1. They are related in the sense that

$$D_{k,\ell} \lesssim \widetilde{D}_{k,\ell} := \log \frac{\det(\mathbf{W}^{(\ell)})}{\det(\mathbf{W}^{(k)})} + \text{Tr} \left( \mathbf{W}^{(k)} \left( (\mathbf{W}^{(\ell)})^{-1} - (\mathbf{W}^{(k)})^{-1} \right) \right) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}}, \quad (69)$$

where  $\widetilde{D}_{k,\ell}$  is defined in the same way as  $\widehat{D}_{k,\ell}$ , except that the coarse models are replaced with the accurate ones; moreover, we have

$$D_{k,\ell} \lesssim \widehat{D}_{k,\ell}, \quad \text{provided that } \epsilon_A \lesssim \sqrt{\frac{W_{\min} D_{k,\ell}}{D_x^2}}, \quad \epsilon_W \lesssim W_{\min} \sqrt{\frac{D_{k,\ell}}{d}}. \quad (70)$$

- (Low-order term is negligible.) With (70) in place, we have

$$\begin{aligned} (A) + (B) - (C) &\geq C_6 T \cdot D_{k,\ell} - C_5 \left[ \sqrt{T} \left\| (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right\|_{\text{F}} + \kappa_w \kappa_{w,\text{cross}} \right] \log \frac{1}{\delta}. \end{aligned} \quad (71)$$

We claim that

$$\text{if } T \gtrsim \left( \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}} + 1 \right) \log^2 \frac{1}{\delta}, \quad \text{then } (A) + (B) - (C) \geq C_7 T \cdot D_{k,\ell} > 0. \quad (72)$$

**Step 3: putting things together.** So far, we have proved that for a short trajectory generated by  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)})$  and for a fixed  $\ell \neq k$ , it holds with probability at least  $1 - \delta$  that  $L(\widehat{\mathbf{A}}^{(\ell)}, \widehat{\mathbf{W}}^{(\ell)}) > L(\widehat{\mathbf{A}}^{(k)}, \widehat{\mathbf{W}}^{(k)})$ , provided that

$$\epsilon_A \lesssim \sqrt{\frac{W_{\min} D_{k,\ell}}{D_x^2}}, \quad \epsilon_W \lesssim W_{\min} \cdot \min \left\{ 1, \sqrt{\frac{D_{k,\ell}}{d}} \right\}, \quad T \gtrsim \begin{cases} \left( \kappa_w^2 + \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}} \right) \log^2 \frac{1}{\delta} & \text{for Case 0,} \\ \left( \kappa_w^2 + \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}} \right) \log^2 \frac{1}{\delta} + \frac{1}{1-\rho} \log(2\kappa_A) & \text{for Case 1.} \end{cases}$$

Plugging in the relation  $D_{k,\ell} \gtrsim \Delta_{A,W}^2 / \kappa_{w,\text{cross}}$  and  $D_x \lesssim \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}$  (for Case 0) or  $D_x \lesssim \kappa_A \sqrt{\Gamma_{\max}(d + \log \frac{T_{\text{total}}}{\delta})}$  (for Case 1), the above conditions become

$$\begin{aligned} \text{For Case 0: } \quad \epsilon_A &\lesssim \sqrt{\frac{W_{\min} \Delta_{A,W}^2}{\Gamma_{\max} \kappa_{w,\text{cross}} (d + \log \frac{T_{\text{total}}}{\delta})}}, \quad \epsilon_W \lesssim W_{\min} \cdot \min \left\{ 1, \frac{\Delta_{A,W}}{\sqrt{\kappa_{w,\text{cross}} d}} \right\}, \\ T &\gtrsim \left( \kappa_w^2 + \frac{\kappa_{w,\text{cross}}^6}{\Delta_{A,W}^2} \right) \log^2 \frac{1}{\delta}; \\ \text{For Case 1: } \quad \epsilon_A &\lesssim \sqrt{\frac{W_{\min} \Delta_{A,W}^2}{\Gamma_{\max} \kappa_{w,\text{cross}} \kappa_A^2 (d + \log \frac{T_{\text{total}}}{\delta})}}, \quad \epsilon_W \lesssim W_{\min} \cdot \min \left\{ 1, \frac{\Delta_{A,W}}{\sqrt{\kappa_{w,\text{cross}} d}} \right\}, \\ T &\gtrsim \left( \kappa_w^2 + \frac{\kappa_{w,\text{cross}}^6}{\Delta_{A,W}^2} \right) \log^2 \frac{1}{\delta} + \frac{1}{1-\rho} \log(2\kappa_A). \end{aligned}$$

Finally, taking the union bound over all  $\ell \neq k$  as well as over all trajectories in  $\mathcal{M}_{\text{classification}}$  finishes the proof of Theorem B.4.

### C.5.1. PROOF OF (65).

Since  $\mathbf{w}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{W}^{(k)})$ , we have

$$\sum_{t=0}^{T-1} \mathbf{w}_t^\top \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) \mathbf{w}_t = \mathbf{z}^\top \mathbf{M} \mathbf{z},$$

where  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_{Td})$ , and  $\mathbf{M} \in \mathbb{R}^{Td \times Td}$  is a block-diagonal matrix with  $\mathbf{Q} := (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \in \mathbb{R}^{d \times d}$  as its diagonal blocks. Therefore, by the Hanson-Wright inequality (Theorem 6.2.1 of (Vershynin, 2018)), we have

$$\mathbb{P} \left( \left| \mathbf{z}^\top \mathbf{M} \mathbf{z} - \mathbb{E}[\mathbf{z}^\top \mathbf{M} \mathbf{z}] \right| \geq u \right) \leq 2 \exp \left( -c \min \left\{ \frac{u^2}{\|\mathbf{M}\|_{\mathbb{F}}^2}, \frac{u}{\|\mathbf{M}\|} \right\} \right),$$

where  $\|\mathbf{M}\|_{\mathbb{F}}^2 = T \|\mathbf{Q}\|_{\mathbb{F}}^2$ ,  $\|\mathbf{M}\| = \|\mathbf{Q}\|$ , and  $\mathbb{E}[\mathbf{z}^\top \mathbf{M} \mathbf{z}] = \text{Tr}(\mathbf{M}) = T \cdot \text{Tr}(\mathbf{Q})$ . Choosing  $u \gtrsim \sqrt{T} \|\mathbf{Q}\|_{\mathbb{F}} \log \frac{1}{\delta}$ , we have with probability at least  $1 - \delta$ ,  $|\mathbf{z}^\top \mathbf{M} \mathbf{z} - T \cdot \text{Tr}(\mathbf{Q})| \leq C_1 \sqrt{T} \|\mathbf{Q}\|_{\mathbb{F}} \log \frac{1}{\delta}$ , which immediately leads to our lower bound (65) for (A).

### C.5.2. PROOF OF (66).

Denote  $\mathbf{u}_t = (\mathbf{W}^{(k)})^{-1/2} \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\mathbf{y}_t = (\mathbf{W}^{(k)})^{1/2} (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t$ . Then we have

$$(B) = \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t + 2 \sum_{t=0}^{T-1} \mathbf{u}_t^\top \mathbf{y}_t.$$

By Lemma D.2, we have with probability at least  $1 - \delta$ ,  $\sum_{t=0}^{T-1} \mathbf{u}_t^\top \mathbf{y}_t \geq -\left(\frac{\lambda}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda} \log \frac{2}{\delta}\right)$  for any fixed  $\lambda > 0$ . This implies that

$$(B) \geq \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t - \lambda \sum_{t=0}^{T-1} \mathbf{y}_t^\top \mathbf{y}_t - \frac{2}{\lambda} \log \frac{2}{\delta}$$

$$= \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - \lambda \cdot (\widehat{\mathbf{W}}^{(\ell)})^{-1} \mathbf{W}^{(k)} (\widehat{\mathbf{W}}^{(\ell)})^{-1} \right) (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t - \frac{2}{\lambda} \log \frac{2}{\delta}.$$

Choosing  $\lambda = 0.05/(\kappa_w \kappa_{w, \text{cross}})$ , we have  $\lambda \cdot (\widehat{\mathbf{W}}^{(\ell)})^{-1} \mathbf{W}^{(k)} (\widehat{\mathbf{W}}^{(\ell)})^{-1} \preceq 0.1 (\widehat{\mathbf{W}}^{(\ell)})^{-1}$ , and thus

$$\begin{aligned} \text{(B)} &\geq 0.9 \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t - 40 \kappa_w \kappa_{w, \text{cross}} \log \frac{2}{\delta} \\ &\geq 0.9 \lambda_{\min}((\widehat{\mathbf{W}}^{(\ell)})^{-1}) \sum_{t=0}^{T-1} \mathbf{x}_t^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)})^\top (\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}) \mathbf{x}_t - 40 \kappa_w \kappa_{w, \text{cross}} \log \frac{2}{\delta}. \end{aligned} \quad (73)$$

Now it remains to lower bound  $\sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta \mathbf{x}_t$ , where  $\Delta := \Delta_A^\top \Delta_A$  and  $\Delta_A := \mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}$ . Since  $\Delta \succcurlyeq \mathbf{0}$ , we have

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta \mathbf{x}_t &\geq \sum_{t=0}^{T-2} \mathbf{x}_{t+1}^\top \Delta \mathbf{x}_{t+1} = \sum_{t=0}^{T-2} (\mathbf{A}^{(k)} \mathbf{x}_t + \mathbf{w}_t)^\top \Delta (\mathbf{A}^{(k)} \mathbf{x}_t + \mathbf{w}_t) \\ &= \underbrace{\sum_{t=0}^{T-2} \mathbf{w}_t^\top \Delta \mathbf{w}_t}_{\text{(i)}} + \underbrace{\sum_{t=0}^{T-2} \mathbf{x}_t^\top \mathbf{A}^{(k)\top} \Delta \mathbf{A}^{(k)} \mathbf{x}_t + 2 \sum_{t=0}^{T-2} \mathbf{w}_t^\top \Delta \mathbf{A}^{(k)} \mathbf{x}_t}_{\text{(ii)}}. \end{aligned} \quad (74)$$

We can lower bound (i) by the Hanson-Wright inequality, similar to our previous proof of (65); the result is that, with probability at least  $1 - \delta$ , one has

$$\text{(i)} \geq T \cdot \text{Tr} \left( (\mathbf{W}^{(k)})^{1/2} \Delta (\mathbf{W}^{(k)})^{1/2} \right) - C_0 \sqrt{T} \| (\mathbf{W}^{(k)})^{1/2} \Delta (\mathbf{W}^{(k)})^{1/2} \|_{\text{F}} \log \frac{1}{\delta}.$$

To lower bound (ii), we apply Lemma D.2, which shows that with probability at least  $1 - \delta$ ,

$$\left| \sum_{t=0}^{T-2} \mathbf{w}_t^\top \Delta \mathbf{A}^{(k)} \mathbf{x}_t \right| \leq \frac{\lambda}{2} \sum_{t=0}^{T-2} \mathbf{x}_t^\top (\mathbf{A}^{(k)})^\top \Delta \mathbf{W}^{(k)} \Delta \mathbf{A}^{(k)} \mathbf{x}_t + \frac{1}{\lambda} \log \frac{2}{\delta},$$

for any fixed  $\lambda > 0$ , hence

$$\text{(ii)} \geq \sum_{t=0}^{T-2} \mathbf{x}_t^\top \mathbf{A}^{(k)\top} (\Delta - \lambda \cdot \Delta \mathbf{W}^{(k)} \Delta) \mathbf{A}^{(k)} \mathbf{x}_t - \frac{2}{\lambda} \log \frac{2}{\delta}.$$

Recall  $\Delta = \Delta_A^\top \Delta_A$ , and thus  $\Delta - \lambda \cdot \Delta \mathbf{W}^{(k)} \Delta = \Delta_A^\top (\mathbf{I}_d - \lambda \cdot \Delta_A \mathbf{W}^{(k)} \Delta_A^\top) \Delta_A \succcurlyeq \mathbf{0}$  if we choose  $\lambda = 1/(\|\mathbf{W}^{(k)}\| \|\Delta_A\|^2) = 1/(\|\mathbf{W}^{(k)}\| \|\Delta\|)$ ; this implies that

$$\text{(ii)} \geq -\frac{2}{\lambda} \log \frac{2}{\delta} = -2 \|\mathbf{W}^{(k)}\| \|\Delta\| \log \frac{2}{\delta}.$$

Putting these back to (74), we have

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta \mathbf{x}_t &\geq \text{(i)} + \text{(ii)} \\ &\geq T \cdot \text{Tr} \left( (\mathbf{W}^{(k)})^{1/2} \Delta (\mathbf{W}^{(k)})^{1/2} \right) - C_0 \sqrt{T} \| (\mathbf{W}^{(k)})^{1/2} \Delta (\mathbf{W}^{(k)})^{1/2} \|_{\text{F}} \log \frac{1}{\delta} - 2 \|\mathbf{W}^{(k)}\| \|\Delta\| \log \frac{2}{\delta} \\ &\geq T \cdot \lambda_{\min}(\mathbf{W}^{(k)}) \text{Tr}(\Delta_A^\top \Delta_A) - C_0 \sqrt{T} \|\mathbf{W}^{(k)}\| \|\Delta_A^\top \Delta_A\|_{\text{F}} \log \frac{1}{\delta} - 2 \|\mathbf{W}^{(k)}\| \|\Delta_A^\top \Delta_A\| \log \frac{2}{\delta} \\ &\geq T \cdot \lambda_{\min}(\mathbf{W}^{(k)}) \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2 \cdot \left( 1 - \frac{C_0}{\sqrt{T}} \kappa(\mathbf{W}^{(k)}) \log \frac{1}{\delta} - \frac{2}{T} \kappa(\mathbf{W}^{(k)}) \log \frac{2}{\delta} \right) \\ &\geq 0.9T \cdot \lambda_{\min}(\mathbf{W}^{(k)}) \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2, \end{aligned}$$

where the last inequality holds if  $T \gtrsim \kappa_w^2 \log^2 \frac{1}{\delta}$ .

Going back to (73), we have

$$\begin{aligned}
 (\text{B}) &\geq 0.9 \lambda_{\min}((\widehat{\mathbf{W}}^{(\ell)})^{-1}) \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta \mathbf{x}_t - 40 \kappa_w \kappa_{w,\text{cross}} \log \frac{2}{\delta} \\
 &\geq 0.81 T \cdot \lambda_{\min}((\widehat{\mathbf{W}}^{(\ell)})^{-1}) \lambda_{\min}(\mathbf{W}^{(k)}) \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\mathbb{F}}^2 - 40 \kappa_w \kappa_{w,\text{cross}} \log \frac{2}{\delta} \\
 &\geq 0.7 T \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} - 40 \kappa_w \kappa_{w,\text{cross}} \log \frac{2}{\delta},
 \end{aligned}$$

which finishes the proof of (66).

### C.5.3. PROOF OF (67).

Denote  $\Delta = \mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}$ ,  $\mathbf{u}_t = (\mathbf{W}^{(k)})^{-1/2} \mathbf{w}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$  and  $\mathbf{y}_t = (\mathbf{W}^{(k)})^{1/2} (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t$ . Then one has

$$(\text{C}) = \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t + 2 \sum_{t=0}^{T-1} \mathbf{w}_t^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t = \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t + 2 \sum_{t=1}^T \mathbf{u}_t^\top \mathbf{y}_t.$$

By Lemma D.2, we have with probability at least  $1 - \delta$ ,  $\sum_{t=0}^{T-1} \mathbf{u}_t^\top \mathbf{y}_t \leq \frac{\lambda}{2} \sum_{t=0}^{T-1} \|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda} \log \frac{2}{\delta}$  for any fixed  $\lambda > 0$ . This implies that

$$\begin{aligned}
 (\text{C}) &\leq \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t + \lambda \sum_{t=0}^{T-1} \mathbf{y}_t^\top \mathbf{y}_t + \frac{2}{\lambda} \log \frac{2}{\delta} \\
 &= \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t + \lambda \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top (\widehat{\mathbf{W}}^{(k)})^{-1} \mathbf{W}^{(k)} (\widehat{\mathbf{W}}^{(k)})^{-1} \Delta \mathbf{x}_t + \frac{2}{\lambda} \log \frac{2}{\delta} \\
 &= \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top \left( (\widehat{\mathbf{W}}^{(k)})^{-1} + \lambda \cdot (\widehat{\mathbf{W}}^{(k)})^{-1} \mathbf{W}^{(k)} (\widehat{\mathbf{W}}^{(k)})^{-1} \right) \Delta \mathbf{x}_t + \frac{2}{\lambda} \log \frac{2}{\delta} \\
 &\leq 1.5 \left( \frac{1}{\lambda_{\min}(\mathbf{W}^{(k)})} + \frac{\lambda \cdot \|\mathbf{W}^{(k)}\|}{\lambda_{\min}(\mathbf{W}^{(k)})^2} \right) \sum_{t=0}^{T-1} \mathbf{x}_t^\top \Delta^\top \Delta \mathbf{x}_t + \frac{2}{\lambda} \log \frac{2}{\delta}.
 \end{aligned}$$

Choosing  $\lambda \asymp 1/\kappa_w$  and recalling  $\|\mathbf{x}_t\|_2 \leq D_x$ , we have (C)  $\lesssim \frac{1}{\widehat{W}_{\min}} T \cdot D_x^2 \|\Delta\|^2 + \kappa_w \log \frac{1}{\delta}$ , which finishes the proof of (67).

### C.5.4. PROOF OF (69).

Denote  $\Delta := \mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}$ . Then the right-hand side of (69) becomes

$$\begin{aligned}
 \widetilde{D}_{k,\ell} &= \log \det \mathbf{W}^{(\ell)} - \log \det \mathbf{W}^{(k)} + \text{Tr} \left( \mathbf{W}^{(k)} (\mathbf{W}^{(\ell)})^{-1} - \mathbf{I}_d \right) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} \\
 &= \log \det \mathbf{W}^{(\ell)} - \log \det (\mathbf{W}^{(\ell)} + \Delta) + \text{Tr} \left( (\mathbf{W}^{(\ell)} + \Delta) (\mathbf{W}^{(\ell)})^{-1} - \mathbf{I}_d \right) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} \\
 &= \log \det \mathbf{W}^{(\ell)} - \log \det \left( (\mathbf{W}^{(\ell)})^{1/2} (\mathbf{I}_d + (\mathbf{W}^{(\ell)})^{-1/2} \Delta (\mathbf{W}^{(\ell)})^{-1/2}) (\mathbf{W}^{(\ell)})^{1/2} \right) \\
 &\quad + \text{Tr} \left( (\mathbf{W}^{(\ell)})^{-1/2} \Delta (\mathbf{W}^{(\ell)})^{-1/2} \right) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} \\
 &= \text{Tr}(\mathbf{X}) - \log \det (\mathbf{I}_d + \mathbf{X}) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}},
 \end{aligned}$$

where we define  $\mathbf{X} := (\mathbf{W}^{(\ell)})^{-1/2} \Delta (\mathbf{W}^{(\ell)})^{-1/2}$ . Notice that  $\mathbf{X}$  is symmetric and satisfies

$$\mathbf{X} + \mathbf{I}_d = (\mathbf{W}^{(\ell)})^{-1/2} \mathbf{W}^{(k)} (\mathbf{W}^{(\ell)})^{-1/2} \succ \mathbf{0},$$

$$\begin{aligned}\|\mathbf{X}\| &\leq \|(\mathbf{W}^{(\ell)})^{-1/2}\|^2 \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\| \leq \frac{2W_{\max}}{W_{\min}} = 2\kappa_{w,\text{cross}}, \\ \|\mathbf{X}\|_{\mathbb{F}}^2 &= \|(\mathbf{W}^{(\ell)})^{-1/2} \Delta (\mathbf{W}^{(\ell)})^{-1/2}\|_{\mathbb{F}}^2 \geq \frac{\|\Delta\|_{\mathbb{F}}^2}{W_{\max}^2}.\end{aligned}$$

Therefore, by Lemma D.3, we have

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) \geq \frac{\|\mathbf{X}\|_{\mathbb{F}}^2}{6\kappa_{w,\text{cross}}} \geq \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\mathbb{F}}^2}{6\kappa_{w,\text{cross}} W_{\max}^2},$$

and thus

$$\tilde{D}_{k,\ell} = \text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} \geq \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\mathbb{F}}^2}{6\kappa_{w,\text{cross}} W_{\max}^2} + \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} \asymp D_{k,\ell},$$

which finishes the proof of (69).

### C.5.5. PROOF OF (70).

Recall the definition

$$\hat{D}_{k,\ell} = \log \frac{\det(\widehat{\mathbf{W}}^{(\ell)})}{\det(\widehat{\mathbf{W}}^{(k)})} + \text{Tr}\left(\mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right)\right) + \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} - \frac{D_x^2 \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}\|_{\mathbb{F}}^2}{W_{\min}}.$$

First, we have

$$\begin{aligned}&\log \frac{\det(\widehat{\mathbf{W}}^{(\ell)})}{\det(\widehat{\mathbf{W}}^{(k)})} + \text{Tr}\left(\mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right)\right) \\ &= \underbrace{\left[ \log \frac{\det(\widehat{\mathbf{W}}^{(\ell)})}{\det(\mathbf{W}^{(k)})} + \text{Tr}\left(\mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\mathbf{W}^{(k)})^{-1} \right)\right) \right]}_{\text{(i)}} \\ &\quad - \underbrace{\left[ \log \frac{\det(\widehat{\mathbf{W}}^{(k)})}{\det(\mathbf{W}^{(k)})} + \text{Tr}\left(\mathbf{W}^{(k)} \left( (\widehat{\mathbf{W}}^{(k)})^{-1} - (\mathbf{W}^{(k)})^{-1} \right)\right) \right]}_{\text{(ii)}}.\end{aligned}$$

We can lower bound (i) by the same idea of our earlier proof for (69), except that we replace  $\mathbf{W}^{(\ell)}$  in that proof with  $\widehat{\mathbf{W}}^{(\ell)}$ ; this gives us

$$\text{(i)} \gtrsim \frac{\|\mathbf{W}^{(k)} - \widehat{\mathbf{W}}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}} W_{\max}^2}.$$

As for (ii), applying Lemma D.3 with  $\mathbf{X} = (\mathbf{W}^{(k)})^{-1/2} (\widehat{\mathbf{W}}^{(k)} - \mathbf{W}^{(k)}) (\mathbf{W}^{(k)})^{-1/2}$ , one has

$$\text{(ii)} = \text{Tr}(\mathbf{X}) - \log \det(\mathbf{X} + \mathbf{I}_d) \leq \|\mathbf{X}\|_{\mathbb{F}}^2 \leq \frac{\epsilon_W^2 d}{W_{\min}^2}.$$

Putting things together, we have

$$\begin{aligned}\hat{D}_{k,\ell} &= \text{(i)} - \text{(ii)} + \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\mathbb{F}}^2}{\kappa_{w,\text{cross}}} - \frac{D_x^2 \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(k)}\|_{\mathbb{F}}^2}{W_{\min}} \\ &\geq \underbrace{C_1 \frac{1}{\kappa_{w,\text{cross}}} \left( \frac{\|\mathbf{W}^{(k)} - \widehat{\mathbf{W}}^{(\ell)}\|_{\mathbb{F}}^2}{W_{\max}^2} + \|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\mathbb{F}}^2 \right)}_{\text{(iii)}} - \underbrace{C_2 \left( \frac{\epsilon_W^2 d}{W_{\min}^2} + \frac{D_x^2 \epsilon_A^2}{W_{\min}} \right)}_{\text{(iv)}}\end{aligned}\quad (75)$$

If  $\epsilon_W, \epsilon_A$  satisfy (70), then (iv)  $\leq c_0 D_{k,\ell}$  for some sufficiently small constant  $c_0 > 0$ . As for (iii), according to the definition of  $D_{k,\ell}$  in (68), there are two possible cases:

- If  $\frac{1}{\kappa_{w,\text{cross}}} \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}}^2}{W_{\max}^2} \geq \frac{D_{k,\ell}}{2}$ , then it is easy to check that (70) implies that

$$\frac{\epsilon_W \sqrt{d}}{W_{\max}} \leq \frac{1}{4} \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}}}{W_{\max}},$$

and hence

$$\begin{aligned} \text{(iii)} &\gtrsim \frac{1}{\kappa_{w,\text{cross}}} \frac{\|\mathbf{W}^{(k)} - \widehat{\mathbf{W}}^{(\ell)}\|_{\text{F}}^2}{W_{\max}^2} \geq \frac{1}{\kappa_{w,\text{cross}}} \left( \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}}}{W_{\max}} - \frac{\epsilon_W \sqrt{d}}{W_{\max}} \right)^2 \\ &\gtrsim \frac{1}{\kappa_{w,\text{cross}}} \frac{\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}}^2}{W_{\max}^2} \gtrsim D_{k,\ell}. \end{aligned}$$

- On the other hand, if  $\frac{1}{\kappa_{w,\text{cross}}} \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}^2 \geq \frac{D_{k,\ell}}{2}$ , then one can check that (70) implies that

$$\epsilon_A \sqrt{d} \leq \frac{1}{4} \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}},$$

and hence

$$\text{(iii)} \gtrsim \frac{\|\mathbf{A}^{(k)} - \widehat{\mathbf{A}}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}} \geq \frac{(\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}} - \epsilon_A \sqrt{d})^2}{\kappa_{w,\text{cross}}} \gtrsim \frac{\|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}^2}{\kappa_{w,\text{cross}}} \gtrsim D_{k,\ell}.$$

In sum, it is always guaranteed that (iii)  $\gtrsim D_{k,\ell}$ . Going back to (75), we claim that  $\widehat{D}_{k,\ell} \geq \text{(iii)} - \text{(iv)} \gtrsim D_{k,\ell}$  as long as  $\epsilon_W$  and  $\epsilon_A$  satisfy (70), which finishes the proof of (70).

#### C.5.6. PROOF OF (72).

Recall the lower bound (71) for (A) + (B) - (C). We want to show that the low-order term is dominated by the order- $T$  term, namely  $T \cdot D_{k,\ell}$ . First, if  $T \gtrsim \frac{\kappa_w \kappa_{w,\text{cross}}}{D_{k,\ell}} \log \frac{1}{\delta}$ , then  $\kappa_w \kappa_{w,\text{cross}} \log \frac{1}{\delta} \lesssim T \cdot D_{k,\ell}$ . Next, we have

$$\begin{aligned} &\left\| (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right\|_{\text{F}} \\ &\leq \|\mathbf{W}^{(k)}\| \cdot \left\| (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right\|_{\text{F}} = \|\mathbf{W}^{(k)}\| \cdot \left\| (\widehat{\mathbf{W}}^{(\ell)})^{-1} (\widehat{\mathbf{W}}^{(k)} - \widehat{\mathbf{W}}^{(\ell)}) (\widehat{\mathbf{W}}^{(k)})^{-1} \right\|_{\text{F}} \\ &\leq \|\mathbf{W}^{(k)}\| \cdot \left\| (\widehat{\mathbf{W}}^{(\ell)})^{-1} \right\| \cdot \left\| (\widehat{\mathbf{W}}^{(k)})^{-1} \right\| \cdot \left( \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} + 2\sqrt{d}\epsilon_W \right) \\ &\leq \frac{2W_{\max}}{W_{\min}^2} \left( \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} + 2\sqrt{d}\epsilon_W \right) = \frac{2\kappa_{w,\text{cross}}}{W_{\min}} \left( \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} + 2\sqrt{d}\epsilon_W \right). \end{aligned}$$

Notice that the definition (68) of  $D_{k,\ell}$  implies that  $\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} \leq \sqrt{\kappa_{w,\text{cross}} W_{\max}^2 D_{k,\ell}}$ , and thus

$$\left\| (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right\|_{\text{F}} \lesssim \frac{\kappa_{w,\text{cross}}}{W_{\min}} \left( \sqrt{\kappa_{w,\text{cross}} W_{\max}^2 D_{k,\ell}} + \sqrt{d}\epsilon_W \right).$$

Now it is easy to checked that

$$\begin{aligned} \text{if } T &\gtrsim \left( \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}} + \left( \frac{\kappa_{w,\text{cross}} \sqrt{d}\epsilon_W}{W_{\min} D_{k,\ell}} \right)^2 \right) \log^2 \frac{1}{\delta}, \\ \text{then } \sqrt{T} &\left\| (\mathbf{W}^{(k)})^{1/2} \left( (\widehat{\mathbf{W}}^{(\ell)})^{-1} - (\widehat{\mathbf{W}}^{(k)})^{-1} \right) (\mathbf{W}^{(k)})^{1/2} \right\|_{\text{F}} \log \frac{1}{\delta} \lesssim T \cdot D_{k,\ell}. \end{aligned}$$

Due to our assumption (70) on  $\epsilon_W$ , we have  $\left( \frac{\kappa_{w,\text{cross}} \sqrt{d}\epsilon_W}{W_{\min} D_{k,\ell}} \right)^2 \lesssim \frac{\kappa_{w,\text{cross}}^2}{D_{k,\ell}} \leq \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}}$ , and thus it suffices to have  $T \gtrsim \left( \frac{\kappa_{w,\text{cross}}^5}{D_{k,\ell}} + 1 \right) \log^2 \frac{1}{\delta}$ , which finishes the proof of (72).

## D. Miscellaneous results

**Proposition D.1.** Consider  $\mathbf{a}_t, \mathbf{b}_t \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_d), 1 \leq t \leq N$ , where  $N \gtrsim d + \log(1/\delta)$ . Then with probability at least  $1 - \delta$ ,

$$\left\| \frac{1}{N} \sum_{t=1}^N \mathbf{a}_t \mathbf{a}_t^\top - \mathbf{I}_d \right\| \lesssim \sqrt{\frac{d + \log \frac{1}{\delta}}{N}}, \quad \left\| \frac{1}{N} \sum_{t=1}^N \mathbf{a}_t \mathbf{b}_t^\top \right\| \lesssim \sqrt{\frac{d + \log \frac{1}{\delta}}{N}}.$$

The proof follows from a standard covering argument (cf. (Vershynin, 2018)), which we skip for brevity.

**Lemma D.2** (Self-normalized concentration, scalar version). Suppose that random vectors  $\{\mathbf{u}_t, \mathbf{y}_t\}_{1 \leq t \leq T}$  and filtrations  $\{\mathcal{F}_t\}_{0 \leq t \leq T-1}$  satisfy  $\mathcal{F}_t = \sigma(\mathbf{u}_i, 1 \leq i \leq t)$ ,  $\mathbf{y}_t \in \mathcal{F}_{t-1}$ , and  $\mathbf{u}_t | \mathcal{F}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_d)$ . Then for any fixed  $\lambda > 0$ , with probability at least  $1 - \delta$ , we have

$$\left| \sum_{t=1}^T \mathbf{u}_t^\top \mathbf{y}_t \right| < \frac{\lambda}{2} \sum_{t=1}^T \|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda} \log \frac{2}{\delta}.$$

*Proof.* Using basic properties of the Gaussian distribution, we have for all fixed  $\lambda \in \mathbb{R}$ ,

$$\begin{aligned} & \mathbb{E} \left[ \exp \left( \sum_{t=1}^T \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \right) \right] \\ &= \mathbb{E} \left[ \exp \left( \sum_{t=1}^{T-1} \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \right) \cdot \underbrace{\mathbb{E} \left[ \exp \left( \lambda \cdot \mathbf{u}_T^\top \mathbf{y}_T - \frac{1}{2} \lambda^2 \|\mathbf{y}_T\|_2^2 \right) | \mathcal{F}_{T-1} \right]}_{\leq 1} \right] \\ &\leq \mathbb{E} \left[ \exp \left( \sum_{t=1}^{T-1} \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \right) \right]. \end{aligned}$$

Continuing this expansion leads to the result  $\mathbb{E}[\exp(\sum_{t=1}^T (\lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2))] \leq 1$ . Now, letting  $z = \log(2/\delta)$  and using Markov's inequality, we have

$$\begin{aligned} \mathbb{P} \left( \sum_{t=1}^T \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \geq z \right) &= \mathbb{P} \left( \exp \left( \sum_{t=1}^T \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \right) \geq \exp(z) \right) \\ &\leq \exp(-z) \cdot \mathbb{E} \left[ \exp \left( \sum_{t=1}^T \left( \lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2 \right) \right) \right] \leq \exp(-z) = \frac{\delta}{2}. \end{aligned}$$

In other words, with probability at least  $1 - \delta/2$ , we have  $\sum_{t=1}^T (\lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2) < z = \log(2/\delta)$ , which implies  $\sum_{t=1}^T \mathbf{u}_t^\top \mathbf{y}_t < \frac{\lambda}{2} \sum_{t=1}^T \|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda} \log \frac{2}{\delta}$  if  $\lambda > 0$ . By a similar argument (but with  $\lambda$  replaced by  $-\lambda$ ), we have with probability at least  $1 - \delta/2$ ,  $\sum_{t=1}^T (-\lambda \cdot \mathbf{u}_t^\top \mathbf{y}_t - \frac{1}{2} \lambda^2 \|\mathbf{y}_t\|_2^2) < \log(2/\delta)$ , which implies  $\sum_{t=1}^T \mathbf{u}_t^\top \mathbf{y}_t > -(\frac{\lambda}{2} \sum_{t=1}^T \|\mathbf{y}_t\|_2^2 + \frac{1}{\lambda} \log \frac{2}{\delta})$  if  $\lambda > 0$ . Finally, taking the union bound finishes the proof of the lemma.  $\square$

**Lemma D.3.** Consider a symmetric matrix  $\mathbf{X} \in \mathbb{R}^{d \times d}$  that satisfies  $\mathbf{I}_d + \mathbf{X} \succ \mathbf{0}$ . If  $\|\mathbf{X}\| \leq B$  for some  $B \geq 2$ , then we have

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) \geq \frac{\|\mathbf{X}\|_{\text{F}}^2}{3B}.$$

On the other hand, if  $\mathbf{X} \succ -\frac{1}{2} \mathbf{I}_d$ , then we have

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) \leq \|\mathbf{X}\|_{\text{F}}^2.$$

*Proof.* Denote  $\{\lambda_i\}_{1 \leq i \leq d}$  as the eigenvalues of  $\mathbf{X}$ , which satisfies  $\lambda_i > -1$  for all  $1 \leq i \leq d$ . Then

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) = \sum_{i=1}^d \lambda_i - \log \prod_{i=1}^d (1 + \lambda_i) = \sum_{i=1}^d \left( \lambda_i - \log(1 + \lambda_i) \right).$$

It can be checked (via elementary calculus) that, for all  $-1 < \lambda \leq B$  where  $B \geq 2$ , it holds that  $\lambda - \log(1 + \lambda) \geq \lambda^2/(3B)$ . Therefore,

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) \geq \sum_{i=1}^d \frac{\lambda_i^2}{3B} = \frac{\|\mathbf{X}\|_{\text{F}}^2}{3B},$$

which completes the proof of our first claim. Similarly, it can be checked that, for all  $\lambda \geq -1/2$ , one has  $\lambda - \log(1 + \lambda) \leq \lambda^2$ , which implies that

$$\text{Tr}(\mathbf{X}) - \log \det(\mathbf{I}_d + \mathbf{X}) \leq \sum_{i=1}^d \lambda_i^2 = \|\mathbf{X}\|_{\text{F}}^2;$$

this completes the proof of our second claim.  $\square$

**Fact D.4.** *In the setting of Section 3.1, it holds that  $\Gamma_{\max} \leq W_{\max} \kappa_A^2 / (1 - \rho)$  and  $\Delta_{\Gamma, Y} \geq \Delta_{A, W} \cdot W_{\max} W_{\min} / (4\kappa_A^2 \Gamma_{\max})$ .*

*Proof.* First, consider  $\Gamma = \sum_{i=0}^{\infty} \mathbf{A}^i \mathbf{W} (\mathbf{A}^i)^{\top}$ , where  $\|\mathbf{W}\| \leq W_{\max}$  and  $\|\mathbf{A}^i\| \leq \kappa_A \rho^i$ . Then we have  $\|\Gamma\| \leq \sum_{i=0}^{\infty} \|\mathbf{A}^i\|^2 \|\mathbf{W}\| \leq W_{\max} \kappa_A^2 \sum_{i=0}^{\infty} \rho^{2i} \leq W_{\max} \kappa_A^2 / (1 - \rho)$ , which proves our upper bound for  $\Gamma_{\max}$ .

Next, let us turn to  $\Delta_{\Gamma, Y}$ . Our proof below can be viewed as a quantitative version of the earlier proof for Fact 2.1. We will show that, if the autocovariance matrices between the  $k$ -th and the  $\ell$ -th models are close (in Frobenius norm), then the models themselves should also be close; our lower bound for  $\Delta_{\Gamma, Y}$  in terms of  $\Delta_{A, W}$  then follows from contraposition.

Consider two LDS models  $(\mathbf{A}^{(k)}, \mathbf{W}^{(k)}) \neq (\mathbf{A}^{(\ell)}, \mathbf{W}^{(\ell)})$  and their autocovariance matrices  $(\Gamma^{(k)}, \mathbf{Y}^{(k)})$ ,  $(\Gamma^{(\ell)}, \mathbf{Y}^{(\ell)})$ . Recall from (17) that  $\mathbf{A}^{(k)} = \mathbf{Y}^{(k)} \Gamma^{(k)^{-1}}$  and  $\mathbf{W}^{(k)} = \Gamma^{(k)} - \mathbf{A}^{(k)} \Gamma^{(k)} \mathbf{A}^{(k)\top} = \Gamma^{(k)} - \mathbf{A}^{(k)} \mathbf{Y}^{(k)\top}$ ;  $\mathbf{A}^{(\ell)}$  and  $\mathbf{W}^{(\ell)}$  can be expressed similarly.

- First, regarding  $\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}$ , one has

$$\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)} = \mathbf{Y}^{(k)} \Gamma^{(k)^{-1}} - \mathbf{Y}^{(\ell)} \Gamma^{(\ell)^{-1}} = (\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}) \Gamma^{(k)^{-1}} + \mathbf{Y}^{(\ell)} (\Gamma^{(k)^{-1}} - \Gamma^{(\ell)^{-1}}),$$

where

$$\mathbf{Y}^{(\ell)} (\Gamma^{(k)^{-1}} - \Gamma^{(\ell)^{-1}}) = \mathbf{Y}^{(\ell)} \Gamma^{(\ell)^{-1}} (\Gamma^{(\ell)} - \Gamma^{(k)}) \Gamma^{(k)^{-1}} = \mathbf{A}^{(\ell)} (\Gamma^{(\ell)} - \Gamma^{(k)}) \Gamma^{(k)^{-1}}.$$

Therefore, we have

$$\begin{aligned} \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}} &\leq \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} \|\Gamma^{(k)^{-1}}\| + \|\mathbf{A}^{(\ell)}\| \|\Gamma^{(\ell)} - \Gamma^{(k)}\|_{\text{F}} \|\Gamma^{(k)^{-1}}\| \\ &\leq \frac{1}{W_{\min}} \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} + \frac{\kappa_A}{W_{\min}} \|\Gamma^{(\ell)} - \Gamma^{(k)}\|_{\text{F}}, \end{aligned} \quad (76)$$

where the last line follows from  $\Gamma^{(k)} \succcurlyeq \mathbf{W}^{(k)} \succcurlyeq W_{\min} \mathbf{I}_d$  and  $\|\mathbf{A}^{(\ell)}\| \leq \kappa_A \rho \leq \kappa_A$ .

- Next, we turn to the analysis of  $\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}$ , which satisfies

$$\begin{aligned} \mathbf{W}^{(k)} - \mathbf{W}^{(\ell)} &= (\Gamma^{(k)} - \Gamma^{(\ell)}) - (\mathbf{A}^{(k)} \mathbf{Y}^{(k)\top} - \mathbf{A}^{(\ell)} \mathbf{Y}^{(\ell)\top}), \\ \|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} &\leq \|\Gamma^{(k)} - \Gamma^{(\ell)}\|_{\text{F}} + \|\mathbf{A}^{(k)} \mathbf{Y}^{(k)\top} - \mathbf{A}^{(\ell)} \mathbf{Y}^{(\ell)\top}\|_{\text{F}}. \end{aligned}$$

Notice that

$$\begin{aligned} \|\mathbf{A}^{(k)} \mathbf{Y}^{(k)\top} - \mathbf{A}^{(\ell)} \mathbf{Y}^{(\ell)\top}\|_{\text{F}} &= \|\mathbf{A}^{(k)} (\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^{\top} + (\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}) \mathbf{Y}^{(\ell)\top}\|_{\text{F}} \\ &\leq \|\mathbf{A}^{(k)} (\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)})^{\top}\|_{\text{F}} + \|(\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}) \mathbf{Y}^{(\ell)\top}\|_{\text{F}} \\ &\leq \kappa_A \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} + \kappa_A \Gamma_{\max} \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}, \end{aligned}$$

where the last line is due to  $\|\mathbf{Y}^{(\ell)}\| = \|\mathbf{A}^{(\ell)} \Gamma^{(\ell)}\| \leq \|\mathbf{A}^{(\ell)}\| \cdot \|\Gamma^{(\ell)}\| \leq \kappa_A \Gamma_{\max}$ . Therefore, we have

$$\|\mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}\|_{\text{F}} \leq \|\Gamma^{(k)} - \Gamma^{(\ell)}\|_{\text{F}} + \kappa_A \|\mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}\|_{\text{F}} + \kappa_A \Gamma_{\max} \|\mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}\|_{\text{F}}. \quad (77)$$



For notational simplify, denote  $\Delta_A := \mathbf{A}^{(k)} - \mathbf{A}^{(\ell)}$ ,  $\Delta_W := \mathbf{W}^{(k)} - \mathbf{W}^{(\ell)}$ ,  $\Delta_\Gamma := \Gamma^{(k)} - \Gamma^{(\ell)}$ ,  $\Delta_Y := \mathbf{Y}^{(k)} - \mathbf{Y}^{(\ell)}$ . From (77), one has

$$\begin{aligned} \|\Delta_A\|_F^2 + \frac{\|\Delta_W\|_F^2}{W_{\max}^2} &\leq \|\Delta_A\|_F^2 + \frac{3}{W_{\max}^2} \left( \|\Delta_\Gamma\|_F^2 + \kappa_A^2 \|\Delta_Y\|_F^2 + \kappa_A^2 \Gamma_{\max}^2 \|\Delta_A\|_F^2 \right) \\ &\leq \frac{3}{W_{\max}^2} \|\Delta_\Gamma\|_F^2 + \frac{3\kappa_A^2}{W_{\max}^2} \|\Delta_Y\|_F^2 + \left( 1 + \frac{3\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \right) \|\Delta_A\|_F^2 \\ &\leq \frac{3}{W_{\max}^2} \|\Delta_\Gamma\|_F^2 + \frac{3\kappa_A^2}{W_{\max}^2} \|\Delta_Y\|_F^2 + \frac{4\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \|\Delta_A\|_F^2, \end{aligned}$$

where the last line follows from  $\kappa_A^2 \Gamma_{\max}^2 / W_{\max}^2 \geq 1$ . Moreover, (76) tells us that  $\|\Delta_A\|_F^2 \leq \frac{2}{W_{\min}^2} \|\Delta_Y\|_F^2 + \frac{2\kappa_A^2}{W_{\min}^2} \|\Delta_\Gamma\|_F^2$ . Putting together, we have

$$\begin{aligned} \|\Delta_A\|_F^2 + \frac{\|\Delta_W\|_F^2}{W_{\max}^2} &\leq \frac{3}{W_{\max}^2} \|\Delta_\Gamma\|_F^2 + \frac{3\kappa_A^2}{W_{\max}^2} \|\Delta_Y\|_F^2 + \frac{4\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \|\Delta_A\|_F^2 \\ &\leq \frac{3}{W_{\max}^2} \|\Delta_\Gamma\|_F^2 + \frac{3\kappa_A^2}{W_{\max}^2} \|\Delta_Y\|_F^2 + \frac{4\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \left( \frac{2}{W_{\min}^2} \|\Delta_Y\|_F^2 + \frac{2\kappa_A^2}{W_{\min}^2} \|\Delta_\Gamma\|_F^2 \right) \\ &= \left( \frac{3}{W_{\max}^2} + \frac{4\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \frac{2\kappa_A^2}{W_{\min}^2} \right) \|\Delta_\Gamma\|_F^2 + \left( \frac{3\kappa_A^2}{W_{\max}^2} + \frac{4\kappa_A^2 \Gamma_{\max}^2}{W_{\max}^2} \frac{2}{W_{\min}^2} \right) \|\Delta_Y\|_F^2 \\ &\leq \frac{11\kappa_A^4 \Gamma_{\max}^2}{W_{\max}^2 W_{\min}^2} \left( \|\Delta_\Gamma\|_F^2 + \|\Delta_Y\|_F^2 \right). \end{aligned}$$

In sum, we have just shown that, if  $\|\Delta_\Gamma\|_F^2 + \|\Delta_Y\|_F^2 < \Delta_{\Gamma,Y}^2$ , then  $\|\Delta_A\|_F^2 + \frac{\|\Delta_W\|_F^2}{W_{\max}^2} < \frac{11\kappa_A^4 \Gamma_{\max}^2}{W_{\max}^2 W_{\min}^2} \Delta_{\Gamma,Y}^2$ . Equivalently (by contraposition), if  $\|\Delta_A\|_F^2 + \frac{\|\Delta_W\|_F^2}{W_{\max}^2} \geq \Delta_{A,W}^2$ , then  $\|\Delta_\Gamma\|_F^2 + \|\Delta_Y\|_F^2 \geq \frac{W_{\max}^2 W_{\min}^2}{11\kappa_A^4 \Gamma_{\max}^2} \Delta_{A,W}^2$ . This proves our lower bound for  $\Delta_{\Gamma,Y}$  in terms of  $\Delta_{A,W}$ .  $\square$

**Example D.5.** It has been known that in our Case 1 (i.e. a single continuous trajectory), the quick switching of multiple LDS models may lead to exponentially large states, even if each individual model is stable (Liberzon, 2003). We give a quick example for completeness. Consider

$$\mathbf{A}^{(k)} = 0.99 \begin{bmatrix} 0 & 2 \\ \frac{1}{2} & 0 \end{bmatrix}, \quad \mathbf{A}^{(\ell)} = 0.99 \begin{bmatrix} 0 & 3 \\ \frac{1}{3} & 0 \end{bmatrix},$$

both satisfying the stability condition in Assumption 3.1 with  $\rho = 0.99 < 1$  and hence  $t_{\text{mix}} \asymp 1/(1-\rho) = 100$ . Suppose that each short trajectory has only a length of 2, and the  $m$ -th (resp.  $(m+1)$ -th) trajectory has label  $k_m = \ell$  (resp.  $k_{m+1} = k$ ). Then  $\mathbf{x}_{m+2,0}$  is equal to  $\mathbf{A}^{(k)} \mathbf{A}^{(\ell)} \mathbf{x}_{m,0}$  plus a mean-zero noise term, where

$$\mathbf{A}^{(k)} \mathbf{A}^{(\ell)} = 0.99^2 \begin{bmatrix} 0 & 2 \\ \frac{1}{2} & 0 \end{bmatrix} \begin{bmatrix} 0 & 3 \\ \frac{1}{3} & 0 \end{bmatrix} = 0.99^2 \begin{bmatrix} \frac{2}{3} & 0 \\ 0 & \frac{3}{2} \end{bmatrix}$$

has spectral radius  $0.99^2 \cdot 3/2 > 1$ ; this will cause the exponential explosion of the states.

## E. Extensions of Algorithm 1

**Different trajectory lengths.** Recall that in Section 2, we assume that all short trajectories within each subset of data  $\mathcal{M}_o$  have the same length  $T_m = T_o$ . If this is not the case, we can easily modify our algorithms in the following ways:

- For subspace estimation, the easiest way to handle different  $T_m$ 's is to simply *truncate* the trajectories in  $\mathcal{M}_{\text{subspace}}$  so that they have the same length  $T_{\text{subspace}} = \min_{m \in \mathcal{M}_{\text{subspace}}} T_m$ , and then apply Algorithm 2 without modification. However, this might waste many samples when some trajectories of  $\mathcal{M}_{\text{subspace}}$  are much longer than others; one way to resolve this is to *manually divide* the longer trajectories into shorter segments of comparable lengths, before doing truncation. A more refined method is to modify Algorithm 2 itself, by *re-defining the index sets*  $\Omega_1, \Omega_2$  separately for each trajectory; moreover, in the definition of  $\widehat{\mathbf{H}}_i$  and  $\widehat{\mathbf{G}}_i$ , one might consider assigning larger weights to longer trajectories, instead of using the uniform weight  $1/|\mathcal{M}_{\text{subspace}}|$ .

- For clustering (or pairwise testing) of  $\mathcal{M}_{\text{clustering}}$ , we can handle various  $T_m$ 's similarly, by either truncating each pair of trajectories to the same length, or modifying Algorithm 3 itself (via re-defining the index sets  $\{\Omega_{g,1}, \Omega_{g,2}\}_{1 \leq g \leq G}$  separately for each trajectory).
- Our methods for model estimation and classification, namely Algorithms 4 and 5, are already adaptive to different  $T_m$ 's in  $\mathcal{M}_{\text{clustering}}$  and  $\mathcal{M}_{\text{classification}}$ , and hence need no modification.

**Unknown parameters.** Next, we show how to handle the case when certain parameters are unknown to the algorithms:

- In Algorithm 2, we set the dimension of the output subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  to be  $K$  (the number of models). If  $K$  is unknown, we might instead examine the eigenvalues of  $\widehat{\mathbf{H}}_i + \widehat{\mathbf{H}}_i^\top$  and  $\widehat{\mathbf{G}}_i + \widehat{\mathbf{G}}_i^\top$ , and pick the subspace dimension that covers most of the energy in the eigenvalues.
- In Algorithm 3, we need to know the separation parameter  $\Delta_{\Gamma, Y}$  (in order to choose the testing threshold  $\tau$  appropriately) and the number of models  $K$  (for clustering). If either  $\Delta_{\Gamma, Y}$  or  $K$  is unknown, we might instead try different values of threshold  $\tau$ , and pick the one that (after permutation) makes the similarity matrix  $\mathbf{S}$  block-diagonal with as few blocks as possible.

## F. Additional experiments

### F.1. Synthetic experiments: clustering and classification

First, we take a closer look at the performance of our clustering method (Algorithm 3) through synthetic experiments. We set the parameters  $d = 40, K = 2, \rho = 0.5, \delta = 0.12$ . The LDS models are generated by  $\mathbf{A}^{(k)} = (\rho \pm \delta)\mathbf{R}$  and  $\mathbf{W}^{(k)} = \mathbf{I}_d$ , where  $\mathbf{R}$  is a random orthogonal matrix. We let  $|\mathcal{M}_{\text{clustering}}| = 5d$ , and vary  $T_{\text{clustering}} \in [10, 60]$ . We run our clustering method on the dataset  $\mathcal{M}_{\text{clustering}}$ , either with or without the assistance of subspace estimation (Algorithm 2) and dimensionality reduction. For the former case, we use the same dataset  $\mathcal{M}_{\text{clustering}}$  for subspace estimation, without sample splitting, which is closer to practice; for the latter, we simply replace the subspaces  $\{\mathbf{V}_i, \mathbf{U}_i\}$  with  $\mathbf{I}_d$ . The numerical results are illustrated in Figure 3 (left), confirming that (1) in both cases, the clustering error decreases as  $T_{\text{clustering}}$  increases, and (2) subspace estimation and dimensionality reduction significantly improve the clustering accuracy.

Next, we examine the performance of our classification method (Algorithm 5) in the same setting as above. We first obtain a coarse model estimation by running Stage 1 of Algorithm 1 on the dataset  $\mathcal{M}_{\text{clustering}}$ , with  $|\mathcal{M}_{\text{clustering}}| = 10d$  and  $T_{\text{clustering}} = 30$ . Then, we run classification on the dataset  $\mathcal{M}_{\text{classification}}$ , with varying  $T_{\text{classification}} \in [4, 50]$ . The numerical results are included in Figure 3 (right), showing that the classification error rapidly decreases to zero as  $T_{\text{classification}}$  grows.

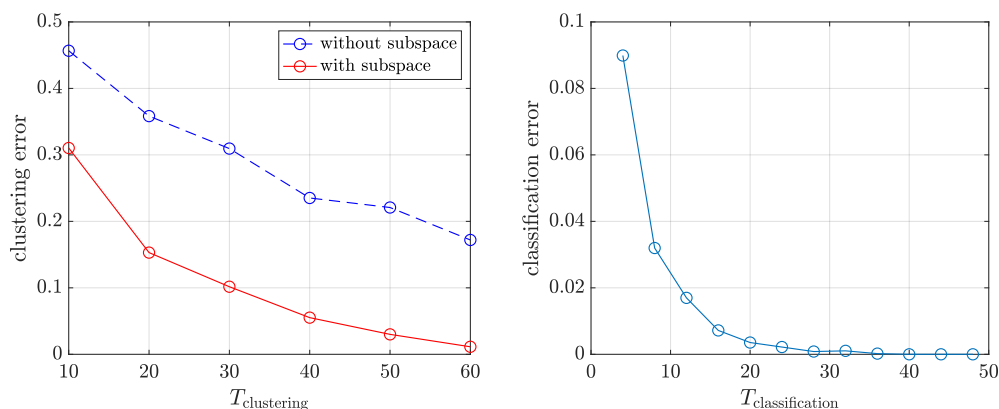


Figure 3. **Left:** mis-clustering rate versus  $T_{\text{clustering}}$ . **Right:** mis-classification rate versus  $T_{\text{classification}}$ .

### F.2. Real-data experiments: MotionSense

To show the practical relevance of the proposed algorithms, we work with the MotionSense dataset (Malekzadeh et al., 2019). This dataset consists of multivariate time series of dimension  $d = 12$ , collected (at a rate of 50Hz) by accelerometer

and gyroscope sensors on a mobile phone while a person performs various activities, such as “jogging”, “walking”, “sitting”, and so on. In our experiments, we break the data into 8-second short trajectories, and treat the human activities as latent variables. Figure 4 (left) illustrates what the data looks like. Notice that the time series do not satisfy the mixing property assumed in our theory, but are rather periodic instead.

As a preliminary attempt to apply our algorithms in the real world, we show that the proposed clustering method (which is one of the most crucial step in our overall approach), without any modification, works reasonably well even for this dataset. To be concrete, we apply Algorithm 3 (without dimensionality reduction, i.e.  $\{\mathbf{V}_i, \mathbf{U}_i\}$  are set to  $\mathbf{I}_d$ ) to a mixture of 12 “jogging” and 12 “walking” trajectories. Figure 4 (right) shows the resulted *distance matrix*, which is defined in the same way as Line 11 of Algorithm 3, but without thresholding. Its clear block structure confirms that, with an appropriate choice of threshold  $\tau$ , Algorithm 3 will return an accurate/exact clustering of the mixed trajectories. These results are strong indication that the proposed algorithms in this work might generalize to much broader settings than what our current theory suggests, and we hope that this will inspire further extensions and applications of the proposed methods.

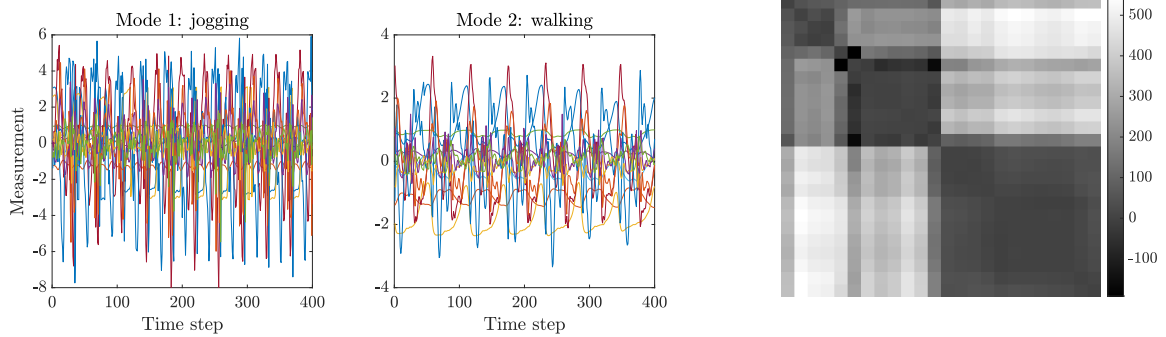


Figure 4. **Left:** examples of “jogging” and “walking” trajectories from the MotionSense dataset. **Right:** the distance matrix constructed by Algorithm 3 for 12 “jogging” and 12 “walking” trajectories.