
Learning Hidden Markov Models from Non-sequence Data via Tensor Decomposition

Tzu-Kuo Huang
Machine Learning Department
Carnegie Mellon University
Pittsburgh, PA 15213
tzukuoh@cs.cmu.edu

Jeff Schneider
Robotics Institute
Carnegie Mellon University
Pittsburgh, PA 15213
schneide@cs.cmu.edu

Abstract

Learning dynamic models from observed data has been a central issue in many scientific studies or engineering tasks. The usual setting is that data are collected sequentially from trajectories of some dynamical system operation. In quite a few modern scientific modeling tasks, however, it turns out that reliable sequential data are rather difficult to gather, whereas out-of-order snapshots are much easier to obtain. Examples include the modeling of galaxies, chronic diseases such as Alzheimer's, or certain biological processes.

Existing methods for learning dynamic model from non-sequence data are mostly based on Expectation-Maximization, which involves non-convex optimization and is thus hard to analyze. Inspired by recent advances in spectral learning methods, we propose to study this problem from a different perspective: moment matching and spectral decomposition. Under that framework, we identify reasonable assumptions on the generative process of non-sequence data, and propose learning algorithms based on the tensor decomposition method [2] to *provably* recover first-order Markov models and hidden Markov models. To the best of our knowledge, this is the first formal guarantee on learning from non-sequence data. Preliminary simulation results confirm our theoretical findings.

1 Introduction

Learning dynamic models from observed data has been a central issue in many fields of study, scientific or engineering tasks. The usual setting is that data are collected sequentially from trajectories of some dynamical system operation, and the goal is to recover parameters of the underlying dynamic model. Although many research and engineering efforts have been devoted to that setting, it turns out that in quite a few modern scientific modeling problems, another situation is more frequently encountered: observed data are out-of-order (or partially-ordered) snapshots rather than full sequential samples of the system operation. As pointed out in [7, 8], this situation may appear in the modeling of celestial objects such as galaxies or chronic diseases such as Alzheimer's, because observations are usually taken from different trajectories (galaxies or patients) at unknown, arbitrary times. Or it may also appear in the study of biological processes, such as cell metabolism under external stimuli, where most measurement techniques are destructive, making it very difficult to repetitively collect observations from the same individual living organisms as they change over time. However, it is much easier to take single snapshots of multiple organisms undergoing the same biological process in a fully asynchronous fashion, hence the lack of timing information. Rabbat et al. [9] noted that in certain network inference problems, the only available data are sets of nodes *co-occurring* in random walks on the network without the order in which they were visited, and the goal is to reconstruct the network structure from such co-occurrence data. This problem is essentially about learning a first-order Markov chain from data lacking sequence information.

As one can imagine, dynamic model learning in a non-sequential setting is much more difficult than in the sequential setting and has not been thoroughly studied. One issue is that the notion of non-sequence data is vague because there can be many different generative processes resulting in non-sequence data. Without any restrictions, one can easily find a case where no meaningful dynamic model can be learnt. It is therefore important to figure out what assumptions on the data and the model would lead to successful learning. However, existing methods for non-sequential settings, e.g., [9, 11, 6, 8], do not shed much light on this issue because they are mostly based on Expectation-Maximization (EM), which require non-convex optimization. Regardless of the assumptions we make, as long as the resulting optimization problem remains non-convex, formal analysis of learning guarantees is still formidable.

We thus propose to take a different approach, based on another long-standing estimation principle: *the method of moments* (MoM). The basic idea of MoM is to find model parameters such that the resulting moments match or resemble the empirical moments. For some estimation problems, this approach is able to give unique and consistent estimates while the maximum-likelihood method gets entangled in multiple and potentially undesirable local maxima. Taking advantage of this property, an emerging area of research in machine learning has recently developed MoM-based learning algorithms *with formal guarantees* for some widely used latent variable models, such as Gaussian mixture models[5], Hidden Markov models [3], Latent Dirichlet Allocation [1, 4], etc. Although many learning algorithms for these models exist, some having been very successful in practice, barely any formal learning guarantee was given until the MoM-based methods were proposed. Such breakthroughs seem surprising, but it turns out that they are mostly based on one crucial property: for quite a few latent variable models, the model parameters can be uniquely determined from *spectral decompositions* of certain low-order moments of observable quantities.

In this work we demonstrate that under the MoM and spectral learning framework, there are reasonable assumptions on the generative process of non-sequence data, under which *the tensor decomposition method* [2], a recent advancement in spectral learning, can provably recover the parameters of *first-order Markov models* and *hidden Markov models*. To the best of our knowledge, ours is the first work that provides formal guarantees for learning from non-sequence data. Interestingly, these assumptions bear much similarity to the usual idea behind *topic modeling*: with the bag-of-words representation which is *invariant to word orderings*, the task of inferring topics is almost impossible given *one single document* (no matter how long it is!), but becomes easier as more documents touching upon various topics become available. For learning dynamic models, what we need in the non-sequence data are *multiple sets* of observations, where each set contains independent samples generated from *its own initial distribution*, and the many different initial distributions together cover the entire (hidden) state space. In some of the aforementioned scientific applications, such as biological studies, this type of assumptions might be realized by running multiple experiments with different initial configurations or amounts of stimuli.

The main body of the paper consists of four sections. Section 2 briefly reviews the essentials of the tensor decomposition framework [2]; Section 3 details our assumptions on non-sequence data, tensor-decomposition based learning algorithms, and theoretical guarantees; Section 4 reports some simulation results confirming our theoretical findings, followed by conclusions in Section 5. Proofs of theoretical results are given in the appendices in the supplementary material.

2 Tensor Decomposition

We mainly follow the exposition in [2], starting with some preliminaries and notations. A real p -th order tensor A is a member of the tensor product space $\bigotimes_{i=1}^p \mathbb{R}^{m_i}$ of p Euclidean spaces. For a vector $\mathbf{x} \in \mathbb{R}^m$, we denote by $\mathbf{x}^{\otimes p} := \mathbf{x} \otimes \mathbf{x} \otimes \cdots \otimes \mathbf{x} \in \bigotimes_{i=1}^p \mathbb{R}^m$ its p -th tensor power. A convenient way to represent $A \in \bigotimes_{i=1}^p \mathbb{R}^{m_i}$ is through a p -way array of real numbers $[A_{i_1 i_2 \dots i_p}]_{1 \leq i_1, i_2, \dots, i_p \leq m_i}$, where $A_{i_1 i_2 \dots i_p}$ denotes the (i_1, i_2, \dots, i_p) -th coordinate of A with respect to a canonical basis. With this representation, we can view A as a multi-linear map that, given a set of p matrices $\{X_i \in \mathbb{R}^{m_i \times m_i}\}_{i=1}^p$, produces another p -th order tensor $A(X_1, X_2, \dots, X_p) \in \bigotimes_{i=1}^p \mathbb{R}^{m_i}$ with the following p -way array representation:

$$A(X_1, X_2, \dots, X_p)_{i_1 i_2 \dots i_p} := \sum_{1 \leq j_1, j_2, \dots, j_p \leq m_i} A_{j_1 j_2 \dots j_p} (X_1)_{j_1 i_1} (X_2)_{j_2 i_2} \cdots (X_p)_{j_p i_p}. \quad (1)$$

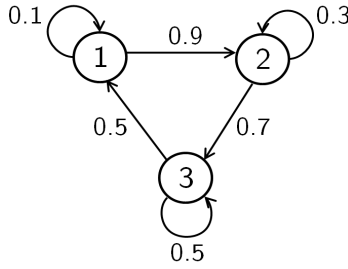


Figure 1: Running example of Markov chain with three states

In this work we consider tensors that are up to the third-order ($p \leq 3$) and, for most of the time, also *symmetric*, meaning that their p -way array representations are invariant under permutations of array indices. More specifically, we focus on second and third-order symmetric tensors in or slightly perturbed from the following form:

$$M_2 := \sum_{i=1}^k \omega_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i, \quad M_3 := \sum_{i=1}^k \omega_i \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i, \quad (2)$$

satisfying the following non-degeneracy conditions:

Condition 1. $\omega_i \geq 0 \forall 1 \leq i \leq k$, $\{\boldsymbol{\mu}_i \in \mathbb{R}^m\}_{i=1}^k$ are linearly independent, and $k \leq m$.

As described in later sections, the core of our learning task involves estimating $\{\omega_i\}_{i=1}^k$ and $\{\boldsymbol{\mu}_i\}_{i=1}^k$ from perturbed or noisy versions of M_2 and M_3 . We solve this estimation problem with the tensor decomposition method recently proposed by Anandkumar et al. [2]. The algorithm and its theoretical guarantee are summarized in Appendix A. The key component of this method is a novel tensor power iteration procedure for factorizing a symmetric orthogonal tensor, which is robust against input perturbation.

3 Learning from Non-sequence Data

We first describe a generative process of non-sequence data for first-order Markov models and demonstrate how to apply tensor decomposition methods to perform consistent learning. Then we extend these ideas to hidden Markov models and provide theoretical guarantees on the sample complexity of the proposed learning algorithm. For notational conveniences we define the following vector-matrix cross product $\otimes_{d \in \{1,2,3\}} : (\mathbf{v} \otimes_1 M)_{ijk} := v_i(M)_{jk}$, $(\mathbf{v} \otimes_2 M)_{ijk} = v_j(M)_{ik}$, $(\mathbf{v} \otimes_3 M)_{ijk} = v_k(M)_{ij}$. For a matrix M we denote by M_i its i -th column.

3.1 First-order Markov Models

Let $P \in [0, 1]^{m \times m}$ be the transition probability matrix of a discrete, first-order, ergodic Markov chain with m states and a unique stationary distribution $\boldsymbol{\pi}$. Let P be of full rank and $\mathbf{1}^\top P = \mathbf{1}^\top$. To give a high-level idea of what makes it possible to learn P from non-sequence data, we use the simple Markov chain with three states shown in Figure 1 as our running example, demonstrating step by step how to extend from a very restrictive generative setting of the data to a reasonably general setting, along with the assumptions made to allow consistent parameter estimation. In the usual setting where we have sequences of observations, say $\{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots\}$ with parenthesized superscripts denoting time, it is straightforward to consistently estimate P . We simply calculate the empirical frequency of consecutive pairs of states:

$$\widehat{P}_{ij} := \frac{\sum_t (\mathbf{x}^{(t+1)} = i, \mathbf{x}^{(t)} = j)}{\sum_t (\mathbf{x}^{(t)} = j)}.$$

Alternatively, suppose for each state j , we have an *i.i.d. sample* of its immediate next state $D_j := \{\mathbf{x}_1^{(1)}, \mathbf{x}_2^{(1)}, \dots \mid \mathbf{x}^{(0)} = j\}$, where subscripts are data indices. Consistent estimation in this case is also easy: the empirical distribution of D_j consistently estimates P_j , the j -th column of P . For

example, the Markov chain in Figure 1 may produce the following three samples, whose empirical distributions estimate the three columns of P respectively:

$$\begin{aligned} D_1 &= \{2, 1, 2, 2, 2, 2, 2, 2, 2\} \Rightarrow \widehat{P}_1 = [0.1 \ 0.9 \ 0.0]^\top, \\ D_2 &= \{3, 3, 2, 3, 2, 3, 3, 2, 3, 3\} \Rightarrow \widehat{P}_2 = [0.0 \ 0.3 \ 0.7]^\top, \\ D_3 &= \{1, 1, 3, 1, 3, 3, 1, 3, 3, 1\} \Rightarrow \widehat{P}_3 = [0.5 \ 0.0 \ 0.5]^\top. \end{aligned}$$

A nice property of these estimates is that, unlike in the sequential setting, they do not depend on any particular ordering of the observations in each set. Nevertheless, such data are not quite non-sequenced because all observations are made at exactly the next time step. We thus consider the following generalization: for each state j , we have $D_j := \{\mathbf{x}_1^{(t_1)}, \mathbf{x}_2^{(t_2)}, \dots \mid \mathbf{x}^{(0)} = j\}$, i.e., independent samples of states drawn at *unknown* future times $\{t_1, t_2, \dots\}$. For example, our data in this setting might be

$$\begin{aligned} D_1 &= \{2, 1, 2, 3, 2, 3, 3, 2, 2, 3\}, \\ D_2 &= \{3, 3, 2, 3, 2, 1, 3, 2, 3, 1\}, \\ D_3 &= \{1, 1, 3, 1, 2, 3, 2, 3, 3, 2\}. \end{aligned} \tag{3}$$

Obviously it is hard to extract information about P from such data. However, if we assume that the unknown times $\{t_i\}$ are i.i.d. random variables following some distribution independent of the initial state j , it can then be easily shown that D_j 's empirical distribution consistently estimates T_j , the j -th column of the the *expected transition probability matrix* $T := \mathbb{E}_t[P^t]$:

$$\begin{aligned} D_1 &= \{2, 1, 2, 3, 2, 3, 3, 2, 2, 3\} \Rightarrow \widehat{T}_1 = [0.1 \ 0.5 \ 0.4]^\top, \\ D_2 &= \{3, 3, 2, 3, 2, 1, 3, 2, 3, 1\} \Rightarrow \widehat{T}_2 = [0.2 \ 0.3 \ 0.5]^\top, \\ D_3 &= \{1, 1, 3, 1, 2, 3, 2, 3, 3, 2\} \Rightarrow \widehat{T}_3 = [0.3 \ 0.3 \ 0.4]^\top. \end{aligned}$$

In general there exist many P 's that result in the same T . Therefore, as detailed later, we make a specific distributional assumption on $\{t_i\}$ to enable unique recovery of the transition matrix P from T (Assumption A.1). Next we consider a further generalization, where the unknowns are not only the time stamps of the observations, but also the initial state j . In other words, we only know each set was generated from the same initial state, but do not know the actual initial state. In this case, the empirical distributions of the sets consistently estimate the columns of T in some *unknown permutation* Π :

$$\begin{aligned} D_{\Pi(3)} &= \{1, 1, 3, 1, 2, 3, 2, 3, 3, 2\} \Rightarrow \widehat{T}_{\Pi(3)} = [0.3 \ 0.3 \ 0.4]^\top, \\ D_{\Pi(2)} &= \{3, 3, 2, 3, 2, 1, 3, 2, 3, 1\} \Rightarrow \widehat{T}_{\Pi(2)} = [0.2 \ 0.3 \ 0.5]^\top, \\ D_{\Pi(1)} &= \{2, 1, 2, 3, 2, 3, 3, 2, 2, 3\} \Rightarrow \widehat{T}_{\Pi(1)} = [0.1 \ 0.5 \ 0.4]^\top. \end{aligned}$$

In order to be able to identify Π , we will again resort to randomness and assume the unknown initial states are random variables following a certain distribution (Assumption A.2) so that the data carry information about Π . Finally, we generalize from a single unknown initial state to an unknown *initial state distribution*, where each set of observations $D := \{\mathbf{x}_1^{(t_1)}, \mathbf{x}_2^{(t_2)}, \dots \mid \pi^{(0)}\}$ consists of independent samples of states drawn at random times from some unknown initial state distribution $\pi^{(0)}$. For example, the data may look like:

$$\begin{aligned} D_{\pi_1^{(0)}} &= \{1, 3, 3, 1, 2, 3, 2, 3, 3, 2\}, \\ D_{\pi_2^{(0)}} &= \{3, 1, 2, 3, 2, 1, 3, 2, 3, 1\}, \\ D_{\pi_3^{(0)}} &= \{2, 1, 2, 3, 3, 3, 3, 1, 2, 3\}, \\ &\vdots \end{aligned}$$

With this final generalization, most would agree that the generated data are non-sequenced and that the generative process is flexible enough to model the real-world situations described in Section 1. However, simple estimation with empirical distributions no longer works because each set may now contain observations from multiple initial states. This is where we take advantage of the tensor

decomposition framework outlined in Section 2, which requires proper assumptions on the initial state distribution $\pi^{(0)}$ (Assumption A.3).

Now we are ready to give the definition of our entire generative process. Assume we have N sets of non-sequence data each containing n observations, and each set of observations $\{\mathbf{x}_i\}_{i=1}^n$ were independently generated by the following:

- Draw an initial distribution $\pi^{(0)} \sim \text{Dirichlet}(\boldsymbol{\alpha})$, (Assumption A.3)
 $\mathbb{E}[\pi^{(0)}] = \boldsymbol{\alpha} / (\sum_{i=1}^m \alpha_i) = \boldsymbol{\pi}$, $\pi_i \neq \pi_j \forall i \neq j$. (Assumption A.2)
- For $i = 1, \dots, n$,
 - Draw a discrete time $t_i \sim \text{Geometric}(r)$, $t_i \in \{1, 2, 3, \dots\}$. (Assumption A.1)
 - Draw an initial state $\mathbf{s}_i \sim \text{Multinomial}(\boldsymbol{\pi}_0)$, $\mathbf{s}_i \in \{0, 1\}^m$.
 - Draw an observation $\mathbf{x}_i \sim \text{Multinomial}(P^{t_i} \mathbf{s}_i)$, $\mathbf{x}_i \in \{0, 1\}^m$.

The above generative process has several properties. First, all the data points in the same set share the same initial state distribution but can have different initial states; the initial state distribution varies across different sets and yet centers at the stationary distribution of the Markov chain. As mentioned in Section 1, this may be achieved in biological studies by running multiple experiments with different input stimuli, so the data collected in the same experiment can be assumed to have the same initial state distribution. Second, each data point is drawn from an independent trajectory of the Markov chain, a similar situation in the modeling of galaxies or Alzheimer’s, and random time steps could be used to compensate for individual variations in speed: a small/large t_i corresponds to a slowly/fast evolving individual object. Finally, the geometric distribution can be interpreted as an overall measure of the magnitude of speed variation: a large success probability r would result in many small t_i ’s, meaning that most objects evolve at similar speeds, while a small r would lead to t_i ’s taking a wide range of values, indicating a large speed variation.

To use the tensor decomposition method in Appendix A, we need the tensor structure (2) in certain low-order moments of observed quantities. The following theorem identifies such quantities:

Theorem 1. *Define the expected transition probability matrix $T := \mathbb{E}_t[P^t] = rP(I - (1-r)P)^{-1}$ and let $\alpha_0 := \sum_i \alpha_i$, $C_2 := \mathbb{E}[\mathbf{x}_1 \mathbf{x}_2^\top]$ and $C_3 := \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3]$. Then the following holds:*

$$\mathbb{E}[\mathbf{x}_1] = \boldsymbol{\pi}, \quad C_2 = \frac{1}{\alpha_0 + 1} T \text{diag}(\boldsymbol{\pi}) T^\top + \frac{\alpha_0}{\alpha_0 + 1} \boldsymbol{\pi} \boldsymbol{\pi}^\top, \quad (4)$$

$$C_3 = \frac{2}{(\alpha_0 + 2)(\alpha_0 + 1)} \sum_i \pi_i T_i^{\otimes 3} + \frac{\alpha_0}{\alpha_0 + 2} \sum_{d=1}^3 \boldsymbol{\pi} \otimes_d C_2 - \frac{2\alpha_0^2}{(\alpha_0 + 2)(\alpha_0 + 1)} \boldsymbol{\pi}^{\otimes 3}, \quad (5)$$

$$M_2 := (\alpha_0 + 1)C_2 - \alpha_0 \boldsymbol{\pi} \boldsymbol{\pi}^\top = T \text{diag}(\boldsymbol{\pi}) T^\top, \quad (6)$$

$$M_3 := \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} C_3 - \frac{(\alpha_0 + 1)\alpha_0}{2} \sum_{d=1}^3 \boldsymbol{\pi} \otimes_d C_2 + \alpha_0^2 \boldsymbol{\pi}^{\otimes 3} = \sum_i \pi_i T_i^{\otimes 3}. \quad (7)$$

The proof is in Appendix B.1, which relies on the special structure in the moments of the Dirichlet distribution (Assumption A.3). It is clear that M_2 and M_3 have the desired tensor structure. Assuming α_0 is known, we can form estimates \widehat{M}_2 and \widehat{M}_3 by computing empirical moments from the data. Note that the \mathbf{x}_i ’s are exchangeable, so we can use all pairs and triples of data points to compute the estimates. Interestingly, these low-order moments have a very similar structure to those in Latent Dirichlet Allocation [1]. Indeed, according to our generative process, we can view a set of non-sequence data points as a document generated by an LDA model with the expected transition matrix T as the topic matrix, the stationary distribution $\boldsymbol{\pi}$ as the topic proportions, and most importantly, the states as *both the words and the topics*. The last property is what distinguishes our generative process from a general LDA model: because both the words and the topics correspond to the states, the topic matrix is no longer invariant to column permutations. Since the tensor decomposition method may return \widehat{T} under any column permutation, we need to recover the correct matching between its rows and columns. Note that the $\widehat{\boldsymbol{\pi}}$ returned by the tensor decomposition method undergoes the same permutation as \widehat{T} ’s columns. Because all π_i ’s have different values by Assumption A.2, we may recover the correct matching by sorting both the returned $\widehat{\boldsymbol{\pi}}$ and the mean $\bar{\boldsymbol{\pi}}$ of all data.

A final issue is estimating P and r from \widehat{T} . This is in general difficult even when the exact T is available because multiple choices of P and r may result in the same T . However, if the true transition matrix P has at least one zero entry, then unique recovery is possible:

Theorem 2. Let P^* , r^* , T^* and π^* denote the true values of the transition probability matrix, the success probability, the expected transition matrix, and the stationary distribution, respectively. Assume that P^* is ergodic and of full rank, and $P_{ij}^* = 0$ for some i and j . Let $\mathcal{S} := \{\lambda/(\lambda - 1) \mid \lambda \text{ is a real negative eigenvalue of } T^*\} \cup \{0\}$. Then the following holds:

- $0 \leq \max(\mathcal{S}) < r^* \leq 1$.
- For all $r \in (0, 1] \setminus \mathcal{S}$, $P(r) := (rI + (1 - r)T^*)^{-1}T^*$ is well-defined and

$$\begin{aligned} \mathbf{1}^\top P(r) &= \mathbf{1}^\top, \quad P(r)\pi^* = \pi^*, \quad P^* = P(r^*), \\ P(r)_{ij} &\geq 0 \quad \forall i, j \quad \iff \quad r \geq r^*. \end{aligned}$$

That is, $P(r)$ is a stochastic matrix if and only if $r \geq r^*$.

The proof is in Appendix C. This theorem indicates that we can determine r^* from T^* by doing bi-section on $(0, 1]$. But this approach fails when we replace T^* by an estimate \hat{T} because even $\hat{P}(r^*)$ might contain negative values. A more practical estimation procedure is the following: for each value of r in a decreasing sequence starting from 1, project $\hat{P}(r) := (rI + (1 - r)\hat{T})^{-1}\hat{T}$ onto the space of stochastic matrices and record the projection distance. Then search in the sequence of projection distances for the first sudden increase¹ starting from 1, and take the corresponding value of r and projected $\hat{P}(r)$ as our estimates.

Assuming the true r and α_0 are known, with the empirical moments being consistent estimators for the true moments and the tensor decomposition method guaranteed to return accurate estimates under small input perturbation, we can conclude that the estimates described above will converge (with high probability) to the true quantities as the sample size N increases. We give sample complexity bound on estimation error in the next section for hidden Markov models.

3.2 Hidden Markov Models

Let P and π now be defined over the hidden discrete state space of size k and have the same properties as the first-order Markov model. The generative process here is almost identical to (and therefore share the same interpretation with) the one in Section 3.1, except for an extra mapping from the discrete hidden state to a continuous observation space:

- Draw a state indicator vector $\mathbf{h}_i \sim \text{Multinomial}(P^{t_i} \mathbf{s}_i)$, $\mathbf{h}_i \in \{0, 1\}^k$.
- Draw an observation: $\mathbf{x}_i = U\mathbf{h}_i + \epsilon_i$, where $U \in \mathbb{R}^{m \times k}$ denotes a rank- k matrix of mean observation vectors for the k hidden states, and the random noise vectors ϵ_i 's are i.i.d satisfying $\mathbb{E}[\epsilon_i] = \mathbf{0}$ and $\text{Var}[\epsilon_i] = \sigma^2 I$.

Note that a spherical covariance² is required for the tensor decomposition method to be applicable. The low-order moments that lead to the desired tensor structure are given in the following:

Theorem 3. Define the expected hidden state transition matrix $T := \mathbb{E}_t[P^t] = rP(I - (1 - r)P)^{-1}$ and let $\alpha_0 := \sum_i \alpha_i$, $V_1 := \mathbb{E}[\mathbf{x}_1]$, $V_2 := \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top]$, $V_3 := \mathbb{E}[\mathbf{x}_1^{\otimes 3}]$, $C_2 := \mathbb{E}[\mathbf{x}_1 \mathbf{x}_2^\top]$ and $C_3 := \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3]$. Then the following holds:

$$\begin{aligned} V_1 &= U\pi, \quad V_2 = U\text{diag}(\pi)U^\top + \sigma^2 I, \quad V_3 = \sum_i \pi_i U_i^{\otimes 3} + \sum_{d=1}^3 V_1 \otimes_d (\sigma^2 I), \\ M_2 &:= V_2 - \sigma^2 I = U\text{diag}(\pi)U^\top, \quad M_3 := V_3 - \sum_{d=1}^3 V_1 \otimes_d (\sigma^2 I) = \sum_i \pi_i U_i^{\otimes 3}, \\ C_2 &= \frac{1}{\alpha_0 + 1} UT\text{diag}(\pi)(UT)^\top + \frac{\alpha_0}{\alpha_0 + 1} V_1 V_1^\top, \\ C_3 &= \frac{2}{(\alpha_0 + 2)(\alpha_0 + 1)} \sum_i \pi_i (UT)_i^{\otimes 3} + \frac{\alpha_0}{\alpha_0 + 2} \sum_{d=1}^3 V_1 \otimes_d C_2 - \frac{2\alpha_0^2}{(\alpha_0 + 2)(\alpha_0 + 1)} V_1^{\otimes 3} \\ M'_2 &:= (\alpha_0 + 1)C_2 - \alpha_0 V_1 V_1^\top = UT\text{diag}(\pi)(UT)^\top, \\ M'_3 &:= \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} C_3 - \frac{(\alpha_0 + 1)\alpha_0}{2} \sum_{d=1}^3 V_1 \otimes_d C_2 + \alpha_0^2 V_1^{\otimes 3} = \sum_i \pi_i (UT)_i^{\otimes 3}. \end{aligned}$$

¹Intuitively the jump should be easier to locate as P gets sparser, but we do not have a formal result.

²We may allow different covariances $\sigma_j^2 I$ for different hidden states. See Section 3.2 of [2] for details.

Algorithm 1 Tensor decomposition method for learning HMM from non-sequence data

input N sets of non-sequence data points, the success probability r , the Dirichlet parameter α_0 , the number of hidden states k , and numbers of iterations L and N .

output Estimates $\hat{\pi}$, \hat{P} and \hat{U} possibly under permutation of state labels.

- 1: Compute empirical averages $\widehat{V}_1, \widehat{V}_2, \widehat{V}_3, \widehat{C}_2, \widehat{C}_3$, and $\widehat{\sigma}^2 := \lambda_{\min}(\widehat{V}_2 - \widehat{V}_1 \widehat{V}_1^\top)$.
- 2: Compute $\widehat{M}_2, \widehat{M}_3, \widehat{M}'_2, \widehat{M}'_3$
- 3: Run Algorithm A.1 (Appendix A) on \widehat{M}_2 and \widehat{M}_3 with the number of hidden states k to obtain a symmetric tensor $\widehat{T} \in \mathbb{R}^{k \times k \times k}$ and a whitening transformation $\widehat{W} \in \mathbb{R}^{m \times k}$.
- 4: Run Algorithm A.2 (Appendix A) k times each with numbers of iterations L and N , the input tensor in the first run set to \widehat{T} and in each subsequent run set to the deflated tensor returned by the previous run, resulting in k pairs of eigenvalue/eigenvector $\{(\hat{\lambda}_i, \hat{\mathbf{v}}_i)\}_{i=1}^k$.
- 5: Repeat Steps 4 and 5 on \widehat{M}'_2 and \widehat{M}'_3 to obtain $\widehat{T}', \widehat{W}'$ and $\{(\hat{\lambda}'_i, \hat{\mathbf{v}}'_i)\}_{i=1}^k$.
- 6: Match $\{(\hat{\lambda}_i, \hat{\mathbf{v}}_i)\}_{i=1}^k$ with $\{(\hat{\lambda}'_i, \hat{\mathbf{v}}'_i)\}_{i=1}^k$ by sorting $\{\hat{\lambda}_i\}_{i=1}^k$ and $\{\hat{\lambda}'_i\}_{i=1}^k$.
- 7: Obtain estimates of HMM parameters:

$$\begin{aligned}\widehat{UT} &:= (\widehat{W}')^\dagger \widehat{V}' \widehat{\Lambda}', & \widehat{U} &:= (\widehat{W}^\top)^\dagger \widehat{V} \widehat{\Lambda}, \\ \widehat{P} &:= (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}, & \widehat{\pi} &:= [\hat{\lambda}_1^{-2} \cdots \hat{\lambda}_k^{-2}]^\top,\end{aligned}$$

where $\widehat{V} := [\widehat{\mathbf{v}}_1 \cdots \widehat{\mathbf{v}}_k]$, $\widehat{\Lambda} := \text{diag}([\hat{\lambda}_1 \cdots \hat{\lambda}_k]^\top)$; \widehat{V}' and $\widehat{\Lambda}'$ are defined in the same way.

- 8: (Optional) Project $\widehat{\pi}$ onto the simplex and \widehat{P} onto the space of stochastic matrices.
-

The proof is in Appendix B.2. This theorem suggests that, unlike first-order Markov models, HMMs require *two* applications of the tensor decomposition methods, one on M_2 and M_3 for extracting the mean observation vectors U , and the other on M'_2 and M'_3 for extracting the matrix product UT . Another issue is that the estimates for M_2 and M_3 require an estimate for the noise variance σ^2 , which is not directly observable. Nevertheless, since M_2 and M_3 are in the form of low-order moments of spherical Gaussian mixtures, we may use the existing result (Theorem 3.2, [2]) to obtain an estimate $\widehat{\sigma}^2 = \lambda_{\min}(\widehat{V}_2 - \widehat{V}_1 \widehat{V}_1^\top)$. The situation regarding permutations of the estimates is also different here. First note that $P = (rU + (1-r)UT)^\dagger UT$, which implies that permuting the columns of U and the columns of UT in the same manner has the effect of permuting both the rows and the columns of P , essentially re-labeling the hidden states. Hence we can only expect to recover P up to some simultaneous row and column permutation. By the assumption that π_i 's are all different, we can sort the two estimates $\widehat{\pi}$ and $\widehat{\pi}'$ to match the columns of \widehat{U} and \widehat{UT} , and obtain \widehat{P} if r is known. When r is unknown, a similar heuristics to the one for first-order Markov models can be used to estimate r , based on the fact that $P = (rU + (1-r)UT)^\dagger UT = (rI + (1-r)T)^{-1}T$, suggesting that Theorem 2 remains true when expressing P by U and UT .

Algorithm 1 gives the complete procedure for learning HMM from non-sequence data. Combining the perturbation bounds of the tensor decomposition method (Appendix A), the whitening procedure (Appendix D.1) and the matrix pseudoinverse [10], and concentration bounds on empirical moments (Appendix D.3), we provide a sample complexity analysis:

Theorem 4. *Suppose the numbers of iterations N and L for Algorithm A.2 satisfy the conditions in Theorem A.1 (Appendix A), and the number of hidden states k , the success probability r , and the Dirichlet parameter α_0 are all given. For any $\eta \in (0, 1)$ and $\epsilon > 0$, if the number of sets N satisfies*

$$N \geq \frac{12 \max(k^2, m) m^3 \nu^3 (\alpha_0 + 2)^2 (\alpha_0 + 1)^2}{\eta} \max \left(\frac{225000}{\delta_{\min}^2}, \frac{4600}{\min(\sigma_k(M'_2), \sigma_k(M_2))^2}, \frac{42000 c^2 \sigma_1(UT)^2 \max(\sigma_1(UT), \sigma_1(U), 1)^2}{\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \min(\sigma_k(UT), \sigma_k(U), 1)^4} \right),$$

where c is some constant, $\nu := \max(\sigma^2 + \max_{i,j}(|U_{ik}|^2), 1)$, $\delta_{\min} := \min_{i,j} |1/\sqrt{\pi_j} - 1/\sqrt{\pi_j}|$, and $\sigma_i(\cdot)$ denotes the i -th largest singular value, then the \widehat{P} and \widehat{U} returned by Algorithm 1 satisfy

$$\text{Prob}(\|P - \widehat{P}\| \leq \epsilon) \geq 1 - \eta \quad \text{and} \quad \text{Prob} \left(\|U - \widehat{U}\| \leq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \geq 1 - \eta,$$

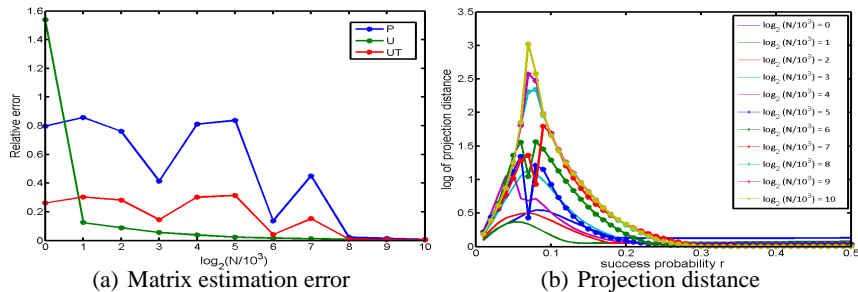


Figure 2: Simulation results

where P and U may undergo label permutation.

The proof is in Appendix E. In this result, the sample size N exhibits a fairly high-order polynomial dependency on m, k, ϵ^{-1} and scales with $1/\eta$ linearly instead of logarithmically, as is common in sample complexity results on spectral learning. This is because we do not impose any constraints on the observation model and simply use the Markov inequality for bounding the deviation in the empirical moments. If we make stronger assumptions such as boundedness or sub-Gaussianity, it is possible to use stronger, exponential tail bounds to obtain better sample complexity. Also worth noting is that δ_{\min}^{-2} acts as a threshold. As shown in our proof, as long as the operator norm of the tensor perturbation is sufficiently smaller than δ_{\min} , which measures the gaps between different π_i 's, we can correctly match the two sets of estimated tensor eigenvalues. Lastly, the lower bound of N , as one would expect, depends on conditions of the matrices being estimated as reflected in the various ratios of singular values. An interesting quantity missing from the sample analysis is the size of each set n . To simplify the analysis we essentially assume $n = 3$, but understanding how n might affect the sample complexity may have a critical impact in practice: when collecting more data, should we collect more sets or larger sets? What is the trade-off between them? This is an interesting direction for future work.

4 Simulation

Our HMM has $m = 40$ and $k = 5$ with Gaussian noise $\sigma^2 = 2$. The mean vectors U were sampled from independent univariate standard normal and then normalized to lie on the unit sphere. The transition matrix P contains one zero entry. For the generative process, we set $\alpha_0 = 1, r = 0.3, n = 1000$, and $N \in 1000\{2^0, 2^1, \dots, 2^{10}\}$. The numbers of iterations for Algorithm A.2 were $N = 200$ and $L = 1000$. Figure 2(a) plots the relative matrix estimation error (in spectral norm) against the sample size N for P, U , and UT obtained by Algorithm 1 given the true r . It is clear that U is the easiest to learn, followed by UT , and P is the most difficult, and that all three errors converge to a very small value for sufficiently large N . Note that in Theorem 4 the bounds for P and U are different. With the model used here, the extra multiplicative factor in the bound for U is less than 0.007, suggesting that U is indeed easier to estimate than P . Figure 2(b) demonstrates the heuristics for determining r , showing projection distances (in logarithm) versus r . As N increases, the take-off point gets closer to the true $r = 0.3$. The large peak indicates a pole (the set S in Theorem 2).

5 Conclusions

We have demonstrated that under reasonable assumptions, tensor decomposition methods can provably learn first-order Markov models and hidden Markov models from non-sequence data. We believe this is the first formal guarantee on learning dynamic models in a non-sequential setting. There are several ways to extend our results. No matter what distribution generates the random time steps, tensor decomposition methods can always learn the expected transition probability matrix T . Depending on the application, it might be better to use some other distribution for the missing time. The proposed algorithm can be modified to learn discrete HMMs under a similar generative process. Finally, applying the proposed methods to real data should be the most interesting future direction.

References

- [1] A. Anandkumar, D. P. Foster, D. Hsu, S. M. Kakade, and Y.-K. Liu. A spectral algorithm for latent dirichlet allocation. *arXiv preprint arXiv:1204.6703v4*, 2013.
- [2] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *arXiv preprint arXiv:1210.7559v2*, 2012.
- [3] A. Anandkumar, D. Hsu, and S. M. Kakade. A method of moments for mixture models and hidden Markov models. *arXiv preprint arXiv:1203.0683*, 2012.
- [4] S. Arora, R. Ge, Y. Halpern, D. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu. A practical algorithm for topic modeling with provable guarantees. *arXiv preprint arXiv:1212.4777*, 2012.
- [5] D. Hsu and S. M. Kakade. Learning mixtures of spherical gaussians: moment methods and spectral decompositions. In *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*, pages 11–20. ACM, 2013.
- [6] T.-K. Huang and J. Schneider. Learning linear dynamical systems without sequence information. In *Proceedings of the 26th International Conference on Machine Learning*, pages 425–432, 2009.
- [7] T.-K. Huang and J. Schneider. Learning auto-regressive models from sequence and non-sequence data. In *Advances in Neural Information Processing Systems 24*, pages 1548–1556. 2011.
- [8] T.-K. Huang, L. Song, and J. Schneider. Learning nonlinear dynamic models from non-sequenced data. In *Proceedings of the 13th International Conference on Artificial Intelligence and Statistics*, 2010.
- [9] M. G. Rabbat, M. A. Figueiredo, and R. D. Nowak. Network inference from co-occurrences. *Information Theory, IEEE Transactions on*, 54(9):4053–4068, 2008.
- [10] G. Stewart. On the perturbation of pseudo-inverses, projections and linear least squares problems. *SIAM review*, 19(4):634–662, 1977.
- [11] X. Zhu, A. B. Goldberg, M. Rabbat, and R. Nowak. Learning bigrams from unigrams. In *the Proceedings of 46th Annual Meeting of the Association for Computational Linguistics: Human Language Technology, Columbus, OH*, 2008.

Supplemental Material for “Learning Hidden Markov Models from Non-sequence Data via Tensor Decomposition”

Tzu-Kuo Huang
 Machine Learning Department
 Carnegie Mellon University
 Pittsburgh, PA 15213
 tzukuoh@cs.cmu.edu

Jeff Schneider
 Robotics Institute
 Carnegie Mellon University
 Pittsburgh, PA 15213
 schneide@cs.cmu.edu

A Tensor Decomposition: Algorithm and Theoretical Guarantee

We start with the population version of the problem: suppose the noiseless M_2 and M_3 defined in (2) of the main text are available, and we want to recover $\omega_{i=1}^k$ and $\{\boldsymbol{\mu}_i\}_{i=1}^k$. First we perform a *whitening step* on them, as outlined in Algorithm A.1, to obtain a whitened, lower-dimensional tensor $\mathcal{T} \in \mathbb{R}^{k \times k \times k}$ and a whitening transformation $W \in \mathbb{R}^{m \times k}$ such that

$$\mathcal{T} := M_3(W, W, W) = \sum_{i=1}^k \omega_i (W^\top \boldsymbol{\mu}_i)^{\otimes 3} = \sum_{i=1}^k \frac{1}{\sqrt{\omega_i}} \tilde{\boldsymbol{\mu}}_i^{\otimes 3},$$

where the vectors $\tilde{\boldsymbol{\mu}}_i := \sqrt{\omega_i} W^\top \boldsymbol{\mu}_i$ form an orthonormal basis of \mathbb{R}^k because $I = W^\top M_2 W = \sum_{i=1}^k W^\top (\sqrt{\omega_i} \boldsymbol{\mu}_i) (\sqrt{\omega_i} \boldsymbol{\mu}_i)^\top W = \sum_{i=1}^k \tilde{\boldsymbol{\mu}}_i \tilde{\boldsymbol{\mu}}_i^\top$. Hence, the symmetric tensor \mathcal{T} has a so-called *orthogonal decomposition*, which may not exist for general symmetric tensors. Then by Theorem 4.3 of [2], which establishes the following results under Condition 1:

1. the set of *robust eigenvectors*¹ of \mathcal{T} correspond exactly to $\{\tilde{\boldsymbol{\mu}}_i\}_{i=1}^k$;
2. the eigenvalue associated with $\tilde{\boldsymbol{\mu}}_i$ is equal to $1/\sqrt{\omega_i}$, $\forall 1 \leq i \leq k$;
3. if (\mathbf{v}, λ) is a pair of robust eigenvector/eigenvalue of \mathcal{T} , then $\boldsymbol{\mu}_i = \lambda(W^\top)^\dagger \mathbf{v}$ for some $1 \leq i \leq k$, where \dagger denotes the Moore-Penrose pseudo inverse;

we can reduce the original problem of recovering the structure in (2) into a robust tensor eigendecomposition problem. Motivated by the power iteration for matrix eigen computation, Anandkumar et al. [2] verify that starting from almost every vector $\boldsymbol{\theta}_0 \in \mathbb{R}^k$, the tensor power iteration $\boldsymbol{\theta}_t := \frac{\mathcal{T}(I, \boldsymbol{\theta}_{t-1}, \boldsymbol{\theta}_{t-1})}{\|\mathcal{T}(I, \boldsymbol{\theta}_{t-1}, \boldsymbol{\theta}_{t-1})\|}$, where $\|\cdot\|$ denotes the vector 2-norm, converges to some robust eigenvector of \mathcal{T} at a quadratic rate, and therefore k successive applications of the tensor power iteration with deflation result in all pairs of robust eigenvectors/eigenvalues.

In practice we almost never have the exact M_2 and M_3 , but only noisy or perturbed versions \widehat{M}_2 and \widehat{M}_3 , which are usually estimates from the data. Perturbation may destroy the nice tensor structure, so the reduced tensor $\widehat{\mathcal{T}}$ resulting from applying Algorithm A.1 to \widehat{M}_2 and \widehat{M}_3 may no longer be orthogonally decomposable, hindering the subsequent robust tensor eigendecomposition. Nevertheless, Anandkumar et al. [2] demonstrate that if the perturbation $E := \widehat{\mathcal{T}} - \mathcal{T}$ is a symmetric tensor with a small operator norm defined as $\|E\| := \sup_{\|\boldsymbol{\theta}\|=1} |E(\boldsymbol{\theta}, \boldsymbol{\theta}, \boldsymbol{\theta})|$, then k successive applications of some *randomized* tensor power iteration coupled with deflation yield accurate estimates of all robust eigenvector/eigenvalue pairs with high probability. More precisely, they propose *the*

¹See Section 4.2 of [2] for the definition of robust eigenvectors/eigenvalues.

Algorithm A.1 Whitening transformation

input A symmetric matrix $M_2 \in \mathbb{R}^{m \times m}$, a symmetric third-order tensor $M_3 \in \mathbb{R}^{m \times m \times m}$, and the target number of dimensions k .

output A reduced third-order tensor $\mathcal{T} \in \mathbb{R}^{k \times k \times k}$ and a whitening transformation $W \in \mathbb{R}^{m \times k}$.

- 1: Compute $W := QD^{-1/2}$, where $Q \in \mathbb{R}^{m \times k}$ denotes the top- k orthonormal eigenvectors of M_2 , and $D \in \mathbb{R}^{k \times k}$ is a diagonal matrix of the corresponding k positive eigenvalues.
 - 2: Compute $\mathcal{T} := M_3(W, W, W)$.
-

Algorithm A.2 Robust tensor power method

input A symmetric tensor $\mathcal{T} \in \mathbb{R}^{k \times k \times k}$, number of iterations L, N .

output the estimated eigenvector/eigenvalue pair; the deflated tensor.

- 1: **for** $\tau = 1$ **to** L **do**
 - 2: Draw $\theta_0^{(\tau)}$ uniformly at random from the unit sphere in \mathbb{R}^k .
 - 3: **for** $t = 1$ **to** N **do**
 - 4:
$$\theta_t^{(\tau)} := \frac{\mathcal{T}(I, \theta_{t-1}^{(\tau)}, \theta_{t-1}^{(\tau)})}{\|\mathcal{T}(I, \theta_{t-1}^{(\tau)}, \theta_{t-1}^{(\tau)})\|}.$$
 - 5: **end for**
 - 6: **end for**
 - 7: Let $\tau^* := \arg \max_{1 \leq \tau \leq L} \mathcal{T}(\theta_N^{(\tau)}, \theta_N^{(\tau)}, \theta_N^{(\tau)})$.
 - 8: Do N power iteration updates (Line 4) starting from $\theta_N^{(\tau^*)}$ to obtain $\hat{\theta}$, and set $\hat{\lambda} := \mathcal{T}(\hat{\theta}, \hat{\theta}, \hat{\theta})$
 - 9: **return** the estimated eigenvector/eigenvalue pair $(\hat{\theta}, \hat{\lambda})$; the deflated tensor $\mathcal{T} - \hat{\lambda} \hat{\theta}^{\otimes 3}$.
-

Robust tensor power method outlined in Algorithm A.2, which employs multiple random restarts, and provide a theoretical guarantee on its robustness against the input perturbation:

Theorem A.1. (Theorem 5.1 of [2]) Let $\hat{\mathcal{T}} = \mathcal{T} + E \in \mathbb{R}^{k \times k \times k}$, where \mathcal{T} is a symmetric tensor with orthogonal decomposition $\mathcal{T} = \sum_{i=1}^k \lambda_i \mathbf{v}_i^{\otimes 3}$ where each $\lambda_i > 0$, $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k\}$ is an orthonormal basis, and E has operator norm $\epsilon := \|E\|$. Define $\lambda_{\min} := \min(\{\lambda_i\}_{i=1}^k)$ and $\lambda_{\max} := \max(\{\lambda_i\}_{i=1}^k)$. There exists universal constants $c_1, c_2, c_3 > 0$ such that the following holds. Pick any $\eta \in (0, 1)$, and suppose

$$\epsilon \leq c_1 \cdot \frac{\lambda_{\min}}{k}, \quad N \geq c_2 \cdot (\log(k) + \log \log(\lambda_{\max}/\epsilon)), \quad \text{and}$$

$$\sqrt{\frac{\ln(L/\log_2(\frac{k}{\eta}))}{\ln(k)}} \cdot \left(1 - \frac{\ln(\ln(L/\log_2(\frac{k}{\eta}))) + c_3}{4 \ln(L/\log_2(\frac{k}{\eta}))} - \sqrt{\frac{\ln(8)}{\ln(L/\log_2(\frac{k}{\eta}))}}\right) \geq 1.02 \left(1 + \sqrt{\frac{\ln(4)}{\ln(k)}}\right).$$

(Note that the condition on L holds with $L = \text{poly}(k) \log(1/\eta)$.) Suppose that Algorithm A.2 is iteratively called k times with numbers of iterations L and N , where the input tensor is $\hat{\mathcal{T}}$ in the first call, and in each subsequent call, the input tensor is the deflated tensor returned by the previous call. Let $(\hat{\mathbf{v}}_1, \hat{\lambda}_1), (\hat{\mathbf{v}}_2, \hat{\lambda}_2), \dots, (\hat{\mathbf{v}}_k, \hat{\lambda}_k)$ be the sequence of estimated eigenvector/eigenvalue pairs returned in these k calls. With probability at least $1 - \eta$, there exists a permutation ρ on $\{1, \dots, k\}$ such that

$$\|\mathbf{v}_{\rho(j)} - \hat{\mathbf{v}}_j\| \leq 8\epsilon/\lambda_{\rho(j)}, \quad |\lambda_{\rho(j)} - \hat{\lambda}_j| \leq 5\epsilon, \quad \forall 1 \leq j \leq k, \quad \text{and} \quad \|\mathcal{T} - \sum_{j=1}^k \hat{\lambda}_j \hat{\mathbf{v}}_j^{\otimes 3}\| \leq 55\epsilon.$$

This result, together with existing perturbation theory regarding the whitening procedure (e.g., Appendix C.1 of [1]), allow us to translate the perturbations in \hat{M}_2 and \hat{M}_3 into the estimation errors in ω_i 's and μ_i 's, guaranteeing accurate estimation under small input perturbation.

B Tensor Structure in Low-order Moments

Here we give proofs of theorems regarding tensor structures in low-order moments.

B.1 Proof of Theorem 1

$$\begin{aligned}\mathbb{E}[\mathbf{x}_1] &= \mathbb{E}_{\pi_0} \mathbb{E}[\mathbf{x}_1 | \pi_0] = \mathbb{E}_{\pi_0} \mathbb{E}[P^{t_1} \mathbf{s}_1 | \pi_0] = \mathbb{E}_{\pi_0} [\mathbb{E}[P^{t_1}] \pi_0] = T \pi = \pi, \\ C_2 &:= \mathbb{E}[\mathbf{x}_1 \mathbf{x}_2^\top] = \mathbb{E}_{\pi_0} \mathbb{E}[P^{t_1} \mathbf{s}_1 \mathbf{s}_2^\top (P^{t_2})^\top | \pi_0] = \mathbb{E}_{\pi_0} [\mathbb{E}[P^{t_1}] \mathbb{E}[\mathbf{s}_1 \mathbf{s}_2^\top | \pi_0] \mathbb{E}[(P^{t_2})^\top]] \\ &= T \mathbb{E}_{\pi_0} [\pi_0 \pi_0^\top] T^\top = T \left(\frac{\text{diag}(\pi)}{\alpha_0 + 1} + \frac{\alpha_0 \pi \pi^\top}{\alpha_0 + 1} \right) T^\top = \frac{T \text{diag}(\pi) T^\top}{\alpha_0 + 1} + \frac{\alpha_0 \pi \pi^\top}{\alpha_0 + 1},\end{aligned}\quad (1)$$

$$\begin{aligned}C_3 &:= \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3] = \mathbb{E}_{\pi_0} \mathbb{E}[(P^{t_1} \mathbf{s}_1) \otimes (P^{t_2} \mathbf{s}_2) \otimes (P^{t_3} \mathbf{s}_3) | \pi_0] \\ &= \mathbb{E}_{\pi_0} [(T \pi_0) \otimes (T \pi_0) \otimes (T \pi_0)] = \frac{\sum_i 2\pi_i T_i \otimes T_i \otimes T_i}{(\alpha_0 + 2)(\alpha_0 + 1)} + \frac{\alpha_0^2 \pi \otimes \pi \otimes \pi}{(\alpha_0 + 2)(\alpha_0 + 1)} \\ &\quad + \frac{\alpha_0 \left(\sum_{ij} (T_i \otimes T_i \otimes T_j + T_i \otimes T_j \otimes T_i + T_j \otimes T_i \otimes T_i) \pi_i \pi_j \right)}{(\alpha_0 + 2)(\alpha_0 + 1)} \\ &= \frac{\sum_i 2\pi_i T_i \otimes T_i \otimes T_i}{(\alpha_0 + 2)(\alpha_0 + 1)} + \frac{\alpha_0 (\pi \otimes_3 C_2 + \pi \otimes_2 C_2 + \pi \otimes_1 C_2)}{\alpha_0 + 2} - \frac{2\alpha_0^2 \pi \otimes \pi \otimes \pi}{(\alpha_0 + 2)(\alpha_0 + 1)},\end{aligned}\quad (2)$$

The second equality in (1) and the two equalities (2) and (3) can be established by using the results on Dirichlet moments in Appendix B.1 of [1].

B.2 Proof of Theorem 3

$$\begin{aligned}V_1 &:= \mathbb{E}[\mathbf{x}_1] = \mathbb{E}[U \mathbf{h}_1 + \epsilon_1] = U \mathbb{E}[P^{t_1} \mathbf{s}_1] = U T \mathbb{E}[\pi_0] = U \pi. \\ V_2 &:= \mathbb{E}[\mathbf{x}_1 \mathbf{x}_1^\top] = \mathbb{E}[(U \mathbf{h}_1 + \epsilon_1)(U \mathbf{h}_1 + \epsilon_1)^\top] \\ &= \mathbb{E}[U \mathbf{h}_1 \mathbf{h}_1^\top U^\top] + \sigma^2 I = U \mathbb{E}[\text{diag}(\mathbf{h}_1)] U^\top + \sigma^2 I \\ &= U \mathbb{E}[\text{diag}(P^{t_1} \mathbf{s}_1)] U^\top + \sigma^2 I = U \mathbb{E}[\text{diag}(T \pi_0)] U^\top + \sigma^2 I \\ &= U \text{diag}(\pi) U^\top + \sigma^2 I. \\ V_3 &:= \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_1 \otimes \mathbf{x}_1] = \mathbb{E}[(U \mathbf{h}_1 + \epsilon_1) \otimes (U \mathbf{h}_1 + \epsilon_1) \otimes (U \mathbf{h}_1 + \epsilon_1)] \\ &= \mathbb{E}[(U \mathbf{h}_1) \otimes (U \mathbf{h}_1) \otimes (U \mathbf{h}_1)] + \mathbb{E}[(U \mathbf{h}_1) \otimes \epsilon_1 \otimes \epsilon_1] + \mathbb{E}[\epsilon_1 \otimes (U \mathbf{h}_1) \otimes \epsilon_1] + \mathbb{E}[\epsilon_1 \otimes \epsilon_1 \otimes (U \mathbf{h}_1)] \\ &= \sum_i \pi_i U_i \otimes U_i \otimes U_i + V_1 \otimes_1 (\sigma^2 I) + V_1 \otimes_2 (\sigma^2 I) + V_1 \otimes_3 (\sigma^2 I). \\ C_2 &:= \mathbb{E}[\mathbf{x}_1 \mathbf{x}_2^\top] = \mathbb{E}[(U \mathbf{h}_1 + \epsilon_1)(U \mathbf{h}_2 + \epsilon_2)^\top] = \mathbb{E}[U \mathbf{h}_1 \mathbf{h}_2^\top U^\top] \\ &= U \mathbb{E}[P^{t_1} \mathbf{s}_1 \mathbf{s}_2^\top (P^{t_2})^\top] U^\top = U T \mathbb{E}[\pi_0 \pi_0^\top] T^\top U^\top = \frac{U T \text{diag}(\pi) (U T)^\top}{\alpha_0 + 1} + \frac{\alpha_0 V_1 V_1^\top}{\alpha_0 + 1}.\end{aligned}\quad (4)$$

$$\begin{aligned}C_3 &:= \mathbb{E}[\mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3] = \mathbb{E}[(U \mathbf{h}_1 + \epsilon_1) \otimes (U \mathbf{h}_2 + \epsilon_2) \otimes (U \mathbf{h}_3 + \epsilon_3)] = \mathbb{E}[(U \mathbf{h}_1) \otimes (U \mathbf{h}_2) \otimes (U \mathbf{h}_3)] \\ &= \mathbb{E}[(U P^{t_1} \mathbf{s}_1) \otimes (U P^{t_2} \mathbf{s}_2) \otimes (U P^{t_3} \mathbf{s}_3)] = \mathbb{E}[(U T \pi_0) \otimes (U T \pi_0) \otimes (U T \pi_0)] \\ &= \frac{\sum_i 2\pi_i (U T)_i \otimes (U T)_i \otimes (U T)_i}{(\alpha_0 + 2)(\alpha_0 + 1)} + \frac{\alpha_0 (V_1 \otimes_3 C_2 + V_1 \otimes_2 C_2 + V_1 \otimes_1 C_2)}{\alpha_0 + 2} - \frac{2\alpha_0^2 V_1 \otimes V_1 \otimes V_1}{(\alpha_0 + 2)(\alpha_0 + 1)},\end{aligned}\quad (5)$$

Again, the last equality in (4) and (5) can be established by using the results on Dirichlet moments in Appendix B.1 of [1].

C Proof of Theorem 2

We first prove the following lemma:

Lemma 1. *If $P(r) := (rI + (1-r)T^*)^{-1}T^*$ exists and is a stochastic matrix for some $r \in (0, 1]$, then $P(r')$ exists and is a stochastic matrix for all $r' \in [r, 1]$.*

Proof. Since $P(r)$ exists we can write $T^* = rP(r)(I - (1-r)P(r))^{-1}$. By assumption P^* is invertible, so T^* is invertible. We then have

$$r'(T^*)^{-1} + (1-r')I = \frac{r'}{r}(P(r)^{-1} - (1-r)I) + (1-r')I = \frac{r'}{r}P(r)^{-1}(I - (1-r/r')P(r)),$$

which is invertible for all $r' \in [r, 1]$. Therefore, we can write

$$P(r') = (r'(T^*)^{-1} + (1-r')I)^{-1} = \frac{r}{r'}P(r)(I - (1-r/r')P(r))^{-1} = \mathbb{E}_t[P(r)],$$

where $t \sim \text{Geometric}(r/r')$, showing that $P(r')$ is a stochastic matrix. \square

To prove Theorem 2 we begin by noting that \mathcal{S} contains all values of r for which $rI + (1-r)T^*$ is singular. Therefore, $P(r)$ is well-defined and invertible for $r \in (0, 1] \setminus \mathcal{S}$. From the identity $T^*\pi^* = \pi^* = (rI + (1-r)T^*)\pi^*$ we have $P(r)\pi^* = \pi^*$, $r \notin \mathcal{S}$. Similarly, the identity $\mathbf{1}^\top T^* = \mathbf{1}^\top = \mathbf{1}^\top(rI + (1-r)T^*)$ and the fact that $(rI + (1-r)T^*)^{-1}T^* = T^*(rI + (1-r)T^*)^{-1}$ imply that $\mathbf{1}^\top P(r) = \mathbf{1}^\top$, $r \notin \mathcal{S}$. It is easy to verify $P(r^*) = P^*$ by plugging in the definition of T^* . Lemma 1 then implies that $\max(\mathcal{S}) < r^*$ and that $P(r')$ is a stochastic matrix for $r' \geq r^*$. To prove the last statement of the theorem we rewrite $P(r)$ by plugging in the definition of T^* :

$$P(r) = \frac{r^*}{r}(I - (1-r^*/r)P^*)^{-1}P^* \quad (6)$$

and consider its first-order derivative w.r.t. r :

$$\frac{\partial P(r)}{\partial r} = -\left(\frac{r}{r^*}I + \left(1 - \frac{r}{r^*}\right)P^*\right)^{-2} \frac{(I - P^*)P^*}{r^*}, \quad (7)$$

which exists for $r \in (0, 1] \setminus \mathcal{S}$. By assumption we have $P_{ij}^* = 0$, and by ergodicity of P^* we can assume $(P^*)_{ij}^2 > 0$ (otherwise there exists $k \neq j$ such that $P_{ik}^* = 0$ and $(P^*)_{ik}^2 > 0$). Then we have

$$\left.\frac{\partial P(r)_{ij}}{\partial r}\right|_{r=r^*} = \frac{(P^*)_{ij}^2}{r^*} > 0, \quad (8)$$

implying that there exists $c > 0$ such that for $r \in [r^* - c, r^*)$, $P(r)_{ij} < P_{ij}^* = 0$. This and Lemma 1 then imply the last statement of the theorem.

D Sample Complexity Analysis

The analyses here mostly follow those in [1]. Let O denote the observation matrix, which can be the T matrix in First-order Markov models, the U matrix or the product UT in Hidden Markov Models. Define

$$\begin{aligned} \tilde{O} &:= O \text{diag}([\sqrt{\pi_1} \sqrt{\pi_2} \cdots \sqrt{\pi_k}]), \\ M_2 &:= O \text{diag}(\pi) O^\top = \tilde{O} \tilde{O}^\top \quad \text{and} \quad M_3 := \sum_{i=1}^k \pi_i O_i \otimes O_i \otimes O_i. \end{aligned}$$

Let $\pi_{\min} := \min_i \pi_i$. We have

$$\sigma_k(O) \sqrt{\pi_{\min}} \leq \sigma_k(\tilde{O}), \quad (9)$$

$$\sigma_1(\tilde{O}) \leq \sigma_1(O), \quad (10)$$

where $\sigma_j(\cdot)$ denotes the j th largest singular value.

Denote by $\|\cdot\|$ the spectral norm of a matrix or the operator norm of a symmetric third-order tensor induced by the vector 2-norm:

$$\|M\| := \sup_{\|\boldsymbol{\theta}\|_2=1} |M(\boldsymbol{\theta}, \boldsymbol{\theta}, \boldsymbol{\theta})|. \quad (11)$$

Suppose

$$\|\widehat{M}_2 - M_2\| = E_2, \quad (12)$$

$$\|\widehat{M}_3 - M_3\| \leq E_3, \quad (13)$$

for some E_2 and E_3 to be determined.

D.1 Perturbation Lemmas

Let $\widehat{M}_{2,k}$ be the best rank k approximation to \widehat{M}_2 in terms of the matrix 2-norm. According to Algorithm A.1, we have

$$\widehat{W}^\top \widehat{M}_{2,k} \widehat{W} = I. \quad (14)$$

Let

$$\widehat{W}^\top M_2 \widehat{W} = ADA^\top \quad (15)$$

be an SVD of $\widehat{W}^\top M_2 \widehat{W}$, where $A \in \mathcal{R}^{k \times k}$. Define

$$W := \widehat{W} A D^{-1/2} A^\top \quad (16)$$

and notice that

$$W^\top M_2 W = A D^{-1/2} A^\top \widehat{W}^\top M_2 \widehat{W} A D^{-1/2} A^\top = I. \quad (17)$$

Let $Q := W^\top \widetilde{O}$ and $\widehat{Q} := \widehat{W}^\top \widetilde{O}$.

Lemma 2. (Lemma C.1 of [1]) Let Π_W be the orthogonal projection onto the range of W and Π be the orthogonal projection onto the range of O . Suppose $E_2 \leq \sigma_k(M_2)/2$. We have the following:

$$\begin{aligned} \|Q\| &= 1, \\ \|\widehat{Q}\| &\leq 2, \\ \|\widehat{W}\| &\leq \frac{2}{\sigma_k(\widetilde{O})}, \\ \|\widehat{W}^\dagger\| &\leq 2\sigma_1(\widetilde{O}), \\ \|W^\dagger\| &\leq 3\sigma_1(\widetilde{O}), \\ \|Q - \widehat{Q}\| &\leq \frac{4E_2}{\sigma_k(\widetilde{O})^2}, \\ \|\widehat{W}^\dagger - W^\dagger\| &\leq \frac{6\sigma_1(\widetilde{O})E_2}{\sigma_k(\widetilde{O})^2}, \\ \|\Pi - \Pi_W\| &\leq \frac{4E_2}{\sigma_k(\widetilde{O})^2}. \end{aligned}$$

Lemma 3. Weyl's Theorem. (Theorem 4.11, p.204 in [4]). Let $A, E \in \mathcal{R}^{m \times n}$ with $m \geq n$ be given. Then

$$\max_{1 \leq i \leq n} |\sigma_i(A + E) - \sigma_i(A)| \leq \|E\|.$$

D.2 Reconstruction Accuracy

Throughout this section we assume that the number of iterations N and L for Algorithm A.2 satisfy the conditions in Theorem A.1.

Lemma 4. Suppose $\max(E_2, E_3) \leq \sigma_k(M_2)/2$. For any $\eta \in (0, 1)$, with probability at least $1 - \eta$ the following holds:

$$\|O - (\widehat{W}^\top)^\dagger \widehat{V} \widehat{\Lambda}\| \leq c \frac{\max(\sigma_1(O), 1)}{\pi_{\min}^{3/2} \min(\sigma_k(O)^2, 1)} \max(E_2, E_3)$$

for some constant $c > 0$.

Proof. By Theorem A.1, the following hold with probability at least $1 - \eta$:

$$\|V - \widehat{V}\|_F = \sqrt{\sum_i \|V_i - \widehat{V}_i\|^2} \leq \sqrt{\sum_i (64E_3^2)/(1/\sqrt{\pi_{\min}})^2} = 8E_3, \quad (18)$$

$$\|\widehat{\Lambda}\| = \max_i 1/\sqrt{\pi_i} \leq \max_i (1/\sqrt{\pi_i} + 5E_3) \leq \pi_{\min}^{-1/2} + 5E_3. \quad (19)$$

With the above two bounds and Lemma 2 we have

$$\begin{aligned}
& \|O - (\widehat{W}^\top)^\dagger \widehat{V} \widehat{\Lambda}\| \leq \|O - \Pi_W O\| + \|\Pi_W O - (\widehat{W}^\top)^\dagger \widehat{V} \widehat{\Lambda}\| \\
& = \|\Pi O - \Pi_W O\| + \|(W^\dagger)^\top V \Lambda - (\widehat{W}^\dagger)^\top \widehat{V} \widehat{\Lambda}\| \\
& \leq \|\Pi - \Pi_W\| \|O\| + \|(W^\dagger)^\top V \Lambda - (W^\dagger)^\top V \widehat{\Lambda}\| + \|(W^\dagger)^\top V \widehat{\Lambda} - (\widehat{W}^\dagger)^\top \widehat{V} \widehat{\Lambda}\| \\
& \leq \|\Pi - \Pi_W\| + \|W^\dagger\| \|V\| \|\Lambda - \widehat{\Lambda}\| + \|(W^\dagger)^\top V \widehat{\Lambda} - (W^\dagger)^\top \widehat{V} \widehat{\Lambda}\| + \|(W^\dagger)^\top \widehat{V} \widehat{\Lambda} - (\widehat{W}^\dagger)^\top \widehat{V} \widehat{\Lambda}\| \\
& \leq \|\Pi - \Pi_W\| + \|W^\dagger\| E_3 + \|W^\dagger\| \|V - \widehat{V}\| \|\widehat{\Lambda}\| + \|W^\dagger - \widehat{W}^\dagger\| \|\widehat{V}\| \|\widehat{\Lambda}\| \\
& \leq \frac{4E_2}{\sigma_k(\widetilde{O})^2} + 3\sigma_1(\widetilde{O})E_3 + 3\sigma_1(\widetilde{O})\|V - \widehat{V}\|_F \|\widehat{\Lambda}\| + \frac{6\sigma_1(\widetilde{O})E_2}{\sigma_k(\widetilde{O})^2} (\|\widehat{V} - V\|_F + 1) \|\widehat{\Lambda}\| \\
& \leq c \left(\left(\frac{24}{\sqrt{\pi_{\min}}} + 3 \right) \sigma_1(O) E_3 + \left(4 + \frac{6\sigma_1(O)}{\sqrt{\pi_{\min}}} \right) \frac{E_2}{\sigma_k(O)^2 \pi_{\min}} \right) \\
& \leq c \left(\frac{27\sigma_1(O)}{\sqrt{\pi_{\min}}} + \frac{10 \max(\sigma_1(O), 1)}{\pi_{\min}^{3/2} \sigma_k(O)^2} \right) \max(E_2, E_3) \\
& \leq c \frac{37 \max(\sigma_1(O), 1)}{\pi_{\min}^{3/2} \min(\sigma_k(O)^2, 1)} \max(E_2, E_3)
\end{aligned}$$

where $c > 0$ is a constant large enough to dominate low-order terms like $E_2 E_3$. \square

Lemma 5. *With a slight abuse of notation, let U denote a column permutation of the true U , UT denote a column permutation of the true UT , and P denote a column-and-row permutation of the true P , where the permutations involved are the same. Suppose*

$$\max(\|U - \widehat{U}\|, \|\widehat{UT} - UT\|) \leq \sigma_k(rU + (1-r)UT)/2.$$

We then have

$$\|P - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| \leq \frac{6\sigma_1(UT)}{\sigma_k(rU + (1-r)UT)^2} \max(\|U - \widehat{U}\|, \|UT - \widehat{UT}\|).$$

Proof. First notice that

$$\begin{aligned}
& (rU + (1-r)UT)^\dagger (UT) \\
& = ((rI + (1-r)T)^\top U^\top (rI + (1-r)T))^{-1} (rI + (1-r)T)^\top U^\top UT \\
& = (rI + (1-r)T)^{-1} T = P.
\end{aligned}$$

Then we have

$$\begin{aligned}
& \|P - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| = \|(rU + (1-r)UT)^\dagger UT - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| \\
& \leq \|(rU + (1-r)UT)^\dagger (UT) - (r\widehat{U} + (1-r)\widehat{UT})^\dagger (UT)\| + \\
& \quad \|(r\widehat{U} + (1-r)\widehat{UT})^\dagger (UT) - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| \\
& \leq \|(rU + (1-r)UT)^\dagger - (r\widehat{U} + (1-r)\widehat{UT})^\dagger\| \|UT\| + \|(r\widehat{U} + (1-r)\widehat{UT})^\dagger\| \|UT - \widehat{UT}\|. \tag{20}
\end{aligned}$$

By Lemma 3 and the assumption of the lemma, we have

$$\sigma_k(rU + (1-r)UT)/2 \leq \sigma_k(r\widehat{U} + (1-r)\widehat{UT}) \leq 3\sigma_k(rU + (1-r)UT)/2,$$

showing that $\text{rank}(r\widehat{U} + (1-r)\widehat{UT}) = k$ and

$$\|(r\widehat{U} + (1-r)\widehat{UT})^\dagger\| = 1/\sigma_k(r\widehat{U} + (1-r)\widehat{UT}) \leq 2/\sigma_k(rU + (1-r)UT).$$

Because $\text{rank}(r\widehat{U} + (1-r)\widehat{UT}) = \text{rank}(rU + (1-r)UT) = k$, Theorem 3.4 in [3] indicates that

$$\begin{aligned}
& \|(rU + (1-r)UT)^\dagger - (r\widehat{U} + (1-r)\widehat{UT})^\dagger\| \\
& \leq \sqrt{2} \|(rU + (1-r)UT)^\dagger\| \|(r\widehat{U} + (1-r)\widehat{UT})^\dagger\| \|r(U - \widehat{U}) + (1-r)(UT - \widehat{UT})\| \\
& \leq \frac{\sqrt{2}(\|rU - \widehat{U}\| + (1-r)\|\widehat{UT} - UT\|)}{\sigma_k(rU + (1-r)UT) \sigma_k(r\widehat{U} + (1-r)\widehat{UT})} \leq \frac{2\sqrt{2}(\|rU - \widehat{U}\| + (1-r)\|UT - \widehat{UT}\|)}{\sigma_k(rU + (1-r)UT)^2}.
\end{aligned}$$

Applying these bounds to (20) then leads to

$$\begin{aligned}
& \|P - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| \\
& \leq \frac{2\sqrt{2}\sigma_1(UT)(r\|U - \widehat{U}\| + (1-r)\|UT - \widehat{UT}\|)}{\sigma_k(rU + (1-r)UT)^2} + \frac{2\|UT - \widehat{UT}\|}{\sigma_k(rU + (1-r)UT)} \\
& = \frac{r2\sqrt{2}\sigma_1(UT)\|U - \widehat{U}\|}{\sigma_k(rU + (1-r)UT)^2} + \frac{((1-r)2\sqrt{2}\sigma_1(UT) + 2\sigma_k(rU + (1-r)UT))\|UT - \widehat{UT}\|}{\sigma_k(rU + (1-r)UT)^2} \\
& \leq \frac{\max(r2\sqrt{2}, (1-r)2\sqrt{2} + 2)\sigma_1(UT)}{\sigma_k(rU + (1-r)UT)^2} \max(\|U - \widehat{U}\|, \|UT - \widehat{UT}\|) \\
& \leq \frac{6\sigma_1(UT)}{\sigma_k(rU + (1-r)UT)^2} \max(\|U - \widehat{U}\|, \|UT - \widehat{UT}\|),
\end{aligned}$$

in which we use the fact $\sigma_1(UT) \geq \sigma_1(rU + (1-r)UT) \geq \sigma_k(rU + (1-r)UT)$. \square

D.3 Concentration of empirical averages

Lemma 6. Let $\{\mathbf{y}_i\}_{i=1}^N$ be N i.i.d. random vectors in \mathcal{R}^m . Let $\boldsymbol{\mu} := \mathbb{E}[\mathbf{y}_i]$, $\Sigma := \text{Var}(\mathbf{y}_i)$ and $\sigma_{\max}^2 := \max_d \Sigma_{dd}$. Let $\bar{\boldsymbol{\mu}} := (\sum_i \mathbf{y}_i)/N$. Then

$$\text{Prob}(\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \geq \epsilon) \leq \frac{m\sigma_{\max}^2}{N\epsilon^2}.$$

Proof. This lemma is a straightforward consequence of the Markov inequality:

$$\text{Prob}(\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2 \geq \epsilon) = \text{Prob}(\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2^2 \geq \epsilon^2) \quad (21)$$

$$\leq \frac{\mathbb{E}[\|\bar{\boldsymbol{\mu}} - \boldsymbol{\mu}\|_2^2]}{\epsilon^2} \quad (22)$$

$$= \frac{\sum_d \mathbb{E}(\bar{\boldsymbol{\mu}}_d - \boldsymbol{\mu}_d)^2}{\epsilon^2} = \frac{\text{Tr}(\Sigma)}{N\epsilon^2} \leq \frac{m\sigma_{\max}^2}{N\epsilon^2}. \quad (23)$$

\square

Lemma 7. Let $\widehat{V}_1, \widehat{V}_2, \widehat{V}_3, \widehat{C}_2, \widehat{C}_3$ denote averages of N independent draws of $\mathbf{x}_1, \mathbf{x}_1\mathbf{x}_1^\top, \mathbf{x}_1 \otimes \mathbf{x}_1 \otimes \mathbf{x}_1, \mathbf{x}_1\mathbf{x}_2^\top, \mathbf{x}_1 \otimes \mathbf{x}_2 \otimes \mathbf{x}_3$ from the generative process in Section 3.2. Let $u_{\max} := \max_{i,j} |U_{ij}|$. Then

$$\text{Prob}(\|\widehat{V}_1 - V_1\|_2 \geq \epsilon) \leq \frac{m(u_{\max}^2 + \sigma^2)}{N\epsilon^2}, \quad (24)$$

$$\text{Prob}(\|\widehat{V}_2 - V_2\|_F \geq \epsilon) \leq \frac{m^2(u_{\max}^2 + \sigma^2)^2}{N\epsilon^2}, \quad (25)$$

$$\text{Prob}(\|\widehat{V}_3 - V_3\|_F \geq \epsilon) \leq \frac{m^3(u_{\max}^2 + \sigma^2)^3}{N\epsilon^2}, \quad (26)$$

$$\text{Prob}(\|\widehat{C}_2 - C_2\|_F \geq \epsilon) \leq \frac{m^2(u_{\max}^2 + \sigma^2)^2}{N\epsilon^2}, \quad (27)$$

$$\text{Prob}(\|\widehat{C}_3 - C_3\|_F \geq \epsilon) \leq \frac{m^3(u_{\max}^2 + \sigma^2)^3}{N\epsilon^2}. \quad (28)$$

Proof. Based on Lemma 6, it suffices to bound σ_{\max}^2 in these five cases:

$$\max_i \text{Var}((\mathbf{x}_1)_i) \leq \max_i \mathbb{E}[(\mathbf{x}_1)_i^2] = \max_i \mathbb{E}_{\mathbf{h}_1}[\sigma^2 + (U\mathbf{h}_1)_i^2] \leq \sigma^2 + \max_{i,k} U_{ik}^2, \quad (29)$$

$$\begin{aligned} \max_{i,j} \text{Var}((\mathbf{x}_1)_i(\mathbf{x}_1)_j) &\leq \max_{i,j} \mathbb{E}[(\mathbf{x}_1)_i^2(\mathbf{x}_1)_j^2] = \max_{i,j} \mathbb{E}_{\mathbf{h}_1}[(\sigma^2 + (U\mathbf{h}_1)_i^2)(\sigma^2 + (U\mathbf{h}_1)_j^2)] \\ &\leq \max_{i,j,l} (\sigma^2 + U_{il}^2)(\sigma^2 + U_{jl}^2) \leq (\sigma^2 + \max_{i,j} U_{ij}^2)^2, \end{aligned} \quad (30)$$

$$\max_{i,j} \text{Var}((\mathbf{x}_1)_i(\mathbf{x}_2)_j) \leq \max_{i,j} \mathbb{E}[(\mathbf{x}_1)_i^2(\mathbf{x}_2)_j^2] = \max_{i,j} \mathbb{E}_{\boldsymbol{\pi}_0}[\mathbb{E}[(\mathbf{x}_1)_i^2|\boldsymbol{\pi}_0]\mathbb{E}[(\mathbf{x}_2)_j^2|\boldsymbol{\pi}_0]] \quad (31)$$

$$\leq \max_{i,j} \sup_{\boldsymbol{\pi}_0} \mathbb{E}[(\mathbf{x}_1)_i^2|\boldsymbol{\pi}_0]\mathbb{E}[(\mathbf{x}_2)_j^2|\boldsymbol{\pi}_0] \leq (\max_{i,j} \sup_{\boldsymbol{\pi}_0} \mathbb{E}[(\mathbf{x}_1)_i^2|\boldsymbol{\pi}_0])^2 \quad (32)$$

$$= (\max_{i,j} \sup_{\boldsymbol{\pi}_0} \sum_k U_{ij}^2 (T\boldsymbol{\pi}_0)_k + \sigma^2)^2 \quad (33)$$

$$= (\max_i \max_{j'} \sum_k U_{ij}^2 T_{jj'} + \sigma^2)^2 \leq (\max_{i,j} U_{ij}^2 + \sigma^2)^2. \quad (34)$$

With similar arguments, we have that

$$\max_{i,j,l} \text{Var}((\mathbf{x}_1)_i(\mathbf{x}_1)_j(\mathbf{x}_1)_l) \leq (\max_{i,j} U_{ij}^2 + \sigma^2)^3, \quad (35)$$

$$\max_{i,j,l} \text{Var}((\mathbf{x}_1)_i(\mathbf{x}_2)_j(\mathbf{x}_3)_l) \leq (\max_{i,j} U_{ij}^2 + \sigma^2)^3. \quad (36)$$

□

Lemma 8. Let $\widehat{M}_2, \widehat{M}_3, \widehat{M}'_2, \widehat{M}'_3$ denote estimates obtained by plugging in empirical averages of independent samples as in Lemma 7 and $\widehat{\sigma}_2 := \lambda_{\min}(\widehat{V}_2 - \widehat{V}_1\widehat{V}_1^\top)$, where $\lambda_{\min}(\cdot)$ denotes the smallest eigenvalue in modulus. Define $\nu := \max(\sigma^2 + u_{\max}^2, 1)$. We then have the following:

$$\begin{aligned} \text{Prob}(\|M_2 - \widehat{M}_2\| \geq \epsilon) &\leq \frac{75m^2\nu^2}{N\epsilon^2}, \\ \text{Prob}(\|M_3 - \widehat{M}_3\| \geq \epsilon) &\leq \frac{1000m^4\nu^3}{N\epsilon^2}, \\ \text{Prob}(\|M'_2 - \widehat{M}'_2\| \geq \epsilon) &\leq \frac{50(\alpha_0 + 1)^2 m^2 \nu^2}{N\epsilon^2}, \\ \text{Prob}(\|M'_3 - \widehat{M}'_3\| \geq \epsilon) &\leq \frac{1100k^2 m^3 (\alpha_0 + 2)^2 (\alpha_0 + 1)^2 \nu^3}{N\epsilon^2}. \end{aligned}$$

Proof. We first note that it is easy to verify $\mathbf{z}^\top (\widehat{V}_2 - \widehat{V}_1\widehat{V}_1^\top) \mathbf{z} \geq 0$ for any real vector \mathbf{z} , so $\widehat{\sigma}_2$ is always non-negative. By Lemma 3, we have

$$\begin{aligned} |\sigma^2 - \widehat{\sigma}_2| &\leq \|V_2 - V_1V_1^\top - (\widehat{V}_2 - \widehat{V}_1\widehat{V}_1^\top)\| \leq \|V_2 - \widehat{V}_2\| + \|V_1V_1^\top - \widehat{V}_1\widehat{V}_1^\top\| \\ &\leq \|V_2 - \widehat{V}_2\| + \|\widehat{V}_1 - V_1\|(\|\widehat{V}_1\| + \|V_1\|) \\ &\leq \|V_2 - \widehat{V}_2\| + 2\|V_1\|\|\widehat{V}_1 - V_1\| + \|V_1 - \widehat{V}_1\|^2. \end{aligned}$$

We also need the following

$$\|V_1\|^2 = \|U\boldsymbol{\pi}\|^2 = \sum_i \left(\sum_j U_{ij}\pi_j \right)^2 \leq \sum_{i,j} \pi_j U_{ij}^2 \leq \sum_i \max_j U_{ij}^2 \leq m u_{\max}^2.$$

Then we have

$$\begin{aligned} \|\widehat{M}_2 - M_2\| &\leq \|\widehat{V}_2 - V_2\| + |\widehat{\sigma}_2 - \sigma^2| \\ &\leq 2\|\widehat{V}_2 - V_2\| + 2\|V_1\|\|\widehat{V}_1 - V_1\| + \|\widehat{V}_1 - V_1\|^2 \\ &\leq 2\|\widehat{V}_2 - V_2\|_F + 2\|V_1\|\|\widehat{V}_1 - V_1\| + \|\widehat{V}_1 - V_1\|^2, \end{aligned}$$

which implies

$$\begin{aligned}
& \text{Prob}(\|\widehat{M}_2 - M_2\| \geq \epsilon) \\
& \leq \text{Prob}(2\|\widehat{V}_2 - V_2\|_F + 2\|V_1\| \|\widehat{V}_1 - V_1\| + \|\widehat{V}_1 - V_1\|^2 \geq \epsilon) \\
& \leq \text{Prob}(2\|\widehat{V}_2 - V_2\|_F \geq \epsilon/3) + \text{Prob}(2\|V_1\| \|\widehat{V}_1 - V_1\| \geq \epsilon/3) + \text{Prob}(\|\widehat{V}_1 - V_1\|^2 \geq \epsilon/3) \\
& \leq \frac{36m^2(u_{\max}^2 + \sigma^2)^2}{N\epsilon^2} + \frac{36\|V_1\|^2 m(u_{\max}^2 + \sigma^2)}{N\epsilon^2} + \frac{3m(u_{\max}^2 + \sigma^2)}{N\epsilon} \\
& \leq \frac{36m^2(u_{\max}^2 + \sigma^2)^2}{N\epsilon^2} + \frac{36m^2 u_{\max}^2 (u_{\max}^2 + \sigma^2)}{N\epsilon^2} + \frac{3m(u_{\max}^2 + \sigma^2)}{N\epsilon} \\
& \leq \frac{75m^2(u_{\max}^2 + \sigma^2)^2}{N\epsilon^2}.
\end{aligned}$$

Similarly, we have

$$\begin{aligned}
\|M_3 - \widehat{M}_3\| & \leq \|V_3 - \widehat{V}_3\|_F + 3\|V_1 \otimes_1 (\sigma^2 I) - \widehat{V}_1 \otimes_1 (\widehat{\sigma}^2 I)\|_F \\
& = \|V_3 - \widehat{V}_3\|_F + 3\sqrt{m} \|\sigma^2 V_1 - \widehat{\sigma}^2 \widehat{V}_1\| \\
& \leq \|V_3 - \widehat{V}_3\|_F + 3\sqrt{m} (\sigma^2 \|V_1 - \widehat{V}_1\| + |\sigma^2 - \widehat{\sigma}^2| (\|V_1\| + \|\widehat{V}_1 - V_1\|)) \\
& \leq \|V_3 - \widehat{V}_3\|_F + \|V_1 - \widehat{V}_1\| 3\sqrt{m} (\sigma^2 + 2mu_{\max}^2) + \|V_2 - \widehat{V}_2\| 3u_{\max} m \\
& \quad + \|V_1 - \widehat{V}_1\|^2 9u_{\max} m + 3\sqrt{m} (\|V_1 - \widehat{V}_1\| \|\widehat{V}_2 - V_2\| + \|V_1 - \widehat{V}_1\|^3),
\end{aligned}$$

implying

$$\begin{aligned}
& \text{Prob}(\|M_3 - \widehat{M}_3\| \geq \epsilon) \\
& \leq \text{Prob}(\|V_3 - \widehat{V}_3\|_F \geq \epsilon/6) + \text{Prob}(\|V_1 - \widehat{V}_1\| \geq \epsilon/(18\sqrt{m}(\sigma^2 + 2mu_{\max}^2))) \\
& \quad + \text{Prob}(\|V_2 - \widehat{V}_2\| \geq \epsilon/(18u_{\max} m)) + \text{Prob}(\|V_1 - \widehat{V}_1\|^2 \geq \epsilon/(54u_{\max} m)) \\
& \quad + \text{Prob}\left(\|V_1 - \widehat{V}_1\| \geq \sqrt{\epsilon/(18\sqrt{m})}\right) + \text{Prob}\left(\|V_2 - \widehat{V}_2\| \geq \sqrt{\epsilon/(18\sqrt{m})}\right) \\
& \quad + \text{Prob}(\|V_1 - \widehat{V}_1\|^3 \geq \epsilon/(18\sqrt{m})) \\
& \leq \frac{36m^3(u_{\max}^2 + \sigma^2)^3}{N\epsilon^2} + \frac{324m^2(\sigma^2 + 2mu_{\max}^2)^2(\sigma^2 + u_{\max}^2)}{N\epsilon^2} + \frac{324u_{\max}^2 m^4(\sigma^2 + u_{\max}^2)^2}{N\epsilon^2} \\
& \quad + \frac{54u_{\max} m^2(\sigma^2 + u_{\max}^2)}{N\epsilon} + \frac{18m^{3/2}(\sigma^2 + u_{\max}^2)}{N\epsilon} + \frac{18m^{5/2}(\sigma^2 + u_{\max}^2)^2}{N\epsilon} \\
& \quad + \frac{36^{1/3} m^{4/3}(\sigma^2 + u_{\max}^2)}{N\epsilon^{2/3}} \\
& \leq \frac{1000m^4(\max(\sigma^2 + u_{\max}^2, 1))^3}{N\epsilon^2}.
\end{aligned}$$

Using similar arguments, we have

$$\begin{aligned}
\|M'_2 - \widehat{M}'_2\| & \leq (\alpha_0 + 1)\|C_2 - \widehat{C}_2\|_F + \alpha_0\|V_1 V_1^\top - \widehat{V}_1 \widehat{V}_1^\top\|_F \\
& \leq (\alpha_0 + 1)\|C_2 - \widehat{C}_2\|_F + 2\alpha_0\|V_1\| \|\widehat{V}_1 - V_1\| + \alpha_0\|\widehat{V}_1 - V_1\|^2,
\end{aligned}$$

and therefore

$$\begin{aligned}
& \text{Prob}(\|M'_2 - \widehat{M}'_2\| \geq \epsilon) \\
& \leq \text{Prob}(\|C_2 - \widehat{C}_2\|_F \geq \frac{\epsilon}{3(\alpha_0 + 1)}) + \text{Prob}(\|\widehat{V}_1 - V_1\| \geq \frac{\epsilon}{6\alpha_0\|V_1\|}) \\
& \quad + \text{Prob}(\|\widehat{V}_1 - V_1\|^2 \geq \frac{\epsilon}{3\alpha_0}) \\
& \leq \frac{9(\alpha_0 + 1)^2 m^2 (\sigma^2 + u_{\max}^2)^2}{N\epsilon^2} + \frac{36\alpha_0^2 m^2 u_{\max}^2 (\sigma^2 + u_{\max}^2)}{N\epsilon^2} + \frac{3\alpha_0 m (\sigma^2 + u_{\max}^2)}{N\epsilon} \\
& \leq \frac{50(\alpha_0 + 1)^2 m^2 (\sigma^2 + u_{\max}^2)^2}{N\epsilon^2}.
\end{aligned}$$

Finally, we have

$$\begin{aligned}
& \|M'_3 - \widehat{M}'_3\| \\
\leq & \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} \|C_3 - \widehat{C}_3\|_F + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1 \otimes_1 C_2 - \widehat{V}_1 \otimes \widehat{C}_2\|_F \\
& + \alpha_0^2 \|V_1 \otimes V_1 \otimes V_1 - \widehat{V}_1 \otimes \widehat{V}_1 \otimes \widehat{V}_1\|_F \\
\leq & \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} \|C_3 - \widehat{C}_3\|_F + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1 - \widehat{V}_1\| \|C_2\|_F + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|\widehat{V}_1\| \|C_2 - \widehat{C}_2\|_F \\
& + 3\alpha_0^2 \|V_1\|^2 \|V_1 - \widehat{V}_1\| + 3\alpha_0^2 \|V_1\| \|V_1 - \widehat{V}_1\|^2 + \alpha_0^2 \|V_1 - \widehat{V}_1\|^3 \\
\leq & \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} \|C_3 - \widehat{C}_3\|_F + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1 - \widehat{V}_1\| \|C_2\|_F + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1\| \|C_2 - \widehat{C}_2\|_F \\
& + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1 - \widehat{V}_1\| \|C_2 - \widehat{C}_2\|_F + 3\alpha_0^2 \|V_1\|^2 \|V_1 - \widehat{V}_1\| + 3\alpha_0^2 \|V_1\| \|V_1 - \widehat{V}_1\|^2 + \alpha_0^2 \|V_1 - \widehat{V}_1\|^3 \\
\leq & \frac{(\alpha_0 + 2)(\alpha_0 + 1)}{2} \|C_3 - \widehat{C}_3\|_F + 5(\alpha_0 + 1)\alpha_0 k m u_{\max}^2 \|V_1 - \widehat{V}_1\| + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1\| \|C_2 - \widehat{C}_2\|_F \\
& + \frac{3(\alpha_0 + 1)\alpha_0}{2} \|V_1 - \widehat{V}_1\| \|C_2 - \widehat{C}_2\|_F + 3\alpha_0^2 \|V_1\| \|V_1 - \widehat{V}_1\|^2 + \alpha_0^2 \|V_1 - \widehat{V}_1\|^3
\end{aligned}$$

using the fact that

$$\|C_2\|_F = \left\| UT \left(\frac{\text{diag} \boldsymbol{\pi} + \alpha_0 \boldsymbol{\pi} \boldsymbol{\pi}^\top}{\alpha_0 + 1} \right) T^\top U^\top \right\|_F \leq \|UT\|_F^2 \leq k m u_{\max}^2,$$

and thus

$$\begin{aligned}
\text{Prob}(\|M'_3 - \widehat{M}'_3\| \geq \epsilon) & \leq \text{Prob} \left(\|C_3 - \widehat{C}_3\|_F \geq \frac{\epsilon}{3(\alpha_0 + 2)(\alpha_0 + 1)} \right) \\
& + \text{Prob} \left(\|V_1 - \widehat{V}_1\|_F \geq \frac{\epsilon}{30(\alpha_0 + 1)\alpha_0 k m u_{\max}^2} \right) + \text{Prob} \left(\|C_2 - \widehat{C}_2\|_F \geq \frac{\epsilon}{9(\alpha_0 + 1)\alpha_0} \right) \\
& + \text{Prob} \left(\|V_1 - \widehat{V}_1\|^2 \geq \frac{\epsilon}{18\alpha_0^2 \|V_1\|} \right) + \text{Prob} \left(\|V_1 - \widehat{V}_1\|^3 \geq \frac{\epsilon}{6\alpha_0^2} \right) \\
\leq & \frac{9m^2(\alpha_0 + 2)^2(\alpha_0 + 1)^2(\sigma^2 + u_{\max}^2)^3}{N\epsilon^2} + \frac{900k^2m^3(\alpha_0 + 1)^2\alpha_0^2u_{\max}^4(\sigma^2 + u_{\max}^2)}{N\epsilon^2} \\
& + \frac{81(\alpha_0 + 1)^2\alpha_0^2m^2(\sigma^2 + u_{\max}^2)^2}{N\epsilon^2} + \frac{18\alpha_0^2m^{3/2}u_{\max}(\sigma^2 + u_{\max}^2)}{N\epsilon} + \frac{6m\alpha_0^{4/3}(\sigma^2 + u_{\max}^2)}{N\epsilon^{2/3}} \\
\leq & \frac{1100k^2m^3(\alpha_0 + 2)^2(\alpha_0 + 1)^2(\sigma^2 + u_{\max}^2)^3}{N\epsilon^2}.
\end{aligned}$$

□

E Proof of Theorem 4

Let \widehat{U} and \widehat{UT} be column-permuted as described in Algorithm 1. Let

$$\delta_{\min} := \min_{i,j} |1/\sqrt{\pi_i} - 1/\sqrt{\pi_j}|.$$

If $\max(E_3, E'_3) \leq \delta_{\min}/15$, Theorem 5.1 of [2] implies that for any $\eta \in (0, 1)$, with probability at least $1 - \eta$, the columns of \widehat{U} and \widehat{UT} are matched to the same permutation of the columns of the true U and UT , respectively. As in Lemma 5, let U, UT , and P denote proper permutations of the true matrices. We then have

$$\begin{aligned}
& \text{Prob} \left(\max(\|U - \widehat{U}\|, \|UT - \widehat{UT}\|) \geq \frac{\epsilon \sigma_k (rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \\
\leq & \text{Prob} \left(\|U - \widehat{U}\| \geq \frac{\epsilon \sigma_k (rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) + \text{Prob} \left(\|UT - \widehat{UT}\| \geq \frac{\epsilon \sigma_k (rU + (1-r)UT)^2}{6\sigma_1(UT)} \right).
\end{aligned}$$

Let the failure probability for the tensor decomposition method be set to $\frac{\eta}{4}$. Then by Lemma 4 we can bound the first term as follows:

$$\begin{aligned} & \text{Prob} \left(\|U - \widehat{U}\| \geq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \\ & \leq \text{Prob} \left(\max(E_2, E_3) \geq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2 \pi_{\min}^{3/2} \min(\sigma_k(U)^2, 1)}{6\sigma_1(UT)c \max(\sigma_1(U), 1)} \right) \\ & \quad + \text{Prob}(\max(E_2, E_3) \geq \sigma_k(M_2)/2) + \frac{\eta}{4} + \text{Prob}(E_3 \geq \delta_{\min}/15), \end{aligned}$$

where the first term in the r.h.s is based on Lemma 4 conditioned on the event that $\max(E_2, E_3) \geq \sigma_k(M_2)/2$ and the tensor decomposition method succeeds, the second and the third terms bound the probability that the event does not occur, and the last term bounds the probability of incorrectly matching the columns of \widehat{U} and U . To continue the bound we use Lemma 8 to have

$$\begin{aligned} & \text{Prob} \left(\max(E_2, E_3) \geq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2 \pi_{\min}^{3/2} \min(\sigma_k(U)^2, 1)}{6\sigma_1(UT)c \max(\sigma_1(U), 1)} \right) \\ & \leq \frac{(2700m^2\nu^2 + 36000m^4\nu^3)\sigma_1(UT)^2 c^2 \max(\sigma_1(U)^2, 1)}{N\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \pi_{\min}^3 \min(\sigma_k(U)^4, 1)} \\ & \leq \frac{39000m^4\nu^3 \sigma_1(UT)^2 c^2 \max(\sigma_1(U)^2, 1)}{N\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \pi_{\min}^3 \min(\sigma_k(U)^4, 1)}, \\ & \text{Prob}(\max(E_2, E_3) \geq \sigma_k(M_2)/2) \leq \frac{300m^2\nu^2 + 4000m^4\nu^3}{N\sigma_k(M_2)^2} \leq \frac{4300m^4\nu^3}{N\sigma_k(M_2)^2}, \\ & \text{Prob}(E_3 \geq \delta_{\min}/15) \leq \frac{225000m^4\nu^3}{N\delta_{\min}^2}. \end{aligned}$$

Thus, by setting the sample size N so that

$$N \geq \frac{12m^4\nu^3}{\eta} \max \left(\frac{225000}{\delta_{\min}^2}, \frac{4300}{\sigma_k(M_2)^2}, \frac{39000\sigma_1(UT)^2 c^2 \max(\sigma_1(U)^2, 1)}{\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \pi_{\min}^3 \min(\sigma_k(U)^4, 1)} \right),$$

we have

$$\text{Prob} \left(\|U - \widehat{U}\| \geq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \leq \frac{\eta}{2}, \quad (37)$$

where the randomness is from both the data and the algorithm. Using similar arguments, we have that for sample size N such that

$$\begin{aligned} N & \geq \frac{12k^2m^3(\alpha_0 + 2)^2(\alpha_0 + 1)^2\nu^3}{\eta} \\ & \quad \max \left(\frac{225000}{\delta_{\min}^2}, \frac{4600}{\sigma_k(M'_2)^2}, \frac{42000\sigma_1(UT)^2(c')^2 \max(\sigma_1(UT)^2, 1)}{\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \pi_{\min}^3 \min(\sigma_k(UT)^4, 1)} \right), \end{aligned}$$

the following holds:

$$\text{Prob} \left(\|UT - \widehat{UT}\| \geq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \leq \frac{\eta}{2}.$$

Combining the two bounds (37) and (E), we have for

$$\begin{aligned} N & \geq \frac{12 \max(k^2, m)m^3\nu^3(\alpha_0 + 2)^2(\alpha_0 + 1)^2}{\eta} \\ & \quad \max \left(\frac{225000}{\delta_{\min}^2}, \frac{4600}{\min(\sigma_k(M'_2), \sigma_k(M_2))^2}, \frac{42000c^2\sigma_1(UT)^2 \max(\sigma_1(UT), \sigma_1(U), 1)^2}{\epsilon^2 \sigma_k(rU + (1-r)UT)^4 \min(\sigma_k(UT), \sigma_k(U), 1)^4} \right), \end{aligned}$$

the following bound holds for any $\epsilon > 0$ and $\eta \in (0, 1)$:

$$\text{Prob} \left(\max(\|U - \widehat{U}\|, \|UT - \widehat{UT}\|) \leq \frac{\epsilon \sigma_k(rU + (1-r)UT)^2}{6\sigma_1(UT)} \right) \geq 1 - \eta,$$

which by Lemma 5 implies that

$$\text{Prob}(\|P - (r\widehat{U} + (1-r)\widehat{UT})^\dagger \widehat{UT}\| \leq \epsilon) \geq 1 - \eta.$$

F Comparison between Proposed Method and Variational EM

For the purpose of comparison, we derive here a variational EM algorithm for learning HMMs from non-sequence data. First notice that the generative process in Section 3.2 specifies only the variance rather than an entire distribution for the observation model. To make it easier to derive a maximum-likelihood based algorithm, here we make an extra assumption of the observation model being Gaussian. The full joint probability of data and latent variables then takes the following form:

$$\begin{aligned}
& f(\{\mathbf{x}_i^j\}, \{\mathbf{h}_i^j\}, \{t_i^j\}, \{\mathbf{s}_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid U, \sigma^2, P, r, \boldsymbol{\alpha}) \\
&= \prod_{j=1}^N \left(\prod_{i=1}^n \left(\prod_{l=1}^k \mathcal{N}(\mathbf{x}_i^j \mid U_l, \sigma^2 I)^{\mathbf{h}_{il}^j} \right) \left(\prod_{l',l} ((P^{t_i^j})_{l'l})^{\mathbf{h}_{i'l'}^j \mathbf{s}_{il}^j} \right) \text{Geometric}(t_i^j \mid r) \left(\prod_l ((\boldsymbol{\pi}_0^j)_l)^{\mathbf{s}_{il}^j} \right) \right) \\
& \quad \text{Dirichlet}(\boldsymbol{\pi}_0^j \mid \boldsymbol{\alpha}),
\end{aligned} \tag{38}$$

in which we use super-script as set indices and sub-scripts as data indices within a set wherever appropriate. The goal is to maximize the marginal likelihood of the data w.r.t to the parameters. We begin by marginalizing over the latent times $\{t_i^j\}$:

$$\begin{aligned}
& f(\{\mathbf{x}_i^j\}, \{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid U, \sigma^2, T, \boldsymbol{\alpha}) \\
&= \prod_{j=1}^N \left(\prod_{i=1}^n \left(\prod_{l=1}^k \mathcal{N}(\mathbf{x}_i^j \mid U_l, \sigma^2 I)^{\mathbf{h}_{il}^j} \right) \left(\prod_{l',l} T_{l'l}^{\mathbf{h}_{i'l'}^j \mathbf{s}_{il}^j} \right) \left(\prod_l ((\boldsymbol{\pi}_0^j)_l)^{\mathbf{s}_{il}^j} \right) \right) \text{Dirichlet}(\boldsymbol{\pi}_0^j \mid \boldsymbol{\alpha}),
\end{aligned} \tag{39}$$

where T is the expected hidden state transition matrix defined in Theorem 3. As in the tensor factorization approach, we recover P and r from the estimated T using the proposed search heuristics. Because the posterior distribution of the remaining latent variables still leads to an intractable E step, we take the following factorized approximation:

$$f(\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid \{\mathbf{x}_i^j\}, U, \sigma^2, T, \boldsymbol{\alpha}) \approx q(\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}) q(\{\boldsymbol{\pi}_0^j\} \mid \{\beta^j\}), \tag{40}$$

where

$$q(\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}) := \prod_{i,j,l',l} ((\Phi_i^j)_{l'l})^{\mathbf{h}_{i'l'}^j \mathbf{s}_{il}^j}, \quad \Phi_i^j \in [0, 1]^{k \times k}, \tag{41}$$

$$q(\{\boldsymbol{\pi}_0^j\} \mid \{\beta^j\}) := \prod_j \text{Dirichlet}(\boldsymbol{\pi}_0^j \mid \beta^j), \tag{42}$$

and obtain the following lower bound on the log marginal likelihood:

$$\begin{aligned}
& g(\{\Phi_i^j\}, \{\beta^j\}, U, \sigma^2, T, \boldsymbol{\alpha}) \\
&:= \mathbb{E}_{\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid \{\beta^j\}} \left[\log \left(\frac{f(\{\mathbf{x}_i^j\}, \{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid U, \sigma^2, T, \boldsymbol{\alpha})}{q(\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}) q(\{\boldsymbol{\pi}_0^j\} \mid \{\beta^j\})} \right) \right] \\
&= \mathbb{E}_{\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid \{\beta^j\}} \left[\log f(\{\mathbf{x}_i^j\}, \{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\}, \{\boldsymbol{\pi}_0^j\} \mid U, \sigma^2, T, \boldsymbol{\alpha}) \right] - \\
& \quad \mathbb{E}_{\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}} \left[\log q(\{\mathbf{h}_i^j\}, \{\mathbf{s}_i^j\} \mid \{\Phi_i^j\}) \right] - \mathbb{E}_{\{\boldsymbol{\pi}_0^j\} \mid \{\beta_0^j\}} \left[\log q(\{\boldsymbol{\pi}_0^j\} \mid \{\beta_0^j\}) \right] \\
&= \sum_{j,i,l,l'} (\Phi_i^j)_{l'l} (\log \mathcal{N}(\mathbf{x}_i^j \mid U_l, \sigma^2 I) + \log T_{l'l}) + \sum_{j,l} \left(\sum_{i,l'} (\Phi_i^j)_{l'l} + \alpha_l - 1 \right) (\psi(\beta_l^j) - \psi(\beta_0^j)) \\
& \quad - N \left(\sum_l \log \Gamma(\alpha_l) - \log \Gamma(\alpha_0) \right) - \sum_{j,i,l,l'} (\Phi_i^j)_{l'l} \log(\Phi_i^j)_{l'l} \\
& \quad - \sum_{j,l} (\beta_l^j - 1) (\psi(\beta_l^j) - \psi(\beta_0^j)) + \sum_j \left(\sum_l \log \Gamma(\beta_l^j) - \log \Gamma(\beta_0^j) \right),
\end{aligned} \tag{43}$$

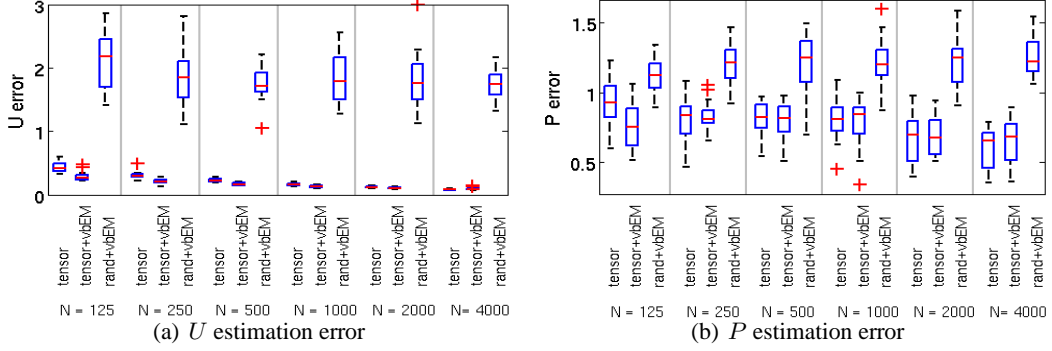


Figure 1: Comparison between Algorithm 1 and Variational EM

where $\psi(\cdot)$ is the digamma function. The variational EM algorithm then amounts to maximizing g iteratively, optimizing over one block of variables at a time while fixing the others:

$$(\Phi_i^j)_{l'w} \propto \mathcal{N}(\mathbf{x}_i^j | U_l, \sigma^2 I) T_{l'w} \exp(\psi(\beta_{l'}^j) - \psi(\beta_0^j)), \quad (44)$$

$$(\beta^j)_l = \sum_{i,l'} (\Phi_i^j)_{l'l} + \alpha_l, \quad (45)$$

$$U_l := \frac{\sum_{j=1}^N \sum_{i=1}^n \sum_{l'=1}^k (\Phi_i^j)_{l'w} \mathbf{x}_i^j}{\sum_{j=1}^N \sum_{i=1}^n \sum_{l'=1}^k (\Phi_i^j)_{l'w}}, \quad (46)$$

$$\sigma^2 := \frac{\sum_{j=1}^N \sum_{i=1}^n \sum_{l,l'} (\Phi_i^j)_{l'w} \|\mathbf{x}_i^j - U_l\|^2}{Nnm}, \quad (47)$$

$$T_{l'w} := \frac{\sum_{j=1}^N \sum_{i=1}^n (\Phi_i^j)_{l'w}}{\sum_{j=1}^N \sum_{i=1}^n \sum_{l=1}^k (\Phi_i^j)_{l'w}}, \quad (48)$$

$$\alpha := \arg \max_{\{\alpha_l \geq 0\}} \sum_{j=1}^N \sum_{l=1}^k (\alpha_l - 1) (\psi(\beta_l^j) - \psi(\beta_0^j)) - N \left(\sum_{l=1}^k \log \Gamma(\alpha_l) - \log \Gamma(\alpha_0) \right) \quad (49)$$

The update for α is a convex optimization problem whose inverse Hessian can be computed in linear time. We iterate these updates until the lower bound stops increasing.

Finally, we have to match the columns of U with the columns of T . Note that the updates imply that the columns of U are aligned with the rows of T , so it suffices to match T 's rows with its columns. Using the assumptions that $\alpha/\alpha_0 = \pi$ and $\pi_i \neq \pi_j \forall i \neq j$, we recover the matching by sorting α/α_0 and $T\alpha/\alpha_0$.

We compare the proposed method and the above Variational EM algorithm on synthetic data generated from the HMM in Section 4 under the same parameter setting, except that the number of sets N takes smaller values $\{125, 250, 500, 1000, 2000, 4000\}$. We repeat the experiment 20 times with different random draws from the generative process. Figure 1 gives the relative estimation errors for U (in spectral norm) and P (in entrywise 1-norm) for three methods: Algorithm 1 (tensor), variational EM initialized with the output of Algorithm 1 (tensor+vbEM), and variational EM initialized with 100 random parameter values (rand+vbEM). Clearly, Algorithm 1 outperforms the randomly initialized variation EM, and there is barely any improvement resulting from combining the two methods, except when N is very small. In terms of computational efficiency, we observe that Algorithm 1 is orders of magnitude faster than the variational EM algorithm. On our platform with 48 cores (2.3 GHz each) and 512GB of memory, Algorithm 1 takes a couple of hours to finish all 20 experiments, but the variational EM method takes days.

References

- [1] A. Anandkumar, D. P. Foster, D. Hsu, S. M. Kakade, and Y.-K. Liu. A spectral algorithm for latent dirichlet allocation. *arXiv preprint arXiv:1204.6703v4*, 2013.
- [2] A. Anandkumar, R. Ge, D. Hsu, S. M. Kakade, and M. Telgarsky. Tensor decompositions for learning latent variable models. *arXiv preprint arXiv:1210.7559v2*, 2012.
- [3] G. Stewart. On the perturbation of pseudo-inverses, projections and linear least squares problems. *SIAM review*, 19(4):634–662, 1977.
- [4] G. W. Stewart and J.-g. Sun. Matrix perturbation theory. 1990.