

ADAPTIVE LOW-NONNEGATIVE-RANK APPROXIMATION FOR STATE AGGREGATION OF MARKOV CHAINS*

Yaqi Duan[†], Mengdi Wang[†], Zaiwen Wen[‡], and Yaxiang Yuan[§]

Abstract. This paper develops a low-nonnegative-rank approximation method to identify the state aggregation structure of a finite-state Markov chain under an assumption that the state space can be mapped into a handful of metastates. The number of metastates is characterized by the nonnegative rank of the Markov transition matrix. Motivated by the success of the nuclear norm relaxation in low-rank minimization problems, we propose an atomic regularizer as a convex surrogate for the nonnegative rank and formulate a convex optimization problem. Because the atomic regularizer itself is not computationally tractable, we instead solve a sequence of problems involving a nonnegative factorization of the Markov transition matrices by using the proximal alternating linearized minimization method. Two methods for adjusting the rank of factorization are developed so that local minima are escaped. One is to append an additional column to the factorized matrices, which can be interpreted as an approximation of a negative subgradient step. The other is to reduce redundant dimensions by means of linear combinations. Overall, the proposed algorithm very likely converges to the global solution. The efficiency and statistical properties of our approach are illustrated on synthetic data. We also apply our state aggregation algorithm on a Manhattan transportation data set and make extensive comparisons with an existing method.

Key words. Markov chain, state aggregation, nonnegative matrix factorization, atomic norm, proximal alternating linearized minimization

AMS subject classifications. 65K05, 90C06, 90C40

DOI. 10.1137/18M1220790

1. Introduction. The Markov chain is a basic model of stochastic processes. As a variant of the Markov chain, the Markov decision process lies at the core of dynamic programming and reinforcement learning [2, 3, 21], and has wide applications in engineering systems, operations research, artificial intelligence, and computer games. The present reinforcement learning algorithms may perform poorly if the state space is of huge ambient dimension. Fortunately, empirical experiences tell us that stochastic systems in the real world can usually be characterized by a handful of key features. Therefore, we consider developing a computational method to identify reduced-dimension representations of a Markov chain, so that decision problems can be solved efficiently with a compressed state space.

In this paper, we are particularly interested in situations where the stochastic system is driven by a latent Markov chain with a much smaller state space. The states in the latent Markov chain are *soft aggregations* of the states in the original system. To be more specific, suppose that the original stochastic system has a state space

*Received by the editors October 15, 2018; accepted for publication (in revised form) by D. J. Higham November 15, 2019; published electronically February 25, 2020.

<https://doi.org/10.1137/18M1220790>

Funding: The work of the third author was partially supported by National Natural Science Foundation of China grants 11831002, 11421101, and by Beijing Academy of Artificial Intelligence (BAAI).

[†]Department of Electrical Engineering, Princeton University, Princeton, NJ 08544 (yaqid@princeton.edu, mengdiw@princeton.edu).

[‡]Beijing International Center for Mathematical Research, Center for Data Science, National Engineering Laboratory for Big Data Analysis and Applications, Peking University, China (wenzw@pku.edu.cn).

[§]State Key Laboratory of Scientific and Engineering Computing, Academy of Mathematics and Systems Science, Chinese Academy of Sciences, China (yyx@lsec.cc.ac.cn).

$\mathcal{S} = \{s_1, s_2, \dots, s_d\}$ with ambient dimension d , and the latent Markov chain has a state space $\tilde{\mathcal{S}} = \{\tilde{s}_1, \tilde{s}_2, \dots, \tilde{s}_r\}$ with intrinsic dimension r , where $r \ll d$. If the system is now at $s_i \in \mathcal{S}$, it first chooses a metastate in $\tilde{\mathcal{S}}$ according to some multinomial distribution $[u_{i1}, u_{i2}, \dots, u_{ir}]$, i.e., \tilde{s}_k is selected with probability u_{ik} . Governed by a transition matrix $\tilde{P} \in \mathbb{R}^{r \times r}$, a one-step transition in $\tilde{\mathcal{S}}$ is then conducted, resulting in a new metastate \tilde{s}_l with probability \tilde{p}_{kl} . From there we finally come back to \mathcal{S} according to another multinomial distribution $[v_{1l}, v_{2l}, \dots, v_{dl}]$. In this way, the transition probabilities of the original system can be expressed by

$$p_{ij} = \sum_{k,l=1}^r u_{ik} \tilde{p}_{kl} v_{jl} \quad \forall s_i, s_j \in \mathcal{S}.$$

In a more compact matrix form,

$$(1.1) \quad P = U \tilde{P} V^T, \quad U, V \in \mathbb{R}^{d \times r}, \tilde{P} \in \mathbb{R}^{r \times r}.$$

Note that the decomposition is not unique.

Based on the state aggregation model above, our goal is then to recover U , V , and \tilde{P} of a Markov chain (up to linear transformation). In many real applications, however, the exact dynamics of the system is never revealed to us, and we only have access to realized trajectories of state-to-state transitions. Our proposed algorithm just requires an approximation of the true probability transition matrix P as input.

1.1. Literature review. Recently, there has been interest in the compression of Markov chains. In [26, 17], the authors propose methods for low-rank approximations of Markov chains and provide theoretical guarantees. It is proved in [26] that the spectral method performs well in the recovery of hard state aggregation structures for “lumpable complex networks.” In this special case, the group membership indicator vectors are mutually orthogonal, so that the singular value decomposition (SVD) coincides with the state aggregation structure. However, for a more general soft state aggregation model (1.1), the orthogonality gets easily violated, because the probability matrices U and V in (1.1) are nonnegative and often dense. This makes spectral methods no longer applicable, and poses additional difficulties in algorithm design and theoretical analysis. A nonconvex estimator based on rank-constrained likelihood maximization is proposed in [17] where a difference of convex functions programming algorithm is used to handle the rank-constrained nonsmooth optimization problem. Statistical upper bounds and convergence results are provided.

As is implied by (1.1), our target matrix P has a low-rank and nonnegative factorization. Recovery of low-rank matrices is intensively studied in the past decade, and matrix completion (MC) is a canonical problem closely related to our work. MC aims to recover a low-rank matrix based on only a fraction of its entries. It is often formulated as finding the matrix with smallest rank satisfying a collection of linear constraints. However, such a nonconvex formulation is NP-hard in general, and the convex surrogate using nuclear norm is a remedy. It is proved by [6, 20] that under mild conditions, the nuclear norm relaxation identifies the exact solution to MC. As is shown in [7], the success of the nuclear norm can be interpreted from the perspective of atomic norms. It inspires us to propose a convex heuristic for state aggregation using similar techniques introduced in [7].

Our treatment of a problem with a low-rank solution through matrix factorization has similar versions in the context of semidefinite programming [12, 5], MC [24, 13],

nuclear norm minimization [19], tensor factorization [11], etc. The scheme presented in our paper is distinct from others since it involves a nonnegative matrix factorization (NMF). The NP-hardness [23] and nonconvexity of NMF make it hard for global optimization, and strategies have been proposed to find local solutions. The algorithms in [16, 18] use alternating (projected) gradient descending minimization with carefully chosen step sizes, while methods in [14, 15] are developed under the alternating nonnegative least square framework. In our state aggregation problem, we adopt another alternating direction method, namely, the proximal alternating linearized method (PALM), whose convergence to critical points is proved in [4].

1.2. Scope of this paper. In this work, we develop a low-nonnegative-rank approximation method to identify an underlying state aggregation structure based on trajectories of a Markov chain. After constructing an empirical transition matrix, we compress the state space via an optimization problem regularized by the nonnegative rank [9], which is similar to the usual notion of rank but has nonnegative requirements on the factorization. Due to the combinatorial nature of nonnegative rank, the regularized optimization problem is nonconvex and hard to solve. Therefore, we propose an atomic regularizer as a convex surrogate. The atomic regularizer is an extension of the atomic norm in [7]. It achieves empirical successes in maintaining low-nonnegative-rank structure, and substantially reduces recovery errors compared to the vanilla empirical transition matrix. We also note that, our atomic regularizer is quite general and not restricted to stochastic matrices. Hence, it can be of independent interest in NMF studies.

For the sake of an efficient algorithm, we reform the convexified problem into a factorized optimization model in terms of $U, V \in \mathbb{R}_+^{d \times s}$, where $X = UV^T$ is an NMF of the variable $X \in \mathbb{R}^{d \times d}$ and s refers to the rank that needs to be adapted. Our factorized problem is more complicated than the usual NMF due to the extra linear constraints on U and V so that $X = UV^T$ is still a Markov transition matrix. After investigating the first-order optimality conditions, we prove that a local minimizer of the factorized optimization model provides a solution to the convex problem if and only if it satisfies a certificate derived from the convex model. Based on this observation, we devise a scheme that involves minimizing a sequence of factorized optimization models until the resulting local minimizer represents a global solution. We develop two strategies for adjusting the rank s so that local minima are escaped. One is to append an additional column to U and V simultaneously, which can be interpreted as an approximation of a negative subgradient step. The other is to reduce the dimensions of U and V together by means of linear combinations so that the rank can be reduced to a proper range. The convergence of our proposed algorithm is very likely ensured since these two methods of rank adjustment guarantee monotone improvement of the objective function. We further introduce a possible application of our state aggregation model on Markov chain simulation. By leveraging the factorization structure, one only needs to generate a Markov chain whose number of states equals the rank s . It largely reduces computational costs.

We illustrate our new approach on synthetic data whose ground-truth solution is known. In situations without sampling noise, the algorithm stops exactly at the intrinsic dimension of the system. Experiments are also conducted on simulated trajectories, where we investigate how the ambient, intrinsic dimensions and the length of trajectory affect the recovery of P^* . We finally apply our model to analyze a Manhattan transportation data set. Compared with an existing method in [25], our zoning results are more aligned to the geometric location without knowing it in advance.

1.3. Outline. The remainder of the paper is organized as follows. In section 2, we formulate state aggregation into an optimization problem regularized by nonnegative rank. In section 3, we introduce the atomic regularizer as a convex surrogate for the nonnegative rank, and derive necessary and sufficient conditions for the global optimality. In section 4, we design a factorized optimization model and look at its KKT conditions that govern the local solutions. In section 5, we develop strategies to escape from local minima and establish an adaptive rank algorithm. In section 6, we propose a method to speed up the simulation of Markov chain. Numerical results are presented in section 7.

1.4. Notation. For a column vector $\mathbf{x} \in \mathbb{R}^p$, we denote its i th entry as x_i , the Euclidean norm as $\|\mathbf{x}\|_2$, and the ℓ_∞ -norm as $\|\mathbf{x}\|_\infty := \max_{1 \leq i \leq p} |x_i|$. For a matrix $A \in \mathbb{R}^{p \times q}$, we denote $(A)_{ij}$ or a_{ij} as the entry in the i th row and j th column of A . We denote $A_j \in \mathbb{R}^p$ as the j th column of A , and $A^i \in \mathbb{R}^q$ as the transpose of the i th row. The inner product of two matrices $A, B \in \mathbb{R}^{p \times q}$ is represented by $\langle A, B \rangle = \sum_{i,j} a_{ij} b_{ij}$. If $A \in \mathbb{R}^{p \times q}$ is entrywisely nonnegative, we write that $A \geq 0$ or $A \in \mathbb{R}_+^{p \times q}$. The nonnegative part of a matrix A is denoted by $[A]_+$, i.e., $([A]_+)_{ij} = \max\{a_{ij}, 0\}$. In this paper, several matrix norms are involved, including the nuclear norm $\|A\|_*$ (the sum of singular values), the spectral norm $\|A\|_2$ (the largest singular value), the Frobenius norm $\|A\|_F = \sqrt{\sum_{i,j} a_{ij}^2}$, and the ℓ_1 -norm $\|A\|_{\ell_1} = \sum_{i,j} |a_{ij}|$. We denote the k th largest singular value of a matrix A by $\sigma_k(A)$. Moreover, $\mathbf{1}_p$ and $\mathbf{0}_p$ are the p -dimensional all-one and all-zero column vectors, respectively. We define $\mathbf{e}_i \in \mathbb{R}^p$ as the vector whose i th entry equals one and the other entries are zeros. $\mathbf{0}_{p \times q}$ represents the p -by- q zero matrix. For a vector $\mathbf{x} \in \mathbb{R}^p$, \mathbf{x}^\perp stands for the orthogonal complement of the subspace spanned by \mathbf{x} , i.e., $\mathbf{x}^\perp = \{\mathbf{y} \in \mathbb{R}^p \mid \mathbf{y}^T \mathbf{x} = 0\}$. We denote by $\text{diag}\{\mathbf{x}\}$ or $\text{diag}\{x_i\}_{i=1}^p$ the p -by- p diagonal matrix with $\mathbf{x} \in \mathbb{R}^p$ on its diagonal. For an index set $\mathcal{I} = \{j_1, j_2, \dots, j_k\} \subseteq \{1, 2, \dots, q\}$, we define a k -by- k diagonal matrix $\text{diag}\{x_i\}_{i \in \mathcal{I}} := \text{diag}\{x_{j_i}\}_{i=1}^k$, and $A_{\mathcal{I}} := [A_{j_1}, A_{j_2}, \dots, A_{j_k}]$, the matrix composed of the j_1, j_2, \dots, j_k th columns in A . The convex hull of a set $\mathcal{C} \subseteq \mathcal{R}^p$ is written as $\text{conv}(\mathcal{C})$.

2. Problem setup. Consider a discrete-time finite-state Markov chain with unknown transition matrix. Suppose that we only have access to observations of state-to-state trajectories. In this work, we are concerned about computational methods to identify state aggregation structures of the Markov chain.

2.1. Markov chains. Suppose that a discrete-time Markov chain has a state space $\mathcal{S} = \{s_1, s_2, \dots, s_d\}$ and is driven by an unknown probability transition matrix $P^* \in \mathbb{R}^{d \times d}$. Let (Y_0, Y_1, Y_2, \dots) be a trajectory generated by the Markov chain. By definition,

$$\mathbb{P}[Y_t = s_j \mid Y_{t-1} = s_i, Y_{t-2}, \dots, Y_0] = p_{ij}^*, \quad t = 1, 2, \dots$$

In our problem, we are given an observed trajectory (Y_0, Y_1, \dots, Y_n) of length $(n+1)$. Based on this sample trajectory, we want to have an initial estimate of the transition probabilities. We assume that the Markov chain is ergodic and has a unique stationary distribution $\xi^* \in \mathbb{R}^d$. One can estimate ξ^* using an empirical stationary distribution $\hat{\xi}^{(n)}$ given by

$$(2.1) \quad \hat{\xi}_j^{(n)} := \frac{1}{n} \sum_{t=1}^n \mathbb{1}_{\{Y_t = s_j\}}, \quad j = 1, 2, \dots, d,$$

where $\mathbb{1}_{\{\cdot\}}$ is an indicator function that takes the value 1 if the subscript event happens or takes 0 otherwise. Similarly, an empirical probability transition matrix $\hat{P}^{(n)}$ is formulated as

$$(2.2) \quad \hat{p}_{ij}^{(n)} := \begin{cases} \frac{\sum_{t=1}^n \mathbb{1}_{\{Y_{t-1}=s_i, Y_t=s_j\}}}{\sum_{t=1}^n \mathbb{1}_{\{Y_{t-1}=s_i\}}} & \text{if } \sum_{t=1}^n \mathbb{1}_{\{Y_{t-1}=s_i\}} \geq 1, \\ \frac{1}{d} & \text{if } \sum_{t=1}^n \mathbb{1}_{\{Y_{t-1}=s_i\}} = 0, \end{cases}$$

which is also the maximum likelihood estimator of P^* [1]. Asymptotically, the ergodic theorem of the Markov chain ensures that $\hat{\xi}^{(n)}$ and $\hat{P}^{(n)}$ are strongly consistent. Finite sampling error bounds are also available; see [26]. However, due to the curse of dimensionality, a large d makes $\hat{P}^{(n)}$ a poor estimator. We next seek for a method to derive a better estimator from $\hat{P}^{(n)}$.

2.2. State aggregation problem. In this part, we leverage the state aggregation model (1.1) to compress the state space of the Markov chain, following the idea of [26]. Based on this additional structure, we expect to recover the ground-truth transition matrix P^* with higher accuracy, compared with the initial estimator $\hat{P}^{(n)}$. A more rigorous definition of the state aggregation structure is shown in Definition 1.

DEFINITION 1 (see [26]). *A d -state Markov chain generated by a probability transition matrix $P^* \in \mathbb{R}^{d \times d}$ admits a state aggregation structure with intrinsic dimension r , if there exist matrices $U^* \in \mathcal{U}^{d \times r} := \{U \in \mathbb{R}_+^{d \times r} \mid U \mathbf{1}_r = \mathbf{1}_d\}$, $V^* \in \mathcal{V}^{d \times r} := \{V \in \mathbb{R}_+^{d \times r} \mid V^T \mathbf{1}_d = \mathbf{1}_r\}$ such that*

$$(2.3) \quad P^* = U^* (V^*)^T.$$

The entries of U^* and V^* refer to aggregation and disaggregation probabilities, respectively.

Definition 1 refines model (1.1). We notice that in expression (1.1), the design of a latent Markov chain \tilde{P} is flexible. Without loss of generality, one can simply take \tilde{P} to be the identity matrix I_r . Therefore, we merge \tilde{P} and U (or V^T) together in (2.3). It is important to note that U^* and V^* here are still not unique. One can only expect to recover the state aggregation structure up to linear transformation.

According to Definition 1, any d -by- d stochastic matrix trivially admits a state aggregation structure with intrinsic dimension d , whereas in practical applications it is desirable to have $r \ll d$. The constraints $U^* \in \mathcal{U}^{d \times r}$ and $V^* \in \mathcal{V}^{d \times r}$ in Definition 1 are proposed so that the aggregation and disaggregation probabilities are meaningful, which leads to an NMF structure in (2.3). In the following, we will denote $\mathcal{U}^{1 \times q}$ by \mathcal{U}^q and $\mathcal{V}^{p \times 1}$ by \mathcal{V}^p for simplicity. If $p = q$, the elements in $\mathcal{U}^{p \times q}$ are stochastic matrices.

In our problem, we assume that the system is implicitly driven by an r -state latent Markov chain as described in section 1 (see also [26, Figure 1]), where $r \ll d$ and is unknown. A state aggregation structure (2.3) with intrinsic dimension r is then taken as an underlying assumption. For the convenience of discussions, we introduce the notion of *nonnegative rank*,

$$\text{rank}_+(A) := \min \{m \mid A = BC^T, B \in \mathbb{R}_+^{d \times m}, C \in \mathbb{R}_+^{d \times m}\} \quad \forall A \in \mathbb{R}^{d \times d}.$$

It can be naturally derived from Definition 1 that the ground-truth transition matrix P^* in our problem satisfies $\text{rank}_+(P^*) \leq r$. We will next propose an optimization

problem to conduct state aggregation by controlling the nonnegative rank of the estimator.

Our goal here is to identify a better estimator $X \in \mathcal{U}^{d \times d}$ of P^* based on the empirical $\hat{P}^{(n)}$. To measure the discrepancy between X and $\hat{P}^{(n)}$, we define

$$(2.4) \quad g(X) := \frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - X) \right\|_F^2,$$

where $\hat{\Xi} = \text{diag}\{\hat{\xi}^{(n)}\}$ is a scaling matrix that assigns more weight to states that have occurred more frequently. An optimization problem is then formulated as

$$(2.5) \quad \text{minimize}_{X \in \mathbb{R}^{d \times d}} \quad g(X) + \chi_{\mathcal{E}}(X) + \lambda \text{rank}_+(X),$$

where $\chi_{\mathcal{E}}$ is a *characteristic function* defined as

$$(2.6) \quad \chi_{\mathcal{E}}(X) := \begin{cases} 0, & X \in \mathcal{E}, \\ +\infty, & \text{otherwise,} \end{cases} \quad \forall X \in \mathbb{R}^{d \times d},$$

$\mathcal{E} := \{X \in \mathbb{R}^{d \times d} \mid X\mathbf{1}_d = \mathbf{1}_d\}$ is the set of row-normalized matrices, and λ refers to the regularization parameter. In (2.5), $g(X)$ indicates the fidelity of data. The implicit constraints $\chi_{\mathcal{E}}(X) < +\infty$ and $\text{rank}_+(X) < +\infty$ imply that $X\mathbf{1}_d = \mathbf{1}_d$ and $X \geq 0$, forcing X to be a stochastic matrix. $\lambda \text{rank}_+(X)$ is regarded as a regularization term that enforces the low-nonnegative-rank property. When λ gets larger, $\text{rank}_+(X)$ tends to be smaller, so that the degree of aggregation increases.

Note that a low-nonnegative-rank solution \hat{X} to problem (2.5) helps recover a state aggregation structure (2.3). In fact, if $\hat{X} \in \mathbb{R}^{d \times d}$ satisfies $\text{rank}_+(\hat{X}) = s \leq d$, then there exists a factorization $\hat{X} = \hat{U}\hat{V}^T$ for some $\hat{U}, \hat{V} \in \mathbb{R}_+^{d \times s}$. Without loss of generality, we assume that \hat{V} does not have a zero column. Taking $\hat{U} := \hat{U} \text{diag}\{\mathbf{1}_d^T \hat{V}\}$ and $\hat{V} := \hat{V}(\text{diag}\{\mathbf{1}_d^T \hat{V}\})^{-1}$, we have $\hat{V}^T \mathbf{1}_d = \mathbf{1}_s$ and $\hat{U} \mathbf{1}_s = \hat{U}(\hat{V}^T \mathbf{1}_d) = \hat{X} \mathbf{1}_d = \mathbf{1}_d$. Hence, $\hat{U} \in \mathcal{U}^{d \times s}$, $\hat{V} \in \mathcal{V}^{d \times s}$, and the factorization $\hat{X} = \hat{U}\hat{V}^T$ identifies a state aggregation structure of \hat{X} .

Due to the combinatorial nature of the nonnegative rank, it is hard to solve (2.5) directly. In the following sections, we will develop computation tools for solving (2.5) with theoretical guarantees.

3. Convexified formulation via atomic regularizers. In this section, we look for a convex surrogate for rank_+ that can be optimized efficiently and in practice yields desirable solutions with low nonnegative rank. We propose an atomic regularizer, inspired by the concept of atomic norm relaxation in [7]. First-order optimality conditions are also developed for the convexified problem.

3.1. Atomic regularizer. We propose a convex surrogate function of rank_+ within the general framework of atomic norm in [7]. The goal of [7] is to represent a vector $\mathbf{x} \in \mathbb{R}^p$ with a few elementary building blocks which we call *atoms*. To be specific, let \mathcal{A} be an *atomic set*; we want $\mathbf{x} \in \mathbb{R}^p$ to be formed as

$$(3.1) \quad \mathbf{x} = \sum_{i=1}^k c_i \mathbf{a}_i, \quad \mathbf{a}_i \in \mathcal{A}, c_i > 0, \quad i = 1, 2, \dots, k.$$

In MC, for instance, one can take the atomic set to be the collection of unit-spectral-norm rank-one matrices, which we denote by \mathcal{A}_* . The affine rank of a matrix is

then interpreted as the smallest number of atoms in \mathcal{A}_* whose conic combination can represent the matrix.

The main idea of [7] is to recover a structure (3.1) with a small k by using a convex program that minimizes the following gauge function:

$$(3.2) \quad \|\mathbf{x}\|_{\mathcal{A}} := \inf\{t > 0 \mid \mathbf{x} \in t \cdot \text{conv}(\mathcal{A})\} \quad \forall \mathbf{x} \in \mathbb{R}^p.$$

If \mathcal{A} is centrally symmetric around the origin, then $\|\cdot\|_{\mathcal{A}}$ is indeed a norm and is called an *atomic norm*. In MC, the atomic set \mathcal{A}_* induces the nuclear norm $\|\cdot\|_*$ [7], which we usually take as a convex surrogate for affine rank. The nuclear norm relaxation consistently yields promising results both theoretically and experimentally [6, 20]. It motivates us to propose a convex relaxation of rank_+ in an analogous way.

It is important to note that the nuclear norm is not suitable for our problem, since it does not guarantee nonnegativity while (2.3) actually requires an NMF. Therefore, we build the following atomic set using nonnegative atoms in \mathcal{A}_* :

$$(3.3) \quad \mathcal{A}_+ := \{\mathbf{u}\mathbf{v}^T \mid \mathbf{u}, \mathbf{v} \in \mathbb{R}_+^d, \|\mathbf{u}\|_2 = 1, \|\mathbf{v}\|_2 = 1\}.$$

In this way, the state aggregation problem (2.5) can be interpreted as approximating $\hat{P}^{(n)}$ with as few atoms in \mathcal{A}_+ as possible. \mathcal{A}_+ is not centrally symmetric. However, we can still derive from \mathcal{A}_+ a convex relaxation of rank_+ that is analogous to (3.2),

$$(3.4) \quad \begin{aligned} \Omega(X) &:= \inf\{t > 0 \mid X \in t \cdot \text{conv}(\mathcal{A}_+)\} \\ &= \inf\left\{\sum_{j=1}^s \|U_j\|_2 \|V_j\|_2 \mid X = UV^T \text{ with } U, V \in \mathbb{R}_+^{d \times s}\right\} \quad \forall X \in \mathbb{R}^{d \times d}. \end{aligned}$$

Since Ω is no longer a norm, we only call it an *atomic regularizer*. For an arbitrary $X \in \mathbb{R}^{d \times d}$, if there exists a factorization $X = UV^T$ with $U, V \in \mathbb{R}_+^{d \times s}$ that achieves the infimum in the definition of Ω , i.e., $\Omega(X) = \sum_{j=1}^s \|U_j\|_2 \|V_j\|_2$, then we say it is an *optimal factorization (with respect to Ω)*.

Replacing $\text{rank}_+(X)$ in (2.5) by the atomic regularizer $\Omega(X)$ yields

$$(3.5) \quad \text{minimize}_{X \in \mathbb{R}^{d \times d}} \quad f_{\lambda}(X) := g(X) + \chi_{\mathcal{E}}(X) + \lambda \Omega(X),$$

which is now a convex optimization problem.

Suppose that $\hat{X} \in \mathbb{R}^{d \times d}$ is a global optimal solution to (3.5). The following Theorem 2 shows that when λ is properly chosen, the convexified problem (3.5) produces a consistent estimator of the ground truth P^* . The deviation bound matches the rate of the central limit theorem (ignoring a logarithmic factor). See Appendix A.1 for the proof.

THEOREM 2 (consistency and convergence rate of \hat{X}). *If we take*

$$\lambda \geq [\Omega(P^*)]^{-1} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F^2,$$

then

$$(3.6) \quad \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F \leq (1 + \sqrt{3}) \sqrt{\lambda \Omega(P^*)}.$$

Additionally, when $\lambda = \gamma \cdot n^{-1} [\log(n)]^2$ for some constant $\gamma > 0$, there exists $C > 0$ such that if $n \geq C$, then with probability at least $1 - n^{-1}$,

$$(3.7) \quad \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F \leq C n^{-1/2} \log(n),$$

where the constant C depends on γ and the ground truth P^ .*

3.2. Global optimality conditions. Since (3.5) is convex over $\mathbb{R}^{d \times d}$, its global solutions can be characterized by first-order information. The objective function f_λ consists of three parts, among which the discrepancy term $g(X)$ is differentiable and the remaining two terms $\chi_\varepsilon(X)$ and $\Omega(X)$ have subdifferentials with explicit form. Next we derive expressions of $\partial\chi_\varepsilon(X)$ and $\partial\Omega(X)$, as well as the sufficient and necessary conditions for global optimality.

The following two lemmas are concerned with $\partial\chi_\varepsilon(X)$ and $\partial\Omega(X)$, whose proofs can be found in Appendices A.2 and A.3.

LEMMA 3. *The subdifferential of characteristic function $\chi_\varepsilon(X)$ is*

$$(3.8) \quad \partial\chi_\varepsilon(X) = \{\mu \mathbf{1}_d^T \mid \mu \in \mathbb{R}^d\}.$$

LEMMA 4. *The subdifferential of the atomic regularizer $\Omega(X)$ is*

$$(3.9) \quad \partial\Omega(X) = \{W \mid \Omega^\circ(W) \leq 1, \langle W, X \rangle = \Omega(X)\},$$

where $\Omega^\circ(W)$ is the support function given by

$$(3.10) \quad \Omega^\circ(W) := \sup_{Z: \Omega(Z) \leq 1} \langle W, Z \rangle.$$

Based on Lemmas 3 and 4, we now establish the sufficient and necessary conditions for global optimality in Theorem 5. The proof can be found in Appendix A.4.

THEOREM 5 (sufficient and necessary conditions for global optimality of (3.5)). *Suppose that $\hat{X} \in \mathbb{R}^{d \times d}$ is factorized as $\hat{X} = \hat{U}\hat{V}^T$, where $\hat{U} \in \mathcal{U}^{d \times s}$, $\hat{V} \in \mathcal{V}^{d \times s}$ for some s , and \hat{U} does not have a zero column. Then \hat{X} is globally optimal for problem (3.5) with an optimal factorization $\hat{X} = \hat{U}\hat{V}^T$ if and only if there exists $\mu \in \mathbb{R}^d$ such that*

$$(3.11a) \quad \left\{ \begin{array}{l} \mathbf{u}^T (\mu \mathbf{1}_d^T - \nabla g(\hat{X})) \mathbf{v} \leq \lambda \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}_+^d \text{ s.t. } \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1, \\ \end{array} \right.$$

$$(3.11b) \quad \left\{ \begin{array}{l} [\mu \mathbf{1}_s^T - \nabla g(\hat{X}) \hat{V}]_+ = \lambda \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^s, \\ \end{array} \right.$$

$$(3.11c) \quad \left\{ \begin{array}{l} [\mathbf{1}_d \mu^T \hat{U} - (\nabla g(\hat{X}))^T \hat{U}]_+ = \lambda \hat{V} \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^s. \\ \end{array} \right.$$

The vector μ is unique if it exists.

In (3.11a), $\mathbf{u}^T (\mu \mathbf{1}_d^T - \nabla g(\hat{X})) \mathbf{v}$ approximates the decrease of $g(X) + \chi_\varepsilon(X)$ in direction $\mathbf{u}\mathbf{v}^T$ at \hat{X} . Intuitively, the inequality means that for a small perturbation along $\mathbf{u}\mathbf{v}^T$, this decrease is dominated by the increase in $\lambda\Omega(X)$. When $(\mathbf{u}, \mathbf{v}) = (\hat{U}_j, \hat{V}_j)$ for $j = 1, 2, \dots, s$, (3.11b) and (3.11c) suggest that the inequality (3.11a) holds as equality. Conditions (3.11a)–(3.11c) provide a certification of global optimality. In order to determine whether $\hat{X} = \hat{U}\hat{V}^T$ is a global solution to (3.5), one can first calculate a vector μ according to (3.11b) and (3.11c), and next verify (3.11a) with the resulting μ .

4. Factorized optimization model. In the previous section, we proposed the convex formulation (3.5) to estimate the state aggregation structure from empirical data. However, due to the implicit form of the atomic regularizer $\Omega(X)$, it is impractical to directly minimize the function $f_\lambda(X)$ over the matrix X . Letting $X = UV^T$, we reformulate (3.5) into a factorized optimization model. Then we investigate its equivalence to the convex model and its KKT conditions.

4.1. Factorized optimization model and its equivalence to convexified model. Although the nonnegative rank of the optimal solution \hat{X} to (3.5) is usually unknown, we can pick an arbitrary s and have an initial guess that $\text{rank}_+(\hat{X}) \leq s$. We recast (3.5) by NMF $X = UV^T$, where $(U, V) \in \mathbb{R}_+^{d \times s} \times \mathbb{R}_+^{d \times s}$. The search space can be further narrowed down to $\mathcal{U}^{d \times s} \times \mathcal{V}^{d \times s}$, since each $(U, V) \in \mathbb{R}_+^{d \times s} \times \mathbb{R}_+^{d \times s}$ with $UV^T \in \mathcal{U}^{d \times d}$ can be mapped into

$$\left(U \text{diag}\{V^T \mathbf{1}_d\}, V (\text{diag}\{V^T \mathbf{1}_d\})^{-1} \right) \in \mathcal{U}^{d \times s} \times \mathcal{V}^{d \times s}$$

under the assumption that V does not have a zero column. This transform does not change the product $X = UV^T$.

We formulate a factorized optimization model as follows:

$$(4.1) \quad \text{minimize}_{U \in \mathcal{U}^{d \times s}, V \in \mathcal{V}^{d \times s}} F_\lambda(U, V) = g(UV^T) + \lambda \sum_{j=1}^s \|U_j\|_2 \|V_j\|_2,$$

where s refers to the rank of the model and is a parameter to adjust. Intuitively, the problem (4.1) not only optimizes $f_\lambda(X)$ but also identifies an optimal factorization of X (with respect to Ω). The equivalence between (4.1) and the convexified problem (3.5) is provided in Lemma 6. See Appendix A.5 for the proof.

LEMMA 6. *When s is sufficiently large, (3.5) and (4.1) are equivalent.*

Lemma 6 indicates that one can solve (3.5) by optimizing a factorized model (4.1) with sufficiently large s . Two strategies for adjusting the rank s are developed in subsection 5.2.

4.2. The KKT conditions of the factorized optimization model. Compared to (3.5), the factorized optimization model (4.1) searches a space of much smaller dimension and provides a tractable way to evaluate the regularization term. However, it is no longer convex.

In order to bridge the gap between a local minimum of (4.1) and a global solution to (3.5), we develop KKT conditions for (4.1) as a counterpart of Theorem 5. The results are summarized in Theorem 7, whose proof can be found in Appendix A.6.

THEOREM 7 (KKT conditions of (4.1)). *Suppose that (\hat{U}, \hat{V}) is a local solution to (4.1), and take $\hat{X} = \hat{U}\hat{V}^T$, $\mathcal{I} = \{j \mid \hat{U}_j \neq 0\}$. Then, there exists $\mu \in \mathbb{R}^d$ such that*

$$(4.2) \quad \begin{cases} \left[\mu \mathbf{1}_{|\mathcal{I}|}^T - \nabla g(\hat{X}) \hat{V}_{\mathcal{I}} \right]_+ = \lambda \hat{U}_{\mathcal{I}} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j \in \mathcal{I}}, \\ \left[\mathbf{1}_d \mu^T \hat{U}_{\mathcal{I}} - \left(\nabla g(\hat{X}) \right)^T \hat{U}_{\mathcal{I}} \right]_+ = \lambda \hat{V}_{\mathcal{I}} \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j \in \mathcal{I}}. \end{cases}$$

The vector μ is unique if existing and refers to the Lagrangian multiplier.

According to (4.2), the existence of a Lagrangian multiplier μ can help determine whether (\hat{U}, \hat{V}) solves (4.1) locally. In addition, given a local solution (\hat{U}, \hat{V}) to (4.1), one readily notices that it satisfies all but one condition in Theorem 5. Comparison between Theorems 5 and 7 therefore provides the following certificate to verify the global optimality of $\hat{X} = \hat{U}\hat{V}^T$.

THEOREM 8 (global optimality certificate). *Suppose that (\hat{U}, \hat{V}) is a local solution to (4.1). Then $\hat{X} = \hat{U}\hat{V}^T$ is globally optimal for (3.5) if and only if*

$$(4.3) \quad \mathbf{u}^T [\mu \mathbf{1}_d^T - \nabla g(\hat{X})] \mathbf{v} \leq \lambda \quad \forall \mathbf{u}, \mathbf{v} \geq 0, \quad \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1,$$

where μ is the Lagrangian multiplier defined in Theorem 7.

5. An adaptive rank estimation algorithm. In this section, we propose an adaptive method that solves a sequence of (4.1) while adjusting values of s until reaching the global optimality condition. Details of the algorithm and convergence properties are provided below.

5.1. A subroutine: PALM. We first deal with the factorized optimization model for a fixed s . In order to solve PALM [4], we rewrite (4.1) into

$$(5.1) \quad \text{minimize}_{U, V \in \mathbb{R}^{d \times s}} \quad \tilde{F}_\lambda(U, V) := F_\lambda(U, V) + \chi_{\mathcal{U}^{d \times s}}(U) + \chi_{\mathcal{V}^{d \times s}}(V),$$

where $\chi_{\mathcal{U}^{d \times s}}$ and $\chi_{\mathcal{V}^{d \times s}}$ are characteristic functions analogous to $\chi_{\mathcal{E}}$. PALM is a Gauss–Seidel-like method that alternately minimizes the linearized objective function with respect to U and V . In the following, we first develop a method to ensure the differentiability of F_λ in each iteration, and next show that the subproblems of PALM can be solved efficiently. We provide two types of step sizes. One is derived from [4], which guarantees the convergence of PALM, in theory. The other combines the Barzilai–Borwein (BB) step sizes with a nonmonotone line search. We use it in numerical experiments.

Under the PALM framework, we generate a sequence $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ to solve (5.1). In order that PALM is well-defined in each iteration, the objective function F_λ is supposed to be differentiable at each (U^k, V^k) . In other words, it is required that U^k does not have any zero columns. To this end, before the $(k+1)$ th iteration, we remove the columns in U^k of which the Euclidean norms are smaller than ε_0 , and normalize U^k so that each row sums to one. Here, ε_0 is a user-defined constant concerning the computation precision. We choose $\varepsilon_0 = 10^{-14}$ by default. The corresponding columns in V^k are also removed. Due to the constraint $U\mathbf{1}_s = \mathbf{1}_d$, the resulting matrices are not empty, and we denote them by \tilde{U}^k and \tilde{V}^k .

After preprocessing, we apply the following modified PALM scheme:

$$(5.2) \quad \begin{aligned} U^{k+1} &\in \arg \min_U \left\{ \left\langle \nabla_U F_\lambda(\tilde{U}^k, \tilde{V}^k), U - \tilde{U}^k \right\rangle + \chi_{\mathcal{U}^{d \times s}}(U) + \frac{1}{2c_k} \|U - \tilde{U}^k\|_F^2 \right\} \\ &= \text{proj}_{\mathcal{U}^{d \times s}} \left(\tilde{U}^k - c_k \nabla_U F_\lambda(\tilde{U}^k, \tilde{V}^k) \right), \end{aligned}$$

$$(5.3) \quad \begin{aligned} V^{k+1} &\in \arg \min_V \left\{ \left\langle \nabla_V F_\lambda(U^{k+1}, \tilde{V}^k), V - \tilde{V}^k \right\rangle + \chi_{\mathcal{V}^{d \times s}}(V) + \frac{1}{2d_k} \|V - \tilde{V}^k\|_F^2 \right\} \\ &= \text{proj}_{\mathcal{V}^{d \times s}} \left(\tilde{V}^k - d_k \nabla_V F_\lambda(U^{k+1}, \tilde{V}^k) \right), \end{aligned}$$

where c_k and d_k are some suitable step sizes, and

$$\text{proj}_{\mathcal{C}}(A) := \arg \min_{X \in \mathcal{C}} \|X - A\|_F \quad \forall A \in \mathbb{R}^{d \times s}$$

for $\mathcal{C} = \mathcal{U}^{d \times s}$ or $\mathcal{V}^{d \times s}$. One can derive the closed-form solutions of $\text{proj}_{\mathcal{U}^{d \times s}}$ and $\text{proj}_{\mathcal{V}^{d \times s}}$ from Lemma 9.

LEMMA 9 (Theorem 2.2 in [8]). For a vector $\mathbf{y} \in \mathbb{R}^p$, let $y^{(1)} \geq y^{(2)} \geq \dots \geq y^{(p)}$ be a rearrangement of the entries in \mathbf{y} , and take

$$l = \max \left\{ j \left| \sum_{k=1}^{j-1} (y^{(k)} - y^{(j)}) < 1, 1 \leq j \leq p \right. \right\} \quad \text{and} \quad \eta = \frac{1}{l} \left(1 - \sum_{k=1}^l y^{(k)} \right).$$

Then, the projection of \mathbf{y} onto $\mathcal{V}^p = \{\mathbf{x} \in \mathbb{R}_+^p \mid \mathbf{x}^T \mathbf{1}_p = 1\}$ is given by

$$\text{proj}_{\mathcal{V}^p}(\mathbf{y}) = [\mathbf{y} + \eta \mathbf{1}_p]_+.$$

Lemma 9 also provides a numerically efficient way to calculate (5.2) and (5.3).

We now consider the step sizes c_k and d_k suggested by [4], which rely on the Lipschitz smoothness of F_λ . The partial gradients of F_λ are

$$\begin{aligned} \nabla_U F_\lambda(U, V) &= -\hat{\Xi}^2(\hat{P}^{(n)} - UV^T)V + \lambda U \text{diag} \left\{ \frac{\|V_j\|_2}{\|U_j\|_2} \right\}_{j=1}^s, \\ \nabla_V F_\lambda(U, V) &= -(\hat{P}^{(n)} - UV^T)^T \hat{\Xi}^2 U + \lambda V \text{diag} \left\{ \frac{\|U_j\|_2}{\|V_j\|_2} \right\}_{j=1}^s, \end{aligned}$$

if existing. Within the feasible set, $\nabla_U F_\lambda(U, V)$ is Lipschitz continuous with constant $L_1(V) = \|\hat{\Xi}\|_F^2 \|V^T V\|_F + \lambda \varepsilon_0^{-1}$ with respect to U for a fixed V , since

$$\begin{aligned} \|\nabla_U F(U', V) - \nabla_U F(U, V)\|_F &\leq \left\| \hat{\Xi}^2(U' - U)V^T V \right\|_F + \lambda \varepsilon_0^{-1} \|U' - U\|_F \\ &\leq \left(\|\hat{\Xi}\|_F^2 \|V^T V\|_F + \lambda \varepsilon_0^{-1} \right) \|U' - U\|_F \quad \forall \text{ feasible } U', U. \end{aligned}$$

Similarly, $\nabla_V F(U, V)$ is also Lipschitz continuous with constant $L_2(U) = \|U^T \hat{\Xi}^2 U\|_F + \lambda \sqrt{d} \|U\|_F$ with respect to V for a fixed U . By taking

$$(5.4) \quad c_k = (\gamma_1 L_1(V^k))^{-1}, \quad d_k = (\gamma_2 L_2(U^k))^{-1}$$

with $\gamma_1 > 1$, $\gamma_2 > 1$, one can prove that the modified PALM converges theoretically.

In fact, as (5.2) and (5.3) iterate, the rank s monotonically decreases and converges to a positive integer. Therefore, when analyzing asymptotic behaviors, we only need to deal with situations where s is fixed and $\|U_j^k\| \geq \varepsilon_0$ for $j = 1, 2, \dots, s$. Under these conditions, the results in [4] can be extended to our proposed scheme. We summarize the convergence properties in Theorem 10, and defer the proof to Appendix A.9.

THEOREM 10 (corollary of Lemma 3 and Theorem 1 in [4]). Let $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ be a sequence generated by (5.2) and (5.3) iteratively and $\|U_j^k\|_2 \geq \varepsilon_0$ for $j = 1, 2, \dots, s$.

1. (Lemma 3(i) in [4]) The sequence $\{F_\lambda(U^k, V^k)\}_{k \in \mathbb{N}}$ is nonincreasing.
2. (Theorem 1(i) in [4]) The sequence $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ has finite length, which means,

$$\sum_{k=1}^{\infty} (\|U^{k+1} - U^k\|_F + \|V^{k+1} - V^k\|_F) < \infty.$$

3. (Theorem 1(ii) in [4]) The sequence $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ converges to a stationary point of problem (5.1).

Theorem 10 shows that the step sizes in (5.4) guarantee monotone improvement of function value and the global convergence of PALM to a stationary point.

In terms of convergence rate, we notice that \tilde{F}_λ is semialgebraic according to the proof of Theorem 10. Remark 6(ii) in [4] suggests that a semialgebraic function satisfies the Łojasiewicz inequality for some $\theta \in [0, 1)$. In other words, there exist $C_0 > 0$ and $k_0 \in \mathbb{N}$ such that

$$\left(\tilde{F}_\lambda(U^k, V^k) - \tilde{F}_\lambda(\hat{U}, \hat{V})\right)^\theta \leq C_0 \inf \left\{ \|G\|_F \mid G \in \partial \tilde{F}_\lambda(U^k, V^k) \right\} \quad \forall k \geq k_0,$$

where (\hat{U}, \hat{V}) is the limit point of $\{(U^k, V^k)\}_{k \in \mathbb{N}}$. If $\theta \leq 1/2$, $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ converges linearly to (\hat{U}, \hat{V}) . Otherwise, we only have theoretical guarantees for a geometric convergence rate.

Aside from (5.4), we also adjust BB step sizes to our problem and adopt the nonmonotone line search technique to enhance numerical performances. Letting

$$c_{k+1,1} = \frac{|\langle S_U^k, H_U^k \rangle|}{\|H_U^k\|_F^2}, \quad c_{k+1,2} = \frac{\|S_U^k\|_F^2}{|\langle S_U^k, H_U^k \rangle|},$$

where $S_U^k = U^{k+1} - \tilde{U}^k$, $H_U^k = \nabla_U f(U^{k+1}, \tilde{V}^k) - \nabla_U f(\tilde{U}^k, \tilde{V}^k)$, and

$$d_{k+1,1} = \frac{|\langle S_V^k, H_V^k \rangle|}{\|H_V^k\|_F^2}, \quad d_{k+1,2} = \frac{\|S_V^k\|_F^2}{|\langle S_V^k, H_V^k \rangle|},$$

where $S_V^k = V^{k+1} - \tilde{V}^k$, $H_V^k = \nabla_V f(U^{k+1}, V^{k+1}) - \nabla_V f(U^{k+1}, \tilde{V}^k)$, we take

$$(5.5) \quad c_k = c_{k,1}\delta^p \text{ or } c_k = c_{k,2}\delta^p, \quad d_k = d_{k,1}\delta^q \text{ or } d_k = d_{k,2}\delta^q$$

for some $\delta \in (0, 1)$, $p, q \in \mathbb{N}$. Here, p and q are, respectively, the smallest integers such that

$$\begin{aligned} f(U^{k+1}, \tilde{V}^k) &\leq \frac{1}{5} \sum_{t=k-4}^k f(\tilde{U}^t, \tilde{V}^t) - \eta c_k \|\nabla_U f(\tilde{U}^k, \tilde{V}^k)\|_F^2, \\ f(U^{k+1}, V^{k+1}) &\leq \frac{1}{5} \sum_{t=k-4}^k f(U^{t+1}, \tilde{V}^t) - \eta d_k \|\nabla_V f(U^{k+1}, \tilde{V}^k)\|_F^2. \end{aligned}$$

In numerical experiments, we say that the subroutine PALM converges locally if the decrease in function value satisfies

$$(5.6) \quad (\bar{f}^k - f^k) / f^k < 10^{-3}, \quad \bar{f}^k = \frac{1}{30} \sum_{t=k-30}^{k-1} f(U^t, V^t).$$

The BB step sizes in (5.5) appear to be efficient in most of our test problems. For instance, when solving a factorized optimization model with $d = 1000$, $s = 10$, and a randomly generated initial point (U^0, V^0) , our scheme converges within 200 steps in most scenarios. See sections 7.2 and 7.3 for further details of the numerical settings.

5.2. Successive refinements to escape from local solutions. When the subroutine PALM converges to a stationary point, one needs to further refine the solution unless it represents a global minimum of the convexified problem (4.1). We propose numerical criteria for the global optimality and develop two methods to escape from local minima without increasing the function value.

5.2.1. Criteria to determine global optimality. Recall that Theorem 8 provides a certificate (4.3) to determine the global optimality of a local solution (\hat{U}, \hat{V}) . We now consider how to verify (4.3) numerically.

Exact stopping rule: We adopt the gradient projection method to solve the following optimization problem:

$$(5.7) \quad \begin{aligned} \sigma := & \underset{\mathbf{u}, \mathbf{v}}{\text{maximize}} && \mathbf{u}^T \left(\mu \mathbf{1}_d^T - \nabla g(\hat{X}) \right) \mathbf{v} \\ \text{s.t.} &&& \|\mathbf{u}\|_2 = 1, \|\mathbf{v}\|_2 = 1, \\ &&& \mathbf{u} \geq 0, \mathbf{v} \geq 0. \end{aligned}$$

The algorithm for (5.7) usually converges in a few steps, and the iterations can be easily calculated. It is safe to say that compared with the PALM process, the time spent on (5.7) is negligible. As (4.3) suggests, the solution $\hat{X} = \hat{U}\hat{V}^T$ is considered to be globally optimal if $\sigma \leq (1 + \varepsilon_{\text{Exa}})\lambda$, where $\varepsilon_{\text{Exa}} > 0$ is a user-defined threshold concerning the precision of the solution.

In some scenarios, it requires a very large rank s to reach a reasonably small ε_{Exa} , and the objective function F_λ improves slowly as s grows. It is a better choice to stop at a smaller s so that we can compress the state space in a more compact way, even if $\hat{X} = \hat{U}\hat{V}^T$ does not solve (3.5) precisely. To this end, we propose an early stopping rule as follows.

Early stopping rule: We define a function

$$\varphi(\mathbf{v}) = \left\| \left[\left(\mu \mathbf{1}_d^T - \nabla g(\hat{X}) \right) \mathbf{v} \right]_+ \right\|_2$$

over the set $\Pi := \{\mathbf{v} \in \mathbb{R}_+^d \mid \|\mathbf{v}\|_2 = 1\}$. The condition (4.3) means that the superlevel set $L_\lambda = \{\mathbf{v} \in \Pi \mid \varphi(\mathbf{v}) > \lambda\}$ is empty. Assume that (\hat{U}, \hat{V}) is a limit point generated by the modified PALM. If $\hat{U}\hat{V}^T$ is close to a global minimum of (3.5), L_λ is supposed to have small measure due to the continuity of φ . We represent the measure of L_λ by probabilities, and make an estimate using a group of independent and identically distributed test vectors $\{\bar{\mathbf{v}}_k\}_{k=1}^N$ uniformly distributed on Π . In our experiments, $N = 5000$. If each $\bar{\mathbf{v}}_k$ satisfies $\varphi(\bar{\mathbf{v}}_k) \leq \lambda$, we say that (\hat{U}, \hat{V}) identifies a global solution. This criterion is a relaxation of condition (4.3) and often results in an early stop in the scheme. If starting from an initial point (U^0, V^0) with a small s , we can usually obtain a low-rank solution. Since $\hat{X} = \hat{U}\hat{V}^T$ is only required to lie in the neighborhood of the genuine minimum, this approximate solution is expected to be more robust to sampling noises and work better in real-world problems.

These two stopping rules are applicable to different situations. When investigating the properties of the state aggregation model (3.5), we prefer the exact stopping rule so that the problem can be solved more accurately. When dealing with real-world applications, the early stopping rule helps identify approximate solutions with much lower nonnegative rank, and compresses the state space more efficiently.

5.2.2. Appending a new column. When condition (4.3) fails to hold, we need to refine the local solution (\hat{U}, \hat{V}) so that PALM can resume from a new initial point. Either by the exact or early stopping rule, one can find a pair of vectors $\bar{\mathbf{u}}, \bar{\mathbf{v}} \in \mathbb{R}^d$ such that

$$(5.8) \quad \bar{\mathbf{u}}^T \left(\mu \mathbf{1}_d^T - \nabla g(\hat{X}) \right) \bar{\mathbf{v}} > \lambda, \quad \bar{\mathbf{u}}, \bar{\mathbf{v}} \geq 0, \quad \|\bar{\mathbf{u}}\|_2 = \|\bar{\mathbf{v}}\|_2 = 1.$$

Intuitively, $\bar{\mathbf{u}}\bar{\mathbf{v}}^T$ approximates the negative subgradient directions of f_λ at $\hat{X} = \hat{U}\hat{V}^T$. In fact, for any $W \in \partial\Omega(\hat{X})$, the matrix $Z = \nabla g(\hat{X}) - \mu\mathbf{1}_d^T + \lambda W \in \partial f_\lambda(\hat{X})$. According to Lemma 3.9, $\Omega^\circ(W) \leq 1$. It implies that

$$\langle Z, \bar{\mathbf{u}}\bar{\mathbf{v}}^T \rangle = \lambda \bar{\mathbf{u}}^T W \bar{\mathbf{v}} - \bar{\mathbf{u}}^T (\mu\mathbf{1}_d^T - \nabla g(\hat{X})) \bar{\mathbf{v}} \leq \lambda - \bar{\mathbf{u}}^T (\mu\mathbf{1}_d^T - \nabla g(\hat{X})) \bar{\mathbf{v}} < 0.$$

Inspired by this, we consider appending scaled $\bar{\mathbf{u}}$ and $\bar{\mathbf{v}}$ as new columns to \hat{U} and \hat{V} , respectively, and propose (5.9) to construct a better solution. See Theorem 11 for details, whose proof can be found in Appendix A.8.

THEOREM 11 (escaping local minima). *Suppose that (\hat{U}, \hat{V}) is a local solution to (4.1), where \hat{U} does not have a zero column, and $\hat{X} = \hat{U}\hat{V}^T$ is not globally optimal for (3.5). According to Theorem 8, there exist $\bar{\mathbf{u}}, \bar{\mathbf{v}} \in \mathbb{R}^d$ satisfying (5.8). Take*

$$(5.9) \quad \bar{U} := [\text{diag}\{\mathbf{1}_d - \kappa\bar{\mathbf{u}}\} \hat{U} \quad \kappa\bar{\mathbf{u}}], \quad \bar{V} := [\hat{V} \quad (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{v}}]$$

for some sufficiently small $\kappa > 0$. Then, (\bar{U}, \bar{V}) is feasible for problem (4.1) and

$$F_\lambda(\bar{U}, \bar{V}) < F_\lambda(\hat{U}, \hat{V}).$$

Theorem 11 suggests that as long as (\hat{U}, \hat{V}) is not globally optimal, one can always reduce the global objective value of (3.5) by appending columns. This will guarantee monotone improvement of the solutions.

In numerical experiments, we adopt backtracking line search to determine the appropriate κ in (5.9). We take $\kappa = 0.5^p \|\bar{\mathbf{u}}\|_\infty^{-1}$, where $p \in \mathbb{N}$ is the smallest integer such that $F_\lambda(\bar{U}, \bar{V}) < (1 - 10^{-5})F_\lambda(\hat{U}, \hat{V})$. If $\kappa < 10^{-8}$, we say that the global optimality is achieved and terminate the algorithm.

5.2.3. Removing redundant dimensions. Suppose that $(\hat{U}, \hat{V}) \in \mathcal{U}^{d \times s} \times \mathcal{V}^{d \times s}$ is a local solution to (4.1). In the case where $\{\hat{U}_j \hat{V}_j^T\}_{j=1}^s$ are linearly dependent, the linear combination $\hat{X} = \hat{U}\hat{V}^T = \sum_{j=1}^s \hat{U}_j \hat{V}_j^T$ can be expressed using less than s nonnegative rank-one matrices. In other words, \hat{X} actually admits a state aggregation structure with less than s metastates. To this end, we hope to identify a factorization $\hat{X} = \tilde{U}\tilde{V}^T$ with rank smaller than s , and (\tilde{U}, \tilde{V}) preserves the objective value.

Recall that $F_\lambda(U, V) = g(UV^T) + \lambda \sum_j \|U_j\|_2 \|V_j\|_2$. Since the product $\tilde{U}\tilde{V}^T = \hat{U}\hat{V}^T$, the first term $g(UV^T)$ is unchanged after the adjustment. As for the second term $\sum_{j=1}^s \|U_j\|_2 \|V_j\|_2$, the following corollary of Theorem 7 suggests that it remains the same under some linear combinations of the rank-one matrices. See Appendix A.7 for the proof.

COROLLARY 12 (of Theorem 7). *Suppose that $(\hat{U}, \hat{V}) \in \mathcal{U}^{d \times s} \times \mathcal{V}^{d \times s}$ is a local solution to (4.1) and there exists $\alpha \in \mathbb{R}^s$ such that $\sum_{j=1}^s \alpha_j \hat{U}_j \hat{V}_j^T = \mathbf{0}_{d \times d}$. Then,*

$$\sum_{j=1}^s \alpha_j \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 = 0.$$

Corollary 12 inspires us to get rid of the “redundant dimensions” in (\hat{U}, \hat{V}) by deleting carefully selected columns and rescaling the rest. To be specific, suppose that the linear equation $\sum_{j=1}^s \alpha_j \hat{U}_j \hat{V}_j^T = \mathbf{0}_{d \times d}$ has s' linearly independent solutions $\alpha^1, \alpha^2, \dots, \alpha^{s'} \in \mathbb{R}^s$, where $1 \leq s' \leq s-1$. One can then construct a vector $\theta \in \mathbb{R}^{s'}$

such that $\sum_{k=1}^{s'} \theta_k \alpha^k \leq \mathbf{1}_s$ and the equality holds on at least s' positions. It can be done, for instance, by solving a corresponding linear program with the simplex method. By taking

$$(5.10) \quad \tilde{U} = \hat{U} \text{diag} \left\{ \mathbf{1}_s - \sum_{k=1}^{s'} \theta_k \alpha^k \right\}, \quad \tilde{V} = \hat{V},$$

and removing all the zero columns in \tilde{U} and the corresponding columns in \tilde{V} , both \tilde{U} and \tilde{V} have at most $s - s'$ columns, which is smaller than s . (\tilde{U}, \tilde{V}) is guaranteed to be feasible to (4.1) and Corollary 12 implies that $F_\lambda(\tilde{U}, \tilde{V}) = F_\lambda(\hat{U}, \hat{V})$. This modification may result in a nonstationary point, from where PALM resumes for a smaller s .

By removing redundant dimensions in local solutions, we have $s \leq d^2 + 1$ throughout the computing process. It also guarantees that the algorithm will terminate at a rank no larger than d^2 .

In numerical experiments, if

$$(5.11) \quad \Delta = \sum_{j=1}^s \sum_{k=1}^{s'} \theta_k \alpha_j^k \hat{U}_j \hat{V}_j^T$$

satisfies $\|\hat{\Xi} \Delta\|_F \leq \varepsilon$, we consider Δ as a negligible component of \hat{X} and apply (5.10) to reduce the rank of factorization. Here, ε refers to the linear dependence threshold that will be determined later.

5.3. An adaptive rank algorithm. Based on the preliminaries, we now propose an adaptive rank algorithm for the convexified problem (3.5). We apply PALM as a subroutine to solve the factorized optimization model (4.1) for fixed s until converging to a stationary point. If the local solution has redundant dimensions as described in Corollary 12, we compress it by (5.10) and continue PALM. Otherwise, we adopt either an exact or early stopping criterion to check condition (4.3). If there exists $(\bar{\mathbf{u}}, \bar{\mathbf{v}})$ violating (4.3), we append columns to the local solution according to Theorem 11, and resume PALM with a new initial point. The full algorithm is summarized as Algorithm 1.

6. Accelerating simulation of Markov chains. In this section, we show that the knowledge of state aggregation models can apply to accelerate the simulation of the large-scale Markov chains. Let $(\hat{U}, \hat{V}) \in \mathbb{R}^{d \times \hat{s}} \times \mathbb{R}^{d \times \hat{s}}$ denote the aggregation and disaggregation probabilities learned from observed data (Y_0, Y_1, \dots, Y_n) by Algorithm 1. One can utilize a well-designed small Markov chain (with \hat{s} states) to generate a random walk trajectory on the d original states. Therefore we are able to simulate the large Markov chain based on limited data and the computational costs are largely reduced.

Specifically, we simulate a new Markov trajectory $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$ in the state space $\mathcal{S} = \{s_1, s_2, \dots, s_d\}$, where the one-step transitions are governed by matrix $\hat{X} = \hat{U} \hat{V}^T$ and the initial state \hat{Y}_0 follows distribution $\xi^\circ \in \mathbb{R}^d$. Note that $\hat{X}^k = \hat{U} (\hat{V}^T \hat{U})^{k-1} \hat{V}^T$ for $k = 1, 2, \dots$, and $\hat{V}^T \hat{U} \mathbf{1}_{\hat{s}} = \hat{V}^T \mathbf{1}_d = \mathbf{1}_{\hat{s}}$. Therefore, $\hat{Q} := \hat{V}^T \hat{U}$ is a stochastic matrix that characterizes the evolution of $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$. It enables us to first produce an \hat{s} -state Markov chain with a transition matrix \hat{Q} and next recover a trajectory in the original \mathcal{S} accordingly.

Algorithm 1 An adaptive rank algorithm for (3.5).

Input: Initial point $(U^0, V^0) \in \mathcal{U}^{d \times s_0} \times \mathcal{V}^{d \times s_0}$.

Termination criterion: exact or early stopping rule.

while the termination criterion is not met **do**

while the local convergence criterion (5.6) is not met **do**

 For each $j = 1, 2, \dots, s$ such that $\|U_j^k\|_2 \leq \varepsilon_0$, remove the j th column from both U^k and V^k . Update s .

 Apply (5.2) and (5.3) to solve (4.1) and obtain (U^{k+1}, V^{k+1}) . $k \leftarrow k + 1$.

if Δ in (5.11) satisfies $\|\hat{\Xi}\Delta\|_F \leq \varepsilon$ **then**

 Compress (U^k, V^k) according to (5.10). Update s .

else

 Check the global optimality of (U^k, V^k) according to the termination criterion.

if the termination criterion is not met **then**

 Determine a step size $\kappa > 0$ in (5.9) with backtracking line search.

if $\kappa \geq 10^{-8}$ **then**

 Append a column to U^k and V^k according to (5.9). $s \leftarrow s + 1$.

else

 Terminate the algorithm.

We propose the following scheme to simulate $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$:

1. Generate a Markov chain $(\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_m)$ in the state space $\tilde{\mathcal{S}} = \{\tilde{s}_1, \tilde{s}_2, \dots, \tilde{s}_{\hat{s}}\}$:

$$(6.1) \quad \begin{cases} \mathbb{P}[\hat{Z}_1 = \tilde{s}_k] = (\hat{U}^T \xi^\circ)_k & \text{for } k = 1, 2, \dots, \hat{s}, \\ \mathbb{P}[\hat{Z}_{t+1} = \tilde{s}_l \mid \hat{Z}_t = \tilde{s}_k] = \hat{q}_{kl} & \text{for } t = 1, 2, \dots, m-1, \quad k, l = 1, 2, \dots, \hat{s}. \end{cases}$$

2. Independently generate $\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m$ in \mathcal{S} according to

$$(6.2) \quad \begin{cases} \mathbb{P}[\hat{Y}_0 = s_i \mid \hat{Z}_1 = \tilde{s}_k] = \frac{\hat{u}_{ik} \xi_i^\circ}{(\hat{U}^T \xi^\circ)_k}, & \mathbb{P}[\hat{Y}_m = s_i \mid \hat{Z}_m = \tilde{s}_k] = \hat{v}_{ik}, \\ \mathbb{P}[\hat{Y}_t = s_i \mid \hat{Z}_t = \tilde{s}_k, \hat{Z}_{t+1} = \tilde{s}_l] = \frac{\hat{v}_{ik} \hat{u}_{il}}{\hat{q}_{kl}} & \text{for } t = 1, 2, \dots, m-1 \end{cases}$$

for $k, l = 1, 2, \dots, \hat{s}$, $i = 1, 2, \dots, d$. $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$ forms a trajectory.

It is worth noting that the random variables $\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m$ in (6.2) are conditionally independent given the trajectory $(\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_m)$, so we can implement (6.2) in a parallel way. To this end, (6.2) is not time consuming compared with (6.1).

The theorem below says that $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$ is a desired Markov chain following the transition matrix \hat{X} and the initial distribution ξ° . See Appendix A.10 for the proof.

THEOREM 13. Suppose that $(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m)$ is generated by (6.1) and (6.2). Then for any trajectory $(s_{i_0}, s_{i_1}, \dots, s_{i_m}) \in \mathcal{S}^{m+1}$, it holds

$$\mathbb{P}[(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m) = (s_{i_0}, s_{i_1}, \dots, s_{i_m})] = \xi_{i_0}^\circ \hat{x}_{i_0, i_1} \hat{x}_{i_1, i_2} \dots \hat{x}_{i_{m-1}, i_m}.$$

Based on Theorem 13, we can safely use (6.1) and (6.2) for Markov chain simulation. The problem is now reduced to generating a Markov trajectory on $\tilde{\mathcal{S}}$. When $\hat{s} \ll d$, it significantly reduces computational time.

7. Numerical experiments. We present some numerical results to illustrate the efficiency of our proposed state aggregation method. When applied to synthetic transition matrices, our scheme stops exactly at the ground-truth intrinsic dimension r when the regularization constant λ is appropriately chosen. We also generate Markovian trajectories with varying dimension d , r , and sampling size n , and investigate how these parameters influence the recovery error. Finally, we use our approach to analyze a real dataset of a Manhattan transportation network, and conduct extensive comparison with an existing method in [25]. The algorithm is implemented in MATLAB and all experiments are performed on a computer with an Intel i5 1.9 GHz and 8 GB of RAM.

7.1. Evaluations of solution quality. We adopt several metrics to evaluate the quality of our solutions. The following *local errors* $relLE1$ and $relLE2$ concern whether the factorized optimization problem (4.1) is solved properly. According to the KKT conditions in Theorem 7, given that the local scheme converges to (\hat{U}, \hat{V}) , we define

$$relLE1 := \frac{\left\| \left[\mu \mathbf{1}_s^T - \nabla g(\hat{X}) \hat{V} \right]_+ - \lambda \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^s \right\|_{\ell_1}}{\left\| \lambda \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^s \right\|_{\ell_1}}$$

and

$$relLE2 := \frac{\left\| \left[\mathbf{1}_d \mu^T \hat{U} - \left(\nabla g(\hat{X}) \right)^T \hat{U} \right]_+ - \lambda \hat{V} \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^s \right\|_{\ell_1}}{\left\| \lambda \hat{V} \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^s \right\|_{\ell_1}}$$

to represent the precision of (\hat{U}, \hat{V}) . Here, the Lagrangian multiplier $\mu \in \mathbb{R}^d$ is estimated by

$$\mu_i = \frac{\sum_{j=1}^s \left(\lambda \hat{u}_{ij} \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} + \left(\nabla g(\hat{X}) \hat{V} \right)_{ij} \right) \mathbb{1}_{\{\hat{u}_{ij} \neq 0\}}}{\sum_{j=1}^s \mathbb{1}_{\{\hat{u}_{ij} \neq 0\}}}, \quad i = 1, 2, \dots, d.$$

As for the global optimality condition the *global error* (GE) is defined as

$$GE := \Omega^\circ(\hat{W}) - 1, \quad \hat{W} = \lambda^{-1}(\mu \mathbf{1}_d^T - \nabla g(\hat{X})),$$

where $\Omega^\circ(\hat{W})$ is calculated by the gradient projection method.

Aside from GE , the duality gap of problem (3.5) serves as another technique to certify global optimality. The Lagrangian dual formulation of (3.5) is given by

$$(7.1) \quad \begin{aligned} \max_{M \in \mathbb{R}^{d \times d}} \quad & -g^*(M) \\ \text{s.t.} \quad & \Omega^\circ(-M) \leq \lambda, \end{aligned}$$

where $g^*(M) = \sup_{X \in \mathcal{E}} \{\langle M, X \rangle - g(X)\}$. After a routine calculation, we have $g^*(M) = \frac{1}{2} \|\hat{\Xi}^{-1} M + \hat{\Xi} \hat{P}^{(n)}\|_F^2 - \frac{1}{2} \|\hat{\Xi} \hat{P}^{(n)}\|_F^2 - \frac{1}{2d} \|\hat{\Xi}^{-1} M \mathbf{1}_d\|_2^2$. The duality gap can be estimated by $g(\hat{X}) + \lambda \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 + g^*(M)$ for some dual candidate M such

that $\Omega^\circ(-M) \leq \lambda$. It is a good choice to take $M = -\frac{\lambda}{\Omega^\circ(\hat{W})}\hat{W}$. We therefore define the *relative duality gap*,

$$relDG := \frac{g(\hat{X}) + \lambda \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 + g^*(M)}{g(\hat{X}) + \lambda \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2}.$$

In situations where the ground truth is known in advance, we can also calculate the *relative recovery error* of a solution \hat{X} ,

$$relRE := \frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2 \bigg/ \frac{1}{2} \left\| \hat{\Xi}P^* \right\|_F^2.$$

After comparing *relRE* with the following *relative sampling error*

$$relSE := \frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F^2 \bigg/ \frac{1}{2} \left\| \hat{\Xi}P^* \right\|_F^2,$$

we can tell if solving problem (3.5) is efficient at revealing the system dynamics.

7.2. Experiments with exact low-nonnegative-rank matrices. In this part, we apply Algorithm 1 to problem (3.5), where $\hat{P}^{(n)}$ and $\hat{\xi}^{(n)}$ are replaced by P^* and ξ^* . We want to check if the algorithm can successfully recover the underlying state aggregation structure when the transition matrix is already seen.

We carry out experiments with $d = 1000, 2000, 5000$, and set the inner dimension $r = 5$. A test Markov chain is created randomly by the following procedure. We first generate two random matrices $U^*, V^* \in \mathbb{R}^{d \times r}$. The rows of U^* and columns of V^* are independent and uniformly distributed on simplexes \mathcal{U}^r and \mathcal{V}^d , respectively. We assemble a Markov transition matrix $P^* = U^*(V^*)^T$ and calculate its stationary distribution ξ^* .

When implementing Algorithm 1, we adopt the exact stopping rule for higher precision. Since the algorithm appends columns one by one, whereas reduces redundant dimensions all at once, it is a better choice to start with a large initial dimension s_0 and let s reduce to a proper range after some computations. To this end, we choose $s_0 = 300$ which is large enough so that s_0 does not affect the terminal dimension \hat{s} .

The numerical results are shown in Tables 1, 2, and 3. We learn that, if the regularization parameter λ is sufficiently small, Algorithm 1 converges to a solution $(\hat{U}, \hat{V}) \in \mathbb{R}^{d \times \hat{s}} \times \mathbb{R}^{d \times \hat{s}}$ with $\hat{s} = r$, the nonnegative rank of P^* . It implies that, in the convexified problem (3.5), the atomic regularizer performs well in identifying low-nonnegative-rank structures. We also notice that, when $\lambda = 10^{-9}$, *relRE* is small ($< (5\%)^2$), and the solutions to (3.5) are close to the ground truth transition matrices. It is safe to say that our proposed method yields reliable results in recovering the state aggregation structures of the ground truth P^* .

7.3. Experiments with simulated data. In this part, we investigate statistical properties of solutions of Algorithm 1 when the input is simulated transition trajectories. Given a fixed trajectory, we solve problem (3.5) with different λ 's, and obtain a series of \hat{X} 's and *relRE*'s. We plot *relRE* against λ and identify the most appropriate λ^* that yields the smallest error *relRE* * . In order to generate a path of solutions corresponding to different values of λ , we employ the *warm-restart* approach. To be specific, we initialize the algorithm with an obtained solution to a problem with slightly larger λ . In each group of our experiments, we investigate how

TABLE 1

The recovery of an exact low-nonnegative-rank matrix P^* with $d = 1000$, $r = 5$, $\sigma_1(\Xi^*P^*) = 1.20 \times 10^{-3}$, $\sigma_r(\Xi^*P^*) = 1.93 \times 10^{-4}$. Algorithm 1 starts with $s_0 = 300$. The linear dependence threshold $\varepsilon = 5 \times 10^{-5}$.

λ	\hat{s}	Obj value	$relRE$	Time (s)	$relLE1$	$relLE2$	GE	$relDG$
1e-06	1	1.14e-06	1.31e-01	13.8	0	6.7e-16	-1.0e-15	2.8e-15
1e-07	5	1.60e-07	2.15e-02	38.2	2.7e-03	5.6e-03	1.6e-03	5.8e-04
1e-08	5	1.80e-08	3.91e-04	77.2	2.3e-03	2.4e-03	1.7e-03	2.6e-03
1e-09	5	1.84e-09	6.76e-06	99.4	5.8e-02	3.5e-02	3.9e-02	4.5e-02

TABLE 2

The recovery of an exact low-nonnegative-rank matrix P^* with $d = 2000$, $r = 5$, $\sigma_1(\Xi^*P^*) = 5.97 \times 10^{-4}$, $\sigma_r(\Xi^*P^*) = 9.81 \times 10^{-5}$. Algorithm 1 starts with $s_0 = 300$. The linear dependence threshold $\varepsilon = 5 \times 10^{-5}$.

λ	\hat{s}	Obj value	$relRE$	Time (s)	$relLE1$	$relLE2$	GE	$relDG$
1e-06	1	1.04e-06	1.86e-01	36.8	2.0e-18	7.6e-16	-1.0e-15	2.4e-15
1e-07	1	1.27e-07	1.08e-01	37.7	7.4e-18	3.5e-06	2.8e-11	7.5e-09
1e-08	5	1.70e-08	4.65e-03	78.7	1.4e-03	2.3e-03	1.8e-03	2.7e-03
1e-09	5	1.80e-09	7.35e-05	140.7	1.7e-02	2.8e-02	3.8e-02	4.7e-02

TABLE 3

The recovery of an exact low-nonnegative-rank matrix P^* with $d = 5000$, $r = 5$, $\sigma_1(\Xi^*P^*) = 2.38 \times 10^{-4}$, $\sigma_r(\Xi^*P^*) = 3.89 \times 10^{-5}$. Algorithm 1 starts with $s_0 = 300$. The linear dependence threshold $\varepsilon = 5 \times 10^{-5}$.

λ	\hat{s}	Obj value	$relRE$	Time (s)	$relLE1$	$relLE2$	GE	$relDG$
1e-06	1	1.01e-06	2.33e-01	75.0	0	5.0e-14	-5.0e-14	1.6e-14
1e-07	1	1.06e-07	1.67e-01	90.1	3.0e-18	3.4e-14	-5.9e-14	1.1e-14
1e-08	5	1.39e-08	8.61e-02	121.8	3.7e-04	1.0e-03	1.5e-03	1.2e-03
1e-09	5	1.75e-09	2.11e-03	181.6	1.3e-02	2.0e-02	1.3e-02	5.3e-03

$relRE^*$ is influenced by one parameter among d , r/d , and n/d^2 . We let the parameter of interest take different values and fix the other two.

We provide some benchmark simulations for the state aggregation problem. A Markov chain P^* with d states and r inherent metastates is created in the same way as in subsection 7.2. After choosing an initial state i_0 under invariant distribution ξ^* , we randomly generate $i_t, t = 1, \dots, n$, step by step, and form a trajectory (i_0, i_1, \dots, i_n) of length $n + 1$.

For the sake of higher precision, we adopt the exact stopping rule when applying Algorithm 1. For the same reasons mentioned before, we also take the initial dimension $s_0 = 300$. $relLE1$, $relLE2$, GE , and $relDG$ are used to measure the optimization errors of Algorithm 1. If these quantities are relatively small, we say that the optimization errors are dominated by the statistical error $relRE$. In this case, it is secure to study the statistical properties of model (3.5) with the computational results. Also, for each group of parameters d , r/d , and n/d^2 , we run Algorithm 1 on 5 independent trajectories so as to reduce the random errors.

Dimension d . Markov chains with $d = 500, 1000, 1200, 1500, 2000, 5000$ and $r/d = 0.01$ are created. For each model, trajectories with sampling size $n/d^2 = 5, 10$, and 20 are generated independently. A plot of the regularization paths is shown in Figure 1. When the sparsity degree r/d and valid sampling size n/d^2 are fixed, the shape of the regularization paths are almost identical, regardless of the dimension d . In Figure 2, we plot the optimal relative recovery error $relRE^*$ against d . We can learn from the results that $relRE^*$ and d have an approximately linear relationship

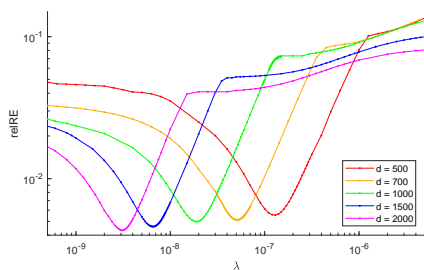


FIG. 1. $relRE$ - λ curves with $r/d = 0.01$, $n/d^2 = 10$, $d = 500, 700, 1000, 1500$, or 2000 .

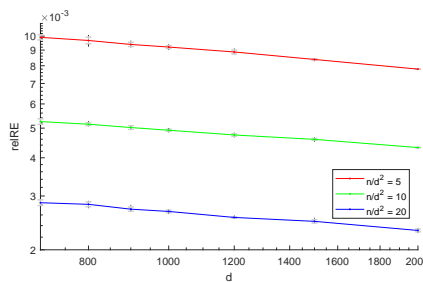


FIG. 2. $relRE^*$ - d curves with $r/d = 0.01$, $n/d^2 = 5, 10$, or 20 . Each black point comes from one independent experiment. The gray bars stand for standard deviations (after logarithm).

TABLE 4
Coefficients in the linear regression $\ln relRE^* = A_1 + B_1 \ln d$.

r/d	n/d^2	Constant A_1	Regression coefficient B_1	Squared correlation coefficient
0.01	5	-3.115	-0.228	0.997
0.01	10	-4.030	-0.186	0.998
0.01	20	-4.536	-0.201	0.996

in the log-log plot. The regression coefficients are shown in Table 4. A summary of computational results is reported in Table 5. With the growth of d , $relSE$ increases whereas $relRE^*$ decreases. Therefore, when r/d and n/d^2 are fixed, one can reconstruct the model more accurately if d is larger.

Sparsity degree r/d . Markov chains with $d = 1000$ and $r = 5, 7, 10, 15$ are created. The sampling size n/d^2 is taken to be 5 or 10. A plot of the regularization paths is shown in Figure 3. We learn that, when d and n/d^2 are fixed, one can recover the model more accurately if r/d is smaller. In Figure 4, we plot $relRE^*$ against r/d . There is also a linear relationship between $relRE^*$ and d in the log-log plot. The linear regression coefficients are shown in Table 6. More detailed computational results are presented in Table 7. When $r \ll d$, as r/d gets smaller, $relSE$ decreases slightly, and $relRE^*$ drops even faster.

Sample size n/d^2 . A Markov chain with $d = 1000$ and $r = 5$ is created. We generate trajectories with $n/d^2 = 5, 10, 20, 50, 100, 200, 500, 1000$. A plot of regularization paths is shown in Figure 5. As the sampling size n grows, the empirical transition matrix $\hat{P}^{(n)}$ gets closer and closer to the ground truth P^* , therefore, λ^* is smaller and $relRE^*$ reduces to zero. In Figure 6, we see the linear relationships between $relRE^*$ and n . The linear regression coefficients are shown in Table 8. By the central limit theorem of the Markov chain, $relSE$ is of order n^{-1} . However, the regression coefficient B_3 is slightly larger than -1 . This results from the bias introduced by the regularization term. Numerical details are summarized in Table 9.

7.4. Experiments with Manhattan taxi data. We use the state aggregation model to partition a Manhattan transportation network into different regions. Our experiment is based on a real dataset of 1.1×10^7 NYC Yellow Cab trips in January 2016 [22]. Each record includes passenger pickup and drop-off information (coordi-

TABLE 5
Numerical results with $r/d = 0.01$, $n/d^2 = 10$. The initial dimension $s_0 = 0.1d$.

(d, λ^*)	(500, 1.6e-07)		(1000, 2.4e-08)		(1200, 7.5e-09)	
ε	1.5e-04	5e-05	1.5e-04	5e-05	1.5e-04	5e-05
\hat{s}	5	13	10	22	13	27
Obj value	4.70e-07	4.70e-07	9.16e-08	9.15e-08	6.09e-08	6.09e-08
$relRE^*$	5.70e-03	5.96e-03	5.45e-03	5.61e-03	5.43e-03	5.77e-03
Time (s)	19.5	19.5	49.4	58.5	90.0	153.0
$relLE1$	1.6e-03	1.8e-03	1.7e-03	5.6e-03	5.7e-04	3.2e-03
$relLE2$	2.3e-03	1.4e-03	4.5e-03	6.9e-03	1.0e-03	9.5e-03
GE	1.3e-01	4.9e-02	2.0e-01	8.8e-02	2.0e-01	5.2e-02
$relDG$	7.1e-02	2.7e-02	8.7e-02	3.9e-02	8.5e-02	2.2e-02
$relSE$	6.30e-02		7.73e-02		7.98e-02	
$\sigma_1(\hat{\Xi}\hat{P}^{(n)})$	2.38e-03		1.10e-03		9.05e-04	
$\sigma_r(\hat{\Xi}\hat{P}^{(n)})$	3.72e-04		8.61e-05		6.18e-05	
$\sigma_{r+1}(\hat{\Xi}\hat{P}^{(n)})$	6.39e-05		2.16e-05		1.63e-05	

(d, λ^*)	(1500, 7.5e-09)		(2000, 3.5e-09)		(5000, 3.0e-10)	
ε	1.5e-04	5e-05	1.5e-04	5e-05	1.5e-04	5e-05
\hat{s}	15	33	20	42	50	124
Obj value	3.53e-08	3.55e-08	1.86e-08	1.86e-08	2.57e-09	2.56e-09
$relRE^*$	4.51e-03	5.21e-03	4.31e-03	4.38e-03	4.58e-03	4.66e-03
Time (s)	256.5	183.8	257.0	370.0	284.1	388.5
$relLE1$	2.6e-03	2.1e-03	8.8e-04	2.9e-03	9.8e-02	2.3e-02
$relLE2$	5.1e-03	7.0e-03	3.4e-03	1.3e-02	1.4e-01	4.7e-02
GE	2.1e-01	9.9e-02	2.6e-01	1.6e-01	2.2e-01	2.9e-02
$relDG$	8.0e-02	3.8e-02	9.5e-02	5.8e-02	6.1e-02	6.9e-03
$relSE$	8.38e-02		8.76e-02		9.29e-02	
$\sigma_1(\hat{\Xi}\hat{P}^{(n)})$	7.10e-04		5.24e-04		2.04e-04	
$\sigma_r(\hat{\Xi}\hat{P}^{(n)})$	4.03e-05		2.23e-05		3.65e-06	
$\sigma_{r+1}(\hat{\Xi}\hat{P}^{(n)})$	7.36e-05		7.36e-06		1.79e-06	

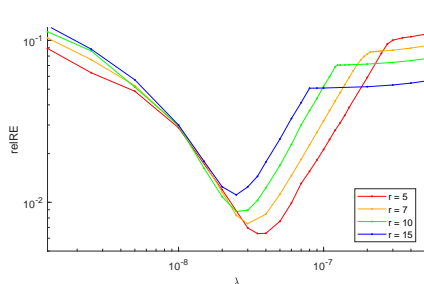


FIG. 3. $relRE$ - λ curves with $d = 1000$, $n/d^2 = 5$, $r = 5, 7, 10$, or 1500 .

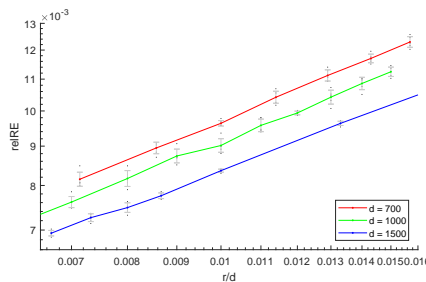


FIG. 4. $relRE^*$ - r/d curves with $n/d^2 = 5$, $d = 700, 1000$, or 1500 . Each black point comes from one independent experiment. The gray bars stand for standard deviations (after logarithm).

TABLE 6
Coefficients in the linear regression $\ln relRE^* = A_2 + B_2 \ln(r/d)$.

d	n/d^2	Constant A_2	Regression coefficient B_2	Squared correlation coefficient
700	5	-0.962	0.525	0.9994
1000	5	-1.026	0.507	0.9990
1500	5	-1.122	0.477	0.9998

TABLE 7
Numerical results with $d = 1000$, $n/d^2 = 10$. The initial dimension $s_0 = 0.1d$.

(r, λ^*)	(5, 3.2e-08)		(10, 2.4e-08)		(15, 2.1e-08)	
ε	1.5e-04	5e-05	1.5e-04	5e-05	1.5e-04	5e-05
\hat{s}	5	14	10	22	16	39
Obj value	1.04e-07	1.04e-07	9.16e-08	9.15e-08	8.42e-08	8.58e-08
$relRE^*$	3.77e-03	3.86e-03	5.45e-03	5.61e-03	6.20e-03	6.67e-03
Time (s)	63.3	68.3	49.4	58.5	65.4	136.4
$relLE1$	5.0e-03	1.8e-03	1.7e-03	5.6e-03	2.3e-04	5.6e-04
$relLE2$	8.6e-03	1.6e-03	4.5e-03	6.9e-03	1.1e-03	1.7e-03
GE	1.9e-01	7.6e-02	2.0e-01	8.8e-02	1.5e-01	3.3e-02
$relDG$	9.3e-02	3.8e-02	8.7e-02	3.9e-02	6.3e-02	1.4e-02
$relSE$	6.31e-02		7.73e-02		8.36e-02	
$\sigma_1(\hat{\Xi}\hat{P}^{(n)})$	1.20e-03		1.10e-03		1.06e-03	
$\sigma_r(\hat{\Xi}\hat{P}^{(n)})$	1.73e-04		8.61e-05		5.74e-05	
$\sigma_{r+1}(\hat{\Xi}\hat{P}^{(n)})$	2.33e-05		2.16e-05		2.09e-05	

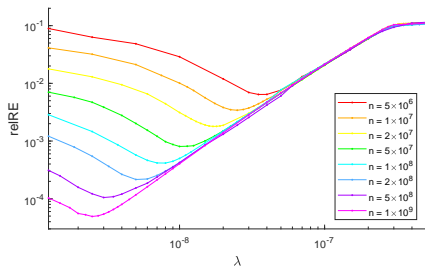


FIG. 5. $relRE$ - λ curves with $d = 1000$, $r/d = 0.005$, $n/d^2 = 5, 10, 20, 50, 100, 200, 500$, or 1500.

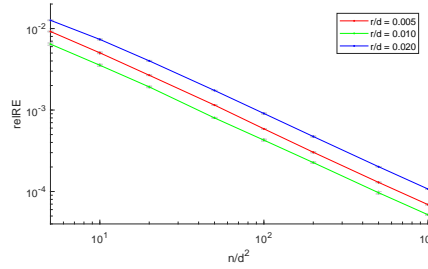


FIG. 6. $relRE$ - n/d^2 curves with $d = 1000$, $r/d = 0.005, 0.010$, or 0.020 . Each black point comes from one independent experiment. The gray bars stand for standard deviations (after log-arithmetic).

TABLE 8
Coefficients in the linear regression $\ln relRE^* = A_3 + B_3 \ln(n/d^2)$.

d	r/d	Constant A_3	Regression coefficient B_3	Squared correlation coefficient
1000	0.005	-1.372	-0.930	0.99987
1000	0.010	-1.539	-0.916	0.99990
1000	0.015	-1.231	-0.910	0.99950

nates, time, etc.) of one trip. We want to construct a stochastic model of the traffic flow. The movements of taxis are nearly memoryless. Therefore, we admit that the stochastic system satisfies the Markov property, and approximate it by a finite-state Markov chain. We divide the map into a fine grid and merge the locations in the same cell into one state. Each trip is a sampled one-step state transition between cells. We formulate an empirical Markov transition matrix \hat{P} and an empirical stationary distribution $\hat{\xi}$ based on the dataset, and seek for a method to simplify the stochastic system.

We further assume that the transportation system is driven by a Markov chain with fewer states that are invisible and are aggregations of the states in the original Markov chain. In order to identify the state aggregation structure, we apply model

TABLE 9
Numerical results with $d = 1000$, $r = 5$. The initial dimension $s_0 = 0.1d$.

(n, λ^*)	(5e+06, 4.1e-08)		(1e+07, 3.2e-08)		(2e+07, 2.4e-08)	
ε	1.5e-04	5e-05	1.5e-04	5e-05	1.5e-04	5e-05
\hat{s}	6	17	5	14	5	14
Obj value	1.67e-07	1.69e-07	1.04e-07	1.04e-07	6.67e-08	6.67e-08
$relRE^*$	5.77e-03	6.40e-03	3.77e-03	3.86e-03	2.33e-03	2.34e-03
Time (s)	87.3	108.8	63.3	68.3	32.5	38.4
$relLE1$	4.0e-03	2.0e-03	5.0e-03	1.8e-03	1.6e-03	1.9e-03
$relLE2$	5.3e-03	1.4e-03	8.6e-03	1.6e-03	2.8e-03	2.4e-03
GE	2.3e-01	8.8e-02	1.9e-01	7.6e-02	1.4e-01	3.9e-02
$relDG$	9.3e-02	3.7e-02	9.3e-02	3.8e-02	8.3e-02	2.4e-02
$relSE$	1.26e-01		6.31e-02		3.15e-02	
$\sigma_1(\hat{\Xi}\hat{P}^{(n)})$	1.20e-03		1.20e-03		1.20e-03	
$\sigma_r(\hat{\Xi}\hat{P}^{(n)})$	1.74e-04		1.73e-04		1.87e-04	
$\sigma_{r+1}(\hat{\Xi}\hat{P}^{(n)})$	3.29e-05		2.33e-05		1.70e-05	

(n, λ^*)	(5e+07, 1.7e-08)		(1e+08, 1.2e-08)		(2e+08, 9.0e-09)	
ε	1.5e-04	5e-05	1.5e-04	5e-05	1.5e-04	5e-05
\hat{s}	5	17	5	18	5	12
Obj value	3.97e-08	3.97e-08	2.61e-08	2.62e-08	1.85e-08	1.85e-08
$relRE^*$	1.18e-03	1.20e-03	6.37e-04	6.83e-04	3.91e-04	3.94e-04
Time (s)	42.0	47.9	42.7	28.6	43.0	45.8
$relLE1$	2.4e-03	2.5e-03	3.2e-03	3.8e-03	7.4e-03	2.0e-02
$relLE2$	3.9e-03	3.7e-03	4.8e-03	2.6e-03	1.2e-02	5.5e-03
GE	6.5e-02	1.2e-02	7.0e-02	1.5e-02	2.7e-02	3.7e-02
$relDG$	4.6e-02	8.4e-03	5.2e-02	1.2e-02	2.1e-02	2.0e-02
$relSE$	1.27e-02		6.37e-03		3.23e-03	
$\sigma_1(\hat{\Xi}\hat{P}^{(n)})$	1.20e-03		1.20e-03		1.20e-03	
$\sigma_r(\hat{\Xi}\hat{P}^{(n)})$	1.72e-04		1.71e-04		1.72e-04	
$\sigma_{r+1}(\hat{\Xi}\hat{P}^{(n)})$	1.04e-05		7.33e-06		5.25e-06	

(3.5) to the estimated \hat{P} and $\hat{\xi}$. The solution \hat{U} and \hat{V} in (3.5) helps embed the states in the original Markov chain into an s -dimensional space. We then cluster the states according to their coordinates, which yields a partition of the transportation network. In this problem, it is desirable to have a small rank s , since clustering algorithms are very likely to fail in a high-dimensional space. For this reason, we choose the early stopping rule and start from $s_0 = 1$ when applying Algorithm 1.

We now introduce some specific settings in our numerical experiments. When preprocessing data, we first delete the records beyond area $80.07^\circ\text{W} \sim 60.92^\circ\text{W}$, $30.69^\circ\text{N} \sim 50.85^\circ\text{N}$, and divide the rectangle into a $0.001^\circ \times 0.001^\circ$ grid. We count the number of times a state appears as a pickup place, and discard the ones with frequency of occurrence lower than 10^{-4} . We end up with 2017 valid states denoted by $\mathcal{S} = \{1, 2, \dots, 2017\}$, and 7.5×10^6 remaining records $(s_t^{\text{pickup}}, s_t^{\text{dropoff}}) \in \mathcal{S}^2$, $t = 1, 2, 3, \dots, T$. The empirical \hat{P} and $\hat{\xi}$ are formulated by

$$\hat{P}_{ij} = \frac{\sum_{t=1}^T \mathbb{1}[s_t^{\text{pickup}} = i, s_t^{\text{dropoff}} = j]}{\sum_{t=1}^T \mathbb{1}[s_t^{\text{pickup}} = i]}, \quad i, j \in \mathcal{S},$$

$$\hat{\xi}_i = \frac{\sum_{t=1}^T \mathbb{1}[s_t^{\text{pickup}} = i]}{T}, \quad i \in \mathcal{S}.$$

We use the MATLAB function `kmeans` to implement a k -means++ algorithm, and specify 500 replicates to help find a lower local minimum.

In order to illustrate the reliability of the state aggregation model, we compare our results with the batch partition procedure in [25], which provides another method to embed the states in a reduced-dimensional space. The network partition algorithm in [25] begins with an optimal rank- s approximation of matrix $\hat{\Xi}\hat{P}$ given by SVD, which we refer to as $\tilde{U}\tilde{D}\tilde{V}^T$. Here, s is the number of parts that the map will be partitioned into. It has been proved in [10] that, for a lumpable Markov chain, one can obtain the optimal partition of state space with the clustering results of $\hat{\Xi}^{-1}\tilde{U}$ or $\hat{\Xi}^{-1}\tilde{V}$. In our numerical experiments, we used the MATLAB function `svds` to implement approximate SVD and `kmeans` with 500 replicates to cluster the states.

The zoning results of \hat{U} in (3.5) and $\hat{\Xi}^{-1}\tilde{U}$ in [25] are shown in Figure 7. Overall, the results of our approach are more aligned to geometric location, while the partition results of the SVD-based method are more scattered. The figures also show that the appropriate regularization parameter λ ranges from 1.5×10^{-10} to 1.8×10^{-10} . Within this interval, Manhattan island is divided into $6 \sim 9$ coherent regions. Note that the regularization parameter λ is small. One possible reason is that the maximal components of \hat{P} are already of the order 10^{-3} . The partition results coincide with the division of lower, midtown, and upper Manhattan and provide abundant information on how the taxi trips are distributed.

8. Conclusion. We propose a convex programming formulation for the state aggregation problem of Markov chains. An atomic regularizer is introduced to control the nonnegative rank of the solutions. The convex formulation is solved by minimizing a sequence of fixed-rank nonnegative factorization models using the PALM. The first-order optimality conditions of the convex and factorized problems are both investigated. They enable us to establish a criterion of whether a stationary point of the factorization model is globally optimal to the convex one. We further develop strategies to adjust the rank of the factorization. By increasing the rank, we can ensure a monotone decrease of the objective function and escape from a local minimum. In situations where the matrix factorization has “redundant ranks,” we can compress the matrices so that the state space is reduced to a proper size. Numerical experiments show that our method always converges to a global solution at a rank much smaller than the ambient dimension. Investigation of statistical properties and a real-world application of the Manhattan transportation data set are also provided.

The performance of our method can be further improved in several aspects, including designing a better regularization than the atomic regularizer and developing a faster method for solving the factorized optimization model. Two particularly important topics of investigation are (i) theoretical analysis of the convergence to the globally optimal solution, and (ii) numerical study for real and practical Markov decision processes.

Appendix A. Proofs.

A.1. Proof of Theorem 2. We first state a preliminary Lemma 14 regarding high-probability deviation bounds of $\hat{\xi}^{(n)}$ and $\hat{P}^{(n)}$.

LEMMA 14. *There exists a constant $C_0 > 0$ such that, if $n \geq C_0$, then with probability at least $1 - n^{-1}$,*

$$\left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F \leq C_0 n^{-1/2} \log(n).$$

Here, the constant C_0 only depends on the ground truth P^ .*

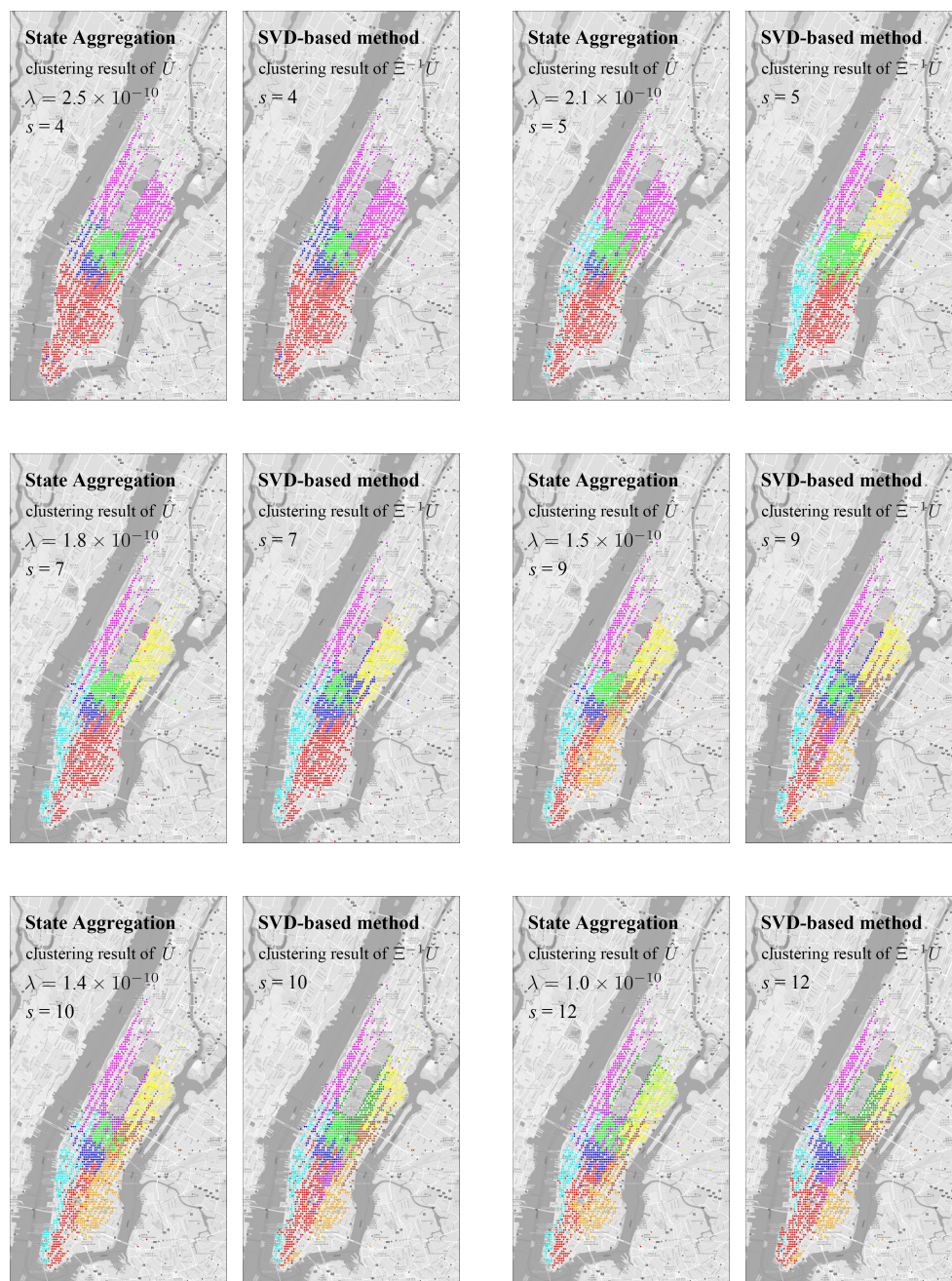


FIG. 7. The partition of the Manhattan transportation network. Each point in the map represents a valid state of the Markov chain. The figures in one pair have exactly the same number of regions, where the left one is produced by the state aggregation model and the right one is provided by the SVD-based method in [25]. In some figures, there are less than s regions appearing on the map, because some points are plotted beyond the boundaries.

Proof of Lemma 14. Lemma 14 is a corollary of the Markov chain concentration inequalities in [26]. Lemma 7 in [26] shows that there exists a constant $C_1 > 0$ dependent on P^* such that, if $n \geq C_1$, then with probability at least $1 - n^{-1}$,

$$(A.1) \quad \left\| \hat{\Xi} \hat{P}^{(n)} - \Xi^* P^* \right\|_2 \leq C_1 n^{-1/2} \log(n) \quad \text{and} \quad \left\| \hat{\xi}^{(n)} - \xi^* \right\|_\infty \leq C_1 n^{-1/2} \log(n).$$

It implies

$$\begin{aligned} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F &\leq \left\| \hat{\Xi} \hat{P}^{(n)} - \Xi^* P^* \right\|_F + \left\| (\hat{\Xi} - \Xi^*) P^* \right\|_F \\ &\leq \sqrt{d} \left\| \hat{\Xi} \hat{P}^{(n)} - \Xi^* P^* \right\|_2 + \left\| \hat{\xi}^{(n)} - \xi^* \right\|_\infty \|P^*\|_F \\ &\leq C_1 (\sqrt{d} + \|P^*\|_F) \cdot n^{-1/2} \log(n). \end{aligned} \quad \square$$

Next we prove Theorem 2.

Proof of Theorem 2. Recall that \hat{X} solves the optimization problem (3.5), therefore, it holds that

$$(A.2) \quad \frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - \hat{X}) \right\|_F^2 + \lambda \Omega(\hat{X}) \leq \frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F^2 + \lambda \Omega(P^*).$$

Note that

$$(A.3) \quad \begin{aligned} &\frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - \hat{X}) \right\|_F^2 \\ &= \frac{1}{2} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F^2 - \left\langle \hat{\Xi}(\hat{P}^{(n)} - P^*), \hat{\Xi}(\hat{X} - P^*) \right\rangle + \frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2. \end{aligned}$$

Plugging (A.3) into (A.2) gives

$$\frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2 \leq \left\langle \hat{\Xi}(\hat{P}^{(n)} - P^*), \hat{\Xi}(\hat{X} - P^*) \right\rangle + \lambda \left(\Omega(P^*) - \Omega(\hat{X}) \right),$$

which implies

$$(A.4) \quad \frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2 \leq \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F + \lambda \Omega(P^*).$$

If $\lambda \geq [\Omega(P^*)]^{-1} \left\| \hat{\Xi}(\hat{P}^{(n)} - P^*) \right\|_F^2$, (A.4) can be further reduced to

$$\frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2 \leq \sqrt{\lambda \Omega(P^*)} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F + \lambda \Omega(P^*),$$

i.e.,

$$(A.5) \quad \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F \leq (1 + \sqrt{3}) \sqrt{\lambda \Omega(P^*)}.$$

Let C_0 be the constant in Lemma 14. We learn from (A.4) and Lemma 14 that, if $n \geq C_0$, then with probability at least $1 - n^{-1}$,

$$\frac{1}{2} \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F^2 \leq C_0 n^{-1/2} \log(n) \left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F + \lambda \Omega(P^*).$$

Let $\lambda = \gamma \cdot n^{-1} [\log(n)]^2$. It follows that

$$\left\| \hat{\Xi}(\hat{X} - P^*) \right\|_F \leq \left(C_0 + \sqrt{C_0^2 + 2\gamma \Omega(P^*)} \right) n^{-1/2} \log(n),$$

which completes the proof of Theorem 2. \square

A.2. Proof of Lemma 3.

Proof. The subgradient of the characteristic function $\chi_{\mathcal{E}}(X)$ is the normal cone of \mathcal{E} , i.e.,

$$\begin{aligned}\partial\chi_{\mathcal{E}}(X) &= \{G \in \mathbb{R}^{d \times d} \mid \langle G, Y - X \rangle \leq 0, \forall Y \in \mathcal{E}\} \\ &= \{G \in \mathbb{R}^{d \times d} \mid \langle G, Z \rangle \leq 0, \forall Z \text{ s.t. } Z\mathbf{1}_d = \mathbf{0}_d\}.\end{aligned}$$

The i th row of G satisfies $\mathbf{z}^T G^i \leq 0$ for all $\mathbf{z} \in \mathbf{1}_d^\perp$. By definition, we also have $-\mathbf{z} \in \mathbf{1}_d^\perp$ so that $\mathbf{z}^T G^i \geq 0$. Hence, $\mathbf{z}^T G^i = 0$ for all $\mathbf{z} \in \mathbf{1}_d^\perp$. It indicates that $G^i = \mu_i \mathbf{1}_d$ for some scalar $\mu_i \in \mathbb{R}$. Combining the rows G^i together gives (3.8). \square

A.3. Proof of Lemma 4.

Proof. Although similar results have been provided in Lemma 13 in [11], we briefly write the proof here in order to keep the paper self-contained. By definition, a matrix $W \in \partial\Omega(X)$ if and only if

$$\Omega(Z) \geq \Omega(X) + \langle W, Z - X \rangle \quad \forall Z \in \mathbb{R}^{d \times d},$$

which is equivalent to

$$\begin{aligned}(\text{A.6}) \quad \langle W, X \rangle - \Omega(X) &\geq \sup_{Z \in \mathbb{R}^{d \times d}} \{\langle W, Z \rangle - \Omega(Z)\} \\ &= \sup_{k \in \mathbb{R}_+} \left\{ k \left(\sup_{Z: \Omega(Z)=1} \langle W, Z \rangle - 1 \right) \right\}.\end{aligned}$$

If $\Omega^\circ(W) > 1$, i.e., $\sup_{Z: \Omega(Z)=1} \langle W, Z \rangle > 1$, one can take $k \rightarrow +\infty$, and the right-hand side of (A.6) goes to infinity. It contradicts the finiteness of the left-hand side. Therefore, one must have $\Omega^\circ(W) \leq 1$. In this case, we set $k = 0$ in order to reach the supremum in (A.6) and arrive at

$$(\text{A.7}) \quad \partial\Omega(X) = \{W \mid \Omega^\circ(W) \leq 1, \langle W, X \rangle \geq \Omega(X)\}.$$

By the definition of $\Omega^\circ(X)$, if $\Omega^\circ(W) \leq 1$, one naturally has $\langle W, X \rangle \leq \Omega(X)$. Hence, (A.7) indicates (3.9). \square

A.4. Proof of Theorem 5.

Proof. A matrix \hat{X} is globally optimal for (3.5) if and only if

$$(\text{A.8}) \quad 0 \in \nabla g(\hat{X}) + \partial\chi_{\mathcal{E}}(\hat{X}) + \lambda \partial\Omega(\hat{X}),$$

where the explicit forms of $\partial\chi_{\mathcal{E}}(\hat{X})$ and $\partial\Omega(\hat{X})$ can be derived from Lemmas 3 and 4. In order to represent condition (A.8) in terms of \hat{U} and \hat{V} , we tailor the expression of $\partial\Omega(\hat{X})$ in (3.9) to the factorization $\hat{X} = \hat{U}\hat{V}^T$,

$$(\text{A.9}) \quad \partial\Omega(\hat{X}) = \left\{ W \mid \Omega^\circ(W) \leq 1, \sum_{j=1}^s \hat{U}_j^T W \hat{V}_j = \Omega(\hat{X}) \right\}.$$

Plugging (3.8) and (A.9) into (A.8), we obtain that the global optimality of \hat{X} is equivalent to the existence of a vector $\mu \in \mathbb{R}^d$ such that

$$\begin{aligned}(\text{A.10}) \quad &\Omega^\circ(\hat{W}) \leq 1, \\ (\text{A.11}) \quad &\sum_{j=1}^s \hat{U}_j^T \hat{W} \hat{V}_j = \Omega(\hat{X}), \quad j = 1, 2, \dots, s,\end{aligned}$$

where the matrix

$$(A.12) \quad \hat{W} = \lambda^{-1} \left[\mu \mathbf{1}_d^T - \nabla g(\hat{X}) \right]$$

is defined by μ . It is also desirable to have $\hat{X} = \hat{U}\hat{V}^T$ as an optimal factorization with respect to Ω . Therefore, we will make some adjustments to (A.11), so that the new condition and (A.10) hold simultaneously if and only if \hat{X} is globally optimal and $\hat{X} = \hat{U}\hat{V}^T$ is an optimal factorization.

Under the condition $\Omega(\hat{X}) = \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2$, it follows from (A.11) that

$$(A.13) \quad \sum_{j=1}^s \hat{U}_j^T \hat{W} \hat{V}_j = \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2.$$

Condition (A.10) also implies $\hat{U}_j^T \hat{W} \hat{V}_j \leq \|\hat{U}_j\|_2 \|\hat{V}_j\|_2$ for $j = 1, 2, \dots, s$. Compared with (A.13), the inequalities need to hold as equalities. Hence, (A.11) can be reduced to

$$(A.14) \quad \hat{U}_j^T \hat{W} \hat{V}_j = \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, 2, \dots, s.$$

Conversely, if there exists $\mu \in \mathbb{R}^d$ satisfying (A.10) and (A.14), then

$$\sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 \stackrel{(A.14)}{=} \sum_{j=1}^s \hat{U}_j^T \hat{W} \hat{V}_j = \langle \hat{W}, \hat{X} \rangle \stackrel{(A.10)}{\leq} \Omega(\hat{X}).$$

The definition of $\Omega(\hat{X})$ suggests that $\Omega(\hat{X}) \leq \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2$. Therefore, $\Omega(\hat{X}) = \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2$ and $\hat{X} = \hat{U}\hat{V}^T$ is an optimal factorization with respect to Ω . Consequently, one can easily derive (A.11) from (A.14), so \hat{X} is a global solution.

To sum up, we have shown that (A.10) and (A.14) are sufficient and necessary conditions for \hat{X} to be a global solution with an optimal factorization $\hat{X} = \hat{U}\hat{V}^T$. In the following paragraphs, we will rewrite these two conditions in a more explicit form.

According to the definition of Ω° , (A.10) can be transformed into

$$(A.15) \quad \mathbf{u}^T \hat{W} \mathbf{v} \leq 1 \quad \forall \mathbf{u}, \mathbf{v} \in \mathbb{R}_+^d \text{ s.t. } \|\mathbf{u}\|_2 = \|\mathbf{v}\|_2 = 1.$$

We next claim that, under (A.10), the following statements (A.16) and (A.17) are both equivalent to (A.14):

$$(A.16) \quad [\hat{W}\hat{V}]_+ = \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^s,$$

$$(A.17) \quad [\hat{W}^T \hat{U}]_+ = \hat{V} \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^s.$$

Due to the symmetry of (A.16) and (A.17), we only need to prove the equivalence between (A.16) and (A.14). The condition $\Omega^\circ(\hat{W}) \leq 1$ implies that, for an arbitrary j , $\frac{\hat{U}_j}{\|\hat{U}_j\|_2}$ optimizes function $\mathbf{u}^T (\hat{W} \hat{V}_j)$ over the set $\{\mathbf{u} \in \mathbb{R}_+^d \mid \|\mathbf{u}\|_2 = 1\}$, i.e.,

$$(A.18) \quad \frac{\hat{U}_j}{\|\hat{U}_j\|_2} = \arg \max_{\mathbf{u} \in \mathbb{R}_+^d, \|\mathbf{u}\|_2=1} \mathbf{u}^T (\hat{W} \hat{V}_j) = \frac{[\hat{W}\hat{V}_j]_+}{\|[\hat{W}\hat{V}_j]_+\|_2}.$$

The second equality in (A.18) holds because $[\hat{W}\hat{V}_j]_+ \neq \mathbf{0}_d$, which is derived by comparing the two sides of (A.14) and considering that $\hat{U}_j \neq \mathbf{0}_d$, $\hat{V}_j \neq \mathbf{0}_d$, and $\hat{U}_j \geq 0$. Additionally, we have

$$\|[\hat{W}\hat{V}_j]_+\|_2 = \left(\frac{[\hat{W}\hat{V}_j]_+}{\|[\hat{W}\hat{V}_j]_+\|_2} \right)^T \hat{W}\hat{V}_j \stackrel{(A.18)}{=} \left(\frac{\hat{U}_j}{\|\hat{U}_j\|_2} \right)^T \hat{W}\hat{V}_j \stackrel{(A.14)}{=} \|\hat{V}_j\|_2,$$

which gives

$$[\hat{W}\hat{V}_j]_+ = \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \hat{U}_j.$$

Combining the vectors together, we arrive at (A.16). On the other hand, when (A.16) holds,

$$\begin{aligned} \hat{U}_j^T \hat{W} \hat{V}_j &= \hat{U}_j^T [\hat{W}\hat{V}_j]_+ \\ &= \hat{U}_j^T \left(\frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \hat{U}_j \right) = \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, \dots, s. \end{aligned}$$

Therefore, (A.14) and (A.16) are equivalent.

Integrating (A.15), (A.16), and (A.17) together, we complete the proof of Theorem 5. We also note that, for an arbitrary $i = 1, 2, \dots, d$, since $(\hat{U}^i)^T \mathbf{1}_s = 1$, there exists $j \in \{1, 2, \dots, s\}$ such that $\hat{u}_{ij} > 0$. The vector $\mu \in \mathbb{R}^d$ has a closed form,

$$(A.19) \quad \mu_i = \left(\nabla g(\hat{X}) \hat{V} \right)_{ij} + \lambda \hat{u}_{ij} \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \quad \text{for any } j \text{ such that } \hat{u}_{ij} > 0.$$

The existence of μ is verified if and only if for any fixed i , the right-hand side of (A.19) takes the same value for all j satisfying $\hat{u}_{ij} > 0$. \square

A.5. Proof of Lemma 6.

Proof. It is obvious that if (U, V) is feasible for (4.1), then $X = UV^T$ is feasible for (3.5). By the definition of $\Omega(X)$, for any $U \in \mathcal{U}^{d \times s}$, $V \in \mathcal{V}^{d \times s}$,

$$f_\lambda(X) = g(UV^T) + \lambda \Omega(X) \leq g(UV^T) + \lambda \sum_{j=1}^s \|U_j\|_2 \|V_j\|_2 = F_\lambda(U, V).$$

Taking the minimum on both sides, we have

$$(A.20) \quad f_\lambda(\hat{X}) \leq F_\lambda(\hat{U}, \hat{V}),$$

where \hat{X} and (\hat{U}, \hat{V}) are, respectively, the optimal solutions to (3.5) and (4.1).

We next aim to show that, given a global solution \hat{X} to (4.1), there exists $(\hat{U}, \hat{V}) \in \mathcal{U}^{d \times s} \times \mathcal{V}^{d \times s}$ such that $F_\lambda(\hat{U}, \hat{V}) = f_\lambda(\hat{X})$ when s is sufficiently large. In fact, the atomic set \mathcal{A}_+ in (3.3) is compact, which implies that its convex hull $\text{conv}(\mathcal{A}_+)$ is also compact. To this end, the infimum in the definition of $\Omega(\hat{X})$ can be achieved, i.e., $\hat{X} \in t\text{conv}(\mathcal{A}_+)$ holds for $t = \Omega(\hat{X})$. According to Carathéodory's theorem, there exist atoms $A_1, A_2, \dots, A_{s_0} \in \mathcal{A}_+$ and parameters $c_1, c_2, \dots, c_{s_0} > 0$, where $s_0 \leq d^2 + 1$, such that

$$\hat{X} = \Omega(\hat{X}) \sum_{j=1}^{s_0} c_j A_j, \quad \sum_{j=1}^{s_0} c_j = 1.$$

Here, each atom A_j can be represented by $A_j = \mathbf{u}_j \mathbf{v}_j^T$ with $\mathbf{u}_j, \mathbf{v}_j \in \mathbb{R}_+^d$, $\|\mathbf{u}_j\|_2 = \|\mathbf{v}_j\|_2 = 1$. When $s \geq s_0$, we take

$$\begin{aligned}\hat{U} &:= \Omega(\hat{X}) \begin{bmatrix} c_1(\mathbf{v}_1^T \mathbf{1}_d) \mathbf{u}_1, & c_2(\mathbf{v}_2^T \mathbf{1}_d) \mathbf{u}_2, & \dots, & c_{s_0}(\mathbf{v}_{s_0}^T \mathbf{1}_d) \mathbf{u}_{s_0}, & \mathbf{0}_{d \times (s-s_0)} \end{bmatrix}, \\ \hat{V} &:= \begin{bmatrix} (\mathbf{v}_1^T \mathbf{1}_d)^{-1} \mathbf{v}_1, & (\mathbf{v}_2^T \mathbf{1}_d)^{-1} \mathbf{v}_2, & \dots, & (\mathbf{v}_{s_0}^T \mathbf{1}_d)^{-1} \mathbf{v}_{s_0}, & V' \end{bmatrix},\end{aligned}$$

where V' is any matrix in $\mathcal{V}^{d \times (s-s_0)}$. In this way, (\hat{U}, \hat{V}) is feasible for (4.1), $\hat{X} = \hat{U} \hat{V}^T$, and $\Omega(\hat{X}) = \sum_{j=1}^s \|\hat{U}_j\|_2 \|\hat{V}_j\|_2$, hence

$$(A.21) \quad F_\lambda(\hat{U}, \hat{V}) = f_\lambda(\hat{X}).$$

It results in the equivalence between (3.5) and (4.1). \square

A.6. Proof of Theorem 7.

Proof. In problem (4.1), the feasible set consists of linear constraints. Hence, the regularity condition for KKT conditions is satisfied. In other words, the KKT conditions are necessary for local optimality.

The Lagrangian function of problem (4.1) is defined as follows:

$$\begin{aligned}\mathcal{L}(U, V, \mu^U, \mu^V, \Theta^U, \Theta^V) &= F_\lambda(U, V) - (\mu^U)^T (U \mathbf{1}_s - \mathbf{1}_d) - \langle \Theta^U, U \rangle \\ &\quad - (\mu^V)^T (V^T \mathbf{1}_d - \mathbf{1}_s) - \langle \Theta^V, V \rangle,\end{aligned}$$

where $\mu^U \in \mathbb{R}^d$, $\mu^V \in \mathbb{R}^s$, $\Theta^U, \Theta^V \in \mathbb{R}_+^{d \times s}$ are dual variables. Given a local solution (\hat{U}, \hat{V}) to (4.1), the KKT conditions are

$$\begin{aligned}(A.22) \quad & \begin{cases} 0 \in \partial_{(U,V)} \mathcal{L} = \partial F_\lambda(\hat{U}, \hat{V}) - \begin{bmatrix} \mu^U \mathbf{1}_s^T \\ \mathbf{1}_d (\mu^V)^T \end{bmatrix} - \begin{bmatrix} \Theta^U \\ \Theta^V \end{bmatrix}, \\ (A.23) \quad & \hat{U} \mathbf{1}_s = \mathbf{1}_d, \quad \langle \Theta^U, \hat{U} \rangle = 0, \hat{U} \geq 0, \Theta^U \geq 0, \\ (A.24) \quad & \hat{V}^T \mathbf{1}_d = \mathbf{1}_s, \quad \langle \Theta^V, \hat{V} \rangle = 0, \hat{V} \geq 0, \Theta^V \geq 0. \end{cases}\end{aligned}$$

For the convenience of notations, we assume that $\hat{U}_j \neq \mathbf{0}_d$ for $j = 1, 2, \dots, s_1$ and $\hat{U}_j = 0$ for $j = s_1 + 1, s_1 + 2, \dots, s$, then decompose \hat{U} and \hat{V} by

$$\hat{U} = \begin{bmatrix} \hat{U}_\alpha & \mathbf{0}_{d \times (s-s_1)} \end{bmatrix}, \quad \hat{V} = \begin{bmatrix} \hat{V}_\alpha & \hat{V}_\beta \end{bmatrix}.$$

Similarly as \hat{U} and \hat{V} , μ^V , Θ^U , and Θ^V are also decomposed as

$$\mu^V = \begin{bmatrix} \mu_\alpha^V \\ \mu_\beta^V \end{bmatrix} \begin{matrix} s_1 \\ s-s_1 \end{matrix}, \quad \Theta^U = \begin{bmatrix} \Theta_\alpha^U & \Theta_\beta^U \end{bmatrix}, \quad \Theta^V = \begin{bmatrix} \Theta_\alpha^V & \Theta_\beta^V \end{bmatrix}.$$

In this way, any subgradient $W \in \partial F_\lambda(\hat{U}, \hat{V})$ can be expressed by

$$(A.25) \quad W = \begin{bmatrix} \nabla g(\hat{X}) \hat{V}_\alpha + \lambda \hat{U}_\alpha \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^{s_1} & \nabla g(\hat{X}) \hat{V}_\beta + \lambda G \text{diag} \left\{ \|\hat{V}_j\|_2 \right\}_{j=s_1+1}^s \\ \left(\nabla g(\hat{X}) \right)^T \hat{U}_\alpha + \lambda \hat{V}_\alpha \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^{s_1} & \mathbf{0}_{d \times (s-s_1)} \end{bmatrix},$$

where G is any d -by- $(s - s_1)$ matrix satisfying $\|G_j\|_2 \leq 1$ for $j = 1, 2, \dots, s - s_1$. Plugging (A.25) into (A.22), we have

$$(A.26) \quad \left\{ \begin{array}{l} \nabla g(\hat{X})\hat{V}_\alpha + \lambda \hat{U}_\alpha \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^{s_1} - \mu^U \mathbf{1}_{s_1}^T - \Theta_\alpha^U = 0, \end{array} \right.$$

$$(A.27) \quad \left\{ \begin{array}{l} \left(\nabla g(\hat{X}) \right)^T \hat{U}_\alpha + \lambda \hat{V}_\alpha \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^{s_1} - \mathbf{1}_d (\mu_\alpha^V)^T - \Theta_\alpha^V = 0. \end{array} \right.$$

Due to the nonnegativity and complementarity of \hat{U}_α and Θ_α^U , (A.26) can be reduced to

$$(A.28) \quad \left[\mu^U \mathbf{1}_{s_1}^T - \nabla g(\hat{X})\hat{V}_\alpha \right]_+ = \lambda \hat{U}_\alpha \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^{s_1}.$$

Similarly, it follows from (A.27) that

$$(A.29) \quad \left[\mathbf{1}_d (\mu_\alpha^V)^T - \left(\nabla g(\hat{X}) \right)^T \hat{U}_\alpha \right]_+ = \lambda \hat{V}_\alpha \text{diag} \left\{ \frac{\|\hat{U}_j\|_2}{\|\hat{V}_j\|_2} \right\}_{j=1}^{s_1}.$$

We next investigate the relationship between μ^U and μ_α^V . Multiplying the j th column in (A.28) by \hat{U}_j^T on the left, we have

$$(A.30) \quad \hat{U}_j^T \left(\mu^U - \nabla g(\hat{X})\hat{V}_j \right) = \lambda \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, 2, \dots, s_1.$$

From (A.29) we also have,

$$(A.31) \quad \hat{V}_j^T \left(\mu_j^V \mathbf{1}_d - \left(\nabla g(\hat{X}) \right)^T \hat{U}_j \right) = \lambda \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, 2, \dots, s_1.$$

Comparing (A.30) and (A.31), and using the fact that $\hat{V}_j^T \mathbf{1}_d = 1$, we arrive at

$$(A.32) \quad \hat{U}_\alpha^T \mu^U = \mu_\alpha^V.$$

Replacing μ_α^V in (A.29) by $\hat{U}_\alpha^T \mu^U$ and combining (A.28) and (A.29) together, we complete the proof of Theorem 7. Similarly to the discussion in the proof of Theorem 5, μ^U can be explicitly evaluated and is unique if existing. \square

A.7. Proof of Corollary 12.

Proof. Because (\hat{U}, \hat{V}) is a local solution to (4.1), the KKT conditions (4.2) imply that

$$\mu^T \hat{U}_j - \hat{U}_j^T \nabla g(\hat{X}) \hat{V}_j^* = \lambda \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, 2, \dots, s,$$

i.e.,

$$\left\langle \mu \mathbf{1}_d^T - \nabla g(\hat{X}), \hat{U}_j \hat{V}_j^T \right\rangle = \lambda \|\hat{U}_j\|_2 \|\hat{V}_j\|_2, \quad j = 1, 2, \dots, s.$$

It follows from the condition $\sum_{j=1}^s \alpha_j \hat{U}_j \hat{V}_j^T = \mathbf{0}_{d \times d}$ that

$$\sum_{j=1}^s \alpha_j \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 = \lambda^{-1} \left\langle \mu \mathbf{1}_d^T - \nabla g(\hat{X}), \sum_{j=1}^s \alpha_j \hat{U}_j \hat{V}_j^T \right\rangle = 0. \quad \square$$

A.8. Proof of Theorem 11.

Proof. It is obvious that when $\kappa > 0$ is sufficiently small, $\bar{U} \in \mathcal{U}^{d \times (s+1)}$, $\bar{V} \in \mathcal{V}^{d \times (s+1)}$, and (\bar{U}, \bar{V}) is feasible. We next analyze the first-order perturbation of F_λ with respect to κ . Define the difference $\Delta := F_\lambda(\bar{U}, \bar{V}) - F_\lambda(\hat{U}, \hat{V})$,

$$(A.33) \quad \begin{aligned} \Delta = & \left[g(\bar{U}\bar{V}^T) - g(\hat{U}\hat{V}^T) \right] + \lambda \|\kappa \bar{\mathbf{u}}\|_2 \left\| (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{v}} \right\|_2 \\ & + \lambda \sum_{j=1}^s \left(\left\| \text{diag}\{\mathbf{1}_d - \kappa \bar{\mathbf{u}}\} \hat{U}_j \right\|_2 \|\hat{V}_j\|_2 - \|\hat{U}_j\|_2 \|\hat{V}_j\|_2 \right). \end{aligned}$$

Because

$$\begin{aligned} \bar{U}\bar{V}^T &= \text{diag}\{\mathbf{1}_d - \kappa \bar{\mathbf{u}}\} \hat{X} + (\kappa \bar{\mathbf{u}}) \left[(\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{v}} \right]^T \\ &= \hat{X} + \kappa \left\{ (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{u}} \bar{\mathbf{v}}^T - \text{diag}\{\bar{\mathbf{u}}\} \hat{X} \right\}, \end{aligned}$$

by the Taylor expansion of g at \hat{X} , when $\kappa \rightarrow 0+$,

$$(A.34) \quad g(\bar{U}\bar{V}^T) - g(\hat{X})$$

$$(A.35) \quad = \kappa \left\langle \nabla g(\hat{X}), (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{u}} \bar{\mathbf{v}}^T - \text{diag}\{\bar{\mathbf{u}}\} \hat{X} \right\rangle + O(\kappa^2)$$

$$(A.36) \quad = \kappa \left\{ (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{u}}^T \nabla g(\hat{X}) \bar{\mathbf{v}} - \left\langle \nabla g(\hat{X}), \text{diag}\{\bar{\mathbf{u}}\} \hat{X} \right\rangle \right\} + O(\kappa^2).$$

Under the condition $\hat{U}_j \neq \mathbf{0}_d$, $j = 1, 2, \dots, s$, the Euclidean norm is differentiable at \hat{U}_j . Therefore, when $\kappa \rightarrow 0+$,

$$(A.37) \quad \begin{aligned} & \left\| \text{diag}\{\mathbf{1}_d - \kappa \bar{\mathbf{u}}\} \hat{U}_j \right\|_2 - \|\hat{U}_j\|_2 = \left\| \hat{U}_j - \kappa \text{diag}\{\bar{\mathbf{u}}\} \hat{U}_j \right\|_2 - \|\hat{U}_j\|_2 \\ & = \left\langle \frac{\hat{U}_j}{\|\hat{U}_j\|_2}, -\kappa \text{diag}\{\bar{\mathbf{u}}\} \hat{U}_j \right\rangle + O(\kappa^2) = -\frac{\kappa}{\|\hat{U}_j\|_2} (\hat{U}_j)^T \text{diag}\{\bar{\mathbf{u}}\} \hat{U}_j + O(\kappa^2). \end{aligned}$$

Plugging (A.34) and (A.37) into (A.33), we obtain

$$(A.38) \quad \begin{aligned} \Delta = & \kappa \left\{ (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{u}}^T \nabla g(\hat{X}) \bar{\mathbf{v}} - \left\langle \nabla g(\hat{X}), \text{diag}\{\bar{\mathbf{u}}\} \hat{X} \right\rangle \right. \\ & \left. + \lambda (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} - \lambda \sum_{j=1}^s \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \hat{U}_j^T \text{diag}\{\bar{\mathbf{u}}\} \hat{U}_j \right\} + O(\kappa^2). \end{aligned}$$

(\hat{U}, \hat{V}) is a local solution to (4.1), so it follows from the results in Theorem 7 that there exists a Lagrangian multiplier $\mu \in \mathbb{R}^d$ such that (4.2) holds. Especially we have

$$\left[\mu \mathbf{1}_s^T - \nabla g(\hat{X}) \hat{V} \right]_+ = \lambda \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\},$$

which implies that

$$\begin{aligned} & \lambda \sum_{j=1}^s \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \hat{U}_j^T \text{diag}\{\bar{\mathbf{u}}\} \hat{U}_j = \left\langle \lambda \hat{U} \text{diag} \left\{ \frac{\|\hat{V}_j\|_2}{\|\hat{U}_j\|_2} \right\}_{j=1}^s, \text{diag}\{\bar{\mathbf{u}}\} \hat{U} \right\rangle \\ & = \left\langle \mu \mathbf{1}_s^T - \nabla g(\hat{X}) \hat{V}, \text{diag}\{\bar{\mathbf{u}}\} \hat{U} \right\rangle = \mu^T \bar{\mathbf{u}} - \left\langle \nabla g(\hat{X}), \text{diag}\{\bar{\mathbf{u}}\} \hat{X} \right\rangle. \end{aligned}$$

Then, (A.38) can be reduced to

$$\begin{aligned}\Delta &= \kappa \left\{ (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \bar{\mathbf{u}}^T \nabla g(\hat{X}) \bar{\mathbf{v}} + \lambda (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} - \bar{\mathbf{u}}^T \mu \right\} + O(\kappa^2) \\ &= \kappa (\bar{\mathbf{v}}^T \mathbf{1}_d)^{-1} \left\{ \lambda - \bar{\mathbf{u}}^T \left[\mu \mathbf{1}_d^T - \nabla g(\hat{X}) \right] \bar{\mathbf{v}} \right\} + O(\kappa^2), \quad (\kappa \rightarrow 0+).\end{aligned}$$

Under the condition $\bar{\mathbf{u}}^T [\mu \mathbf{1}_d^T - \nabla g(\hat{X})] \bar{\mathbf{v}} > \lambda$, $\Delta < 0$ for some sufficiently small $\kappa > 0$. \square

A.9. Proof of Theorem 10.

Proof. Recall that we rewrite the factorization (4.1) into (5.1). According to Lemma 3 and Theorem 1 in [4], it suffices to show that

1. the objective function in (5.1) is a KL function,
2. Assumptions 1 and 2 in [4] hold for our problem,
3. $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ is bounded.

Because the feasible set of (4.1) is bounded, condition 3 naturally holds. Since $F_\lambda \in C^2$ and $\chi_{\mathcal{U}^{d \times s}}, \chi_{\mathcal{V}^{d \times s}}$ are lower semicontinuous, Assumption 1 in [4] is also satisfied.

We next verify the KL property of \tilde{F}_λ . $g(UV^T)$ is a polynomial function, so it is semialgebraic. The Euclidean norm $\|\cdot\|_2$ is a semialgebraic function, and the finite sums and products of semialgebraic functions are also semialgebraic, so $\sum_{j=1}^s \|U_j\|_2 \|V_j\|_2$ is semialgebraic. $\mathcal{U}^{d \times s}$ and $\mathcal{V}^{d \times s}$ consist of linear constraints, so they are semialgebraic sets. Hence, their characteristic functions $\chi_{\mathcal{U}^{d \times s}}$ and $\chi_{\mathcal{V}^{d \times s}}$ are semialgebraic. It can be concluded that \tilde{F}_λ is semialgebraic. According to Theorem 3 in [4], it satisfies the KL property.

As for Assumption 2, the objective function in (5.1) is nonnegative, so it is bounded below. Assumption 2(i) holds. Additionally, we have shown that, when V is fixed, the partial gradient $\nabla_U F_\lambda(U, V)$ is globally Lipschitz with moduli $L_1(V)$. Due to the boundedness of V and the continuity of $L_1(V)$, $L_1(V)$ has positive lower and upper bounds. Likewise, $\nabla_V F_\lambda(U, V)$ also has positive lower and upper bounds for any fixed U . For this reason, Assumptions 2(ii) and (iii) in [4] are satisfied. Because $F_\lambda(U, V) \in C^2$, Assumption 2(iv) holds.

According to Theorem 1 in [4], the sequence $\{(U^k, V^k)\}_{k \in \mathbb{N}}$ has the finite length property and globally converges to a stationary point. \square

A.10. Proof of Theorem 13.

Proof. In the following, we fix a trajectory $(s_{i_0}, s_{i_1}, \dots, s_{i_m}) \in \mathcal{S}^{m+1}$ and calculate the probability

$$I = \mathbb{P}[(\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_m) = (\tilde{s}_{k_1}, \tilde{s}_{k_2}, \dots, \tilde{s}_{k_m}), (\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m) = (s_{i_0}, s_{i_1}, \dots, s_{i_m})]$$

for each possible $(\tilde{s}_{k_1}, \tilde{s}_{k_2}, \dots, \tilde{s}_{k_m}) \in \tilde{\mathcal{S}}^m$.

According to (6.2), the random variables $\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m$ are conditionally independent given $(\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_m)$. Therefore,

(A.39)

$$I = \mathbb{P}[(\hat{Z}_1, \dots, \hat{Z}_m) = (\tilde{s}_{k_1}, \dots, \tilde{s}_{k_m})] \prod_{t=0}^m \mathbb{P}[\hat{Y}_t = s_{i_t} \mid (\hat{Z}_1, \dots, \hat{Z}_m) = (\tilde{s}_{k_1}, \dots, \tilde{s}_{k_m})].$$

Substituting (6.1) and (6.2) into (A.39), we have
(A.40)

$$\begin{aligned} I &= \mathbb{P}[\hat{Z}_1 = \tilde{s}_{k_1}] \prod_{t=1}^{m-1} \mathbb{P}[\hat{Z}_{t+1} = \tilde{s}_{k_{t+1}} \mid \hat{Z}_t = \tilde{s}_{k_t}] \prod_{t=1}^{m-1} \mathbb{P}[\hat{Y}_t = s_{i_t} \mid \hat{Z}_t = \tilde{s}_{k_t}, \hat{Z}_{t+1} = \tilde{s}_{k_{t+1}}] \\ &\quad \cdot \mathbb{P}[\hat{Y}_0 = s_{i_0} \mid \hat{Z}_1 = \tilde{s}_{k_1}] \mathbb{P}[\hat{Y}_m = s_{i_m} \mid \hat{Z}_m = \tilde{s}_{k_m}] \\ &= (\hat{U}^T \xi^\circ)_{k_1} \hat{q}_{k_1, k_2} \hat{q}_{k_2, k_3} \cdots \hat{q}_{k_{m-1}, k_m} \cdot \left(\prod_{t=1}^{m-1} \frac{\hat{v}_{i_t, k_t} \hat{u}_{i_t, k_{t+1}}}{\hat{q}_{k_t, k_{t+1}}} \right) \cdot \frac{\hat{u}_{i_0, k_1} \xi_{i_0}^\circ}{(\hat{U}^T \xi^\circ)_{k_1}} \cdot \hat{v}_{i_m, k_m} \\ &= \xi_{i_0}^\circ \hat{u}_{i_0, k_1} \cdot \left(\prod_{t=1}^{m-1} \hat{v}_{i_t, k_t} \hat{u}_{i_t, k_{t+1}} \right) \cdot \hat{v}_{i_m, k_m} = \xi_{i_0}^\circ \prod_{t=1}^m \hat{u}_{i_{t-1}, k_t} \hat{v}_{i_t, k_t}. \end{aligned}$$

In (A.40), taking the summation over all $(\tilde{s}_{k_1}, \tilde{s}_{k_2}, \dots, \tilde{s}_{k_m})$ gives

$$\begin{aligned} &\mathbb{P}[(\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m) = (s_{i_0}, s_{i_1}, \dots, s_{i_m})] \\ &= \sum_{k_1, \dots, k_m=1}^{\hat{s}} \mathbb{P}[(\hat{Z}_1, \hat{Z}_2, \dots, \hat{Z}_m) = (\tilde{s}_{k_1}, \tilde{s}_{k_2}, \dots, \tilde{s}_{k_m}), \\ &\quad (\hat{Y}_0, \hat{Y}_1, \dots, \hat{Y}_m) = (s_{i_0}, s_{i_1}, \dots, s_{i_m})] \\ &= \sum_{k_1, \dots, k_m=1}^{\hat{s}} \left(\xi_{i_0}^\circ \prod_{t=1}^m \hat{u}_{i_{t-1}, k_t} \hat{v}_{i_t, k_t} \right) \\ &= \xi_{i_0}^\circ \prod_{t=1}^m \left(\sum_{k_t=1}^{\hat{s}} \hat{u}_{i_{t-1}, k_t} \hat{v}_{i_t, k_t} \right) = \xi_{i_0}^\circ \prod_{t=1}^m \hat{x}_{i_{t-1}, i_t}, \end{aligned}$$

where we used $\hat{X} = \hat{U} \hat{V}^T$. □

REFERENCES

- [1] T. W. ANDERSON AND L. A. GOODMAN, *Statistical inference about Markov chains*, Ann. Math. Stat., 28 (1957), pp. 89–110, <http://www.jstor.org/stable/2237025>.
- [2] D. P. BERTSEKAS, *Dynamic Programming and Optimal Control*, Vols. I and II, Athena Scientific, Belmont, MA, 1995.
- [3] D. P. BERTSEKAS, *Abstract Dynamic Programming*, Athena Scientific, Belmont, MA, 2018.
- [4] J. BOLTE, S. SABACH, AND M. TEBoulLE, *Proximal alternating linearized minimization for nonconvex and nonsmooth problems*, Math. Program., 146 (2014), pp. 459–494, <https://doi.org/10.1007/s10107-013-0701-9>.
- [5] S. BURER AND R. D. MONTEIRO, *A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization*, Math. Program., 95 (2003), pp. 329–357.
- [6] E. J. CANDÈS AND B. RECHT, *Exact matrix completion via convex optimization*, Found. Comput. Math., 9 (2009), 717.
- [7] V. CHANDRASEKARAN, B. RECHT, P. A. PARRILO, AND A. S. WILLSKY, *The convex geometry of linear inverse problems*, Found. Comput. Math., 12 (2012), pp. 805–849, <https://doi.org/10.1007/s10208-012-9135-7>.
- [8] Y. CHEN AND X. YE, *Projection Onto a Simplex*, preprint, <https://arxiv.org/abs/1101.6081>, 2011.
- [9] J. E. COHEN AND U. G. ROTHBLUM, *Nonnegative ranks, decompositions, and factorizations of nonnegative matrices*, Linear Algebra Appl., 190 (1993), pp. 149–168, [https://doi.org/10.1016/0024-3795\(93\)90224-C](https://doi.org/10.1016/0024-3795(93)90224-C).
- [10] W. E, T. LI, AND E. VANDEN-ELJNDEN, *Optimal partition and effective dynamics of complex networks*, Proc. Natl. Acad. Sci. USA, 105 (2008), pp. 7907–7912.
- [11] B. D. HAEFFLE AND R. VIDAL, *Global Optimality in Tensor Factorization, Deep Learning, and Beyond*, preprint, <https://arxiv.org/abs/1506.07540>, 2015.

- [12] M. JOURNÉE, F. BACH, P.-A. ABSIL, AND R. SEPULCHRE, *Low-rank optimization on the cone of positive semidefinite matrices*, SIAM J. Optim., 20 (2010), pp. 2327–2351.
- [13] R. H. KESHAVAN AND S. OH, *A Gradient Descent Algorithm on the Grassman Manifold for Matrix Completion*, preprint, <https://arxiv.org/abs/0910.5260>, 2009.
- [14] H. KIM AND H. PARK, *Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method*, SIAM J. Matrix Anal. Appl., 30 (2008), pp. 713–730.
- [15] J. KIM AND H. PARK, *Fast nonnegative matrix factorization: An active-set-like method and comparisons*, SIAM J. Sci. Comput., 33 (2011), pp. 3261–3281.
- [16] D. D. LEE AND H. S. SEUNG, *Algorithms for non-negative matrix factorization*, Adv. Neural Inf. Process. Syst., 13 (2001), pp. 556–562. <http://papers.nips.cc/paper/1861-algorithms-for-non-negative-matrix-factorization.pdf>.
- [17] X. LI, M. WANG, AND A. ZHANG, *Estimation of Markov Chain via Rank-constrained Likelihood*, Proc. Mach. Learn. Res., 80 (2018), pp. 3033–3042.
- [18] C.-J. LIN, *Projected gradient methods for nonnegative matrix factorization*, Neural Comput., 19 (2007), pp. 2756–2779, <https://doi.org/10.1162/neco.2007.19.10.2756>.
- [19] B. MISHRA, G. MEYER, F. BACH, AND R. SEPULCHRE, *Low-rank optimization with trace norm penalty*, SIAM J. Optim., 23 (2013), pp. 2124–2149.
- [20] B. RECHT, M. FAZEL, AND P. A. PARRILO, *Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization*, SIAM Rev., 52 (2010), pp. 471–501, <https://doi.org/10.1137/070697835>.
- [21] R. S. SUTTON AND A. G. BARTO, *Reinforcement Learning: An Introduction*, MIT Press, Cambridge, MA, 2018.
- [22] N. TLC, *NYC Taxi and Limousine Commission (TLC) Trip Record Data*, http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml, 2018.
- [23] S. A. VAVASIS, *On the complexity of nonnegative matrix factorization*, SIAM J. Optim., 20 (2010), pp. 1364–1377, <https://doi.org/10.1137/070709967>.
- [24] Z. WEN, W. YIN, AND Y. ZHANG, *Solving a low-rank factorization model for matrix completion by a nonlinear successive over-relaxation algorithm*, Math. Program. Comput., 4 (2012), pp. 333–361.
- [25] L. F. YANG, V. BRAVERMAN, T. ZHAO, AND M. WANG, *Online Factorization and Partition of Complex Networks From Random Walks*, Proceedings of the Conference on Uncertainty in Artificial Intelligence, <http://auai.org/uai2019/proceedings/papers/299.pdf>, 2017.
- [26] A. ZHANG AND M. WANG, *Spectral State Compression of Markov Processes*, IEEE Trans. Inform. Theory, to appear, 2019.