

GLOBAL IDENTIFIABILITY OF OVERCOMPLETE DICTIONARY LEARNING VIA L_1 AND VOLUME MINIMIZATION

Yuchen Sun & Kejun Huang

Department of Computer and Information Science and Engineering
 University of Florida
 Gainesville, FL 32611, USA
 {yuchen.sun, kejun.huang}@ufl.edu

ABSTRACT

We propose a novel formulation for dictionary learning with an overcomplete dictionary, i.e., when the number of atoms is larger than the dimension of the dictionary. The proposed formulation consists of a weighted sum of ℓ_1 norms of the rows of the sparse coefficient matrix plus the log of the matrix volume of the dictionary matrix. The main contribution of this work is to show that this novel formulation guarantees global identifiability of the overcomplete dictionary, under a mild condition that the sparse coefficient matrix satisfies a strong scattering condition in the hypercube. Furthermore, if every column of the coefficient matrix is sparse and the dictionary guarantees ℓ_1 recovery, then the coefficient matrix is identifiable as well. This is a major breakthrough for not only dictionary learning but also general matrix factorization models as identifiability is guaranteed even when the latent dimension is higher than the ambient dimension. We also provide a probabilistic analysis and show that if the sparse coefficient matrix is generated from the widely adopted sparse-Gaussian model, then the $m \times k$ overcomplete dictionary is globally identifiable if the sample size is bigger than a constant times $(k^2/m) \log(k^2/m)$ with overwhelming probability. Finally, we propose an algorithm based on alternating minimization to solve the new proposed formulation.

1 INTRODUCTION

Dictionary learning (DL) amounts to factor a data matrix as $\mathbf{X} = \mathbf{A}\mathbf{S}$ where \mathbf{S} is sparse (Tošić & Frossard, 2011), which may also be known as sparse coding (Olshausen & Field, 1997) or sparse component analysis (Georgiev et al., 2005) in various fields. Treating $\mathbf{X} \in \mathbb{R}^{m \times n}$ or $\mathbb{C}^{m \times n}$ as a collection of data samples as its columns, this factorization means that each sample is a *sparse* combination of the columns of \mathbf{A} , or in other words atoms of the dictionary. Unlike the task of compressive sensing or sparse vector recovery, in which case the dictionary matrix \mathbf{A} is given, dictionary learning tries to find both \mathbf{A} and \mathbf{S} , therefore the problem is a lot more challenging. Depending on the shape of the dictionary matrix \mathbf{A} , we may seek to find a complete dictionary if \mathbf{A} is square or an overcomplete dictionary if \mathbf{A} is wide. In this paper we focus on overcomplete dictionary learning, therefore $\mathbf{A} \in \mathbb{R}^{m \times k}$ and $\mathbf{S} \in \mathbb{R}^{k \times n}$ (or $\mathbb{C}^{m \times k}$ and $\mathbb{C}^{k \times n}$, respectively) where $m < k$.

Dictionary learning has found numerous applications in signal denoising (Elad & Aharon, 2006), audio coding (Plumbley et al., 2009), and medical imaging (Tošić et al., 2010), to name just a few. On the theory side, most of the existing works have focused on algorithm design. Famous algorithms include k -SVD (Aharon et al., 2006a) and online dictionary learning (Mairal et al., 2009), among numerous other algorithms based on generic nonconvex algorithm design with guarantee of convergence to a stationary point. More recently, there has appeared a line of research that attempts to show global optimality for dictionary learning under more restrictive assumptions, such as (Spielman et al., 2012; Agarwal et al., 2016; Arora et al., 2014; 2015; Sun et al., 2016a;b; Rambhatla et al., 2019; Bai et al., 2019; Zhai et al., 2020a;b; Shen et al., 2020; Tolooshams & Ba, 2022).

1.1 PRIOR WORK ON IDENTIFIABILITY OF DL

A matrix factorization model without any additional assumptions on the latent factors is known to be not unique, since we can always “insert” an invertible matrix \mathbf{W} and \mathbf{W}^{-1} as $\mathbf{X} = \tilde{\mathbf{A}}\tilde{\mathbf{S}}$ where $\tilde{\mathbf{A}} = \mathbf{A}\mathbf{W}^{-1}$ and $\tilde{\mathbf{S}} = \mathbf{W}\mathbf{S}$, and one cannot distinguish whether \mathbf{S} or $\tilde{\mathbf{S}}$ are the groundtruth sources. Furthermore, if the dictionary \mathbf{A} is overcomplete with $m < k$, columns of $\tilde{\mathbf{S}}$ could include additional components that are orthogonal to the row space of \mathbf{A} , i.e., $\tilde{\mathbf{S}} = \mathbf{W}(\mathbf{S} + \mathbf{B})$ where $\mathbf{A}\mathbf{B} = \mathbf{0}$, while we still have $\mathbf{X} = \mathbf{A}\mathbf{S} = \mathbf{A}\mathbf{W}^{-1}\mathbf{W}(\mathbf{S} + \mathbf{B}) = \tilde{\mathbf{A}}\tilde{\mathbf{S}}$. If, however, a learning criterion $q(\mathbf{A}, \mathbf{S})$ is imposed so that the resulting ambiguities can only be permutations and scaling, then we say the model is identifiable, as is formalized as follows. Notice that we differentiate the identifiability of just the dictionary \mathbf{A} and the whole factorization model, which is indeed not equivalent when the dictionary is overcomplete.

Definition 1 (Identifiability). Consider the generative model $\mathbf{X} = \mathbf{A}_\natural\mathbf{S}_\natural$, where \mathbf{A}_\natural and \mathbf{S}_\natural are the groundtruth latent factors. Let $(\mathbf{A}_\star, \mathbf{S}_\star)$ be optimal for an identification criterion q

$$(\mathbf{A}_\star, \mathbf{S}_\star) = \arg \min_{\mathbf{X}=\mathbf{A}\mathbf{S}} q(\mathbf{A}, \mathbf{S}).$$

If \mathbf{A}_\natural and/or \mathbf{S}_\natural satisfy some condition such that for any $(\mathbf{A}_\star, \mathbf{S}_\star)$, there exist a permutation matrix $\mathbf{\Pi}$ and a diagonal matrix \mathbf{D} such that $\mathbf{A}_\natural = \mathbf{A}_\star\mathbf{D}\mathbf{\Pi}$, then we say \mathbf{A}_\natural is essentially identifiable, up to permutation and scaling, under that condition; if we further have that $\mathbf{S}_\natural = \mathbf{\Pi}^T\mathbf{D}^{-1}\mathbf{S}_\star$, then we say that the matrix factorization model is essentially identifiable, up to permutation and scaling, under that condition.

When dictionary learning was first proposed, a common learning criterion $q(\mathbf{A}, \mathbf{S})$ is simply the total number of nonzeros in \mathbf{S} , sometimes also called the ℓ_0 (pseudo-)norm $\|\mathbf{S}\|_0$. If every column of \mathbf{S}_\natural is s -sparse, then via some combinatorial calculation, it has been shown that the ℓ_0 -norm minimization criterion guarantees identifiability if the spark of \mathbf{A} is at least $2s$ and the sample size n is $O((s+1)\binom{k}{s})$ (Aharon et al., 2006b; Hillar & Sommer, 2015; Garfinkle & Hillar, 2019). The main drawback is that the required sample size n is usually too large to be practical. Cohen & Gillis (2019) reduced the sample complexity down to $O(k^3/(k-s)^2)$, but also restricted the dictionary to be complete, i.e., $m \geq k$.

Another famous learning criterion for DL, inspired by the success of compressive sensing (Donoho, 2006; Candès & Wakin, 2008), is the following formulation with $q(\mathbf{A}, \mathbf{S})$ being the summation of the absolute values of \mathbf{S} plus indicator functions that columns of \mathbf{A} have bounded ℓ_2 norms:

$$\underset{\mathbf{A}, \mathbf{S}}{\text{minimize}} \|\mathbf{S}\|_1 \quad \text{subject to } \mathbf{X} = \mathbf{A}\mathbf{S}, \|\mathbf{A}\mathbf{e}_c\|_2 \leq 1, c = 1, \dots, k. \quad (1)$$

Identifiability results based on the ℓ_1 norm formulation have been predominantly local, meaning the model is identifiable within a neighborhood of the groundtruth factors $(\mathbf{A}^\natural, \mathbf{S}^\natural)$, while the dictionary is restricted to be complete (and incoherent) (Gribonval & Schnass, 2010; Wu & Yu, 2017; Wang et al., 2020), with the sole exception of (Geng & Wright, 2014) for overcomplete DL. The advantage is that the sample size requirement is typically down to $O(k \log k)$ and allows the existence of dense outliers. Global identifiability is achieved by Hu & Huang (2023a); Sun & Huang (2024) by using a matrix volume criterion $|\det \mathbf{A}|$ while constraining the ℓ_1 norms of the rows of \mathbf{S} with the same sample complexity, although as the criterion suggests it only applies to complete dictionaries.

1.2 THIS PAPER

In this paper, we propose the following novel formulation for overcomplete dictionary learning, and show that global identifiability can be achieved under mild conditions:

$$\underset{\mathbf{A}, \mathbf{S}}{\text{minimize}} \frac{1}{2} \log \det \mathbf{A}\mathbf{A}^T + \max_{\|d\|_2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^T \mathbf{S}\|_1 \quad \text{subject to } \mathbf{X} = \mathbf{A}\mathbf{S} \quad (2)$$

This means the learning criterion $q(\mathbf{A}, \mathbf{S})$ consists of two parts: a weighted sum of the ℓ_1 norms of the rows of \mathbf{S} and a term that is proportional to the “volume” of the dictionary matrix \mathbf{A} —for an overcomplete dictionary, the volume is defined as $\det \mathbf{A}\mathbf{A}^T$ (Ben-Israel, 1992). Our contributions are as follows:

1. We give a deterministic characterization of global identifiability of overcomplete dictionary learning via solving (2). Our analysis shows that a sufficient condition is that 1) every column of \mathbf{S}_\natural is at most s -sparse and \mathbf{A} is a dictionary that guarantees exact recovery of all s -sparse vectors via ℓ_1 minimization, and 2) the cellular hull of \mathbf{S}_\natural is m -sufficiently scattered in the k -hypercube $[-1, 1]^k$. The resulting identifiability condition is almost minimal: the first condition is obviously necessary as otherwise \mathbf{S}_\natural would not be identifiable even if the overcomplete \mathbf{A}_\natural is correctly recovered; the second condition is a slightly stronger condition than that of complete dictionary learning. It is appealing to see that no other conditions are needed to guarantee global identifiability.
2. We further provide a probabilistic characterization of when a randomly generated factor matrix satisfies the aforementioned identifiability conditions. Since we only require \mathbf{A}_\natural to guarantee exact recovery of s -sparse vectors via ℓ_1 minimization, for which there exist numerous work on this topic (such as when an i.i.d. Gaussian matrix satisfies the restricted isometry property with high probability), we will be focusing on studying the sample complexity of \mathbf{S} . We adopt the sparse-Gaussian model, i.e., every column contains at most s nonzero values that are drawn from i.i.d. standard normal and show that the resulting \mathbf{S} satisfies the m -strongly scattered in the k -hypercube $[-1, 1]^k$ with overwhelming probability if $k < m$ and n is $O((k^2/m) \log(k^2/m))$. Notice that it is again a sharp generalization of complete dictionary learning with sample complexity $O(k \log(k))$ (Hu & Huang, 2023a), and a factor of k better than the works that focus on global optimality guarantees of overcomplete DL (Agarwal et al., 2016; Rambhatla et al., 2019).
3. We propose an alternating minimization algorithm for the novel identification criterion of overcomplete DL. The formulation is first modified slightly by moving the exact factorization constraint as a data fidelity term in the objective function, then the overcomplete dictionary and the sparse coefficient matrices are updated alternately via a gradient-type step. As the problem is NP-hard, no known algorithm is able to guarantee convergence to a global optimum. The proposed algorithm is applied to synthetically generated data to demonstrate that global identifiability can indeed be guaranteed via solving (2).

2 IDENTIFIABILITY ANALYSIS

In this section, we provide analysis on when solving (2) guarantees the exact recovery of the overcomplete dictionary \mathbf{A}_\natural and/or the sparse coefficient matrix \mathbf{S}_\natural . In Definition 1 we mentioned that it is acceptable to recover \mathbf{A}_\natural up to column permutation and scaling. While column permutation does not affect the objective value of (2), column scaling does. Therefore, we first study the optimal scaling that is induced from solving (2), which provides important insights into the subsequent analysis of identifiability. We provide both a deterministic condition and a probabilistic generative model that guarantees identifiability with overwhelming probability.

2.1 OPTIMAL SCALING

Lemma 1. *Let $(\mathbf{A}_\star, \mathbf{S}_\star)$ be an optimal solution of (2), then*

$$\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \sqrt{\left[\mathbf{A}_\star^\top (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star \right]_{cc}} = d_{\star c}, \quad c = 1, \dots, k. \quad (3)$$

where $d_{\star c}$ are the optimal weights that reach the maximum of $\sum_c d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1$.

Proof. If $(\mathbf{A}_\star, \mathbf{S}_\star)$ is feasible for (2), then so is $(\mathbf{A}_\star \Psi, \Psi^{-1} \mathbf{S}_\star)$ where Ψ is a diagonal matrix with c th diagonal entry ψ_c . Plugging $(\mathbf{A}_\star \Psi, \Psi^{-1} \mathbf{S}_\star)$ into the objective of (2) and optimize with respect to Ψ while fixing $(\mathbf{A}_\star, \mathbf{S}_\star)$, then $\Psi = \mathbf{I}$ should be an optimal solution. Taking the derivative with respect to ψ_c and setting it equal to zero, we get

$$\mathbf{a}_c^\top \left(\mathbf{A}_\star \Psi^2 \mathbf{A}_\star^\top \right)^{-1} \mathbf{a}_c \psi_c - d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 / \psi_c^2 = 0,$$

where \mathbf{a}_c^\top is the c th row of \mathbf{A}_\star . If $\psi_c = 1$ is optimal, then we must have

$$d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \left[\mathbf{A}_\star^\top (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star \right]_{cc}. \quad (4)$$

To maximize $\sum_c d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1$ subject to $\|\mathbf{d}\|_2^2 = m$, we know from the Cauchy-Schwarz inequality that \mathbf{d}_\star should be chosen as some scalar α times $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1$ such that

$$\|\mathbf{d}_\star\|_2^2 = \sum_{c=1}^k \alpha^2 \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1^2 = m.$$

Plugging $\mathbf{d}_\star = \alpha \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1$ into (4) and sum over $c = 1, \dots, k$ shows

$$\sum_{c=1}^k \alpha \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1^2 = \sum_{c=1}^k \left[\mathbf{A}_\star^\top (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star \right]_{cc} = \text{Tr} \mathbf{A}_\star^\top (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star = \text{Tr} (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star \mathbf{A}_\star^\top = m.$$

This means $\alpha = 1$, and therefore (3) holds. \square

2.2 IDENTIFIABILITY ANALYSIS

Assumption 1. The columns of \mathbf{A}_\natural and rows of \mathbf{S}_\natural are scaled and counter-scaled to satisfy:

$$\|\mathbf{e}_c^\top \mathbf{S}_\natural\|_1 = \sqrt{\left[\mathbf{A}_\natural^\top (\mathbf{A}_\natural \mathbf{A}_\natural^\top)^{-1} \mathbf{A}_\natural \right]_{cc}}, \quad c = 1, \dots, k. \quad (5)$$

Assumption 2. Rows of \mathbf{A}_\natural and \mathbf{S}_\natural are both linearly independent. Matrix \mathbf{A}_\natural does not contain zero columns.

Assumption 3 (m -strongly scattered in the k -hypercube). Let C_k denote the k -hypercube $C_k = \{\mathbf{x} \in \mathbb{R}^k \mid \|\mathbf{x}\|_\infty \leq 1\}$. Define \mathcal{B}_m as the following set

$$\mathcal{B}_m = \left\{ \text{Diag}(\|\mathbf{q}_1\|_2, \dots, \|\mathbf{q}_k\|_2)^\dagger \mathbf{Q} \mathbf{p} \mid \forall \mathbf{Q} \in \mathbb{R}^{k \times m} : \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}, \mathbf{p} \in \mathbb{R}^m : \|\mathbf{p}\|_2 = 1 \right\},$$

where \mathbf{q}_c denotes the c th row of \mathbf{Q} . A set $\mathcal{S} \in \mathbb{R}^k$ is m -strongly scattered in the k -hypercube if:

1. $\mathcal{B}_m \subseteq \mathcal{S} \subseteq C_k$;
2. $\partial \mathcal{B}_m \cap \partial \mathcal{S} = \{\text{Diag}(\|\mathbf{q}_1\|_2, \dots, \|\mathbf{q}_k\|_2)^\dagger \mathbf{Q} \mathbf{q} / \|\mathbf{q}\|_2 \mid \mathbf{q} \text{ are rows of } \mathbf{Q} \text{ with } \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}\}$, where ∂ denotes the boundary of the set.

The definition of the set \mathcal{B}_m involves all $k \times m$ matrices with orthonormal columns \mathbf{Q} and normalized m -vectors \mathbf{p} . To see that \mathcal{B}_m is indeed a subset of C , notice that the rows of $\mathbf{P} \mathbf{Q}$ are either zero or unit norms by construction, so all elements of the vector $\mathbf{P} \mathbf{Q} \mathbf{p}$ are in $[-1, 1]$ using Cauchy-Schwarz inequality. Assumption 3 is equivalent to the sufficiently scattered condition proposed by Hu & Huang (2023a) when $m = k$, and more restrictive when $m < k$. As we will see, such restriction will be useful to establish identifiability for overcomplete DL. The sufficiently scattered condition has many variations for various identifiable unsupervised learning models, such as nonnegative matrix factorization (Huang et al., 2013; Fu et al., 2015; Huang et al., 2016; 2018), simplicial representation learning (Fu et al., 2015; Lin et al., 2015; Huang & Fu, 2019), and bounded component analysis (Tatli & Erdogan, 2021; Hu & Huang, 2023b; 2024). However, to the best of our knowledge, this is the first variation that is capable of guaranteeing identifiability when the latent dimension is *higher* than the ambient dimension. An illustration of 2-strongly scattered in the 3-hypercube is shown on the right of Figure 1, compared with 3-strongly scattered on the left, which is equivalent to the sufficiently scattered condition presented by Hu & Huang (2023a).

Assumption 3 will be imposed on the cellular hull of \mathbf{S}_\natural , which is defined as follows:

Definition 2 (Cellular hull). The cellular hull of a finite set of vectors $\{\mathbf{s}_1, \dots, \mathbf{s}_n\}$, stacked as the columns of the matrix \mathbf{S} , is

$$\text{cell}(\mathbf{S}) = \left\{ \mathbf{S} \boldsymbol{\theta} \mid \|\boldsymbol{\theta}\|_\infty \leq 1 \right\}.$$

Consider the groundtruth sparse coefficient matrix \mathbf{S}_\natural , if we rescale its rows to have unit ℓ_1 norms, denoted as $\tilde{\mathbf{S}}_\natural$, then $\text{cell}(\tilde{\mathbf{S}}_\natural) \subseteq C_k$ due to Hölder’s inequality $|\mathbf{a}^\top \mathbf{b}| \leq \|\mathbf{a}\|_1 \|\mathbf{b}\|_\infty$. For identifiability of overcomplete DL, we would require $\text{cell}(\tilde{\mathbf{S}}_\natural)$ to be m -strongly scattered in the k -hypercube, as formally stated as follows:

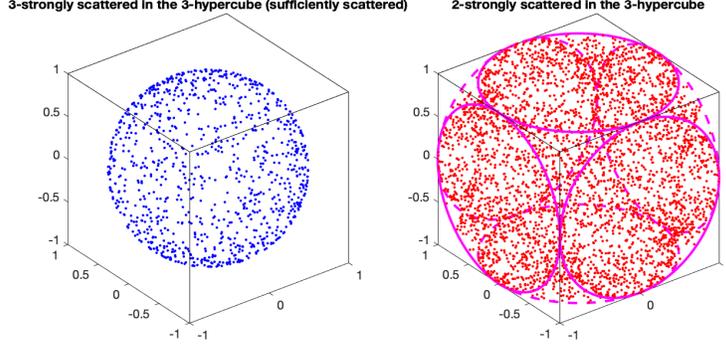


Figure 1: An illustration of \mathcal{B}_2 in \mathbb{R}^3 on the right. In comparison, \mathcal{B}_3 in \mathbb{R}^3 is the Euclidean ball illustrated on the left. While \mathcal{B}_3 touches the boundary of $[-1, 1]^3$ at only 6 points $\pm e_1, \pm e_2$, and $\pm e_3$, \mathcal{B}_2 touches each face of $[-1, 1]^3$ at a circle with radius 1, as shown in magenta on the right. In the context of dictionary learning, a 3×3 complete dictionary is identifiable if $\text{cell}(\mathbf{S}_\dagger) \supseteq \mathcal{B}_3$, shown on the left, while a 2×3 overcomplete dictionary is identifiable if $\text{cell}(\mathbf{S}_\dagger) \supseteq \mathcal{B}_2$, shown on the right.

Theorem 1. Consider the overcomplete DL model $\mathbf{X} = \mathbf{A}_\dagger \mathbf{S}_\dagger$, where $\mathbf{A}_\dagger \in \mathbb{R}^{m \times k}$ is the groundtruth mixing matrix and $\mathbf{S}_\dagger \in \mathbb{R}^{k \times n}$ is the groundtruth sparse coefficient matrix. Suppose \mathbf{A}_\dagger and \mathbf{S}_\dagger satisfies Assumptions 1 and 2. Furthermore, let $\tilde{\mathbf{S}}_\dagger$ denote the matrix obtained from rescaling the rows of \mathbf{S}_\dagger to have unit ℓ_1 norms, and assume that $\text{cell}(\tilde{\mathbf{S}}_\dagger)$ is m -strongly scattered in the k -hypercube. Then for any solution of (2), denoted as $(\mathbf{A}_\star, \mathbf{S}_\star)$, there exist a permutation matrix $\mathbf{\Pi}$ and a diagonal matrix \mathbf{D} such that $\mathbf{A}_\dagger = \mathbf{A}_\star \mathbf{D} \mathbf{\Pi}$. In other words, an overcomplete dictionary \mathbf{A}_\dagger is identifiable if the groundtruth \mathbf{A}_\dagger and \mathbf{S}_\dagger satisfies Assumptions 1, 2, and $\text{cell}(\tilde{\mathbf{S}}_\dagger)$ satisfies Assumption 3.

Proof sketch. Assumption 2 asserts that rows of \mathbf{A}_\dagger are linearly independent, so there exists a $k \times m$ matrix \mathbf{Q} with orthonormal columns that spans the row space of \mathbf{A}_\dagger , then $\mathbf{Q}^\top \mathbf{Q} = \mathbf{I}$ and

$$\mathbf{Q} \mathbf{Q}^\top = \mathbf{A}_\dagger^\top (\mathbf{A}_\dagger \mathbf{A}_\dagger^\top)^{-1} \mathbf{A}_\dagger.$$

Using the diagonal matrix \mathbf{D}_\dagger , which is defined as

$$[\mathbf{D}_\dagger]_{cc} = \sqrt{\left[\mathbf{A}_\dagger^\top (\mathbf{A}_\dagger \mathbf{A}_\dagger^\top)^{-1} \mathbf{A}_\dagger \right]_{cc}}, \quad c = 1, \dots, k, \quad (6)$$

and \mathbf{q}_c^\top as rows of \mathbf{Q} defined in Assumption 3, we have

$$[\mathbf{D}_\dagger]_{cc} = \|\mathbf{e}_c^\top \mathbf{S}_\dagger\|_1 = \sqrt{[\mathbf{Q} \mathbf{Q}^\top]_{cc}} = \|\mathbf{q}_c\|_2. \quad (7)$$

For two wide $k \times n$ matrices \mathbf{S}_\dagger and \mathbf{S}_\star , there exist a $k \times k$ matrix \mathbf{W} and $k \times n$ matrix \mathbf{B} such that

$$\mathbf{S}_\star = \mathbf{W} \mathbf{S}_\dagger + \mathbf{B},$$

where $\mathbf{W} = \mathbf{S}_\dagger^\dagger \mathbf{S}_\star$ and rows of \mathbf{B} are in the null space of \mathbf{S}_\dagger , i.e., $\mathbf{S}_\dagger \mathbf{B}^\top = 0$. Denote \mathbf{w}_c^\top and \mathbf{b}_c^\top as the c th row of \mathbf{W} and \mathbf{B} , respectively, then

$$\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}_\dagger + \mathbf{b}_c^\top\|_1 \geq \mathbf{w}_c^\top \mathbf{S}_\dagger \boldsymbol{\theta} + \mathbf{b}_c^\top \boldsymbol{\theta} = \mathbf{w}_c^\top \mathbf{D}_\dagger \tilde{\mathbf{S}}_\dagger \boldsymbol{\theta} + \mathbf{b}_c^\top \boldsymbol{\theta}, \quad \forall \|\boldsymbol{\theta}\|_\infty \leq 1.$$

If $\text{cell}(\tilde{\mathbf{S}}_\dagger)$ is m -strongly scattered in the k -hypercube, then for all $\|\mathbf{p}\|_2 = 1$ we can find a $\|\boldsymbol{\theta}\|_\infty$ such that $\mathbf{D}_\dagger^{-1} \mathbf{Q} \mathbf{p} = \tilde{\mathbf{S}}_\dagger \boldsymbol{\theta}$. In Appendix A we show that $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}_\dagger\|_1$, as a result,

$$\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}_\dagger\|_1 = \|\mathbf{w}_c^\top \mathbf{D}_\dagger \tilde{\mathbf{S}}_\dagger\|_1 \geq \mathbf{w}_c^\top \mathbf{D}_\dagger \mathbf{D}_\dagger^{-1} \mathbf{Q} \mathbf{p} = \mathbf{w}_c^\top \mathbf{Q} \mathbf{p}, \quad \forall \|\mathbf{p}\|_2 = 1.$$

Choose $\mathbf{p} = \mathbf{Q}^\top \mathbf{w}_c / \|\mathbf{Q}^\top \mathbf{w}_c\|_2$, we have $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 \geq \|\mathbf{Q}^\top \mathbf{w}_c\|_2$. Square both sides and sum over $c = 1, \dots, k$, we get

$$\|\mathbf{Q}^\top \mathbf{W}^\top\|_F^2 = \sum_{c=1}^k \|\mathbf{Q}^\top \mathbf{w}_c\|_2^2 \leq \sum_{c=1}^k \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1^2 = \text{Tr} \mathbf{A}_\star^\top (\mathbf{A}_\star \mathbf{A}_\star^\top)^{-1} \mathbf{A}_\star = m. \quad (8)$$

Since both $(\mathbf{A}_h, \mathbf{S}_h)$ and $(\mathbf{A}_\star, \mathbf{S}_\star)$ are feasible, $\mathbf{X} = \mathbf{A}_h \mathbf{S}_h = \mathbf{A}_\star \mathbf{S}_\star = \mathbf{A}_\star (\mathbf{W} \mathbf{S}_h + \mathbf{B})$. Multiplying both sides by \mathbf{S}_h^\dagger gives

$$\mathbf{A}_h = \mathbf{A}_h \mathbf{S}_h \mathbf{S}_h^\dagger = \mathbf{A}_\star (\mathbf{W} \mathbf{S}_h + \mathbf{B}) \mathbf{S}_h^\dagger = \mathbf{A}_\star \mathbf{W}.$$

Then we have

$$\log \det \mathbf{A}_\star \mathbf{A}_\star^\top \geq \log \det \mathbf{A}_h \mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{A}_h^\top = \log \det \mathbf{A}_h \mathbf{A}_h^\top + \log \det \mathbf{Q}^\top \mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{Q} \quad (9)$$

where the first inequality is shown in Appendix A. Regarding the second term, we have

$$\log \det \mathbf{Q}^\top \mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{Q} \geq -\log \det \mathbf{Q}^\top \mathbf{W}^\top \mathbf{W} \mathbf{Q} \quad (10a)$$

$$\geq -m \log \frac{1}{m} \|\mathbf{Q}^\top \mathbf{W}^\top\|_F^2, \quad (10b)$$

where (10a) and (10b) are shown Appendix A. Combining (8), (9), and (10) shows that

$$\log \det \mathbf{A}_\star \mathbf{A}_\star^\top \geq \log \det \mathbf{A}_h \mathbf{A}_h^\top. \quad (11)$$

On the other hand, since $(\mathbf{A}_\star, \mathbf{S}_\star)$ is optimal for (2), we have

$$\frac{1}{2} \log \det \mathbf{A}_\star \mathbf{A}_\star^\top + \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 \leq \frac{1}{2} \log \det \mathbf{A}_h \mathbf{A}_h^\top + \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_h\|_1.$$

Lemma 1 and Assumption 1 shows that

$$\max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_h\|_1 = m,$$

therefore

$$\log \det \mathbf{A}_\star \mathbf{A}_\star^\top \leq \log \det \mathbf{A}_h \mathbf{A}_h^\top \quad (12)$$

Combining (11) and (12) shows that \mathbf{A}_h when scaled according to Assumption 1, or any of its column permutation and/or sign flips, is optimal for (2). In Appendix A, we complete the proof that the second requirement of m -strongly scattered in the k -hypercube guarantees that every solution must satisfy $\mathbf{A}_h = \mathbf{A}_\star \mathbf{D} \mathbf{\Pi}$, where \mathbf{D} is a diagonal matrix with only ± 1 on the diagonal, hence the overcomplete dictionary is identifiable. \square

Our analysis so far has not explicitly mentioned the sparsity of \mathbf{S}_h , which may seem somewhat counter-intuitive. The explanation is two-fold: in the next subsection, we will show that if \mathbf{S}_h follows a sparse generative model, then $\text{cell}(\tilde{\mathbf{S}}_h)$ will be m -strongly scattered in the k -hypercube with very high probability, thus sparsity is implicitly implied in Assumption 3. In particular, our analysis shows that it is necessary that every column of \mathbf{S}_h contains no more than m nonzeros. On the other hand, Theorem 1 only shows that the dictionary \mathbf{A}_h is identifiable, but for an overcomplete dictionary it does not necessarily mean that the sparse coefficient \mathbf{S}_h is identifiable. Fortunately, with the knowledge of the dictionary, the identifiability of the sparse coefficients has been studied extensively (Donoho, 2006; Candès & Wakin, 2008). Here we provide a general result.

Assumption 4. Every column of \mathbf{S}_h contains at most s nonzeros. In addition, \mathbf{A}_h is a dictionary such that for every \mathbf{s}_0 with no more than s nonzeros, \mathbf{s}_0 is the unique solution to the following optimization problem

$$\underset{\mathbf{s}}{\text{minimize}} \quad \|\mathbf{s}\|_1 \quad \text{subject to} \quad \mathbf{A}_h \mathbf{D}_h^{-1} \mathbf{s} = \mathbf{A}_h \mathbf{D}_h^{-1} \mathbf{s}_0,$$

where \mathbf{D}_h is a diagonal matrix with c th diagonal defined in (6).

Corollary 1. Consider the overcomplete DL model $\mathbf{X} = \mathbf{A}_h \mathbf{S}_h$, where $\mathbf{A}_h \in \mathbb{R}^{m \times k}$ is the groundtruth mixing matrix and $\mathbf{S}_h \in \mathbb{R}^{k \times n}$ is the groundtruth sparse coefficient matrix. Suppose \mathbf{A}_h and \mathbf{S}_h satisfies Assumptions 1–4. Then for any solution of (2), denoted as $(\mathbf{A}_\star, \mathbf{S}_\star)$, there exist a permutation matrix $\mathbf{\Pi}$ and a diagonal matrix \mathbf{D} such that $\mathbf{A}_h = \mathbf{A}_\star \mathbf{D} \mathbf{\Pi}$ and $\mathbf{S}_h = \mathbf{\Pi}^\top \mathbf{D}^{-1} \mathbf{S}_\star$.

Proof. Theorem 1 shows that \mathbf{A}_\dagger is identifiable if Assumption 1–3 are satisfied. Assumption 4 guarantees that \mathbf{S}_\dagger is uniquely determined if \mathbf{A}_\dagger is given. To see this, we fix $\mathbf{A} = \mathbf{A}_\dagger$ in (2). From Lemma 1, we know that the optimal \mathbf{d} should be the diagonal of \mathbf{D}_\dagger as defined in Assumption 4. Then optimizing (2) with respect to \mathbf{S} is equivalent to the following problem with a change-of-variable $\tilde{\mathbf{S}} = \mathbf{D}_\dagger \mathbf{S}$

$$\underset{\tilde{\mathbf{S}}}{\text{minimize}} \quad \|\tilde{\mathbf{S}}\|_1 \quad \text{subject to} \quad \mathbf{X} = \mathbf{A}_\dagger \mathbf{S}_\dagger = \mathbf{A}_\dagger \mathbf{D}_\dagger^{-1} \tilde{\mathbf{S}}.$$

If every column of \mathbf{S}_\dagger is at most s -sparse, then the optimal $\tilde{\mathbf{S}}_\star = \mathbf{D}_\dagger \mathbf{S}_\dagger$, therefore $\mathbf{S}_\star = \mathbf{S}_\dagger$. The result still holds if columns of \mathbf{A}_\dagger are permuted and/or multiplied with ± 1 . \square

2.3 SAMPLE COMPLEXITY ANALYSIS

Theorem 1 states that an overcomplete dictionary \mathbf{A}_\dagger is identifiable if Assumptions 1–3 are satisfied. Assumptions 1 and 2 are quite easy to satisfy, as it is very reasonable to assume that rows of wide matrices \mathbf{A}_\dagger and \mathbf{S}_\dagger are linearly independent, and given any \mathbf{A}_\dagger and \mathbf{S}_\dagger one can always find the scaling to satisfy (5). The most crucial condition is Assumption 3, or the fact that $\text{cell}(\tilde{\mathbf{S}}_\dagger)$ is m -strongly scattered in the k -hypercube. In this section we assume that a sparse coefficient matrix \mathbf{S} is generated from the sparse-Gaussian model, which has appeared in (Wu & Yu, 2017; Wang et al., 2020), and show that in this case Assumption 3 is satisfied with high probability. This is a different generative model than prior works that also use a volume criterion for complete DL (Hu & Huang, 2023a; Sun & Huang, 2024), in which a Bernoulli-Gaussian model is considered. This is because such a generative model cannot guarantee that every column of \mathbf{S} is at least s -sparse, thus Corollary 1 cannot be invoked to identify both \mathbf{A}_\dagger and \mathbf{S}_\dagger .

Assumption 5 (Sparse-Gaussian model). The matrix $\mathbf{S} \in \mathbb{R}^{k \times n}$ is generated from a sparse-Gaussian model with parameter $s < k$, denoted as $\mathbf{S} \sim \mathcal{SG}(s)$, if every column of \mathbf{S} is independently and identically distributed from the following process: a subset \mathcal{I} of size s is uniformly drawn from all size- s subsets of $\{1, \dots, k\}$, let $\mathbf{s} \in \mathbb{R}^k$ be such that $s_i = 0$ if $i \in \mathcal{I}$ and $s_i \sim \mathcal{N}(0, 1)$ if $i \notin \mathcal{I}$, where $\mathcal{N}(0, 1)$ stands for a standard normal distribution.

To check whether $\mathcal{B}_m \subseteq \text{cell}(\tilde{\mathbf{S}})$, where $\tilde{\mathbf{S}}$ is obtained from scaling its rows to have unit ℓ_1 norms, it is easier to equivalently check its polar version $\text{cell}(\tilde{\mathbf{S}})^\circ \subseteq \mathcal{B}_m^\circ$, where the polar of set \mathcal{S} is defined as $\mathcal{S}^\circ = \{\mathbf{x} \mid \mathbf{x}^\top \mathbf{y} \leq 1, \forall \mathbf{y} \in \mathcal{S}\}$. For $\text{cell}(\tilde{\mathbf{S}})$, its polar has a relatively simple form

$$\text{cell}(\tilde{\mathbf{S}})^\circ = \{\mathbf{w} \mid \|\mathbf{w}^\top \tilde{\mathbf{S}}\|_1 \leq 1\}.$$

The polar for \mathcal{B}_m has a more complicated form

$$\mathcal{B}_m^\circ = \{\mathbf{w} \mid \|\mathbf{Q}^\top \text{Diag}(\|\mathbf{q}_1\|_2, \dots, \|\mathbf{q}_k\|_2)^\dagger \mathbf{w}\|_2 \leq 1, \forall \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}\}.$$

Therefore, checking whether $\text{cell}(\tilde{\mathbf{S}})^\circ \subseteq \mathcal{B}_m^\circ$ is equivalent to checking whether the optimal value of the following problem equals 1:

$$\underset{\mathbf{Q}, \mathbf{w}}{\text{maximize}} \quad \|\mathbf{Q}^\top \mathbf{D}^\dagger \mathbf{w}\|_2^2 \quad \text{subject to} \quad \|\mathbf{w}^\top \tilde{\mathbf{S}}\|_1 \leq 1, \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}, \quad (13)$$

where $\mathbf{D} = \text{Diag}(\|\mathbf{q}_1\|_2, \dots, \|\mathbf{q}_k\|_2)$.

Theorem 2. Suppose $\mathbf{S} \in \mathbb{R}^{k \times n}$ is generated from the sparse-Gaussian model $\mathcal{SG}(s)$, where $s < m$, and $\tilde{\mathbf{S}}$ is obtained by scaling its rows to have unit ℓ_1 norm. Then

$$\Pr \left[\sup_{\substack{\|\mathbf{w}^\top \tilde{\mathbf{S}}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^\dagger \mathbf{w}\| > 1 \right] \leq 4 \exp \left(\frac{k}{2} \log \frac{k^2}{m} - n \frac{s^2 m}{k^3} \right). \quad (14)$$

The probability goes to zero exponentially fast as

$$n = O \left(\frac{k^2}{m} \log \frac{k^2}{m} \right).$$

The proof is relegated to Appendix B. Comparing this result to prior work on complete DL (Hu & Huang, 2023a), in which case the sample complexity is $O(k \log k)$, we see that the bounds agree when $m = k$, which is a good sign that the bound is tight. On the other hand, there is one step in the proof that shows that for overcomplete DL, it is necessary that every column of \mathbf{S} is at most s -sparse, where $s < m$; this is not required for complete DL. This shows the necessity of adopting the sparse-Gaussian model rather than the Bernoulli-Gaussian model, even if identifiability of \mathbf{S}_t is not required. In fact, even the most relaxed condition on sparse recovery would require $s < m/2$, so assuming $\mathbf{S} \sim \mathcal{SG}(s)$ with $s < m$ is a very reasonable assumption in practice.

3 ALGORITHM VIA ALTERNATING MINIMIZATION

We will now design an algorithm for the novel formulation (2) for overcomplete DL, whose global correctness in recovering the ground-truth dictionary has been theoretically established above. The main idea is similar to the (inexact) alternating optimization framework that most DL algorithms adopt. First of all, since most practical applications admit approximate factorization, we move the constraint $\mathbf{X} = \mathbf{A}\mathbf{S}$ in (2) as a penalty term as follows

$$\underset{\mathbf{A}, \mathbf{S}}{\text{minimize}} \quad \frac{\lambda}{2} \|\mathbf{A}\mathbf{S} - \mathbf{X}\|_F^2 + \frac{1}{2} \log \det \mathbf{A}\mathbf{A}^\top + \max_{\|\mathbf{d}\|_2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}\|_1, \quad (15)$$

where the hyper-parameter λ balances data fidelity and the identifiability criterion. The rest of this section will focus on designing an iterative algorithm for solving (15). We denote $(\mathbf{A}_t, \mathbf{S}_t)$ as the updates obtained at the t th iteration. In an alternating fashion, \mathbf{A}_{t+1} is obtained by fixing $\mathbf{S} = \mathbf{S}_t$, and \mathbf{S}_{t+1} is obtained by fixing $\mathbf{A} = \mathbf{A}_{t+1}$.

3.1 UPDATE OF \mathbf{A}

Fixing $\mathbf{S} = \mathbf{S}_t$, the update of \mathbf{A} amounts to solving a log-determinant regularized least squares, which is nonconvex optimization problem. We propose two types of updates:

- Gradient descent. As the gradient of $(1/2) \log \det \mathbf{A}\mathbf{A}^\top$ at \mathbf{A}_t is $(\mathbf{A}_t^\dagger)^\top$, a simple choice of update is to move along the negative gradient direction with step size γ

$$\mathbf{A}_{t+1} \leftarrow \mathbf{A}_t - \gamma \left(\lambda (\mathbf{A}\mathbf{S}_t - \mathbf{X})\mathbf{S}_t^\top + (\mathbf{A}_t^\dagger)^\top \right)$$

- Majorization minimization. Notice that the log determinant of a positive definite matrix is concave, therefore

$$\log \det \mathbf{A}\mathbf{A}^\top \leq \log \det \mathbf{A}_t \mathbf{A}_t^\top + \text{Tr}[(\mathbf{A}_t \mathbf{A}_t^\top)^{-1} (\mathbf{A}\mathbf{A}^\top - \mathbf{A}_t \mathbf{A}_t^\top)].$$

This defines a quadratic majorization function for $(1/2) \log \det \mathbf{A}\mathbf{A}^\top$. Minimizing this term plus the fitting term with respect to \mathbf{A} amounts to solving the following linear equation:

$$\lambda (\mathbf{A}\mathbf{S}_t - \mathbf{X})\mathbf{S}_t^\top + (\mathbf{A}_t \mathbf{A}_t^\top)^{-1} \mathbf{A} = 0.$$

This is Sylvester's equation, which can be solved by taking the eigen-decomposition of $\mathbf{S}_t \mathbf{S}_t^\top$ and $\mathbf{A}_t \mathbf{A}_t^\top$.

For simplicity, we resort to the gradient descent update in the rest of this paper.

3.2 UPDATE OF \mathbf{S}

If \mathbf{d} is fixed, then the subproblem of \mathbf{S} is a ℓ_1 regularized least squares problem, which has been extensively studied. It is known to have no closed-form solutions, which is not preferable as one step of an iterative algorithm. Thus, we propose to update \mathbf{S} with a proximal gradient step. Since it amounts to getting a linear approximation at \mathbf{S}_t , we also set \mathbf{d} as the optimal choice with \mathbf{S}_t , which is easy to obtain from Cauchy-Schwarz:

$$d_c = \sqrt{\frac{m}{\sum_{j=1}^k \|\mathbf{e}_j^\top \mathbf{S}_t\|_1^2}} \|\mathbf{e}_c^\top \mathbf{S}_t\|_1. \quad (16)$$

As a result, the proximal gradient update of \mathbf{S} with step size γ take the form

$$\mathbf{S}_{t+1} \leftarrow \mathcal{T}_{\gamma d} (\mathbf{S}_t - \gamma \lambda \mathbf{A}_{t+1}^\top (\mathbf{A}_{t+1} \mathbf{S}_t - \mathbf{X})),$$

where $\mathcal{T}_{\gamma d}(\cdot)$ is the soft-thresholding operator on a matrix with k rows, and the threshold for components on the c th row is d_c/γ .

3.3 SUMMARY AND EXPERIMENTAL DEMONSTRATION

The proposed algorithm based on alternating minimization is summarized in

Algorithm 1 Solving (14) via alternating minimization

```

1: initialize  $\mathbf{A}_0$  and  $\mathbf{S}_0$ 
2: for  $t = 0, 1, 2, \dots$  until convergence do
3:    $\mathbf{A}_{t+1} \leftarrow \mathbf{A}_t - \gamma (\lambda (\mathbf{A} \mathbf{S}_t - \mathbf{X}) \mathbf{S}_t^\top + (\mathbf{A}_t^\dagger)^\top)$ 
4:   for  $c = 1, \dots, k$  do
5:      $d_c = \sqrt{\frac{m}{\sum_{j=1}^k \|e_j^\top \mathbf{S}_t\|_1^2}} \|e_c^\top \mathbf{S}_t\|_1$ 
6:   end for
7:    $\mathbf{S}_{t+1} \leftarrow \mathcal{T}_{\gamma d} (\mathbf{S}_t - \gamma \lambda \mathbf{A}_{t+1}^\top (\mathbf{A}_{t+1} \mathbf{S}_t - \mathbf{X}))$ 
8: end for

```

As (14) is nonconvex and NP-hard, no known algorithm is able to guarantee convergence to a global optimum. In the following, we provide a brief demonstration of the performance of the proposed algorithm. Admittedly, the proposed formulation with identifiability guarantees opens up a new direction for research on algorithm design for dictionary learning, which is a challenging task in itself as it involves several terms that are nontrivial to handle.

We synthetically generate random problems with $s = 5$, $m = 10$, $k = 20$, and $n = 200$. A groundtruth sparse coefficient matrix $\mathbf{S}_\dagger \in \mathbb{R}^{k \times n}$ is generated from the sparse-Gaussian model $\mathcal{SG}(s)$, while the groundtruth dictionary $\mathbf{A}_\dagger \in \mathbb{R}^{m \times k}$ is simply generated from a standard normal distribution. The data matrix is then generated as $\mathbf{X} = \mathbf{A}_\dagger \mathbf{S}_\dagger$. Algorithm 1 runs on \mathbf{X} with $\lambda = 1000$ (since we want the data fidelity term to be almost zero) and $\gamma = 10^{-5}$. At the end, both columns of \mathbf{A}_\dagger and \mathbf{A}_\star are scaled to unit ℓ_2 norm, and the Hungarian algorithm (Kuhn, 1955) is used to find the best column matching. The resulting estimation error $\|\mathbf{A}_\dagger - \mathbf{A}_\star\|_F^2$ remains approximately 10^{-8} over multiple runs. Considering the \mathbf{A}_\dagger is 10×20 with unit column norms this is a satisfactory initial result.

4 CONCLUSION

In this paper, we provide perhaps the first identifiability analysis of a matrix factorization model when the latent dimension is *higher* than the ambient dimension, namely the overcomplete dictionary learning problem. Classical works on this problem rely on combinatorial mathematics, which in turn requires the sample size to be factorial to the latent dimension. Our work is based on a novel formulation that uses a hybrid of weighted ℓ_1 norm of the sparse coefficients and the volume of the overcomplete dictionary as the identification criterion, and we show that identifiability of the overcomplete dictionary of size $m \times k$ can be guaranteed if a geometric condition is satisfied, namely the cellular hull of the sparse coefficient matrix is m -strongly scattered in the k -hypercube. If the sparse coefficient matrix is generated from the sparse-Gaussian model, then such identifiability condition can be satisfied with very high probability if the sample size is $O((k^2/m) \log(k^2/m))$, which is a huge improvement compared to prior work with factorial complexity. We also propose an algorithm for the novel overcomplete dictionary learning formulation. The proposed novel formulation for overcomplete DL with a global identifiability guarantee leaves much room for faster and more efficient algorithm design.

ACKNOWLEDGMENTS

This work is supported in part by NSF ECCS-2237640 and NIH R01LM014027.

REFERENCES

- Alekh Agarwal, Animashree Anandkumar, Prateek Jain, and Praneeth Netrapalli. Learning sparsely used overcomplete dictionaries via alternating minimization. *SIAM Journal on Optimization*, 26(4):2775–2799, 2016.
- Michal Aharon, Michael Elad, and Alfred Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing*, 54(11):4311–4322, 2006a.
- Michal Aharon, Michael Elad, and Alfred M Bruckstein. On the uniqueness of overcomplete dictionaries, and a practical way to retrieve them. *Linear algebra and its applications*, 416(1):48–67, 2006b.
- Sanjeev Arora, Rong Ge, and Ankur Moitra. New algorithms for learning incoherent and overcomplete dictionaries. In *Conference on Learning Theory*, pp. 779–806. PMLR, 2014.
- Sanjeev Arora, Rong Ge, Tengyu Ma, and Ankur Moitra. Simple, efficient, and neural algorithms for sparse coding. In *Conference on Learning Theory*, pp. 113–149. PMLR, 2015.
- Yu Bai, Qijia Jiang, and Ju Sun. Subgradient descent learns orthogonal dictionaries. In *International Conference on Learning Representations*, 2019.
- Adi Ben-Israel. A volume associated with $m \times n$ matrices. *Linear algebra and its applications*, 167:87–111, 1992.
- George Bennett. Probability inequalities for the sum of independent random variables. *Journal of the American Statistical Association*, 57(297):33–45, 1962.
- Emmanuel J Candès and Michael B Wakin. An introduction to compressive sampling. *IEEE signal processing magazine*, 25(2):21–30, 2008.
- Jeremy E Cohen and Nicolas Gillis. Identifiability of complete dictionary learning. *SIAM Journal on Mathematics of Data Science*, 1(3):518–536, 2019.
- David L Donoho. Compressed sensing. *IEEE Transactions on information theory*, 52(4):1289–1306, 2006.
- Michael Elad and Michal Aharon. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing*, 15(12):3736–3745, 2006.
- Xiao Fu, Wing-Kin Ma, Kejun Huang, and Nicholas D Sidiropoulos. Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain. *IEEE Transactions on Signal Processing*, 63(9):2306–2320, 2015.
- Charles J Garfinkle and Christopher J Hillar. On the uniqueness and stability of dictionaries for sparse representation of noisy signals. *IEEE Transactions on Signal Processing*, 67(23):5884–5892, 2019.
- Quan Geng and John Wright. On the local correctness of ℓ_1 -minimization for dictionary learning. In *2014 IEEE International Symposium on Information Theory*, pp. 3180–3184. IEEE, 2014.
- Pando Georgiev, Fabian Theis, and Andrzej Cichocki. Sparse component analysis and blind source separation of underdetermined mixtures. *IEEE transactions on neural networks*, 16(4):992–996, 2005.
- Rémi Gribonval and Karin Schnass. Dictionary identification—sparse matrix-factorization via ℓ_1 -minimization. *IEEE Transactions on Information Theory*, 56(7):3523–3539, 2010.
- Christopher J Hillar and Friedrich T Sommer. When can dictionary learning uniquely recover sparse data from subsamples? *IEEE Transactions on Information Theory*, 61(11):6290–6297, 2015.
- Jingzhou Hu and Kejun Huang. Global identifiability of ℓ_1 -based dictionary learning via matrix volume optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 36, 2023a.

- Jingzhou Hu and Kejun Huang. Identifiable bounded component analysis via minimum volume enclosing parallelotope. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023b.
- Jingzhou Hu and Kejun Huang. Complex bounded component analysis: Identifiability and algorithm. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024.
- Kejun Huang and Xiao Fu. Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm. In *International Conference on Machine Learning*, pp. 2859–2868, 2019.
- Kejun Huang, Nicholas D Sidiropoulos, and Ananthram Swami. Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition. *IEEE Transactions on Signal Processing*, 62(1):211–224, 2013.
- Kejun Huang, Xiao Fu, and Nikolaos D Sidiropoulos. Anchor-free correlated topic modeling: Identifiability and algorithm. *Advances in Neural Information Processing Systems*, 29, 2016.
- Kejun Huang, Xiao Fu, and Nicholas D. Sidiropoulos. Learning hidden Markov models from pairwise co-occurrences with application to topic modeling. In *International Conference on Machine Learning (ICML)*, pp. 2073–2082, 2018.
- Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955.
- Chia-Hsiang Lin, Wing-Kin Ma, Wei-Chiang Li, Chong-Yung Chi, and ArulMurugan Ambikapathi. Identifiability of the simplex volume minimization criterion for blind hyperspectral unmixing: The no-pure-pixel case. *IEEE Transactions on Geoscience and Remote Sensing*, 53(10):5530–5546, 2015.
- Julien Mairal, Francis Bach, Jean Ponce, and Guillermo Sapiro. Online dictionary learning for sparse coding. In *Proceedings of the 26th annual international conference on machine learning*, pp. 689–696, 2009.
- Bruno A Olshausen and David J Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision Research*, 37(23):3311–3325, 1997.
- Gilles Pisier. *The volume of convex bodies and Banach space geometry*, volume 94. Cambridge University Press, 1999.
- Mark D Plumbley, Thomas Blumensath, Laurent Daudet, Rémi Gribonval, and Mike E Davies. Sparse representations in audio and music: from coding to source separation. *Proceedings of the IEEE*, 98(6):995–1005, 2009.
- Sirisha Rambhatla, Xingguo Li, and Jarvis Haupt. Noodl: Provable online dictionary learning and sparse coding. In *International Conference on Learning Representations*, 2019.
- Yifei Shen, Ye Xue, Jun Zhang, Khaled Letaief, and Vincent Lau. Complete dictionary learning via lp-norm maximization. In *Conference on Uncertainty in Artificial Intelligence*, pp. 280–289. PMLR, 2020.
- Daniel A. Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In Shie Mannor, Nathan Srebro, and Robert C. Williamson (eds.), *Proceedings of the 25th Annual Conference on Learning Theory*, volume 23 of *Proceedings of Machine Learning Research*, pp. 37.1–37.18, Edinburgh, Scotland, 25–27 Jun 2012. PMLR.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere i: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 63(2):853–884, 2016a.
- Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere ii: Recovery by riemannian trust-region method. *IEEE Transactions on Information Theory*, 63(2):885–914, 2016b.

- Yuchen Sun and Kejun Huang. Improved identifiability and sample complexity analysis of complete dictionary learning. In *2024 IEEE 13rd Sensor Array and Multichannel Signal Processing Workshop (SAM)*, pp. 1–5. IEEE, 2024.
- Gokcan Tatli and Alper T Erdogan. Polytopic matrix factorization: Determinant maximization based criterion and identifiability. *IEEE Transactions on Signal Processing*, 69:5431–5447, 2021.
- Bahareh Tolooshams and Demba E Ba. Stable and interpretable unrolled dictionary learning. *Transactions on Machine Learning Research*, 2022.
- Ivana Tošić and Pascal Frossard. Dictionary learning. *IEEE Signal Processing Magazine*, 28(2): 27–38, 2011.
- Ivana Tošić, Ivana Jovanović, Pascal Frossard, Martin Vetterli, and Neb Durić. Ultrasound tomography with learned dictionaries. In *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 5502–5505. IEEE, 2010.
- Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*, volume 47. Cambridge university press, 2018.
- Yu Wang, Siqi Wu, and Bin Yu. Unique sharp local minimum in l_1 -minimization complete dictionary learning. *Journal of Machine Learning Research*, 21(63):1–52, 2020.
- Siqi Wu and Bin Yu. Local identifiability of l_1 -minimization dictionary learning: a sufficient and almost necessary condition. *The Journal of Machine Learning Research*, 18(1):6121–6176, 2017.
- Yuexiang Zhai, Hermish Mehta, Zhengyuan Zhou, and Yi Ma. Understanding l_4 -based dictionary learning: Interpretation, stability, and robustness. In *International Conference on Learning Representations*, 2020a.
- Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. Complete dictionary learning via l_4 -norm maximization over the orthogonal group. *Journal of Machine Learning Research*, 21(165):1–68, 2020b.

A PROOF OF THEOREM 1

In the proof sketch of Theorem 1, we showed that

$$\frac{1}{2} \log \det \mathbf{A}_\dagger \mathbf{A}_\dagger^\top + \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_\dagger\|_1 = \frac{1}{2} \log \det \mathbf{A}_\star \mathbf{A}_\star^\top + \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{e}_c^\top \mathbf{S}_\star\|_1,$$

so $(\mathbf{A}_\dagger, \mathbf{S}_\dagger)$ is at least one candidate solution for (2). In this section, we first complete the proof by showing that if the second requirement of Assumption 3 is satisfied, then any optimal solution $(\mathbf{A}_\star, \mathbf{S}_\star)$ must satisfy that $\mathbf{A}_\star = \mathbf{A}_\dagger \mathbf{\Pi} \mathbf{D}$, where $\mathbf{\Pi}$ is a permutation matrix and \mathbf{D} is a diagonal matrix with ± 1 on the diagonal.

In (8) we proved that $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 \geq \|\mathbf{Q}^\top \mathbf{w}_c\|_2$ by choosing $\|\boldsymbol{\theta}\|_\infty \leq 1$ such that

$$\tilde{\mathbf{S}} \boldsymbol{\theta} = \mathbf{D}_\dagger^{-1} \mathbf{Q} \mathbf{Q}^\top \mathbf{w}_c / \|\mathbf{Q}^\top \mathbf{w}_c\|_2.$$

Suppose $\mathbf{Q}^\top \mathbf{w}_c / \|\mathbf{Q}^\top \mathbf{w}_c\|_2 \neq \pm \mathbf{q}_c / \|\mathbf{q}_c\|_2$ for all $c = 1, \dots, k$, then according to the second requirement of Assumption 3, it is in the interior of $\text{cell}(\tilde{\mathbf{S}}_\dagger)$, which means there exists $\alpha > 1$ and $\|\tilde{\boldsymbol{\theta}}\|_\infty \leq 1$ such that

$$\tilde{\mathbf{S}} \tilde{\boldsymbol{\theta}} = \alpha \mathbf{D}_\dagger^{-1} \mathbf{Q} \mathbf{Q}^\top \mathbf{w}_c / \|\mathbf{Q}^\top \mathbf{w}_c\|_2.$$

Therefore

$$\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}\|_1 \geq \mathbf{w}_c^\top \mathbf{D}_\dagger \tilde{\mathbf{S}} \tilde{\boldsymbol{\theta}} = \alpha \|\mathbf{Q}^\top \mathbf{w}_c\|_2 > \|\mathbf{Q}^\top \mathbf{w}_c\|_2.$$

This means equality of $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 \geq \|\mathbf{Q}^\top \mathbf{w}_c\|_2$ is only attained when $\mathbf{Q}^\top \mathbf{w}_c = \pm \mathbf{q}_c$ for some $c = 1, \dots, k$.

Finally, it is proven in (10b) that

$$\log \det \mathbf{Q}^\top \mathbf{W}^\top \mathbf{W} \mathbf{Q} \leq m \log \frac{1}{m} \|\mathbf{Q}^\top \mathbf{W}\|_F^2,$$

and if $\det \mathbf{A}_\dagger \mathbf{A}_\dagger^\top = \det \mathbf{A}_\star \mathbf{A}_\star^\top$ the above inequality must hold as an equality. In the proof of (10b) we showed that equality holds only if all the eigenvalues of $\mathbf{Q}^\top \mathbf{W}^\top \mathbf{W} \mathbf{Q}$ are equal, meaning columns of $\mathbf{W} \mathbf{Q}$ are orthonormal. Since \mathbf{Q} itself is orthonormal, it is possible only if $\mathbf{W} \mathbf{Q} = \mathbf{D} \mathbf{\Pi} \mathbf{Q}$, where \mathbf{D} is a diagonal matrix with ± 1 on the diagonals and $\mathbf{\Pi}$ is a permutation matrix. This shows that

$$\mathbf{A}_\star = \mathbf{A}_\dagger \mathbf{\Pi} \mathbf{D}.$$

Q.E.D.

The remaining of this section shows some key equalities and inequalities in the proof sketch that are skipped for clarity.

Proof that $\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}_\dagger\|_1$. As we argued in the proof sketch of Theorem 1, the constraint $\mathbf{X} = \mathbf{A} \mathbf{S}$ is equivalent to $\mathbf{A} \mathbf{W} = \mathbf{A}_\dagger$ and $\mathbf{S} = \mathbf{W} \mathbf{S}_\dagger + \mathbf{B}$ where $\mathbf{S}_\dagger \mathbf{B}^\top = 0$. Substituting them into (2) eliminates the constraint $\mathbf{X} = \mathbf{A} \mathbf{S}$:

$$\underset{\mathbf{A}, \mathbf{W}, \mathbf{B}}{\text{minimize}} \quad \frac{1}{2} \log \det \mathbf{A} \mathbf{A}^\top + \max_{\|\mathbf{d}\|_2^2=m} \sum_{c=1}^k d_c \|\mathbf{w}_c^\top \mathbf{S}_\dagger + \mathbf{b}_c^\top\|_1 \quad \text{subject to } \mathbf{A} \mathbf{W} = \mathbf{A}_\dagger, \mathbf{S}_\dagger \mathbf{B}^\top = 0.$$

Taking the Clarke generalized derivative of the Lagrange function with respect to \mathbf{b}_c and setting it equal to zero gives

$$d_c \boldsymbol{\theta}_c + \mathbf{S}_\dagger^\top \boldsymbol{\mu}_c = 0,$$

where $\boldsymbol{\theta}_c = \text{sign}(\mathbf{S}_\dagger^\top \mathbf{w}_c + \mathbf{b}_c)$ and $\boldsymbol{\mu}_c$ is the Lagrange multiplier for the c th column of $\mathbf{S}_\dagger \mathbf{B}^\top$. This shows that $\boldsymbol{\theta}_c$ is a linear combination of rows of \mathbf{S}_\dagger , which means it is orthogonal to \mathbf{b}_c . As a result,

$$\|\mathbf{e}_c^\top \mathbf{S}_\star\|_1 = \|\mathbf{w}_c^\top \mathbf{S}_\dagger + \mathbf{b}_c^\top\|_1 = \mathbf{w}_c^\top \mathbf{S}_\dagger \boldsymbol{\theta}_c + \mathbf{b}_c^\top \boldsymbol{\theta}_c = \mathbf{w}_c^\top \mathbf{S}_\dagger \boldsymbol{\theta}_c = \|\mathbf{w}_c^\top \mathbf{S}_\dagger\|_1.$$

□

Proof of inequality (9). The matrix $\mathbf{W}\mathbf{W}^\dagger$ defines a projection matrix, which is symmetric, with the following properties:

$$\mathbf{W}\mathbf{W}^\dagger \preceq \mathbf{I}, \quad \mathbf{W}\mathbf{W}^\dagger = \mathbf{W}\mathbf{W}^\dagger\mathbf{W}\mathbf{W}^\dagger = \mathbf{W}\mathbf{W}^\dagger(\mathbf{W}\mathbf{W}^\dagger)^\top = \mathbf{W}\mathbf{W}^\dagger(\mathbf{W}^\dagger)^\top\mathbf{W}^\top.$$

As a result,

$$\det \mathbf{A}_\star \mathbf{A}_\star^\top \geq \det \mathbf{A}_\star \mathbf{W}\mathbf{W}^\dagger \mathbf{A}_\star^\top = \det \mathbf{A}_\star \mathbf{W}\mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{W}^\top \mathbf{A}_\star^\top = \det \mathbf{A}_\dagger \mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{A}_\dagger^\top,$$

where the last step is because $\mathbf{A}_\dagger = \mathbf{A}_\star \mathbf{W}$. \square

Proof of inequality (10a). Since $\mathbf{A}_\dagger = \mathbf{A}_\star \mathbf{W}$, rows of \mathbf{A}_\dagger are in the row space of \mathbf{W} , and so are columns of \mathbf{Q} . Therefore

$$\mathbf{W}\mathbf{W}^\dagger \mathbf{Q} = \mathbf{Q}.$$

This means $\mathbf{W}\mathbf{Q}$ has linearly independent columns, therefore $\mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q}$ is invertible and

$$(\mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q})^{-1} = (\mathbf{W}\mathbf{Q})^\dagger (\mathbf{Q}^\top \mathbf{W}^\top)^\dagger.$$

On the other hand, since $\mathbf{Q}^\top \mathbf{W}\mathbf{W}^\dagger \mathbf{Q} = \mathbf{I}$, this means

$$\mathbf{Q}^\top \mathbf{W}^\dagger = (\mathbf{W}\mathbf{Q})^\dagger + \mathbf{C},$$

where rows of \mathbf{C} are orthogonal to columns of $\mathbf{W}\mathbf{Q}$, i.e., $\mathbf{C}\mathbf{W}\mathbf{Q} = 0$. Columns of $\mathbf{W}\mathbf{Q}$ and rows of $(\mathbf{W}\mathbf{Q})^\dagger$ span the same subspace, so we also have $\mathbf{C}^\top (\mathbf{W}\mathbf{Q})^\dagger = 0$. As a result,

$$\begin{aligned} \det \mathbf{Q}^\top \mathbf{W}^\dagger (\mathbf{W}^\dagger)^\top \mathbf{Q} &= \det \left((\mathbf{W}\mathbf{Q})^\dagger + \mathbf{C} \right) \left((\mathbf{W}\mathbf{Q})^\dagger + \mathbf{C} \right)^\top \\ &= \det \left((\mathbf{W}\mathbf{Q})^\dagger (\mathbf{Q}^\top \mathbf{W}^\top)^\dagger + \mathbf{C}\mathbf{C}^\top \right) \\ &\geq \det (\mathbf{W}\mathbf{Q})^\dagger (\mathbf{Q}^\top \mathbf{W}^\top)^\dagger = \det (\mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q})^{-1}. \end{aligned}$$

Taking the log on both sides shows (10a). \square

Proof of inequality (10b). Denote the eigenvalues of $\mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q}$ as $\lambda_1, \dots, \lambda_m$, which are all nonnegative, then

$$(\det \mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q})^{1/m} = \left(\prod_{j=1}^m \lambda_j \right)^{1/m} \leq \frac{1}{m} \sum_{j=1}^m \lambda_j = \text{Tr } \mathbf{Q}^\top \mathbf{W}^\top \mathbf{W}\mathbf{Q} = \|\mathbf{Q}^\top \mathbf{W}^\top\|_{\text{F}}^2,$$

where the inequality in the middle is the geometric-arithmetic mean inequality. Taking the log on both sides and rearranging gives (10b). Notice that equality holds only if $\lambda_1 = \dots = \lambda_m$, i.e., columns of $\mathbf{W}\mathbf{Q}$ are orthonormal. \square

B PROOF OF THEOREM 2

We assume that $\mathbf{S} \sim \mathcal{S}\mathcal{G}(s)$, and the first thing we do is rescale its rows to have unit ℓ_1 norms and use it for the problem 13. To simplify the analysis, we can instead directly maximize $\|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\|_1^2$ subject to $\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1$, and compare it with the largest ℓ_1 norm of the rows of \mathbf{S} . The complement of the intended probability can be bounded as

$$\Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| \leq 1 \right] \geq \Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| \leq \alpha \cap \max_j \|\mathbf{S}_{j,:}\|_1 \geq \alpha \right],$$

with an arbitrary choice of α . Conversely,

$$\begin{aligned} \Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| > 1 \right] &\leq \Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{w}\| > \alpha \cup \max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \\ &\leq \Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| > \alpha \right] + \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \quad (17) \end{aligned}$$

where the second inequality is obtained from the union bound. The rest of this section is dedicated to bounding the above two terms. Both of these results rely on the following version of the Bernstein inequality (Bennett, 1962):

Theorem 3 (Bernstein’s inequality). *Let Z_1, \dots, Z_n be independent random variables with $\mathbb{E}[Z_i^2] \leq v^2$ and there exists some constant c such that for all integer $d > 2$*

$$\mathbb{E}[|Z_i|^d] \leq \frac{1}{2} d! v^2 c^{d-2}. \quad (18)$$

Then

$$\Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > \epsilon \right] \leq 2 \exp \left(-\frac{\epsilon^2}{2(nv^2 + c\epsilon)} \right)$$

Lemma 2 (Bounding the second term in (17)). *Suppose $\mathbf{S} \in \mathbb{R}^{k \times n}$ is generated from the sparse-Gaussian model $\mathcal{SG}(s)$. Then*

$$\Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) \right] \leq 2k \exp \left(-\frac{n(s/k)\epsilon^2}{2 + \sqrt{2}\epsilon} \right)$$

Proof. Let $\mathbf{s} = (s_1, \dots, s_n)$ be one row of \mathbf{S} generated from $\mathcal{SG}(s)$, then each s_i has probability s/k to be standard normal and probability $1 - s/k$ to be zero. We will use Bernstein’s inequality with $Z_i = |s_i|$. Let g denote a standard normal random variable, then $|g|$ follows a Chi-distribution of degree 1, so its moments are

$$\mathbb{E}[|g|^d] = 2^{d/2} \frac{\Gamma\left(\frac{d+1}{2}\right)}{\Gamma\left(\frac{1}{2}\right)}. \quad (19)$$

As a result, we have $\mathbb{E}[Z_i] = (s/k)\sqrt{2/\pi}$ and $\mathbb{E}[Z_i^2] = s/k$. Using the recurrence relation for the Gamma function $\Gamma(t+1) = t\Gamma(t)$ and $\sqrt{2}/\Gamma(1/2) = \sqrt{2/\pi} < 1$ we can bound the rest of the moments with $d > 2$ as

$$\mathbb{E}[|Z_i|^d] \leq (s/k) \frac{d!}{2^{d/2}}. \quad (20)$$

Therefore the moments satisfy (18) with $c = 1/\sqrt{2}$, results in

$$\begin{aligned} \Pr \left[\|\mathbf{s}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) \right] &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n(s/k)\epsilon \right] \\ &\leq 2 \exp \left(-\frac{n^2(s/k)^2\epsilon^2}{2(n(s/k) + n(s/k)\epsilon/\sqrt{2})} \right) \\ &= 2 \exp \left(-\frac{n(s/k)\epsilon^2}{2 + \sqrt{2}\epsilon} \right). \end{aligned}$$

Finally, using the union bound

$$\begin{aligned} \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) \right] &\leq k \Pr \left[\|\mathbf{S}_{j,:}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) \right] \\ &\leq 2k \exp \left(-\frac{n(s/k)\epsilon^2}{2 + \sqrt{2}\epsilon} \right). \end{aligned}$$

□

We now proceed to bound the first term in (17). First, we note the following equivalence:

$$\Pr \left[\sup_{\substack{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| > \alpha \right] = \Pr \left[\inf_{\substack{\|\mathbf{Q}^\top \mathbf{D}^{-1} \mathbf{w}\| = 1 \\ \mathbf{Q}^\top \mathbf{Q} = \mathbf{I}}} \|\mathbf{S}^\top \mathbf{w}\|_1 < 1/\alpha \right] \quad (21)$$

We are also going to use the following notion of δ -cover from convex geometry (Pisier, 1999):

Definition 3 (δ -cover). A finite δ -cover of a set \mathcal{S} in \mathbb{R}^k is a finite set $\mathcal{N}(\mathcal{S}, \delta)$ of points on \mathcal{S} such that any point on \mathcal{S} is within δ away from an element in $\mathcal{N}(\mathcal{S}, \delta)$, i.e.

$$\min_{\mathbf{w}_i \in \mathcal{N}(\mathcal{S}, \delta)} \|\mathbf{w} - \mathbf{w}_i\| < \delta, \quad \forall \mathbf{w} \in \mathcal{S}.$$

A well-known result is that the δ -cover of the unit sphere $\mathcal{B}_k = \{\mathbf{w} \mid \|\mathbf{w}\|_2 = 1\}$ has bounded cardinality (Vershynin, 2018):

$$|\mathcal{N}(\mathcal{B}_k, \delta)| \leq \left(\frac{2}{\delta} + 1\right)^k.$$

Notice that \mathcal{B}_k is a special case of \mathcal{B}_m defined in Assumption 3 with $k = m$, implying that $\mathcal{B}_k \subseteq \mathcal{B}_m$ and thus $\mathcal{B}_m^\circ \subseteq \mathcal{B}_k^\circ$. We will use this result to prove a (rather loose) bound for the δ -cover of \mathcal{B}_m° .

Lemma 3. *The δ -cover of \mathcal{B}_m° satisfies*

$$|\mathcal{N}(\mathcal{B}_m, \delta)|_2 \leq \left(\frac{2}{\delta} + 1\right)^k$$

Proof. It is easy to see that \mathcal{B}_k is self-polar, i.e., $\mathcal{B}_k^\circ = \mathcal{B}_k$, so $\mathcal{B}_m^\circ \subseteq \mathcal{B}_k$. Since the superset \mathcal{B}_k satisfies $|\mathcal{N}(\mathcal{B}_k, \delta)| \leq (1 + 2/\delta)^k$, so is the subset \mathcal{B}_m° . \square

Lemma 4. *Let $\mathcal{N}(\mathcal{B}_m, \delta) = \{\mathbf{w}_i\}$ be a δ -cover for \mathcal{B}_m° in \mathbb{R}^k . Assume that we have both the lowerbound*

$$\|\mathbf{S}^\top \mathbf{w}_i\|_1 \geq \beta, \quad \forall \mathbf{w}_i \in \mathcal{N}(\mathcal{B}_m, \delta)$$

and the upperbound

$$\|\mathbf{S}^\top\|_1 = \sup_{\|\mathbf{w}\|_1 \leq 1} \|\mathbf{S}^\top \mathbf{w}\|_1 \leq \gamma.$$

Then

$$\inf_{\mathbf{w} \in \mathcal{B}_m^\circ} \|\mathbf{S}^\top \mathbf{w}\|_1 \geq \beta - \gamma \delta \sqrt{k}$$

Proof. By definition of the δ -cover, for all $\mathbf{w} \in \mathcal{B}_m^\circ$ we can find $\mathbf{w}_i \in \mathcal{N}(\mathcal{B}_m^\circ, \delta)$ with $\|\mathbf{w} - \mathbf{w}_i\| < \delta$. Therefore

$$\begin{aligned} \|\mathbf{S}^\top \mathbf{w}\|_1 &\geq \|\mathbf{S}^\top \mathbf{w}_i\|_1 - \|\mathbf{S}^\top(\mathbf{w} - \mathbf{w}_i)\|_1 \geq \beta - \|\mathbf{S}^\top\|_1 \|\mathbf{w} - \mathbf{w}_i\|_1 \\ &\geq \beta - \|\mathbf{S}^\top\|_1 \|\mathbf{w} - \mathbf{w}_i\|_2 \sqrt{k} \geq \beta - \gamma \delta \sqrt{k}. \end{aligned}$$

\square

Lemma 5 (Bounding the first term in (17)). *Suppose $\mathbf{S} \in \mathbb{R}^{k \times n}$ is generated from the sparse-Gaussian model $\mathcal{SG}(p)$. Then*

$$\Pr \left[\inf_{\mathbf{w} \in \mathcal{B}_m^\circ} \|\mathbf{S}^\top \mathbf{w}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) - \delta \sqrt{k} n(s/k)(\sqrt{2/\pi} + \epsilon) \right] \leq \left(\left(\frac{2}{\delta}\right)^k + 2 \right) 2 \exp \left(-\frac{n(s/k)^2 \epsilon^2}{2 + \sqrt{2}\epsilon} \right),$$

where $\delta \in (0, 1)$ represents any choice of δ -cover for \mathcal{B}_m° .

Proof. Following Lemma 4, we have

$$\Pr \left[\inf_{\mathbf{w} \in \mathcal{B}_m^\circ} \|\mathbf{S}^\top \mathbf{w}\|_1 < \beta - \gamma \delta \sqrt{k} \right] \leq \sum_{\mathbf{w}_i \in \mathcal{N}(\mathcal{B}_m^\circ, \delta)} \Pr \left[\|\mathbf{S}^\top \mathbf{w}_i\|_1 < \beta \right] + \Pr \left[\|\mathbf{S}^\top\|_1 > \gamma \right], \quad (22)$$

where $|\mathcal{N}(\mathcal{B}_m^\circ, \delta)| < (1 + 2/\delta)^k$ according to Lemma 3.

The bound to the first term in (22) is almost identical to Lemma 2. Dropping the subscript of \mathbf{w}_i , we write

$$\|\mathbf{S}^\top \mathbf{w}\|_1 = \sum_{i=1}^n \left| \sum_{j=1}^k s_{ij} w_j \right| := \sum_{i=1}^n |Z_i|.$$

Without the absolute value, Z_i is normally distributed with zero mean and variance $\sigma^2 = \sum_{j \in \mathcal{I}_i} w_j^2$, where \mathcal{I}_i is the index set of nonzero elements in \mathbf{s}_i with $|\mathcal{I}_i| = s < m$. Then $\sigma^2 \leq 1$. To see this, suppose without loss of generality that $\mathcal{I}_i = \{1, \dots, s\}$ and let $\mathbf{Q} = [\mathbf{I} \ 0]^\top$, then $\mathbf{w} \in \mathcal{B}_m^\circ$ implies that $\sum_{j \in \mathcal{I}_i} w_j^2 = \|\mathbf{Q}\mathbf{w}\|_2^2 \leq 1$. In other words, if $\mathbf{w} \in \mathcal{B}_m^\circ$, the squared sum of no more than m elements of \mathbf{w} must be ≤ 1 . Therefore for $d \geq 2$

$$\mathbb{E}[|Z_i|^d] \leq \mathbb{E}[|Z_i|^2] \leq 1,$$

while for $d = 1$ we have

$$\mathbb{E}[|Z_i|] = \frac{s}{k} \sqrt{\frac{2}{\pi}}.$$

Again, the moments satisfy (18) with $c = 1/\sqrt{2}$, results in

$$\begin{aligned} \Pr \left[\|\mathbf{S}^\top \mathbf{w}\|_1 < n(s/k)(\sqrt{2/\pi} - \epsilon) \right] &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n(s/k)\epsilon \right] \\ &\leq 2 \exp \left(-\frac{n^2(s/k)^2 \epsilon^2}{2(n + n(s/k)\epsilon/\sqrt{2})} \right) \\ &= 2 \exp \left(-\frac{n(s/k)^2 \epsilon^2}{2 + \sqrt{2}\epsilon} \right). \end{aligned} \quad (23)$$

To bound $\|\mathbf{S}^\top\|_1$, we recall that this is the ℓ_1 induced norm for matrix \mathbf{S}^\top , which is shown to be the maximum of the ℓ_1 norms of the columns of \mathbf{S}^\top . This means we can use similar arguments used in Lemma 2 (but applied to the other direction) to have

$$\begin{aligned} \Pr \left[\|\mathbf{S}^\top\|_1 > np(\sqrt{2/\pi} + \epsilon) \right] &= \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 > n(s/k)(\sqrt{2/\pi} + \epsilon) \right] \\ &\leq \Pr \left[\|\mathbf{S}_{j,:}\|_1 > n(s/k)(\sqrt{2/\pi} + \epsilon) \right] \\ &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n(s/k)\epsilon \right] \\ &\leq 2 \exp \left(-\frac{n(s/k)\epsilon^2}{2 + \sqrt{2}\epsilon} \right) \leq 2 \exp \left(-\frac{n(s/k)^2 \epsilon^2}{2 + \sqrt{2}\epsilon} \right), \end{aligned} \quad (24)$$

where we pick an arbitrary $j \in [k]$ in the second line since this event implies that the maximum ℓ_1 norm of the rows is lowerbounded, and in the third line each Z_i satisfies (20). The proof is complete by combining (22), (23), and (24) with $\beta = n(s/k)(\sqrt{2/\pi} - \epsilon)$ and $\gamma = n(s/k)(\sqrt{2/\pi} + \epsilon)$. \square

Proof of Theorem 2. We first instantiate Lemma 5 with

$$\delta = \frac{n^2(s/k)^2(\sqrt{2/\pi} - \epsilon)^2 - 1}{n^2(s/k)^2(\sqrt{2/\pi} - \epsilon)(\sqrt{2/\pi} + \epsilon)\sqrt{k}},$$

which satisfies $\delta > 0$ if

$$\epsilon < \sqrt{\frac{2}{\pi}} - \frac{1}{n(s/k)}.$$

Then we have

$$\Pr \left[\inf_{\mathbf{w} \in \mathcal{B}_m^\circ} \|\mathbf{S}^\top \mathbf{w}\|_1 < 1/n(s/k)^2(\sqrt{2/\pi} - \epsilon) \right] \leq \left(\left(\frac{2}{\delta} \right)^k + 2 \right) 2 \exp \left(-\frac{n(s/k)\epsilon^2}{2 + \sqrt{2}\epsilon} \right),$$

Combining (17), (21), and Lemma 2 with $\alpha = n(s/k)(\sqrt{2/\pi} - \epsilon)$, we obtain

$$\begin{aligned} \Pr \left[\sup_{\|\mathbf{w}^\top \hat{\mathbf{S}}\|_1 \leq 1} \|\mathbf{w}\| > 1 \right] &\leq \Pr \left[\inf_{\|\mathbf{w}\|_1=1} \|\mathbf{S}^\top \mathbf{w}\|_1 < 1/\alpha \right] + \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \\ &\leq 2 \left(k + \left(\frac{2}{\delta} \right)^k + 2 \right) \exp \left(-\frac{n(s/k)^2 \epsilon^2}{2 + \sqrt{2}\epsilon} \right). \end{aligned} \quad (25)$$

Further, with

$$\epsilon < \frac{((m/k)^{1/4} - 1)\sqrt{2/\pi}}{(m/k)^{1/4} + 1},$$

where the right hand side is obviously positive, we have

$$\delta = \frac{n^2(s/k)^2(\sqrt{2/\pi} - \epsilon)^2 - 1}{n^2(s/k)^2(\sqrt{2/\pi} - \epsilon)(\sqrt{2/\pi} + \epsilon)\sqrt{k}} > \frac{n^2(s/k)^2(\sqrt{2/\pi} - \epsilon)^2}{n^2(s/k)^2(\sqrt{2/\pi} + \epsilon)^2\sqrt{k}} > \frac{\sqrt{m}}{k}.$$

We can further relax (25) to

$$\begin{aligned} \Pr \left[\sup_{\mathbf{w} \in \mathcal{B}_m^{\circ}} \|\mathbf{w}\| > 1 \right] &\leq 2 \left(k + (2k/\sqrt{m})^k + 2 \right) \exp \left(-\frac{n(s/k)^2\epsilon^2}{2 + \sqrt{2}\epsilon} \right) \\ &\leq 4 \left(2k/\sqrt{m} \right)^k \exp \left(-\frac{n(s/k)^2\epsilon^2}{2 + \sqrt{2}\epsilon} \right) \\ &\leq 4 \exp \left((k/2) \log(k^2/m) - n(s/k)^2(m/k) \right), \end{aligned}$$

where in the last inequality we simply picked a small enough ϵ so that

$$\frac{\epsilon^2}{2 + \sqrt{2}\epsilon} < \frac{m}{k}.$$

This completes the proof. □