

Improved Identifiability and Sample Complexity Analysis of Complete Dictionary Learning

Yuchen Sun and Kejun Huang

Department of Computer and Information Science and Engineering

University of Florida

Gainesville, Florida 32611

Email: (yuchen.sun, kejun.huang)@ufl.edu

Abstract—It has been recently shown that dictionary learning is identifiable with a complete dictionary, with relatively loose assumptions that the dictionary is invertible and that a geometric body obtained from the sparse coefficient matrix is sufficiently scattered in the hypercube. In addition, if the sparse coefficients are generated from the Bernoulli-Gaussian model, then the number of samples required to guarantee identifiability is $O(k \log(k))$. In this paper, we further improve the identifiability results in two-folds. First, we extend the deterministic identifiability results from the real domain in the previous work to the complex domain, hence significantly expand its applications. Second, we consider the generative model for the sparse coefficients from the more specific (also widely adopted) Bernoulli-Gaussian to a much wider class where the nonzero entries can be drawn from any subexponential distributions, and show that the same sample complexity results still holds. The recovery performance is confirmed in a Monte-Carlo simulation where the sparse coefficients are drawn from Bernoulli-Laplacian model.

Index Terms—dictionary learning, sufficiently scattered, Bernoulli-subexponential.

I. INTRODUCTION

Dictionary learning (DL) amounts to factor a data matrix as $\mathbf{X} = \mathbf{A}\mathbf{S}$ where \mathbf{S} is sparse [1]. Treating $\mathbf{X} \in \mathbb{R}^{k \times n}$ or $\mathbb{C}^{k \times n}$ 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 complete dictionary learning, therefore $\mathbf{A} \in \mathbb{R}^{k \times k}$ and $\mathbf{S} \in \mathbb{R}^{k \times n}$ (or $\mathbb{C}^{k \times k}$ and $\mathbb{C}^{k \times n}$, respectively).

Identifiability of dictionary learning has been an open question ever since the problem was posed. As an instance of matrix factorization, inherent ambiguities of permutation and scaling of the columns of \mathbf{A} (and in turn same permutation and counter-scaling of the rows of \mathbf{S}) are unavoidable and inconsequential in practice, and the factorization is essentially unique or identifiable if all admissible factorization are permutation and scaling of each other. Earlier works focus on directly enforcing sparsity constraints on \mathbf{S} , which show that over-complete DL is identifiable if $n > O((k+1)\binom{k}{s})$ [2], [3] and complete DL is identifiable with $O(k^3/(k-s)^2)$ samples. Notice that these work require *all* columns of \mathbf{S} to be strictly s -sparse—no dense outliers is allowed for the analysis to work. More recently, there have been work on identifiability of DL with ℓ_1 norm as the regularization to promote sparsity, although they are only able to show identifiability in a local region [4]–[6] if the dictionary \mathbf{A} is incoherent and the sparse coefficients \mathbf{S} follows some sparse generative model such as the Bernoulli-Gaussian model. The benefit is not only that algorithm

design is easier with ℓ_1 norm regularization, but also that the number of samples n required is usually $O(k \log k)$, while a few outlying dense columns of \mathbf{S} are allowed.

The best known result regarding the identifiability of DL is in [7], which shows that complete dictionary learning is *globally* identifiable if: 1) the complete dictionary $\mathbf{A} \in \mathbb{R}^{k \times k}$ is invertible (hence no incoherence is necessary for complete dictionaries), and 2) the sparse coefficients $\mathbf{S} \in \mathbb{R}^{k \times n}$ satisfies some sufficiently scattered condition, which we will elaborate in the next section. It was also shown that if \mathbf{S} is generated from the Bernoulli-Gaussian model, it will satisfy the sufficiently scattered condition with overwhelming probability as long as the sample size n is larger than $O(k \log(k))$, which confirms the sample complexity analysis in [4]–[6] but identifiability is guaranteed globally, not just in a local region near the groundtruth factors. This new framework for identifiable DL is based on a new problem formulation by constraining the ℓ_1 norm of each row of \mathbf{S} while minimizing the matrix volume of the dictionary, i.e., $|\det \mathbf{A}|$.

In this paper, we aim at further improving the identifiability result for complete dictionary learning. First, noticing that all the results in [7] applies only to real dictionary learning, we extend the identifiability result from the real case to the complex case, i.e., allowing both the dictionary matrix \mathbf{A} and the sparse coefficients matrix \mathbf{S} to have complex values. A complex matrix \mathbf{S} is sparse if a significant portion of its values are equal to zero, but for its nonzero values we do not care if their real or imaginary parts are zero—they are simply treated as nonzeros. Second and more importantly, we extend the sample analyses to a much wider class of generative models for \mathbf{S} called Bernoulli-subexponential models; compared to the Bernoulli-Gaussian model that is widely adopted in DL research, we now allow the nonzero values of \mathbf{S} to be drawn from any subexponential distributions, which includes not only Gaussian but also a large variety of distributions such as Laplacian, Cauchy, chi, and any bounded ones [8].

II. IDENTIFIABILITY OF DICTIONARY LEARNING

We denote the groundtruth factorization of the DL model as $\mathbf{X} = \mathbf{A}_\dagger \mathbf{S}_\dagger$. Suppose the rows of \mathbf{X} are linearly independent, which is a valid assumption since $k > n$, then the row space of \mathbf{X} is the same as that of \mathbf{S}_\dagger . Consider another factorization $\mathbf{X} = \mathbf{A}\mathbf{S}$, then the same rule applies, which means that there exists an invertible matrix \mathbf{W} such that $\mathbf{S} = \mathbf{W}\mathbf{S}_\dagger$. Furthermore, both \mathbf{S} and \mathbf{S}_\dagger are right-invertible, and they are related as $\mathbf{S}^\dagger = \mathbf{S}_\dagger^\dagger \mathbf{W}^{-1}$; multiplying it from the right on both sides of $\mathbf{A}\mathbf{S} = \mathbf{A}_\dagger \mathbf{S}_\dagger$ shows that $\mathbf{A} = \mathbf{A}_\dagger \mathbf{W}^{-1}$. This shows that matrix factorization in general is not identifiable without any structural constraints. On the other hand, people have discovered a plethora of structural constraints that do guarantee identifiability, i.e., restricting the choice of \mathbf{W} to be the product of a permutation matrix

Supported in part by NSF ECCS-2237640 and NIH R01LM014027.

and a diagonal matrix, such as nonnegative [9]–[11], simplicial [12], [13], or bounded [14]–[16].

A breakthrough result regarding identifiability of DL has been shown in [7] using the following formulation:

$$\begin{aligned} & \underset{A, S}{\text{minimize}} \quad |\det A| \\ & \text{subject to} \quad X = AS, \|S_{j,:}\|_1 \leq 1, j = 1, \dots, k, \end{aligned} \quad (1)$$

In the remainder of this section, we present the identifiability result in the complex domain, although the original result in [7] was in the real domain, and show that similar identifiability result can be obtained in a much more general complex case.

First we introduce some definitions.

Definition 1 (Cellular hull). The cellular hull of a finite set of vectors $\{s_1, \dots, s_n\}$, stacked as the columns of the matrix S , is

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

Notice that the definition works in both real and complex cases: in the real domain, it restricts the absolute values of the coefficients θ_i , and in the complex domain, it restricts the magnitudes of the complex coefficients θ_i .

Assumption 1 (Sufficiently scattered in the complex hypercube). Let \mathcal{B} denote the complex Euclidean ball $\mathcal{B} = \{x \in \mathbb{C}^k \mid \|x\| \leq 1\}$ and \mathcal{C} denote the complex hypercube $\mathcal{C} = \{x \in \mathbb{C}^k \mid \|x\|_\infty \leq 1\}$. A set S is sufficiently scattered in the complex hypercube if:

- 1) $\mathcal{B} \subseteq S \subseteq \mathcal{C}$;
- 2) $\partial \mathcal{B} \cap \partial S = \{\alpha e_i \mid |\alpha| = 1, i = 1, \dots, k\}$, where ∂ denotes the boundary of the set, and e_1, \dots, e_k are the k unit vectors in \mathbb{R}^k .

Our main result on the identifiability of complex DL is presented as follows:

Theorem 1. Consider the complex DL model $X = A_{\mathfrak{h}} S_{\mathfrak{h}}$, where $A_{\mathfrak{h}} \in \mathbb{C}^{d \times k}$ is the groundtruth mixing matrix and $S_{\mathfrak{h}}$ is the groundtruth sparse coefficient matrix. Let $\tilde{S}_{\mathfrak{h}}$ denote the matrix obtained from rescaling the rows of $S_{\mathfrak{h}}$ to have unit ℓ_1 norms. If $\text{rank}(A_{\mathfrak{h}}) = k$ and $\text{cell}(\tilde{S}_{\mathfrak{h}})$ is sufficiently scattered in the complex hypercube as in Assumption 1, then for any solution of (1), denoted as (A_\star, S_\star) , there exist a permutation matrix Π and a complex diagonal matrix D such that

$$A_{\mathfrak{h}} = A_\star D \Pi \quad \text{and} \quad S_{\mathfrak{h}} = \Pi^\top D^{-1} S_\star.$$

In other words, complex DL is identifiable if the groundtruth $A_{\mathfrak{h}}$ has full column rank and the cellular hull of $S_{\mathfrak{h}}$ is sufficiently scattered.

Proof. Let $\tilde{S}_{\mathfrak{h}}$ denote the matrix obtained from rescaling the rows of $S_{\mathfrak{h}}$ to have unit ℓ_1 norms, and counter-scale the columns of $A_{\mathfrak{h}}$ to obtain $\tilde{A}_{\mathfrak{h}}$. Since both $(\tilde{A}_{\mathfrak{h}}, \tilde{S}_{\mathfrak{h}})$ and (A_\star, S_\star) are feasible for (1), we immediately have that $|\det A_\star| \leq |\det \tilde{A}_{\mathfrak{h}}|$. Define $W = A_\star^{-1} \tilde{A}_{\mathfrak{h}}$, then

$$|\det W| = |\det A_\star^{-1} \tilde{A}_{\mathfrak{h}}| = |\det \tilde{A}_{\mathfrak{h}}| / |\det A_\star| \geq 1. \quad (2)$$

On the other hand, since $W \tilde{S}_{\mathfrak{h}} = A_\star^{-1} \tilde{A}_{\mathfrak{h}} \tilde{S}_{\mathfrak{h}} = S_\star$, we also have that $\|W \tilde{S}_{\mathfrak{h}}\|_{j,:} \leq 1$. Let w^H be any row of W , then $\|w^H \tilde{S}_{\mathfrak{h}}\|_1 \leq 1$. This is equivalent to a set of linear inequalities

$$w^H \tilde{S}_{\mathfrak{h}} \theta \leq 1, \quad \forall \|\theta\|_\infty \leq 1.$$

Now we invoke the assumption that $\text{cell}(\tilde{S}_{\mathfrak{h}})$ is sufficiently scattered, then for any v with unit norm $\|v\| = 1$, there must exist θ with $\|\theta\|_\infty \leq 1$ such that $v = \tilde{S}_{\mathfrak{h}} \theta$. Therefore

$$|w^H v| = |w^H \tilde{S}_{\mathfrak{h}} \theta| \leq 1. \quad (3)$$

Since (3) holds for every v with unit norm, we must have $\|w\| \leq 1$ as well (otherwise we would let $v = w/\|w\|$ and have that $|w^H v| = \|w\| > 1$, contradicting (3)), which means every row of W has norm no greater than 1. This gives us

$$|\det W| \leq \prod_{j=1}^k \|w_j\| \leq 1, \quad (4)$$

where the first inequality is due to the Hadamard inequality. Combining (2) and (4) shows that $(A_{\mathfrak{h}}, S_{\mathfrak{h}})$ or any of their column permutation and unimodular scaling is in the set of optimal solutions of (1). The second requirement of Assumption 1 further ensures that they are the only possible solutions. This shows the identifiability of complete dictionary learning in the complex domain. \square

III. IMPROVED SAMPLE COMPLEXITY ANALYSIS

The analysis given in the previous section gives an exact characterization of when dictionary learning is identifiable using the proposed formulation (1) with a matrix volume identification criterion. The analysis is inspired by the line of work from nonnegative matrix factorization and simplicial component analysis [9], [11], [13] using geometric interpretations, but with the introduction of cellular hulls in the complex domain, the geometric interpretation becomes harder to visualize. Nevertheless, there is a simple algebraic representation: $\text{cell}(\tilde{S})$ is sufficiently scattered in the complex hypercube if the solution set of the following optimization problem is $\{\alpha e_j \mid |\alpha| = 1, j = 1, \dots, k\}$:

$$\begin{aligned} & \underset{w}{\text{maximize}} \quad \|w\|^2 \\ & \text{subject to} \quad \|\tilde{S}^H w\|_1 \leq 1. \end{aligned} \quad (5)$$

Solving (5) exactly is NP-hard. In this section, we assume that the sparse coefficient matrix S is generated from a probabilistic generative model and show that it would satisfy guarantee identifiability with high probability, provided that the number of data points n is $O(k \log k)$. What differs from the previous work [7] as well as several earlier work on local identifiability is that here we consider a much broader class of generative model called Bernoulli-Subexponential model as introduced here.

Assumption 2 (Bernoulli-Subexponential model). The matrix $S \in \mathbb{R}^{k \times n}$ or $\mathbb{C}^{k \times n}$ is generated from a Bernoulli-Subexponential model with parameter $p \in (0, 1)$, denoted as $S \sim \mathcal{BE}(p)$, if its elements are i.i.d. with $S_{ij} = b_{ij} g_{ij}$, where $b_{ij} \in \{0, 1\}$ are i.i.d. Bernoulli random variables with $\Pr[b_{ij} = 1] = p$ (thus $\Pr[b_{ij} = 0] = 1 - p$) and g_{ij} are i.i.d. random variables drawn from a subexponential distribution with subexponential norm v .

First let us review the definition of subexponential distributions. As the name suggests, a random variable is called subexponential if its tail distribution can be upperbounded by an exponential function. There are several equivalent characterizations of subexponential distributions [8]. One of the more explicit definitions, which involves the notion of subexponential norms in Assumption 2, is as follows:

Definition 2 (Subexponential random variables). A random variable $Z \in \mathbb{R}$ or \mathbb{C} is called a sub-exponential random variable if the moment-generating function of $|Z|$ satisfies $E[\exp(|Z|/t)] \leq 2$ for some $t > 0$. The smallest t that satisfies such inequality is called the subexponential norm of Z , denoted as v in Assumption 2, i.e.,

$$v = \inf\{t > 0 : E[\exp(|Z|/t)] \leq 2\}.$$

A wide variety of random variables can be categorized as subexponential, including not only common ones such as Gaussian and

bounded, but also heavy-tailed ones such as Laplacian, Cauchy, and chi. By allowing the nonzero entries of \mathbf{S} to be drawn from any subexponential distribution, we greatly expand the applicability of the sample complexity analysis given in [7].

One of the nice properties of subexponential random variables is that sums of independent ones are concentrated around their mean, as shown in Bernstein's inequality:

Theorem 2 (Bernstein's inequality [8]). *Let Z_1, \dots, Z_n be i.i.d. subexponential random variables with mean μ and subexponential norm v . Then, for any $t > 0$, we have*

$$\Pr \left[\left| \frac{1}{n} \sum_{i=1}^n Z_i - \mu \right| > \epsilon \right] \leq 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{v^2}, \frac{\epsilon}{v} \right) \right),$$

where $c > 0$ is an absolute constant.

We are now ready to present the main theorem about improved sample complexity for identifiable dictionary learning:

Theorem 3. *Suppose $\mathbf{S} \in \mathbb{R}^{k \times n}$ or $\mathbb{C}^{k \times n}$ is generated from the Bernoulli-Subexponential model $\mathcal{BE}(p)$, and $\tilde{\mathbf{S}}$ is obtained by scaling its rows to have unit ℓ_1 norm. Then*

$$\Pr \left[\sup_{\|\mathbf{w}^\top \tilde{\mathbf{S}}\|_1 \leq 1} \|\mathbf{w}\| > 1 \right] \leq 4 \exp \left(k \log(3\sqrt{k}) - cn(\mu/v)^2 p(1-p) \right), \quad (6)$$

where μ and v are parameters of the subexponential distribution, and c is an absolute constant (c.f. Theorem 2).

Proof. We assume that $\mathbf{S} \sim \mathcal{BE}(p)$, and the first thing we do is rescaling its rows to have unit ℓ_1 norms and use it for the problem (5). To simplify the analysis, we can instead directly maximize $\|\mathbf{w}\|^2$ 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

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

with arbitrary choice of α . Conversely,

$$\begin{aligned} \Pr \left[\sup_{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > 1 \right] &\leq \Pr \left[\sup_{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > \alpha \cup \max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \\ &\leq \Pr \left[\sup_{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > \alpha \right] + \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \end{aligned} \quad (7)$$

where the second inequality is obtained from the union bound.

a) Bounding the second term in (7): Let $\mathbf{s} = (s_1, \dots, s_n)$ be generated from $\mathcal{BG}(p)$, i.e., each $s_i = b_i g_i$ with b_i Bernoulli with probability p and g_i subexponential with norm v , we will use Bernstein's inequality with $Z_i = b_i |g_i|$. If $\mathbb{E}[g_i] = \mu$, then $\mathbb{E} Z_i = p\mu$, and Z_i is also subexponential with norm $p v$. Therefore

$$\begin{aligned} \Pr [\|\mathbf{s}\|_1 < n(p\mu - \epsilon)] &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n\epsilon \right] \\ &\leq 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{p v} \right) \right) \end{aligned}$$

This puts a bound on the probability that one row of \mathbf{S} has bounded ℓ_1 norm. The second term in (7) requires all k rows to be bounded, which clearly has an even smaller probability, therefore

$$\Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \leq 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{p v} \right) \right). \quad (8)$$

b) Bounding the first term in (7): First we note the following equivalence:

$$\Pr \left[\sup_{\|\mathbf{w}^\top \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > \alpha \right] = \Pr \left[\inf_{\|\mathbf{w}\|_1 = 1} \|\mathbf{S}^\top \mathbf{w}\|_1 < 1/\alpha \right] \quad (9)$$

We are also going to use the following notion of δ -cover from convex geometry [17] that holds for all ℓ_p norm balls, but we are only to instantiate the Euclidean ball:

Lemma 1 (δ -cover). *A finite δ -cover of the unit sphere in \mathbb{R}^k is a finite set C_δ of points with unit ℓ_2 norm such that any point on the unit sphere is within ϵ away from an element in C_δ , i.e.*

$$\min_{\mathbf{w}_i \in C_\delta} \|\mathbf{w} - \mathbf{w}_i\| < \delta, \quad \forall \|\mathbf{w}\| = 1.$$

For $\delta \in (0, 1)$ there always exists a δ -cover C_δ with cardinality $|C_\delta| < (3/\delta)^k$. In the complex case, the ℓ_2 norm essentially treats the real and imaginary parts separately, therefore the δ -cover is $|C_\delta| < (3/\delta)^{2k}$

Let $C_\delta = \{\mathbf{w}_i\}$ be an δ -cover for the sphere in \mathbb{R}^k or \mathbb{C}^k . Assume that we have both the lowerbound

$$\|\mathbf{S}^\top \mathbf{w}_i\|_1 \geq \beta, \quad \forall \mathbf{w}_i \in C_\delta$$

and the upperbound

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

Then

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

Therefore

$$\inf_{\|\mathbf{w}\| \leq 1} \|\mathbf{S}^\top \mathbf{w}\|_1 \geq \beta - \gamma \delta \sqrt{k}.$$

As a result, we have

$$\begin{aligned} \Pr \left[\inf_{\|\mathbf{w}\| \leq 1} \|\mathbf{S}^\top \mathbf{w}\|_1 < \beta - \gamma \delta \sqrt{k} \right] \\ \leq \sum_{\mathbf{w}_i \in C_\delta} \Pr [\|\mathbf{S}^\top \mathbf{w}_i\|_1 < \beta] + \Pr [\|\mathbf{S}^\top\|_1 > \gamma], \end{aligned} \quad (10)$$

where C_δ is a δ -cover of the unit sphere with cardinality $|C_\delta| < (3/\delta)^k$ in \mathbb{R}^k or $|C_\delta| < (3/\delta)^{2k}$ in \mathbb{C}^k .

The bound to the first term in (10) is almost identical to (7). Dropping the subscript of \mathbf{w}_i , we write

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

Each Z_i is the absolute value of the weighted sum of k subexponentials with norm $p v$, so it is also subexponential with norm less than $p v$. Denote $\mu_w = \mathbb{E}[Z_i]$, then

$$\begin{aligned} \Pr [\|\mathbf{S}^\top \mathbf{w}\|_1 < n(\mu_w - \epsilon)] &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n\epsilon \right] \\ &\leq 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{p v} \right) \right). \end{aligned} \quad (11)$$

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}^H . This means we can use similar arguments used in (8) (but applied to the other direction) to have

$$\begin{aligned} \Pr \left[\|\mathbf{S}^H\|_1 > n(p\mu + \epsilon) \right] &= \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 > n(p\mu + \epsilon) \right] \\ &\leq \Pr \left[\|\mathbf{S}_{j,:}\|_1 > n(p\mu + \epsilon) \right] \\ &\leq \Pr \left[\left| \sum_{i=1}^n (Z_i - \mathbb{E}[Z_i]) \right| > n\epsilon \right] \\ &\leq 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{pv} \right) \right), \end{aligned} \quad (12)$$

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 are lowerbounded.

Combining (11), (11), and (12) with $\beta = n(\mu_w - \epsilon)$, $\gamma = n(p\mu + \epsilon)$, and

$$\delta = \frac{n^2(\mu_w - \epsilon)(p\mu - \epsilon) - 1}{n^2(p\mu + \epsilon)(p\mu - \epsilon)\sqrt{k}},$$

which satisfies $0 < \delta < 1$ for small enough ϵ , then we have

$$\begin{aligned} \Pr \left[\inf_{\|\mathbf{w}\|=1} \|\mathbf{S}^T \mathbf{w}\|_1 < 1/(n(p\mu - \epsilon)) \right] \\ \leq ((3/\delta)^k + 1) 2 \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{pv} \right) \right). \end{aligned} \quad (13)$$

Combining (13) and (8) with $\alpha = n(p\mu - \epsilon)$ into (7) gives us

$$\Pr \left[\sup_{\|\mathbf{w}^H \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > 1 \right] \quad (14)$$

$$\begin{aligned} &\leq \Pr \left[\inf_{\|\mathbf{w}\|=1} \|\mathbf{S}^H \mathbf{w}\|_1 < 1/\alpha \right] + \Pr \left[\max_j \|\mathbf{S}_{j,:}\|_1 < \alpha \right] \\ &\leq 2 \left(\left(\frac{3}{\delta} \right)^k + 2 \right) \exp \left(-cn \min \left(\frac{\epsilon^2}{p^2 v^2}, \frac{\epsilon}{pv} \right) \right). \end{aligned} \quad (15)$$

For any $0 < \epsilon < p\mu$ we have $\delta > 1/\sqrt{k}$; if we further require $\epsilon^2/p^2 v^2 < \epsilon/pv$, i.e. $\epsilon < pv$, then (14) becomes

$$\begin{aligned} \Pr \left[\sup_{\|\mathbf{w}^H \mathbf{S}\|_1 \leq 1} \|\mathbf{w}\| > 1 \right] &\leq 2 \left((3\sqrt{k})^k + 2 \right) \exp \left(-\frac{cn\epsilon^2}{p^2 v^2} \right) \\ &\leq 4 (3\sqrt{k})^k \exp \left(-\frac{cn\epsilon^2}{p^2 v^2} \right) \\ &\leq 4 \exp \left(k \log(3\sqrt{k}) - cn(\mu/v)^2 p(1-p) \right). \end{aligned}$$

□

Theorem 3 shows that if \mathbf{S} is generated from a Bernoulli-Subexponential model, then the probability that it does not satisfy the sufficiently scattered condition can be upperbounded by (6), which approaches zero at an exponential rate as $n \gg k \log(2\sqrt{k})/p(1-p)$. As a result, we have the following corollary that directly bounds the probability that such a Bernoulli-Subexponential \mathbf{S} can be uniquely identified via complex dictionary learning.

Corollary 1. Consider the generative model $\mathbf{X} = \mathbf{A}_\mathfrak{h} \mathbf{S}_\mathfrak{h}$, where $\mathbf{A}_\mathfrak{h} \in \mathbb{C}^{k \times k}$ is the groundtruth dictionary and \mathbf{S} is the groundtruth sparse coefficients. If $\text{rank}(\mathbf{A}_\mathfrak{h}) = k$ and the matrix $\mathbf{S}_\mathfrak{h} \in \mathbb{C}^{k \times n}$ is generated from the Bernoulli-Subexponential model $\mathcal{BE}(p)$, then $(\mathbf{A}_\mathfrak{h}, \mathbf{S}_\mathfrak{h})$ are globally identifiable via optimizing (1) with probability at least

$$1 - 4 \exp \left(k \log(3\sqrt{k}) - cn(\mu/v)^2 p(1-p) \right).$$

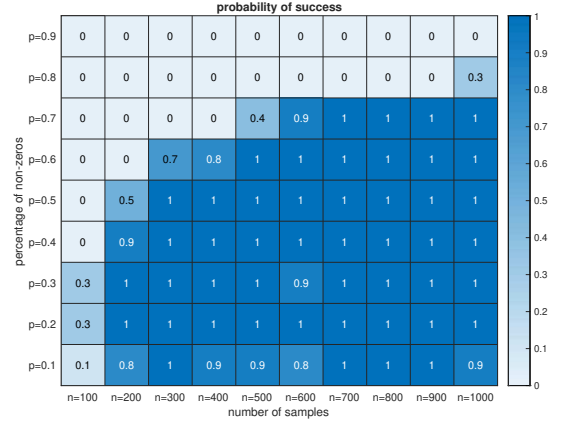


Fig. 1. Groundtruth $\mathbf{S}_\mathfrak{h}$ is generated from a Bernoulli-Laplacian model. The probability of exactly recovering the dictionary for various sample size n and sparsity p (probability of nonzeros in $\mathbf{S}_\mathfrak{h}$). The dictionary size is fixed with $k = 20$.

IV. EXPERIMENTS

We justify the improved sample complexity analysis of identifiable dictionary learning in the following numerical experiment conducted in MATLAB. The algorithm to approximately solve formulation (1) is the L-ADMM algorithm proposed in [7]. Different from the simulations done in [7], the groundtruth dictionary $\mathbf{S}_\mathfrak{h}$ is generated from a Bernoulli-Laplacian model; the Laplacian distribution has a much heavier tail distribution compared to Gaussian, but is still in the category of subexponential distributions, and we demonstrate that similar identifiability results still hold here, as confirmed by Theorem 3.

We fix the dictionary size $k = 20$, and change the sample size n from 100 to 1000 and the probability p in the Bernoulli-Subexponential model from 0.1 to 0.9. To ensure definitively that the dictionary is recovered, we will normalize the column and use the Hungarian algorithm [18] to find the best column matching, and then calculate the estimation error. We declare success if the estimation error is smaller than $1e-5$. In each setting the experiments are repeated 10 times and the percentage of exact recovery are recorded. The results are shown in Figure 1, which agrees with the bound in Theorem 3.

V. CONCLUSION

In this work we greatly improve the identifiability result of complete dictionary learning in [7] in two ways. First, the result is generalized to the complex domain, while [7] was only shown in the real domain. Second, we show that the same sample complexity analysis could be extended to a wider class of generative model in which the nonzeros of the sparse coefficients can be drawn from any subexponential distributions, as opposed to only Gaussian models in the previous work. We demonstrate the recovery performance in Monte-Carlo simulation where the nonzeros of the sparse coefficients are drawn from the Laplacian distributions, and the results show that the identifiability agrees with the developed theory.

REFERENCES

- [1] I. Tošić and P. Frossard, "Dictionary learning," *IEEE Signal Processing Magazine*, vol. 28, no. 2, pp. 27–38, 2011.
- [2] M. Aharon, M. Elad, and A. M. Bruckstein, "On the uniqueness of overcomplete dictionaries, and a practical way to retrieve them," *Linear algebra and its applications*, vol. 416, no. 1, pp. 48–67, 2006.

- [3] C. J. Hillar and F. T. Sommer, "When can dictionary learning uniquely recover sparse data from subsamples?" *IEEE Transactions on Information Theory*, vol. 61, no. 11, pp. 6290–6297, 2015.
- [4] R. Gribonval and K. Schnass, "Dictionary identification—sparse matrix-factorization via ℓ_1 -minimization," *IEEE Transactions on Information Theory*, vol. 56, no. 7, pp. 3523–3539, 2010.
- [5] S. Wu and B. Yu, "Local identifiability of ℓ_1 -minimization dictionary learning: a sufficient and almost necessary condition," *The Journal of Machine Learning Research*, vol. 18, no. 1, pp. 6121–6176, 2017.
- [6] Y. Wang, S. Wu, and B. Yu, "Unique sharp local minimum in ℓ_1 -minimization complete dictionary learning," *Journal of Machine Learning Research*, vol. 21, no. 63, pp. 1–52, 2020.
- [7] J. Hu and K. Huang, "Global identifiability of ℓ_1 -based dictionary learning via matrix volume optimization," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 36, 2023.
- [8] R. Vershynin, *High-dimensional Probability: An Introduction with Applications in Data Science*. Cambridge university press, 2018, vol. 47.
- [9] K. Huang, N. D. Sidiropoulos, and A. Swami, "Non-negative matrix factorization revisited: Uniqueness and algorithm for symmetric decomposition," *IEEE Transactions on Signal Processing*, vol. 62, no. 1, pp. 211–224, 2013.
- [10] K. Huang, X. Fu, and N. D. Sidiropoulos, "Anchor-free correlated topic modeling: Identifiability and algorithm," *Advances in Neural Information Processing Systems*, vol. 29, 2016.
- [11] X. Fu, K. Huang, and N. D. Sidiropoulos, "On identifiability of nonnegative matrix factorization," *IEEE Signal Processing Letters*, vol. 25, no. 3, pp. 328–332, 2018.
- [12] X. Fu, W.-K. Ma, K. Huang, and N. D. Sidiropoulos, "Blind separation of quasi-stationary sources: Exploiting convex geometry in covariance domain," *IEEE Transactions on Signal Processing*, vol. 63, no. 9, pp. 2306–2320, 2015.
- [13] K. Huang and X. Fu, "Detecting overlapping and correlated communities without pure nodes: Identifiability and algorithm," in *International Conference on Machine Learning*, 2019, pp. 2859–2868.
- [14] G. Tatli and A. T. Erdogan, "Polytopic matrix factorization: Determinant maximization based criterion and identifiability," *IEEE Transactions on Signal Processing*, vol. 69, pp. 5431–5447, 2021.
- [15] J. Hu and K. Huang, "Identifiable bounded component analysis via minimum volume enclosing parallelotope," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2023.
- [16] —, "Complex bounded component analysis: Identifiability and algorithm," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024.
- [17] G. Pisier, *The volume of convex bodies and Banach space geometry*. Cambridge University Press, 1999, vol. 94.
- [18] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.