

A Provable Approach for Double-Sparse Coding

Thanh V. Nguyen

ECE Department
Iowa State University
thanhnv@iastate.edu

Raymond K. W. Wong

Statistics Department
Texas A&M University
raywong@tamu.edu

Chinmay Hegde *

ECE Department
Iowa State University
chinmay@iastate.edu

Abstract

Sparse coding is a crucial subroutine in algorithms for various signal processing, deep learning, and other machine learning applications. The central goal is to learn an overcomplete dictionary that can sparsely represent a given dataset. However, storage, transmission, and processing of the learned dictionary can be untenably high if the data dimension is high. In this paper, we consider the double-sparsity model introduced by Rubinstein, Zibulevsky, and Elad (2010) where the dictionary itself is the product of a fixed, known basis and a data-adaptive sparse component. First, we introduce a simple algorithm for double-sparse coding that can be amenable to efficient implementation via neural architectures. Second, we theoretically analyze its performance and demonstrate asymptotic sample complexity and running time benefits over existing (provable) approaches for sparse coding. To our knowledge, our work introduces the first computationally efficient algorithm for double-sparse coding that enjoys rigorous statistical guarantees. Finally, we support our analysis via several numerical experiments on simulated data, confirming that our method can indeed be useful in problem sizes encountered in practical applications.

Introduction

We consider the problem of *dictionary learning* (also known as sparse coding), a common and powerful technique in unsupervised feature learning. The high-level idea of sparse coding is to represent a set of data vectors in terms of *sparse* linear combinations of atoms from a learned basis (or dictionary). Sparse coding has a rich history in diverse fields such as image processing, machine learning, and neuroscience (Krim et al. 1999; Elad and Aharon 2006; Rubinstein, Bruckstein, and Elad 2010; Mairal et al. 2009). Sparse coding forms a core component of several neural learning systems, both biological (Olshausen and Field 1997) and artificial (Gregor and LeCun 2010; Boureau et al. 2010; Mazumdar and Rawat 2017).

Formally, suppose we are given p data samples $Y = [y^{(1)}, y^{(2)}, \dots, y^{(p)}] \in \mathbb{R}^{n \times p}$. We wish to find a dictionary $D \in \mathbb{R}^{n \times m}$ (with $n < m$) and corresponding sparse code

vectors $X = [x^{(1)}, x^{(2)}, \dots, x^{(p)}] \in \mathbb{R}^{m \times p}$ such that the representation DX fits the data samples as well as possible. The typical approach is to pose the dictionary (and codes) as the solution to the constrained optimization problem:

$$\begin{aligned} \min_{D, X} \mathcal{L}(D, X) &= \frac{1}{2} \sum_{j=1}^p \|y^{(j)} - Dx^{(j)}\|_2^2, \\ \text{s.t. } \sum_{j=1}^p \mathcal{S}(x^{(j)}) &\leq S, \end{aligned} \quad (1)$$

where $\mathcal{S}(\cdot)$ is some sparsity-inducing penalty function, such as the ℓ_1 -norm. However, even a cursory attempt at solving (1) reveals several conceptual obstacles:

Theoretical challenges. The constrained optimization problem (1) involves a non-convex (bilinear) objective function, as well as potentially non-convex constraints depending on the function \mathcal{S} . Therefore, design and analysis of *provably* correct algorithms for this problem can be difficult. Indeed, the vast majority of practical approaches for sparse coding are based on heuristics; barring a few recent papers from the learning theory community (Spielman, Wang, and Wright 2012; Agarwal et al. 2014; Arora et al. 2015; Sun, Qu, and Wright 2015), very few methods come equipped with global correctness guarantees.

Challenges in applications. Even if we ignore theoretical correctness issues and somehow are able to learn good enough sparse codes, we often find that applications using such learned sparse codes encounter *memory* and *running-time* issues. Indeed, in the overcomplete case, merely storing the learned dictionary D incurs $mn = \Omega(n^2)$ memory cost, which is prohibitive when n is large. In practical applications such as image analysis, one typically resorts to chopping the data into smaller blocks (e.g., partitioning image data into patches) to make the problem manageable.

An approach used to resolve those practical computational difficulties is to assume some type of structure in the (learned) dictionary D ; e.g., the dictionary is assumed to be either separable, or obey a convolutional structure. One such variant is *double-sparse* coding (Rubinstein, Zibulevsky, and Elad 2010; Sulam et al. 2016) where the dictionary D itself exhibits a sparse structure. More precisely,

$$D = \Phi A,$$

*This work is supported in part by the National Science Foundation under the grants CCF-1566281 and DMS-1612985. Copyright © 2018, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

| Setting | Reference | Sample complexity (w/o noise) | Sample complexity (w/ noise) | Upper bound on running time | Expt |
|---------------|---|-------------------------------|--|-----------------------------|--------------|
| Regular | MOD (Engan, Aase, and Husoy 1999) | \times | \times | \times | \checkmark |
| | K-SVD (Aharon, Elad, and Bruckstein 2006) | \times | \times | \times | \checkmark |
| | (Spielman, Wang, and Wright 2012) | $O(n^2 \log n)$ | \times | $\tilde{\Omega}(n^4)$ | \checkmark |
| | (Arora, Ge, and Moitra 2014) | $\tilde{O}(m^2/k^2)$ | \times | $\tilde{O}(np^2)$ | \times |
| | (Gribonval, Jenatton, and Bach 2015) | $O(nm^3)$ | $O(nm^3)$ | \times | \times |
| | (Arora et al. 2015) | $\tilde{O}(mk)$ | \times | $\tilde{O}(mn^2p)$ | \times |
| Double Sparse | Double Sparsity (Rubinstein, Zibulevsky, and Elad 2010) | \times | \times | \times | \checkmark |
| | (Gribonval et al. 2015) | $\tilde{O}(mr)$ | $\tilde{O}(mr)$ | \times | \times |
| | Trainlets (Sulam et al. 2016) | \times | \times | \times | \checkmark |
| | This paper | $\tilde{O}(mr)$ | $\tilde{O}(mr + \sigma_\varepsilon^2 \frac{mnr}{k})$ | $\tilde{O}(mnp)$ | \checkmark |

Table 1: Comparison of various sparse coding techniques. Expt: whether numerical experiments have been conducted. \times in all other columns indicates no provable guarantees. Here, n is the signal dimension, and m is the number of atoms. The sparsity levels for A and x are r and k respectively, and p is the sample size.

with a known “base dictionary” $\Phi \in \mathbb{R}^{n \times n}$ and a learned column-sparse “synthesis” matrix $A \in \mathbb{R}^{n \times m}$. The base dictionary Φ is typically any orthonormal basis (such as the canonical or wavelet basis) chosen according to domain knowledge, while the synthesis matrix A is *column-wise sparse* and is learned from the data. Such a double-sparsity assumption is appealing conceptually, since it lets us combine the knowledge of good dictionaries Φ to synthesize new representations tailored to specific data families. Moreover, the sparse structure of the synthesis matrix produces more interpretable features than the regular model; see (Rubinstein, Zibulevsky, and Elad 2010; Sulam et al. 2016) for an extensive discussion and illustration. If the columns of A has only (say) $r \ll n$ non-zero elements, then the overall memory burden of storing and transmitting A is $O(mr)$, which is much lower than that for general unstructured dictionaries. Moreover, this approach performs comparably with (regular) sparse coding approaches, as demonstrated by the extensive empirical evaluations in (Sulam et al. 2016). However, two obstacles remain: the associated *training* algorithms used to learn double-sparse codes incur significant running time, and no rigorous theoretical analysis of their performance has been reported in the literature.

Our Contributions

We provide a new algorithmic approach to double-sparse coding. To the best of our knowledge, our approach is the first method that enjoys *provable* statistical and algorithmic guarantees for the double-sparse coding problem. In addition, our approach enjoys three benefits: (i) our method is *tractable* as well as *neurally plausible*, i.e., its execution can plausibly be achieved using a neural network architecture; (ii) our method enjoys *noise robustness* guarantees; (iii) we demonstrate *practical relevance* via several simulations.

Inspired by the aforementioned recent theoretical papers in sparse coding, we assume a learning-theoretic setup where the data samples arise from a ground-truth generative

model. Informally, suppose there exists a true (but unknown) synthesis matrix $A^* \in \mathbb{R}^{n \times m}$ whose columns have only r non-zero elements, and the i^{th} data sample is generated as:

$$y^{(i)} = \Phi A^* x^{*(i)} + \text{noise}, \quad i = 1, 2, \dots, p,$$

where the code vector $x^{*(i)}$ is independently drawn from a distribution supported on the set of k -sparse vectors. We desire to learn the matrix A^* . We suppose that the synthesis matrix A^* is *incoherent* (the columns of A^* are sufficiently close to orthogonal) and has bounded spectral norm, that m is at most a constant multiple of n , and that the noise is sub-Gaussian. All these assumptions are standard¹.

First, we propose and analyze an algorithm that produces a coarse estimate of the synthesis matrix that is sufficiently close to the ground truth A^* . Our method builds upon the method of *spectral initialization* that have recently gained popularity in non-convex machine learning (Zhang et al. 2016; Wang, Zhang, and Gu 2016).

Second, given such a coarse estimate of the synthesis matrix A^* , we propose and analyze a gradient descent-style algorithm to refine this estimate. This algorithm is simpler than previously studied double-sparse coding algorithms that rely on alternating minimization (such as the Trainlets approach of (Sulam et al. 2016)), while still giving good statistical performance.

Put together, the above constitutes the first provably polynomial-time method for double-sparse coding. In particular, in the absence of noise, we prove that $p = \Omega(mr \text{ polylog } n)$ samples are sufficient to obtain a good enough estimate in the initialization, and also to obtain guaranteed linear convergence during descent to provably recover A^* . See Table 1 for the summary and a comparison with the existing work. Indeed, our sample complexity result matches with what is achieved in (Gribonval et al. 2015). Nevertheless, we provide a practical polynomial-time algorithm for learning the sparse dictionary whereas Gribonval

¹We clarify the generative model in concrete terms below.

et al. only study properties of a theoretical estimator. Also, our approach results in strict improvement in sample complexity, as well as running time over rigorous methods for (regular) sparse coding, such as (Arora et al. 2015).

We analyze our approach in a more realistic setting with the presence of additive noise, and demonstrate its stability. While our analysis mainly consists of sufficiency results and involves several (absolute) unspecified constants, in practice we have found that these constants are reasonable. We justify our observations by reporting a suite of numerical experiments on synthetic test datasets.

Techniques

The remainder of the paper is fairly technical; therefore, for clarity let us provide some non-rigorous intuition for our approach. At a high level, our method extends the neural sparse coding approach of (Arora et al. 2015) to the double-sparse case. A major barrier in the analysis of sparse coding algorithms is that the gradient of \mathcal{L} in (1) with respect to D inherently depends on the codes of the training samples (i.e., the columns of X), but these codes are unknown *a priori*. However, the main insight in (Arora et al. 2015) is that within a small enough neighborhood of the true dictionary, an approximated version of X^* can be estimated, and therefore the overall method is similar to performing *approximate* gradient descent towards the population parameter A^* . Regarding the actual algorithm as its noisy variation allows us to overcome the finite-sample variability of the loss, and obtain a descent property directly related to A^* .

The descent stage of our approach leverages this intuition. However, instead of standard gradient descent, we perform *approximate projected* gradient descent so that the column-wise r -sparsity property is enforced in each new estimate of A^* . This extra projection step is critical in showing sample complexity improvement over (regular) sparse coding methods. The key novelty is in figuring out how to perform the projection in each gradient iteration. For this purpose, we develop a novel initialization algorithm that identifies the locations of the non-zeroes in A^* even before commencing the descent phase. This is non-trivially different from previous rigorous methods for sparse coding, and the analysis is somewhat more involved.

In (Arora et al. 2015), (the principal eigenvector of) the weighted covariance matrix of y , given by a suitable weighted average of outer products $y_i y_i^T$, is shown to provide a coarse estimate of a given dictionary atom. We leverage this idea and rigorously show that the diagonal of the weighted covariance matrix serves as a good indicator of the support of a column in A^* . The success relies on the concentration of the diagonal vector with dimension n , instead of the covariance matrix with dimensions $n \times n$. With the support selected, our scheme only utilizes a *truncated* weighted covariance matrix with dimensions $r \times r$. This initialization scheme enables us to effectively reduce the dimension of the problem, and therefore leads to significant improvement in sample complexity and running time over previous (provable) sparse coding methods when the data representation sparsity k is much smaller than m .

Further, we rigorously analyze the proposed algorithms

in the presence of noise with a bounded expected norm. Our analysis shows that our method is stable, and in the case of i.i.d. Gaussian noise with bounded expected ℓ_2 -norms, is at least a polynomial factor better than previous polynomial time algorithms for sparse coding in terms of running time. Our analysis of the descent stage follows from (Arora et al. 2015), where the descent property is first shown under an ideal algorithm which uses the expectation of the noisy (approximate) gradient, and is later established to the practical case via a concentration argument. Our novel initialization algorithm allows an accurate determination of the support of A^* , and therefore, for each column of A^* , we can focus on an r -dimensional subvector of the noisy (approximate) gradient vector, rather than the full n -dimensional vector. This allows us to sharpen the sample complexity beyond what has been established in the earlier work.

Setup and Definitions

Notation. Let $[m] \triangleq \{1, 2, \dots, m\}$ for some integer m . For any vector $x = [x_1, x_2, \dots, x_m]^T \in \mathbb{R}^m$, let $\text{supp}(x) \triangleq \{i \in [m] : x_i \neq 0\}$. Given any subset $S \subseteq [m]$, x_S corresponds to the sub-vector of x indexed by the elements of S . For any matrix $A \in \mathbb{R}^{n \times m}$, we use $A_{\bullet i}$ and $A_{j \bullet}^T$ to represent the i^{th} column and the j^{th} row respectively. For some appropriate sets R and S , let $A_{R \bullet}$ (respectively, $A_{\bullet S}$) be the submatrix of A with rows (respectively columns) indexed by the elements in R (respectively S). For the i^{th} column $A_{\bullet i}$, use $A_{R, i}$ to denote the sub-vector indexed by the elements of R . Use $A_{R \bullet}^T$ to indicate $(A_{R \bullet})^T$. Let \circ and $\text{sgn}(\cdot)$ represent the (element-wise) Hadamard operator and sign function. Further, $\text{threshold}_K(x)$ is a thresholding operator that replaces any elements of x with magnitude less than K by zero.

The ℓ_2 -norm $\|x\|$ for a vector x and the spectral norm $\|A\|$ for a matrix A are used extensively in this paper. In some cases, we also utilize the Frobenius norm $\|A\|_F$ and a special matrix operator norm $\|A\|_{1,2} \triangleq \max_{\|x\|_1 \leq 1} \|Ax\|$.

For clarity purposes, we adopt big-Oh notation extensively. The symbols $\tilde{\Omega}(\cdot)$ and $\tilde{O}(\cdot)$ represent $\Omega(\cdot)$ and $O(\cdot)$ up to a multiplicative poly-logarithmic factor of n respectively. Throughout the paper, we use the phrase “with high probability” (abbreviated to w.h.p.) to describe an event with failure probability of $O(n^{-\omega(1)})$. In addition, $g(n) = O^*(f(n))$ means $g(n) \leq Kf(n)$ for some small enough constant K .

Model. Suppose that the observed samples are given by

$$y^{(i)} = Dx^{*(i)} + \varepsilon, \quad i = 1, \dots, p;$$

i.e., we are given p samples of y generated from a fixed (but unknown) dictionary D where the sparse code x^* and the error ε are drawn from a joint distribution \mathcal{D} specified below. In the double-sparse setting, the dictionary is assumed to follow a decomposition $D = \Phi A^*$, where $\Phi \in \mathbb{R}^{n \times n}$ is a known orthonormal basis matrix and A^* is an unknown, ground truth synthesis matrix. Our approach relies upon the following assumptions on the synthesis dictionary A^* :

- A1** The dimensions of A^* obey $m = O(n)$.
- A2** A^* is μ -incoherent, i.e., for $i \neq j$, $|\langle A_{\bullet i}^*, A_{\bullet j}^* \rangle| \leq \mu/\sqrt{n}$.

A3 $A_{\bullet i}^*$ has exactly r non-zero elements, and is normalized such that $\|A_{\bullet i}^*\| = 1$ for all i . Moreover, $|A_{ij}^*| \geq \tau$ for $A_{ij}^* \neq 0$ and $\tau = \Omega(1/\sqrt{r})$.

A4 A^* has bounded spectral norm: $\|A^*\| \leq O(\sqrt{m/n})$.

These assumptions are standard. In Assumption **A2**, the incoherence μ is $O(1/\log n)$ with high probability for a normal random matrix (Arora, Ge, and Moitra 2014). Assumption **A3** is a common assumption for sparse signal recovery². Assumption **A4** is also standard (Arora et al. 2015). In addition to Assumptions **A1-A4**, we make the following distributional assumptions on \mathcal{D} :

B1 The support $S = \text{supp}(x^*)$ is of size at most k ; its indices are uniformly drawn without replacement from $[m]$.

B2 The nonzero entries x_S^* are pairwise independent and sub-Gaussian conditioned on the support S , with $\mathbb{E}[x_i^* | i \in S] = 0$ and $\mathbb{E}[x_i^{*2} | i \in S] = 1$.

B3 For $i \in S$, $|x_i^*| \geq C$ where $0 < C \leq 1$.

B4 The additive noise ε has i.i.d. Gaussian entries with variance σ_ε^2 with $\sigma_\varepsilon = O(1/\sqrt{n})$.

Similar sub-Gaussian models for \mathcal{D} have been previously considered in (Jenatton, Gribonval, and Bach 2012).

For the rest of the paper, for notational simplicity we set $\Phi = I_n$, i.e., the identity matrix. This does not affect anything, since one can study the equivalent problem:

$$y' = A^* x^* + \varepsilon',$$

where $y' = \Phi^T y$ and $\varepsilon' = \Phi^T \varepsilon$. Due to the Gaussian assumption on ε , it follows that ε' also has independent Gaussian entries. The analysis can be extended to sub-Gaussian noise with several minor (but tedious) changes.

Our goal is to devise an algorithm that produces an provably “good” estimate of A^* . For this, we need to define a suitable measure of “goodness”. We use the following notion of distance that measures the maximal column-wise difference in ℓ_2 -norm under a suitable transformation.

Definition 1 ((δ, κ) -nearness). *A is said to be δ -close to A^* if there is a permutation $\pi : [m] \rightarrow [m]$ and a sign flip $\sigma : [m] \rightarrow \{\pm 1\}$ such that $\|\sigma(i)A_{\bullet \pi(i)} - A_{\bullet i}^*\| \leq \delta$ for every i . In addition, A is said to be (δ, κ) -near to A^* if $\|A_{\bullet \pi} - A^*\| \leq \kappa \|A^*\|$ also holds.*

For notational simplicity, in our theorems we simply replace π and σ in Definition 1 with the identity permutation $\pi(i) = i$ and the positive sign $\sigma(\cdot) = +1$ while keeping in mind that in reality we are referring to finding one element of the equivalence class of all permutations and sign flip transforms of A^* .

We will also need some technical tools from (Arora et al. 2015) to analyze gradient descent-style methods. Consider an iterative algorithm that looks for a desired solution $z^* \in \mathbb{R}^n$ to optimize some function $f(z)$. Suppose that the

²The requirement of exactly r non-zero elements is merely for simplicity and there is no technical difficulty to extend our algorithms and corresponding analyses to the case with at most r non-zero elements.

algorithm produces a sequence of estimates z^1, \dots, z^s via the update rule:

$$z^{s+1} = z^s - \eta g^s,$$

for some vector g^s and scalar step size η . The goal is to characterize “good” directions g^s such that the sequence converges to z^* under the Euclidean distance. The following gives one such sufficient condition for g^s .

Definition 2. *A vector g^s at the s^{th} iteration is $(\alpha, \beta, \gamma_s)$ -correlated with a desired solution z^* if*

$$\langle g^s, z^s - z^* \rangle \geq \alpha \|z^s - z^*\|^2 + \beta \|g^s\|^2 - \gamma_s.$$

We know from convex optimization that if f is 2α -strongly convex and $1/2\beta$ -smooth, and g^s is chosen as the gradient $\nabla_z f(z)$, then g^s is $(\alpha, \beta, 0)$ -correlated with z^* . In our setting, the desired solution corresponds to A^* , the ground-truth synthesis matrix. In (Arora et al. 2015), it is shown that $g^s = \mathbb{E}_y[(A^s x - y)\text{sgn}(x)^T]$, where $x = \text{threshold}_{C/2}((A^s)^T y)$ indeed satisfies Definition 2. This g^s is a population quantity and not explicitly available, but one can estimate such g^s using an empirical average. The corresponding estimator \hat{g}^s is a random variable, so we also need a related *correlated-with-high-probability* condition:

Definition 3. *A direction \hat{g}^s at the s^{th} iteration is $(\alpha, \beta, \gamma_s)$ -correlated-w.h.p. with a desired solution z^* if, w.h.p.,*

$$\langle \hat{g}^s, z^s - z^* \rangle \geq \alpha \|z^s - z^*\|^2 + \beta \|\hat{g}^s\|^2 - \gamma_s.$$

From Definition 2, one can establish a form of descent property in each update step, as shown in Theorem 1.

Theorem 1 (Convergence of approximate gradient descent). *Suppose that g^s satisfies the condition described in Definition 2 for $s = 1, 2, \dots, T$. Moreover, $0 < \eta \leq 2\beta$ and $\gamma = \max_{s=1}^T \gamma_s$. Then, the following holds for all s :*

$$\|z^{s+1} - z^*\|^2 \leq (1 - 2\alpha\eta) \|z^s - z^*\|^2 + 2\eta\gamma_s.$$

In particular, the above update converges geometrically to z^ with an error γ/α . That is,*

$$\|z^{s+1} - z^*\|^2 \leq (1 - 2\alpha\eta)^s \|z^0 - z^*\|^2 + 2\gamma/\alpha.$$

We can obtain a similar result for Definition 3 except that $\|z^{s+1} - z^*\|^2$ is replaced with its expectation.

Armed with the above tools, we are now ready to introduce our method. As discussed above, our approach consists of two stages: an initialization algorithm that produces a coarse estimate of A^* , and a descent-style algorithm that refines this estimate to accurately recover A^* .

Stage 1: Initialization

The first stage of our approach iteratively estimates the columns of A^* (up to sign flips) in a manner similar to (Arora et al. 2015). However, their initialization algorithm incurs severe computational costs in terms of running time. More precisely, the expected value of the running time is $\tilde{\Omega}(mn^2p)$, which is unrealistic for large m and n .

In contrast, we leverage the double-sparsity assumption in our generative model to obtain a more efficient approach.

Algorithm 1 Truncated Pairwise Reweighting

Initialize $L = \emptyset$

Randomly divide p samples into two disjoint sets \mathcal{P}_1 and \mathcal{P}_2 of sizes p_1 and p_2 respectively

While $|L| < m$. Pick u and v from \mathcal{P}_1 at random
For every $l = 1, 2, \dots, n$, compute

$$\hat{e}_l = \frac{1}{p_2} \sum_{i=1}^{p_2} \langle y^{(i)}, u \rangle \langle y^{(i)}, v \rangle (y_l^{(i)})^2$$

Sort $\hat{e}_1, \hat{e}_2, \dots, \hat{e}_n$ in descending order

If $\hat{e}_{(r)} \geq \Omega(k/mr) \wedge \hat{e}_{(r+1)}/\hat{e}_{(r)} < O^*(r/\log^2 n)$

Let \hat{R} be set of the r largest entries of \hat{e}

$$\widehat{M}_{u,v} = \frac{1}{p_2} \sum_{i=1}^{p_2} \langle y^{(i)}, u \rangle \langle y^{(i)}, v \rangle y_{\hat{R}}^{(i)} (y_{\hat{R}}^{(i)})^T$$

$\delta_1, \delta_2 \leftarrow$ top singular values of $\widehat{M}_{u,v}$

$z_{\hat{R}} \leftarrow$ top singular vector of $\widehat{M}_{u,v}$

If $\delta_1 \geq \Omega(k/m)$ and $\delta_2 < O^*(k/m \log n)$

If $\text{dist}(\pm z, l) > 1/\log n$ for any $l \in L$

Update $L = L \cup \{z\}$

Return $A^0 = (L_1, \dots, L_m)$

The key ingredient of our method is a novel spectral procedure that gives us an estimate of the column supports purely from the observed samples. The full algorithm, that we call *Truncated Pairwise Reweighting*, is listed in pseudocode form as Algorithm 1.

We first state a theoretical result characterizing the performance of Algorithm 1.

Theorem 2. *Suppose that Assumptions B1-B4 hold and Assumptions A1-A4 hold with parameters $\mu = O^*(\frac{\sqrt{n}}{k \log^3 n})$, $k = O^*(\frac{\sqrt{n}}{\log n})$ and $r = o(\log^2 n)$. Then, with high probability, Algorithm 1 returns an initial estimate A^0 whose columns share the same support as A^* and is $(\delta, 2)$ -near to A^* with $\delta = O^*(1/\log n)$ if $p_1 = \tilde{\Omega}(m)$ and $p_2 = \tilde{\Omega}(mr)$.*

The formal proof is available in our extended version (Nguyen, Wong, and Hegde 2017). To provide some intuition about the working of the algorithm (and proof of Theorem 2), let us consider the setting where we have access to infinitely many samples. Of course, this setting is fictional. However, the analysis of this case is much simpler since we can deal with expected values rather than empirical averages. Moreover, the analysis reveals several key lemmas, which we will reuse extensively for proving Theorem 2.

First, we give some intuition behind the definition of the “scores”, \hat{e}_l . Fix a sample $y = A^* x^* + \varepsilon_y$ from the available training set, and consider two other samples

$$u = A^* \alpha + \varepsilon_u, \quad v = A^* \alpha' + \varepsilon_v.$$

Consider the (very coarse) estimate for the sparse code of u with respect to A^* :

$$\beta = A^{*T} u = A^{*T} A^* \alpha + A^{*T} \varepsilon_u.$$

As long as A^* is incoherent enough and $\varepsilon_u, \varepsilon_y$ is small, the estimate β “looks” like α in the following sense:

$$\langle y, u \rangle \approx \langle x^*, \beta \rangle \approx \langle x^*, \alpha \rangle.$$

Moreover, the above inner products are large only if α and x^* share some elements in their supports; else, they are likely to be small. Likewise, the weight $\langle y, u \rangle \langle y, v \rangle$ is large only when x^* shares common elements with both α and α' . The following lemma leverages this intuition; given sufficiently many samples, \hat{e}_l gives an indicator of how large the “overlap” between α and α' is.

Lemma 1. *Fix samples u and v . Suppose that $y = A^* x^* + \varepsilon$ is a random sample independent of u and v , whose codes α and α' have supports U and V respectively. Then*

$$e_l \triangleq \mathbb{E}[\langle y, u \rangle \langle y, v \rangle y_l^2] = \sum_{i \in U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} + E,$$

where $q_i = \mathbb{P}[i \in S]$, $q_{ij} = \mathbb{P}[i, j \in S]$ and $c_i = \mathbb{E}[x_i^4 | i \in S]$. Also, E has absolute value $O^*(k/m \log^2 n)$ w.h.p.

Now, suppose for a moment that u and v share exactly one common atom in their codes, i.e., $U \cap V = \{i\}$. Lemma 1 suggests that e_l is proportional to A_{li}^{*2} ; therefore, the scores e_l corresponding to the r largest coefficients of the shared atom will dominate the rest. This lets us isolate the support, R , of the shared atom. We still need a mechanism to estimate its non-zero coefficients. This is handled in the following two Lemmas, which shows that the spectrum of a certain (truncated) weighted covariance matrix reveals this information. This step is reminiscent of covariance-thresholding methods for sparse PCA (Johnstone and Lu 2004; Deshpande and Montanari 2014), and distinguishes our approach from that in (Arora et al. 2015).

Lemma 2. *The truncated re-weighting matrix obeys:*

$$\begin{aligned} M_{u,v}^R &\triangleq \mathbb{E}[\langle y, u \rangle \langle y, v \rangle y_R y_R^T] \\ &= \sum_{i \in U \cap V} q_i c_i \beta_i \beta'_i A_{R,i}^* A_{R,i}^{*T} + E', \end{aligned}$$

where E' have spectrum norm at most $O^*(k/m \log n)$ w.h.p.

Lemma 3. *If $U \cap V = \{i\}$, then the r largest entries of e_l are of magnitude at least $\Omega(k/mr)$ and are supported on R . Moreover, the top singular vector of $M_{u,v}^R$ is δ -close to $A_{R,i}^*$ for $\delta = O^*(1/\log n)$.*

Using the same argument for bounding E in Lemma 1, we can see that $M_0 \triangleq q_i c_i \beta_i \beta'_i A_{R,i}^* A_{R,i}^{*T}$ has norm at least $\Omega(k/m)$ when u and v share a unique element i . Therefore, the spectral norm of M_0 dominates those of the perturbation term E' . Thus, given R , we can use the first singular vector of $M_{u,v}^R$ as an estimate of $A_{R,i}^*$.

The question remains when and how we can certify that u and v share a unique single element in the support of their code vectors. Fortunately, this condition can be confirmed by checking the decay of the singular values of the (truncated) covariance matrix. This is quantified as follows.

Lemma 4. *If the top singular value of $M_{u,v}$ is at least $\Omega(k/m)$ and the second largest one is at most $O^*(k/m \log n)$, then u and v share a unique dictionary element with high probability.*

Algorithm 2 Double-Sparse Coding Descent Algorithm

Initialize A^0 is $(\delta, 2)$ -near to A^* . $H = (h_{ij})_{n \times m}$ where $h_{ij} = 1$ if $i \in \text{supp}(A_{\bullet j}^0)$ and 0 otherwise.

Repeat for $s = 0, 1, \dots, T$

Encode: $x^{(i)} = \text{threshold}_{C/2}((A^s)^T y^{(i)})$

Update: $A^{s+1} = \mathcal{P}_H(A^s - \eta \hat{g}^s) = A^s - \eta \mathcal{P}_H(\hat{g}^s)$,

where $\hat{g}^s = \frac{1}{p} \sum_{i=1}^p (A^s x^{(i)} - y^{(i)}) \text{sgn}(x^{(i)})^T$

and $\mathcal{P}_H(G) = H \circ G$

The above discussion assumes infinitely many available samples. However, we can derive analogous finite-sample lemmas which hold w.h.p. via concentration arguments. See the appendix for details. Similar to (Arora et al. 2015), our algorithm requires $\tilde{O}(m)$ iterations to estimate all the atoms, and hence the expected running time is $\tilde{O}(mnp)$.

Lemma 3 indicates that the support, as well as a coarse (δ -close) estimate, of each column of A^* can be estimated using our proposed initialization method. We now show how to refine this estimate using a descent-style method.

Stage 2: Descent

We adapt the neural sparse coding approach of (Arora et al. 2015) to obtain an improved estimate of A^* . As mentioned earlier, at a high level the algorithm is akin to performing approximate gradient descent. The insight is that within a small enough neighborhood (in the sense of δ -closeness) of the true A^* , an estimate of the ground-truth code vectors, X^* , can be constructed using a neurally plausible algorithm. It can be used to construct a noisy approximate gradient \hat{g}^s .

The innovation, in our case, is the double-sparsity model since we know *a priori* that A^* is itself sparse. Under sufficiently many samples, the support of A^* can be deduced from the initialization stage; therefore we perform an extra *projection* step in each iteration of gradient descent. In this sense, our method is non-trivially different from (Arora et al. 2015). The full algorithm is presented as Algorithm 2.

As discussed in the Setup section above, convergence of noisy approximate gradient descent can be achieved as long as \hat{g}^s is correlated-w.h.p. with the true solution. However, an analogous convergence result for projected gradient descent does not exist in the literature. We fill this gap via a careful analysis. Due to the projection, we only require the correlated-w.h.p. property for a *part* of \hat{g}^s with A^* when it is restricted to some support set. The descent property is still achieved via Theorem 3. Due to the various perturbation terms, \hat{g} is only a biased estimate of $\nabla_A \mathcal{L}(A, X)$; therefore, we can only refine the estimate of A^* until the column-wise error is of the order of $O(\sqrt{k/n})$. The performance of Algorithm 2 can be characterized via the following theorem.

Theorem 3. *Suppose that the initial estimate A^0 has the correct column supports and is $(\delta, 2)$ -near to A^* with $\delta = O^*(1/\log n)$. If Algorithm 2 is provided with $p = \tilde{\Omega}(m + \sigma_\epsilon^2 \frac{mnr}{k})$ samples at each step and $\eta = \Theta(m/k)$, then*

$$\mathbb{E}[\|A_{\bullet i}^s - A_{\bullet i}^*\|^2] \leq (1 - \rho)^s \|A_{\bullet i}^0 - A_{\bullet i}^*\|^2 + O(k/n)$$

for some $0 < \rho < 1/2$ and for $s = 1, 2, \dots, T$. Consequently, A^s converges to A^* geometrically until column-wise error is $O(\sqrt{k/n})$.

The formal proof of Theorem 3 is available in our extended version (Nguyen, Wong, and Hegde 2017). Here, we shed some light on the analysis techniques by studying the case of infinite samples. Therefore, the estimate \hat{g}^s can be replaced by its expectation,

$$g^s \triangleq \mathbb{E}[(A^s x - y) \text{sgn}(x)^T].$$

Let us focus on the i^{th} column. Given the knowledge of the support R of $A_{\bullet i}^*$, we only have to restrict our focus to $g_{R,i}^s$. A key component is to establish the $(\alpha, \beta, \gamma_s)$ -correlation of $g_{R,i}^s$ with $A_{R,i}^*$ so as to obtain a descent property, similar to Theorem 3, for infinite number of samples. To this end, we establish the following lemma, using the same strategy as in (Arora et al. 2015).

Lemma 5. *Suppose that the initial estimate A^0 has the correct column supports and is $(\delta, 2)$ -near to A^* with $\delta = O^*(1/\log n)$. The update is of the form $g_{R,i}^s = p_i q_i (\lambda_i^s A_{R,i}^s - A_{R,i}^* + \xi_i^s \pm \zeta)$ where $R = \text{supp}(A_{\bullet i}^*)$ and*

$$\xi_i^s = A_{R,-i}^s \text{diag}(q_{ij}) (A_{\bullet -i}^s)^T A_{\bullet i}^* / q_i$$

and $\lambda_i^s = \langle A_{\bullet i}^s, A_{\bullet i}^* \rangle$. In addition, $\|\xi_i^s\| \leq O(k/n)$ and ζ is negligible.

Intuitively, Lemma 5 suggests that $g_{R,i}^s$ is almost equal to $p_i q_i (A_{R,i}^s - A_{R,i}^*)$ (since $\lambda_i^s \approx 1$), which is a desired direction. Then, we can prove the correlation and descent results accordingly:

Lemma 6. *If A^s is $(2\delta, 2)$ -near to A^* with $\delta = O^*(1/\log n)$ and $R = \text{supp}(A_{\bullet i}^*)$, then $2g_{R,i}^s$ is $(\alpha, 1/2\alpha, \epsilon^2/\alpha)$ -correlated with $A_{R,i}^*$ by*

$$\langle 2g_{R,i}^s, A_{R,i}^s - A_{R,i}^* \rangle \geq \alpha \|A_{R,i}^s - A_{R,i}^*\|^2 + 1/(2\alpha) \|g_{R,i}^s\|^2 - \epsilon^2/\alpha$$

where $\epsilon = O(k^2/mn)$. In particular, $g_{R,i}^s$ is (α, β, γ) -correlated with $A_{R,i}^*$ for $\alpha = \Omega(k/m)$, $\beta = \Omega(m/k)$ and $\gamma = O(k^3/mn^2)$.

After the results under infinite samples are achieved, we study the concentration of the empirical average \hat{g}^s to its mean. Again, due to the knowledge of the column supports, for each column of \hat{g}^s , we only have to establish such concentration over a r -dimensional sub-vector. This helps to achieve a better sample complexity especially when r is small. To sum up, the respective sample complexities for the descent and the initialization stage are $\tilde{O}(m + \sigma_\epsilon^2 \frac{mnr}{k})$ and $\tilde{O}(mr)$. Overall, the sample complexity $\tilde{O}(mr + \sigma_\epsilon^2 \frac{mnr}{k})$ sufficiently guarantees the success of our approach.

Regarding the running time, the running time per iteration of Algorithm 2 is $O(m \max(k, r)p)$ due to the sparsity of both A and x . The main bottleneck is at the initialization stage with the expected running time is $\tilde{O}(mnp)$. Consequently, the total computational complexity of our approach is $\tilde{O}(mnp)$.

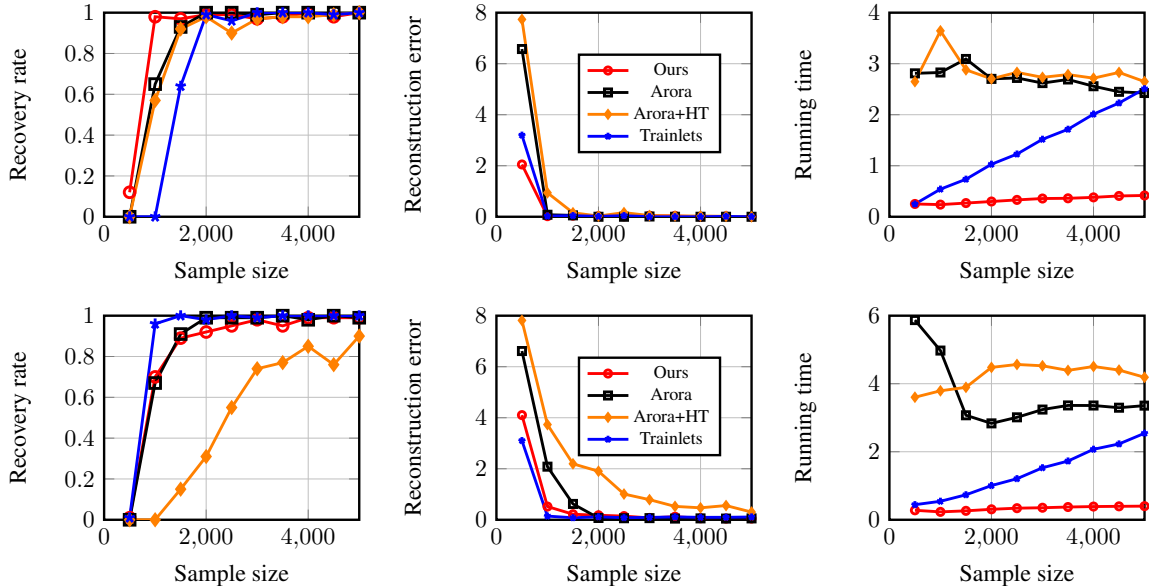


Figure 1: (top) The performance of four methods on three metrics (recovery rate, reconstruction error and running time) in sample size in the noiseless case. (bottom) The same metrics are measured for the noisy case.

Empirical Study

We compare our method with three different methods for both standard sparse and double-sparse coding. For the standard approach, we implement the algorithm proposed in (Arora et al. 2015), which currently is the best theoretically sound method for provable sparse coding. However, since their method does not explicitly leverage the double-sparsity model, we also implement a heuristic modification that performs a hard thresholding (HT)-based post-processing step in the initialization and learning procedures (which we dub *Arora + HT*). The final comparison is the *Trainlets* approach of (Sulam et al. 2016).

We generate a synthetic training dataset according to the model described in the Setup. The base dictionary Φ is the identity matrix of size $n = 64$ and the square synthesis matrix A^* is a block diagonal matrix with 32 blocks. Each 2×2 block is of form $\begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix}$ (i.e., the column sparsity $r = 2$). The support of x^* is drawn uniformly over all 6-dimensional subsets of $[m]$, and the nonzero coefficients are randomly set to ± 1 with equal probability. In our simulations with noise, we add Gaussian noise ε with entrywise variance $\sigma_\varepsilon^2 = 0.01$ to each of those above samples. For all the approaches except *Trainlets*, we use $T = 2000$ iterations for the initialization procedure, and set the number of steps in the descent stage to 25. Since *Trainlets* does not have a specified initialization procedure, we initialize it with a random Gaussian matrix upon which column-wise sparse thresholding is then performed. The learning step of *Trainlets*³ is executed for 50 iterations, which tolerates its initialization deficiency. For each Monte Carlo trial, we uniformly draw p samples, feed these samples to the four different algorithms, and observe

their ability to reconstruct A^* .

We evaluate these approaches on three metrics as a function of the number of available samples: (i) fraction of trials in which each algorithm successfully recovers the ground truth A^* ; (ii) reconstruction error; and (iii) running time. The synthesis matrix is said to be “successfully recovered” if the Frobenius norm of the difference between the estimate \hat{A} and the ground truth A^* is smaller than a threshold which is set to 10^{-4} in the noiseless case, and to 0.5 in the other. All three metrics are averaged over 100 Monte Carlo simulations. As discussed above, the Frobenius norm is only meaningful under a suitable permutation and sign flip transformation linking \hat{A} and A^* . We estimate this transformation using a simple maximum weight matching algorithm. Specifically, we construct a weighted bipartite graph with nodes representing columns of A^* and \hat{A} and adjacency matrix defined as $G = |A^{*T} \hat{A}|$, where $|\cdot|$ is taken element-wise. We compute the optimal matching using the Hungarian algorithm, and then estimate the sign flips by looking at the sign of the inner products between the matched columns.

The results of our experiments are shown in Figure 1 with the top and bottom rows respectively for the noiseless and noisy cases. The two leftmost figures suggest that all algorithms exhibit a “phase transitions” in sample complexity that occurs in the range of 500-2000 samples. In the noiseless case, our method achieves the phase transition with the fewest number of samples. In the noisy case, our method nearly matches the best sample complexity performance (next to *Trainlets*, which is a heuristic and computationally expensive). Our method achieves the best performance in terms of (wall-clock) running time in all cases.

³We utilize *Trainlets*’s implementation provided at <http://jsulam.cswp.cs.technion.ac.il/home/software/>.

References

- Agarwal, A.; Anandkumar, A.; Jain, P.; Netrapalli, P.; and Tandon, R. 2014. Learning sparsely used overcomplete dictionaries. In *Conference on Learning Theory*, 123–137.
- Aharon, M.; Elad, M.; and Bruckstein, A. 2006. k -svd: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on signal processing* 54(11):4311–4322.
- Arora, S.; Ge, R.; Ma, T.; and Moitra, A. 2015. Simple, efficient, and neural algorithms for sparse coding. In *Conference on Learning Theory*, 113–149.
- Arora, S.; Ge, R.; and Moitra, A. 2014. New algorithms for learning incoherent and overcomplete dictionaries. In *Conference on Learning Theory*, 779–806.
- Boureau, Y.-L.; Bach, F.; LeCun, Y.; and Ponce, J. 2010. Learning mid-level features for recognition. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, 2559–2566. IEEE.
- Deshpande, Y., and Montanari, A. 2014. Sparse pca via covariance thresholding. In *Advances in Neural Information Processing Systems*, 334–342.
- Elad, M., and Aharon, M. 2006. Image denoising via sparse and redundant representations over learned dictionaries. *IEEE Transactions on Image processing* 15(12):3736–3745.
- Engan, K.; Aase, S. O.; and Husoy, J. H. 1999. Method of optimal directions for frame design. In *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, volume 5, 2443–2446. IEEE.
- Gregor, K., and LeCun, Y. 2010. Learning fast approximations of sparse coding. In *Proceedings of the 27th International Conference on Machine Learning (ICML)*, 399–406.
- Gribonval, R.; Jenatton, R.; Bach, F.; Kleinstueber, M.; and Seibert, M. 2015. Sample complexity of dictionary learning and other matrix factorizations. *IEEE Transactions on Information Theory* 61(6):3469–3486.
- Gribonval, R.; Jenatton, R.; and Bach, F. 2015. Sparse and spurious: dictionary learning with noise and outliers. *IEEE Transactions on Information Theory* 61(11):6298–6319.
- Jenatton, R.; Gribonval, R.; and Bach, F. 2012. Local stability and robustness of sparse dictionary learning in the presence of noise. *arXiv preprint arXiv:1210.0685*.
- Johnstone, I. M., and Lu, A. Y. 2004. Sparse principal components analysis. *Unpublished manuscript* 7.
- Krim, H.; Tucker, D.; Mallat, S.; and Donoho, D. 1999. On denoising and best signal representation. *IEEE Transactions on Information Theory* 45(7):2225–2238.
- Mairal, J.; Bach, F.; Ponce, J.; and Sapiro, G. 2009. On-line dictionary learning for sparse coding. In *Proceedings of the 26th International Conference on Machine Learning (ICML)*, 689–696.
- Mazumdar, A., and Rawat, A. S. 2017. Associative memory using dictionary learning and expander decoding. In *Proc. Conf. American Assoc. Artificial Intelligence (AAAI)*, 267–273.
- Nguyen, T.; Wong, R. K. W.; and Hegde, C. 2017. A provable approach for double-sparse coding. *arXiv preprint arXiv:1711.03638*.
- Oja, E. 1992. Principal components, minor components, and linear neural networks. *Neural networks* 5(6):927–935.
- Olshausen, B. A., and Field, D. J. 1997. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research* 37(23):3311–3325.
- Rubinstein, R.; Bruckstein, A. M.; and Elad, M. 2010. Dictionaries for sparse representation modeling. *Proceedings of the IEEE* 98(6):1045–1057.
- Rubinstein, R.; Zibulevsky, M.; and Elad, M. 2010. Double sparsity: Learning sparse dictionaries for sparse signal approximation. *IEEE Transactions on Signal Processing* 58(3):1553–1564.
- Spielman, D. A.; Wang, H.; and Wright, J. 2012. Exact recovery of sparsely-used dictionaries. In *Conference on Learning Theory*, 37–1.
- Sulam, J.; Ophir, B.; Zibulevsky, M.; and Elad, M. 2016. Trainlets: Dictionary learning in high dimensions. *IEEE Transactions on Signal Processing* 64(12):3180–3193.
- Sun, J.; Qu, Q.; and Wright, J. 2015. Complete dictionary recovery using nonconvex optimization. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2351–2360.
- Wang, L.; Zhang, X.; and Gu, Q. 2016. A unified computational and statistical framework for nonconvex low-rank matrix estimation. *arXiv preprint arXiv:1610.05275*.
- Zhang, Y.; Chen, X.; Zhou, D.; and Jordan, M. I. 2016. Spectral methods meet em: A provably optimal algorithm for crowdsourcing. *Journal of Machine Learning Research* 17(1):3537–3580.

Auxiliary Lemma

Claim 1 (Maximal row ℓ_1 -norm). *Given that $\|A^*\|_F^2 = m$ and $\|A^*\| = O(\sqrt{m/n})$, then $\|A^{*T}\|_{1,2} = \Theta(\sqrt{m/n})$.*

Proof. Recall the definition of the operator norm:

$$\|A^{*T}\|_{1,2} = \sup_{x \neq 0} \frac{\|A^T x\|}{\|x\|_1} \leq \sup_{x \neq 0} \frac{\|A^T x\|}{\|x\|} = \|A^{*T}\| = O(\sqrt{m/n}).$$

Since $\|A^*\|_F^2 = m$, $\|A^{*T}\|_{1,2} \geq \|A^*\|_F / \sqrt{n} = \sqrt{m/n}$. Combining with the above, we have $\|A^{*T}\|_{1,2} = \Theta(\sqrt{m/n})$. \square

Along with Assumptions **A1** and **A3**, the above claim implies the number of nonzero entries in each row is $O(r)$. This Claim is an important ingredient in our analysis of our initialization algorithm shown in Section ‘‘Stage 1: Initialization’’.

Analysis of Initialization Algorithm

Proof of Lemma 1

The proof of Lemma 1 can be divided into three steps: 1) we first establish useful properties of β with respect to α ; 2) we then explicitly derive e_l in terms of the generative model parameters and β ; and 3) we finally bound the error terms in E based on the first result and appropriate assumptions.

Claim 2. *In the generative model, $\|x^*\| \leq \tilde{O}(\sqrt{k})$ and $\|\varepsilon\| \leq \tilde{O}(\sigma_\varepsilon \sqrt{n})$ with high probability.*

Proof. The claim directly follows from the fact that x^* is a k -sparse random vector whose nonzero entries are independent sub-Gaussian with variance 1. Meanwhile, ε has n independent Gaussian entries of variance σ_ε^2 . \square

Despite its simplicity, this claim will be used in many proofs throughout the paper. Note also that in this section we will calculate the expectation over y and often refer probabilistic bounds (w.h.p.) under the randomness of u and v .

Claim 3. *Suppose that $u = A^* \alpha + \varepsilon_u$ is a random sample and $U = \text{supp}(\alpha)$. Let $\beta = A^{*T} u$, then, w.h.p., we have (a) $|\beta_i - \alpha_i| \leq \frac{\mu k \log n}{\sqrt{n}} + \sigma_\varepsilon \log n$ for each i and (b) $\|\beta\| \leq \tilde{O}(\sqrt{k} + \sigma_\varepsilon \sqrt{n})$.*

Proof. The proof mostly follows from Claim 36 of (Arora et al. 2015), with an additional consideration of the error ε_u . Write $W = U \setminus \{i\}$ and observe that

$$|\beta_i - \alpha_i| = |A_{\bullet i}^{*T} A_{\bullet W}^* \alpha_W + A_{\bullet i}^{*T} \varepsilon_u| \leq |\langle A_{\bullet W}^{*T} A_{\bullet i}^*, \alpha_W \rangle| + |\langle A_{\bullet i}^*, \varepsilon_u \rangle|.$$

Since A^* is μ -incoherence, then $\|A_{\bullet i}^{*T} A_{\bullet W}^*\| \leq \mu \sqrt{k/n}$. Moreover, α_W has $k-1$ independent sub-Gaussian entries of variance 1, therefore $|\langle A_{\bullet W}^{*T} A_{\bullet i}^*, \alpha_W \rangle| \leq \frac{\mu k \log n}{\sqrt{n}}$ with high probability. Also recall that ε_u has independent Gaussian entries of variance σ_ε^2 , then $A_{\bullet i}^{*T} \varepsilon_u$ is Gaussian with the same variance ($\|A_{\bullet i}^*\| = 1$). Hence $|A_{\bullet i}^{*T} \varepsilon| \leq \sigma_\varepsilon \log n$ with high probability. Consequently, $|\beta_i - \alpha_i| \leq \frac{\mu k \log n}{\sqrt{n}} + \sigma_\varepsilon \log n$, which is the first part of the claim.

Next, in order to bound $\|\beta\|$, we express β as

$$\|\beta\| = \|A^{*T} A_{\bullet U}^* \alpha_U + A^{*T} \varepsilon_u\| \leq \|A^*\| \|A_{\bullet U}^*\| \|\alpha_U\| + \|A^*\| \|\varepsilon_u\|.$$

Using Claim 2 to get $\|\alpha_U\| \leq \tilde{O}(\sqrt{k})$ and $\|\varepsilon_u\| \leq \tilde{O}(\sigma_\varepsilon \sqrt{n})$ w.h.p., and further noticing that $\|A_{\bullet U}^*\| \leq \|A^*\| \leq O(1)$, we complete the proof for the second part. \square

Claim 3 suggests that the difference between β_i and α_i is bounded above by $O^*(1/\log^2 n)$ w.h.p. if $\mu = O^*(\frac{\sqrt{n}}{k \log^3 n})$. Therefore, w.h.p., $C - o(1) \leq |\beta_i| \leq |\alpha_i| + o(1) \leq O(\log m)$ for $i \in U$ and $|\beta_i| \leq O^*(1/\log^2 n)$ otherwise. On the other hand, under Assumption **B4**, $\|\beta\| \leq \tilde{O}(\sqrt{k})$ w.h.p. We will use these results multiple times in the next few proofs.

Proof of Lemma 1. We decompose e_l into small parts so that the stochastic model \mathcal{D} is made use.

$$\begin{aligned} e_l &= \mathbb{E}[\langle y, u \rangle \langle y, v \rangle y_l^2] = \mathbb{E}[\langle A^* x^* + \varepsilon, u \rangle \langle A^* x^* + \varepsilon, v \rangle (\langle A_{l \bullet}^*, x^* \rangle + \varepsilon)^2] \\ &= \mathbb{E}[\{ \langle x^*, \beta \rangle \langle x^*, \beta' \rangle + x^{*T} (\beta v^T + \beta' u^T) \varepsilon + u^T \varepsilon \varepsilon^T v \} \{ \langle A_{l \bullet}^*, x^* \rangle^2 + 2 \langle A_{l \bullet}^*, x^* \rangle \varepsilon_l + \varepsilon_l \}] \\ &= E_1 + E_2 + \dots + E_9. \end{aligned}$$

where the terms are

$$\begin{aligned}
E_1 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \langle A_{l\bullet}^*, x^* \rangle^2] \\
E_2 &= 2\mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \langle A_{l\bullet}^*, x^* \rangle \varepsilon_l] \\
E_3 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \varepsilon_l^2] \\
E_4 &= \mathbb{E}[\langle A_{l\bullet}^*, x^* \rangle^2 x^{*T} (\beta v^T + \beta' u^T) \varepsilon] \\
E_5 &= \mathbb{E}[\langle A_{l\bullet}^*, x^* \rangle x^{*T} (\beta v^T + \beta' u^T) \varepsilon \varepsilon_l] \\
E_6 &= \mathbb{E}[(\beta v^T + \beta' u^T) \varepsilon \varepsilon_l^2] \\
E_7 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v \langle A_{l\bullet}^*, x^* \rangle^2] \\
E_8 &= 2\mathbb{E}[u^T \varepsilon \varepsilon^T v \langle A_{l\bullet}^*, x^* \rangle \varepsilon_l] \\
E_9 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_l^2].
\end{aligned} \tag{2}$$

Because x^* and ε are independent and have zero mean, E_2 and E_4 are clearly zero. $E_6 = (\beta v^T + \beta' u^T) \mathbb{E}[\varepsilon \varepsilon_l^2] = 0$ due to the fact that $\mathbb{E}[\varepsilon_j \varepsilon_l^2] = 0$, for $j \neq l$, and $\mathbb{E}[\varepsilon_l^3] = 0$; and $E_8 = A_{l\bullet}^{*T} \mathbb{E}[x^*] \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_l] = 0$. We bound the remaining terms separately in the following claims.

Claim 4. *In the decomposition (2), E_1 is of the form*

$$E_1 = \sum_{i \in U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} + \sum_{i \notin U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} + \sum_{j \neq i} q_{ij} (\beta_i \beta'_i A_{lj}^{*2} + 2\beta_i \beta'_j A_{li}^* A_{lj}^*),$$

where all those terms except $\sum_{i \in U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2}$ have magnitude at most $O^*(k/m \log^2 n)$ w.h.p.

Proof. Using the generative model in Assumptions **B1-B4**, we have

$$\begin{aligned}
E_1 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \langle A_{l\bullet}^*, x^* \rangle^2] \\
&= \mathbb{E}_S[\mathbb{E}_{x^*|S}[\sum_{i \in S} \beta_i x_i^* \sum_{i \in S} \beta'_i x_i^* (\sum_{i \in S} A_{li}^* x_i^*)^2]] \\
&= \sum_{i \in [m]} q_i c_i \beta_i \beta'_i A_{li}^{*2} + \sum_{i, j \in [m], j \neq i} q_{ij} (\beta_i \beta'_i A_{lj}^{*2} + 2\beta_i \beta'_j A_{li}^* A_{lj}^*) \\
&= \sum_{i \in U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} + \sum_{i \notin U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} + \sum_{j \neq i} q_{ij} (\beta_i \beta'_i A_{lj}^{*2} + 2\beta_i \beta'_j A_{li}^* A_{lj}^*),
\end{aligned}$$

where we have used $q_i = \mathbb{P}[i \in S]$, $q_{ij} = \mathbb{P}[i, j \in S]$ and $c_i = \mathbb{E}[x_i^4 | i \in S]$ and Assumptions **B1-B4**. We now prove that the last three terms are upper bounded by $O^*(k/m \log^2 n)$. The key observation is that all these terms typically involve a quadratic form of the l -th row $A_{l\bullet}^{*T}$ whose norm is bounded by $O(1)$ (by Claim 1). Moreover, $|\beta_i \beta'_i|$ is relatively small for $i \notin U \cap V$ while $q_{ij} = O(k^2/m^2)$. For the second term, we apply Claim 3 for $i \in [m] \setminus (U \cap V)$ to get $|\beta_i \beta'_i| \leq O^*(\frac{1}{\log^4 n})$ w.h.p. for $\mu = O^*(\frac{\sqrt{n}}{k \log^3 n})$ and use the bound $q_i c_i = \Theta(k/m)$, then, w.h.p.,

$$\left| \sum_{i \notin U \cap V} q_i c_i \beta_i \beta'_i A_{li}^{*2} \right| \leq \max_i |q_i c_i \beta_i \beta'_i| \sum_{i \notin U \cap V} A_{li}^{*2} \leq \max_i |q_i c_i \beta_i \beta'_i| \|A^*\|_{1,2}^2 \leq O^*(k/m \log^4 n).$$

For the third term, we make use of the bounds on $\|\beta\|$ and $\|\beta'\|$ from the previous claim where $\|\beta\| \|\beta'\| \leq \tilde{O}(k)$ w.h.p., and on $q_{ij} = \Theta(k^2/m^2)$. More precisely, w.h.p.,

$$\begin{aligned}
\left| \sum_{j \neq i} q_{ij} \beta_i \beta'_j A_{lj}^{*2} \right| &= \left| \sum_i \beta_i \beta'_i \sum_{j \neq i} q_{ij} A_{lj}^{*2} \right| \leq \sum_i |\beta_i \beta'_i| \left(\sum_{j \neq i} q_{ij} A_{lj}^{*2} \right) \\
&\leq (\max_{i \neq j} q_{ij}) \sum_i |\beta_i \beta'_i| \left(\sum_j A_{lj}^{*2} \right) \leq (\max_{i \neq j} q_{ij}) \|\beta\| \|\beta'\| \|A^*\|_{1,2}^2 \leq \tilde{O}(k^3/m^2),
\end{aligned}$$

where the second last inequality follows from the Cauchy-Schwarz inequality. For the last term, we write it in a matrix form as $\sum_{j \neq i} q_{ij} \beta_i \beta'_j A_{li}^* A_{lj}^* = A_{l\bullet}^{*T} Q_\beta A_{l\bullet}^*$ where $(Q_\beta)_{ij} = q_{ij} \beta_i \beta'_j$ for $i \neq j$ and $(Q_\beta)_{ij} = 0$ for $i = j$. Then

$$|A_{l\bullet}^{*T} Q_\beta A_{l\bullet}^*| \leq \|Q_\beta\| \|A_{l\bullet}^*\|^2 \leq \|Q_\beta\|_F \|A^*\|_{1,2}^2,$$

where $\|Q_\beta\|_F^2 = \sum_{i \neq j} q_{ij}^2 \beta_i^2 (\beta'_j)^2 \leq (\max_{i \neq j} q_{ij}^2) \sum_i \beta_i^2 \sum_j (\beta'_j)^2 \leq (\max_{i \neq j} q_{ij}^2) \|\beta\|^2 \|\beta'\|^2$. Finally,

$$\left| \sum_{j \neq i} q_{ij} \beta_i \beta'_j A_{li}^* A_{lj}^* \right| \leq (\max_{i \neq j} q_{ij}) \|\beta\| \|\beta'\| \|A^*\|_{1,2} \leq \tilde{O}(k^3/m^2).$$

Under Assumption $k = O^*(\frac{\sqrt{n}}{\log n})$ and hence $\tilde{O}(k^3/m^2) \leq O^*(k/m \log^2 n)$, we have the same bound $O^*(k/m \log^2 n)$ for those last two terms w.h.p. Therefore, we have completed the proof of the claim. \square

Claim 5. In the decomposition (2), $|E_3|$, $|E_5|$, $|E_7|$ and $|E_9|$ are at most $O^*(k/m \log^2 n)$.

Proof. Recall that $\mathbb{E}[x_i^2 | S] = 1$ and $q_i = \mathbb{P}[i \in S] = \Theta(k/m)$ for $S = \text{supp}(x^*)$, then

$$\begin{aligned} E_3 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \varepsilon_l^2] = \sigma_\varepsilon^2 \mathbb{E}_S [\mathbb{E}_{x^* | S} [\sum_{i,j \in S} \beta_i \beta'_j x_i^* x_j^*]] \\ &= \sigma_\varepsilon^2 \mathbb{E}_S [\sum_{i \in S} \beta_i \beta'_i] = \sum_i \sigma_\varepsilon^2 q_i \beta_i \beta'_i. \end{aligned}$$

Write $Q = \text{diag}(q_1, q_2, \dots, q_m)$. It is verified that $|E_3| = |\sigma_\varepsilon^2 \langle Q\beta, \beta' \rangle| \leq \sigma_\varepsilon^2 \|Q\| \|\beta\| \|\beta'\| \leq \tilde{O}(\sigma_\varepsilon^2 k^2/m) = \tilde{O}(k^3/mn)$ where we have used $\|\beta\| \leq \tilde{O}(\sqrt{k})$ and $\sigma_\varepsilon = O(1/\sqrt{n})$. For convenience, we handle the seventh term before E_5 :

$$E_7 = \mathbb{E}[u^T \varepsilon \varepsilon^T v \langle A_{l\bullet}^*, x^* \rangle^2] = \mathbb{E}[\langle A_{l\bullet}^*, x^* \rangle^2] u^T \mathbb{E}[\varepsilon \varepsilon^T] v = \sum_i \sigma_\varepsilon^2 \langle u, v \rangle q_i A_{li}^2 = \sigma_\varepsilon^2 \langle u, v \rangle A_{l\bullet}^T Q A_{l\bullet}.$$

To bound this term, we use Claim 9 in Appendix Section ‘‘Sample Complexity’’ below to get $\|u\| = \|A^* \alpha + \varepsilon_u\| \leq \tilde{O}(k)$ w.h.p. and $\langle u, v \rangle \leq \|u\| \|v\| \leq \tilde{O}(\sqrt{k})$ w.h.p. Consequently, $|E_7| \leq \sigma_\varepsilon^2 \|Q\| \|A_{l\bullet}\|^2 |\langle u, v \rangle| \leq \tilde{O}(k^3/mn)$ because $\|A_{l\bullet}\|^2 \leq O(m/n)$. Now, we work on the fifth term E_5 as follows:

$$\begin{aligned} E_5 &= \mathbb{E}[\langle A_{l\bullet}^*, x^* \rangle x^{*T} (\beta v^T + \beta' u^T) \varepsilon \varepsilon_l] \\ &= A_{l\bullet}^{*T} \mathbb{E}[x^* x^{*T}] (\beta v^T + \beta' u^T) \mathbb{E}[\varepsilon \varepsilon_l] \\ &= \sigma_\varepsilon^2 A_{l\bullet}^{*T} Q (v_l \beta + u_l \beta'), \end{aligned}$$

and $|E_5| \leq \sigma_\varepsilon^2 \|A_{l\bullet}^{*T}\| \|Q\| \|v_l \beta + u_l \beta'\| \leq \sigma_\varepsilon^2 \|A_{l\bullet}^{*T}\| \|Q\| \|v_l \beta + u_l \beta'\|$. To show that E_5 is bounded by $\tilde{O}(k^3/mn)$, it suffices to show that $\|v_l \beta + u_l \beta'\| \leq 2\|u\| \|\beta\| \leq \tilde{O}(k)$ w.h.p. using the result $\|u\| \leq \tilde{O}(k)$, which follows from Claim 9. The last term

$$E_9 = \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_l^2] = u^T \mathbb{E}[\varepsilon \varepsilon^T \varepsilon_l^2] v = 9\sigma_\varepsilon^4 \langle u, v \rangle,$$

due to the independent entries of ε and $\mathbb{E}[\varepsilon_l^4] = 9\sigma_\varepsilon^4$. Therefore, $|E_9| \leq 9\sigma_\varepsilon^4 \|u\| \|v\| \leq O(k^3/n^2)$. Since $m = O(n)$ and $k \leq O^*(\frac{\sqrt{n}}{\log n})$, we have the same bound $\tilde{O}(k/m \log^2 n)$ for all $|E_3|$, $|E_5|$, $|E_7|$ and $|E_9|$, and conclude the proof of the claim. \square

Combining the bounds from Claims 4 and 5 for every single term in (2), we go to finish the proof for Lemma 1. \square

Proof of Lemma 2

We prove this lemma by using the same strategy used to prove Lemma 1.

$$\begin{aligned} M_{u,v} &\triangleq \mathbb{E}[\langle y, u \rangle \langle y, v \rangle y_R y_R^T] = \mathbb{E}[\langle A^* x^* + \varepsilon, u \rangle \langle A^* x^* + \varepsilon, v \rangle (A_{R\bullet}^* x^* + \varepsilon_R) (A_{R\bullet}^* x^* + \varepsilon_R)^T] \\ &= \mathbb{E}[\{\langle x^*, \beta \rangle \langle x^*, \beta' \rangle + x^{*T} (\beta v^T + \beta' u^T) \varepsilon + u^T \varepsilon \varepsilon^T v\} \{A_{R\bullet}^* x^* x^{*T} A_{R\bullet}^{*T} + A_{R\bullet}^* x^* \varepsilon_R^T + \varepsilon_R x^{*T} A_{R\bullet}^{*T} + \varepsilon_R \varepsilon_R^T\}] \\ &= M_1 + \dots + M_8, \end{aligned}$$

in which only nontrivial terms are kept in place, including

$$\begin{aligned} M_1 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle A_{R\bullet}^* x^* x^{*T} A_{R\bullet}^{*T}] \\ M_2 &= \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \varepsilon_R \varepsilon_R^T] \\ M_3 &= \mathbb{E}[x^{*T} (\beta v^T + \beta' u^T) \varepsilon A_{R\bullet}^* x^* \varepsilon_R^T] \\ M_4 &= \mathbb{E}[x^{*T} (\beta v^T + \beta' u^T) \varepsilon \varepsilon_R x^{*T} A_{R\bullet}^{*T}] \\ M_5 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v A_{R\bullet}^* x^* x^{*T} A_{R\bullet}^{*T}] \\ M_6 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v A_{R\bullet}^* x^* \varepsilon_R^T] \\ M_7 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_R^T x^{*T} A_{R\bullet}^{*T}] \\ M_8 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_R \varepsilon_R^T]. \end{aligned} \tag{3}$$

By swapping inner product terms and taking advantage of the independence, we can show that $M_6 = \mathbb{E}[A_{R_\bullet}^* x^* u^T \varepsilon \varepsilon^T v \varepsilon_R^T] = 0$ and $M_7 = \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_R^T x^* A_{R_\bullet}^{*T}] = 0$. The remaining are bounded in the next claims.

Claim 6. In the decomposition (3),

$$M_1 = \sum_{i \in U \cup V} q_i c_i \beta_i \beta'_i A_{R_\bullet, i}^* A_{R_\bullet, i}^{*T} + E'_1 + E'_2 + E'_3,$$

where $E'_1 = \sum_{i \notin U \cup V} q_i c_i \beta_i \beta'_i A_{R_\bullet, i}^* A_{R_\bullet, i}^{*T}$, $E'_2 = \sum_{i \neq j} q_{ij} \beta_i \beta'_i A_{R_\bullet, j}^* A_{R_\bullet, j}^{*T}$ and $E'_3 = \sum_{i \neq j} q_{ij} (\beta_i A_{R_\bullet, i}^* \beta'_j A_{R_\bullet, j}^{*T} + \beta'_i A_{R_\bullet, i}^* \beta_j A_{R_\bullet, j}^{*T})$ have norms bounded by $O^*(k/m \log n)$ w.h.p.

Proof. The derivation of M_1 is similar to the way E_1 is derived in the proof of Lemma 1, we state and use the expression of M_1 without proof. To prove the claim, we bound all the terms with respect to the spectral norm of $A_{R_\bullet}^*$ and use the spectral norm bound (Assumption **A4**) to obtain the exact upper bounds.

For the first term E'_1 , rewrite $E'_1 = A_{R_\bullet, S}^* D_1 A_{R_\bullet, S}^{*T}$ where $S = [m] \setminus (U \cup V)$ and D_1 is a diagonal matrix whose entries are $q_i c_i \beta_i \beta'_i$. Clearly, $\|D_1\| \leq \max_{i \in S} |q_i c_i \beta_i \beta'_i| \leq O^*(k/m \log^2 n)$ as shown in Claim 4, then

$$\|E'_1\| \leq \max_{i \in S} |q_i c_i \beta_i \beta'_i| \|A_{R_\bullet, S}^*\|^2 \leq \max_{i \in S} |q_i c_i \beta_i \beta'_i| \|A_{R_\bullet}^*\|^2,$$

where $\|A_{R_\bullet, S}^*\| \leq \|A_{R_\bullet}^*\|$. The second term E'_2 is a sum of positive semidefinite matrices, then

$$E'_2 = \sum_{i \neq j} q_{ij} \beta_i \beta'_i A_{R_\bullet, j}^* A_{R_\bullet, j}^{*T} \preceq \max_{i \neq j} q_{ij} \left(\sum_i \beta_i \beta'_i \right) \left(\sum_j A_{R_\bullet, j}^* A_{R_\bullet, j}^{*T} \right) \preceq (\max_{i \neq j} q_{ij}) \|\beta\| \|\beta'\| A_{R_\bullet}^* A_{R_\bullet}^{*T},$$

which implies that $\|E'_2\| \leq (\max_{i \neq j} q_{ij}) \|\beta\| \|\beta'\| \|A_{R_\bullet}^*\|^2$. Observe that E'_3 has the same form as the last term in Claim 4. Effectively, $E'_3 = A_{R_\bullet}^{*T} Q_\beta A_{R_\bullet}^*$, then

$$\|E'_3\| \leq \|Q_\beta\| \|A_{R_\bullet}^*\|^2 \leq (\max_{i \neq j} q_{ij}) \|\beta\| \|\beta'\| \|A_{R_\bullet}^*\|^2.$$

By Claim 3, we have $\|\beta\|$ and $\|\beta'\|$ are bounded by $O(\sqrt{k \log n})$ w.h.p. In addition, $\|A_{R_\bullet}^*\| \leq \|A^*\| \leq O(1)$ and note that $k \leq O^*(\sqrt{n}/\log n)$, then we complete the proof for Lemma 6. \square

Claim 7. In the decomposition (3), M_2, M_3, M_4, M_5 and M_8 have norms bounded by $O^*(k/m \log n)$ w.h.p.

Proof. Recall the definition of Q in Claim 5 and use the fact that $\mathbb{E}[x^* x^{*T}] = Q$, we can get $M_2 = \mathbb{E}[\langle x^*, \beta \rangle \langle x^*, \beta' \rangle \varepsilon_R \varepsilon_R^T] = \sum_i \sigma_\varepsilon^2 q_i \beta_i \beta'_i I_r$. Then, $\|M_2\| \leq \sigma_\varepsilon^2 \max_i q_i \|\beta\| \|\beta'\| \leq \tilde{O}(k^3/mn)$. The next three terms all involve $A_{R_\bullet}^*$ whose norm is bounded, so we work on them at the same time.

$$\begin{aligned} M_3 &= \mathbb{E}[x^{*T} (\beta v^T + \beta' u^T) \varepsilon A_{R_\bullet}^* x^* \varepsilon_R^T] = \mathbb{E}[A_{R_\bullet}^* x^* x^{*T} (\beta v^T + \beta' u^T) \varepsilon \varepsilon_R^T] \\ &= A_{R_\bullet}^* \mathbb{E}[x^* x^{*T}] (\beta v^T + \beta' u^T) \mathbb{E}[\varepsilon \varepsilon_R^T] \\ &= A_{R_\bullet}^* Q (\beta v^T + \beta' u^T) \mathbb{E}[\varepsilon \varepsilon_R^T], \end{aligned}$$

and

$$\begin{aligned} M_4 &= \mathbb{E}[x^{*T} (\beta v^T + \beta' u^T) \varepsilon \varepsilon_R x^{*T} A_{R_\bullet}^{*T}] = \mathbb{E}[\varepsilon_R \varepsilon^T (v \beta^T + u \beta'^T) x^* x^{*T} A_{R_\bullet}^{*T}] \\ &= \mathbb{E}[\varepsilon_R \varepsilon^T] (v \beta^T + u \beta'^T) \mathbb{E}[x^* x^{*T}] A_{R_\bullet}^{*T} \\ &= \mathbb{E}[\varepsilon_R \varepsilon^T] (v \beta^T + u \beta'^T) Q A_{R_\bullet}^{*T}, \end{aligned}$$

and the fifth term $M_5 = \mathbb{E}[u^T \varepsilon \varepsilon^T v A_{R_\bullet}^* x^* x^{*T} A_{R_\bullet}^{*T}] = \sigma_\varepsilon^2 u^T v A_{R_\bullet}^* \mathbb{E}[x^* x^{*T}] A_{R_\bullet}^{*T} = \sigma_\varepsilon^2 u^T v A_{R_\bullet}^* Q A_{R_\bullet}^{*T}$. We already have $\|\mathbb{E}[\varepsilon \varepsilon^T]\| = \sigma_\varepsilon^2$, $\|Q\| \leq O(k/m)$ and $|u^T v| \leq \tilde{O}(k)$ (Claim 9). The remaining work is to bound $\|\beta v^T + \beta' u^T\|$, from which the bound of $v \beta^T + u \beta'^T$ directly follows. We have $\|\beta v^T\| = \|A^* u v^T\| \leq \|A^*\| \|u\| \|v\| \leq \tilde{O}(k)$, hence all three terms M_3, M_4 and M_5 are bounded in norm by $\tilde{O}(\sigma_\varepsilon^2 k^2/m) \leq \tilde{O}(k^3/mn)$ w.h.p.

The remaining term is

$$\begin{aligned} M_8 &= \mathbb{E}[u^T \varepsilon \varepsilon^T v \varepsilon_R \varepsilon_R^T] = \mathbb{E}[(\sum_{i,j} u_i v_j \varepsilon_i \varepsilon_j) \varepsilon_R \varepsilon_R^T] \\ &= \mathbb{E}[(\sum_{i \in R} u_i v_i \varepsilon_i^2 \varepsilon_R \varepsilon_R^T)] + \mathbb{E}[(\sum_{i \neq j} u_i v_j \varepsilon_i \varepsilon_j) \varepsilon_R \varepsilon_R^T] \\ &= \sigma_\varepsilon^4 u_R v_R^T, \end{aligned}$$

where $u_R = A_{R_\bullet}^* \alpha + (\varepsilon_u)_R$ and $v_R = A_{R_\bullet}^* \alpha' + (\varepsilon_v)_R$. We can see that $\|u_R\| \leq \|A_{R_\bullet}^*\| \|\alpha\| + \|(\varepsilon_u)_R\| \leq O(\sqrt{k} \log n)$, so $\|M_8\| \leq \tilde{O}(\sigma_\varepsilon^4 k) = \tilde{O}(k^3/n^2)$. Since $m = O(n)$ and $k \leq O^*(\frac{\sqrt{n}}{\log n})$, we can bound all the above terms by $O^*(k/m \log n)$ and finish the proof of Claim 7. \square

Combining the results of Claims 6 and 7, we complete the proof of Lemma 2.

Proof of Lemma 3

The recovery of $A_{\bullet,i}^*$'s support directly follows from Lemma 1. For the latter part, recall from Lemma 2 that

$$M_{u,v} = q_i c_i \beta_i \beta_i' A_{R,i}^* A_{R,i}^{*T} + \text{perturbation terms.}$$

The perturbation terms have norms bounded by $O^*(1/\log n)$ w.h.p. On the other hand, the first term has norm at least $\Omega(k/m)$ since $\|A_{R,i}^*\| = 1$ for the correct support R and $|q_i c_i \beta_i \beta_i'| \geq \Omega(k/m)$. Then applying Wedin's Theorem to $M_{u,v}$, we can conclude that the top singular vector must be $O^*(k/m \log n)/\Omega(k/m) = O^*(1/\log n)$ -close to $A_{R,i}^*$. \square

Proof of Lemma 4

The proof follows from that of Lemma 37 in (Arora et al. 2015). The main idea is to separate the possible cases of how u and v share support and to use Lemma 2 with the bounded perturbation terms to conclude when u and v share exactly one. We note that due to the condition where the ordered statistics $e_{(r)} \geq \Omega(k/mr)$ and $e_{(r+1)}/e_{(r)} \leq O^*(r/\log^2 n)$, then it must be the case that u and v share only one atom or share more than one atoms with the same support. When their supports overlap more than one, then the first singular value cannot dominate the second one, and hence u and v can only share a unique element. \square

Analysis of Main Algorithm

Simple Encoding

Lemma 7. Assume that A^s is δ -close to A^* for $\delta = O^*(1/\log n)$ and $\mu \leq \frac{\sqrt{n}}{2k}$, and $k \geq \Omega(\log m)$. Then with high probability over random samples $y = A^* x^* + \varepsilon$,

$$\text{sgn}(\text{threshold}_{C/2}((A^s)^T y)) = \text{sgn}(x^*). \quad (4)$$

Proof. We follow the same proof strategy from (Arora et al. 2015) (Lemmas 16 and 17) to prove a more general version in which the noise ε is taken into account. Write $S = \text{supp}(x^*)$ and skip the superscript s on A^s for the readability. What we need is to show $S = \{i \in [m] : \langle A_{\bullet,i}, y \rangle \geq C/2\}$ and then $\text{sgn}(\langle A_{\bullet,i}^s, y \rangle) = \text{sgn}(x_i^*)$ for each $i \in S$ with high probability. Following the same argument of (Arora et al. 2015), we prove in below a stronger statement that, even conditioned on the support S , $S = \{i \in [m] : \langle A_{\bullet,i}, y \rangle \geq C/2\}$ with high probability.

Rewrite

$$\langle A_{\bullet,i}, y \rangle = \langle A_{\bullet,i}, A^* x^* + \varepsilon \rangle = \langle A_{\bullet,i}, A_{\bullet,i}^* \rangle x_i^* + \sum_{j \neq i} \langle A_{\bullet,i}, A_{\bullet,j}^* \rangle x_j^* + \langle A_{\bullet,i}, \varepsilon \rangle,$$

and observe that, due to the closeness of $A_{\bullet,i}$ and $A_{\bullet,i}^*$, the first term is either close to x_i^* or equal to 0 depending on whether or not $i \in S$. Meanwhile, the rest are small due to the incoherence and the concentration in the weighted average of noise. We will show that both $Z_i = \sum_{S \setminus \{i\}} \langle A_{\bullet,i}, A_{\bullet,j}^* \rangle x_j^*$ and $\langle A_{\bullet,i}, \varepsilon \rangle$ are bounded by $C/8$ with high probability.

The cross-term $Z_i = \sum_{S \setminus \{i\}} \langle A_{\bullet,i}, A_{\bullet,j}^* \rangle x_j^*$ is a sum of zero-mean independent sub-Gaussian random variables, which is another sub-Gaussian random variable with variance $\sigma_{Z_i}^2 = \sum_{S \setminus \{i\}} \langle A_{\bullet,i}, A_{\bullet,j}^* \rangle^2$. Note that

$$\langle A_{\bullet,i}, A_{\bullet,j}^* \rangle^2 \leq 2(\langle A_{\bullet,i}^*, A_{\bullet,j}^* \rangle^2 + \langle A_{\bullet,i} - A_{\bullet,i}^*, A_{\bullet,j}^* \rangle^2) \leq 2\mu^2/n + 2\langle A_{\bullet,i} - A_{\bullet,i}^*, A_{\bullet,j}^* \rangle^2,$$

where we use Cauchy-Schwarz inequality and the μ -incoherence of A^* . Therefore,

$$\sigma_{Z_i}^2 \leq 2\mu^2 k/n + 2\|A_{\bullet,S}^{*T}(A_{\bullet,i} - A_{\bullet,i}^*)\|_F^2 \leq 2\mu^2 k/n + 2\|A_{\bullet,S}^*\|^2 \|A_{\bullet,i} - A_{\bullet,i}^*\|^2 \leq O(1/\log n),$$

under the condition $\mu \leq \frac{\sqrt{n}}{2k}$ and $k = \Omega(\log m)$. Applying Bernstein's inequality, we get $|Z_i| \leq C/8$ with high probability. What remains is to bound the noise term $\langle A_{\bullet,i}, \varepsilon \rangle$. In fact, $\langle A_{\bullet,i}, \varepsilon \rangle$ is a sum of n Gaussian random variables, which is a sub-Gaussian with variance σ_ε^2 . It is easy to see that $|\langle A_{\bullet,i}, \varepsilon \rangle| \leq \sigma_\varepsilon \log n$ with high probability. Notice that $\sigma_\varepsilon = O(1/\sqrt{n})$.

Finally, we combine these bounds to have $|Z_i + \langle A_{\bullet,i}, \varepsilon \rangle| \leq C/4$. Therefore, for $i \in S$, then $|\langle A_{\bullet,i}, y \rangle| \geq C/2$ and negligible otherwise. Using union bound for every $i = 1, 2, \dots, m$, we finish the proof of the lemma. \square

Approximate Gradient in Expectation

Proof of Lemma 5. Having the result from Lemma 7, we are now able to study the expected update direction $g^s = \mathbb{E}[(A^s x - y) \text{sgn}(x)^T]$. Recall that A^s is the update at the s -th iteration and $x \triangleq \text{threshold}_{C/2}((A^s)^T y)$. Based on the generative model, denote $p_i = \mathbb{E}[x_i^* \text{sgn}(x_i^*) | i \in S]$, $q_i = \mathbb{P}[i \in S]$ and $q_{ij} = \mathbb{P}[i, j \in S]$. Throughout this section, we will use ζ to denote any vector whose norm is negligible although they can be different across their appearances. A_{-i} denotes the sub-matrix of A whose i -th column is removed. To avoid overwhelming appearance of the superscript s , we skip it from A^s for neatness. Write \mathcal{F}_{x^*} as the event under which the support of x is the same as that of x^* , and $\bar{\mathcal{F}}_{x^*}$ is its complement. In other words, $\mathbf{1}_{\mathcal{F}_{x^*}} = \mathbf{1}[\text{sgn}(x) = \text{sgn}(x^*)]$ and $\mathbf{1}_{\mathcal{F}_{x^*}} + \mathbf{1}_{\bar{\mathcal{F}}_{x^*}} = 1$. Note that

$$g_{\bullet,i}^s = \mathbb{E}[(Ax - y) \text{sgn}(x_i)] = \mathbb{E}[(Ax - y) \text{sgn}(x_i) \mathbf{1}_{\mathcal{F}_{x^*}}] \pm \zeta.$$

Using the fact that $y = A^*x^* + \varepsilon$ and that under \mathcal{F}_{x^*} , we have $Ax = A_{\bullet S}x_S = A_{\bullet S}A_{\bullet S}^T y = A_{\bullet S}A_{\bullet S}^T A^*x^* + A_{\bullet S}A_{\bullet S}^T \varepsilon$. Using the independence of ε and x^* to get rid of the noise term, we get

$$\begin{aligned} g_{\bullet i}^s &= \mathbb{E}[(A_{\bullet S}A_{\bullet S}^T - I_n)A^*x^*\mathbf{1}_{\mathcal{F}_{x^*}}] + \mathbb{E}[(A_{\bullet S}A_{\bullet S}^T - I_n)\varepsilon \text{sgn}(x_i)\mathbf{1}_{\mathcal{F}_{x^*}}] \pm \zeta \\ &= \mathbb{E}[(A_{\bullet S}A_{\bullet S}^T - I_n)A^*x^*\text{sgn}(x_i)\mathbf{1}_{\mathcal{F}_{x^*}}] \pm \zeta \quad (\text{Independence of } \varepsilon \text{ and } x^* \text{'s}) \\ &= \mathbb{E}[(A_{\bullet S}A_{\bullet S}^T - I_n)A^*x^*\text{sgn}(x_i^*)(1 - \mathbf{1}_{\mathcal{F}_{x^*}})] \pm \zeta \quad (\text{Under } \mathcal{F}_{x^*} \text{ event}) \\ &= \mathbb{E}[(A_{\bullet S}A_{\bullet S}^T - I_n)A^*x^*\text{sgn}(x_i^*)] \pm \zeta. \end{aligned}$$

Recall from the generative model assumptions that $S = \text{supp}(x^*)$ is random and the entries of x^* are pairwise independent given the support, so

$$\begin{aligned} g_{\bullet i}^s &= \mathbb{E}_S \mathbb{E}_{x^*|S} [(A_{\bullet S}A_{\bullet S}^T - I_n)A^*x^*\text{sgn}(x_i^*)] \pm \zeta \\ &= p_i \mathbb{E}_{S, i \in S} [(A_{\bullet S}A_{\bullet S}^T - I_n)A_{\bullet i}^*] \pm \zeta \\ &= p_i \mathbb{E}_{S, i \in S} [(A_{\bullet i}A_{\bullet i}^T - I_n)A_{\bullet i}^*] + p_i \mathbb{E}_{S, i \in S} \left[\sum_{l \in S, l \neq i} A_{\bullet l}A_{\bullet l}^T A_{\bullet i}^* \right] \pm \zeta \\ &= p_i q_i (A_{\bullet i}A_{\bullet i}^T - I_n)A_{\bullet i}^* + p_i \sum_{l \in [m], l \neq i} q_{il} A_{\bullet l}A_{\bullet l}^T A_{\bullet i}^* \pm \zeta \\ &= p_i q_i (\lambda_i A_{\bullet i} - A_{\bullet i}^*) + p_i A_{\bullet -i} \text{diag}(q_{ij}) A_{\bullet -i}^T A_{\bullet i}^* \pm \zeta, \end{aligned}$$

where $\lambda_i^s = \langle A_{\bullet i}^s, A_{\bullet i}^* \rangle$. Let $\xi_i^s = A_{R, -i} \text{diag}(q_{ij}) A_{\bullet -i}^T A_{\bullet i}^* / q_i$ for $j = 1, \dots, m$. We now have the full expression of the expected approximate gradient

$$g_{R, i}^s = p_i q_i (\lambda_i A_{R, i}^s - A_{R, i}^* + \xi_i^s) \pm \zeta_R. \quad (5)$$

What remains is to bound the norms of ξ_i^s and ζ . It follows from the nearness that $\|A_{R, -i}^s\| \leq O(\sqrt{m/n})$ and $\|A_{-i}^s\| \leq O(\sqrt{m/n})$. Then, along with the fact that $\|A_i^*\| = 1$, we can see that

$$\|\xi_i^s\| \leq \|A_{R, -i}^s\| \max_{j \neq i} \frac{q_{ij}}{q_i} \|A_{-i}^s\| \leq O(k/n). \quad (6)$$

Next, we show that ζ is negligible in norm. Since \mathcal{F}_{x^*} happens with very high probability, it suffices to bound norm of $(Ax - y)\text{sgn}(x_i)$ which can easily be done using Lemma 12 and Lemma 11 in Section ‘‘Sample Complexity’’. This concludes the proof for Lemma 5. \square

$(\alpha, \beta, \gamma_s)$ -Correlation

Proof of Lemma 6. Throughout the proof, we omit the superscript s for simplicity and denote $2\alpha = p_i q_i$. First, we rewrite $g_{\bullet i}^s$ as a combination of the true direction $A_{\bullet i}^s - A_{\bullet i}^*$ and a term with small norm:

$$g_{R, i} = 2\alpha(A_{R, i} - A_{R, i}^*) + v, \quad (7)$$

where $v = 2\alpha[(\lambda_i - 1)A_{\bullet i} + \epsilon_i]$ with norm bounded. In fact, since $A_{\bullet i}$ is δ -close to $A_{\bullet i}^*$, and both have unit norm, then $\|2\alpha(\lambda_i - 1)A_{\bullet i}\| = \alpha\|A_{\bullet i} - A_{\bullet i}^*\|^2 \leq \alpha\|A_{\bullet i} - A_{\bullet i}^*\|$ and $\|\xi_i\| \leq O(k/n)$ from (6). Therefore,

$$\|v\| = \|2\alpha(\lambda_i - 1)A_{R, i} + 2\alpha\xi_i\| \leq \alpha\|A_{R, i} - A_{R, i}^*\| + \epsilon,$$

where $\epsilon = O(k^2/mn)$. Now, we make use of (7) to show the first part of Lemma 6:

$$\langle 2g_{R, i}, A_{R, i} - A_{R, i}^* \rangle = 4\alpha\|A_{R, i} - A_{R, i}^*\|^2 + \langle 2v, A_{R, i} - A_{R, i}^* \rangle. \quad (8)$$

We want to lower bound the inner product term in $\|g_{R, i}\|^2$ and $\|A_{R, i} - A_{R, i}^*\|^2$. Effectively, from (7)

$$\begin{aligned} 4\alpha\langle v, A_{\bullet i} - A_{\bullet i}^* \rangle &= \|g_{R, i}\|^2 - 4\alpha^2\|A_{R, i} - A_{R, i}^*\|^2 - \|v\|^2 \\ &\geq \|g_{R, i}\|^2 - 6\alpha^2\|A_{R, i} - A_{R, i}^*\|^2 - 2\epsilon^2, \end{aligned} \quad (9)$$

where the last step is due to Cauchy-Schwarz inequality: $\|v\|^2 \leq 2(\alpha^2\|A_{R, i} - A_{R, i}^*\|^2 + \epsilon^2)$.

Substitute $2\langle v, A_{\bullet i} - A_{\bullet i}^* \rangle$ in (8) for the right hand side of (9), we get the first result:

$$\langle 2g_{R, i}, A_{R, i} - A_{R, i}^* \rangle \geq \alpha\|A_{R, i} - A_{R, i}^*\|^2 + \frac{1}{2\alpha}\|g_{R, i}\|^2 - \frac{\epsilon^2}{\alpha}.$$

The second part directly follows from the correlation and Theorem 1. Moreover, we have $p_i = \Theta(k/m)$ and $q_i = \Theta(1)$, then $\alpha = \Theta(k/m)$, $\beta = \Theta(m/k)$ and $\gamma_s = O(k^3/mn^2)$. Then $g_{R, i}^s$ is $(\Omega(k/m), \Omega(m/k), O(k^3/mn^2))$ -correlated with the true solution. \square

Nearness

Lemma 8. *If A^s is $(2\delta, 2)$ -near to A^* , then $\|A^{s+1} - A^*\| \leq 2\|A^*\|$.*

Proof. From Lemma 5 we have $g_{\bullet_i}^s = p_i q_i (\lambda_i A_{\bullet_i}^s - A_{\bullet_i}^*) + A_{\bullet_{-i}} \text{diag}(q_{ij}) A_{\bullet_{-i}}^T A_{\bullet_i}^* \pm \zeta$. Denote $\bar{R} = [n] \setminus R$, then it is obvious that $g_{\bar{R},i}^s = A_{\bar{R},-i} \text{diag}(q_{ij}) A_{\bar{R},-i}^T A_{\bullet_i}^* \pm \zeta$ is bounded by $O(k^2/m^2)$. Then we follow the proof of Lemma 24 in (Arora et al. 2015) for the nearness with full $g^s = g_{R,i}^s + g_{\bar{R},i}^s$ to finish the proof for this lemma. \square

Sample Complexity

In previous sections, we rigorously analyzed both initialization and learning algorithms as if the expectations g^s , e and $M_{u,v}$ were given. Here we show that corresponding estimates based on empirical means are sufficient for the algorithms to succeed, and identify how many samples are required. Technically, this requires the study of their concentrations around their expectations. Having had these concentrations, we are ready to prove Theorems 2 and 3.

The entire section involves a variety of concentration bounds. Here we make heavy use of Bernstein's inequality for different types of random variables (including scalar, vector and matrix). The Bernstein's inequality is stated as follows.

Lemma 9 (Bernstein's Inequality). *Suppose that $Z^{(1)}, Z^{(2)}, \dots, Z^{(p)}$ are p i.i.d. samples from some distribution \mathcal{D} . If $\mathbb{E}[Z] = 0$, $\|Z^{(j)}\| \leq \mathcal{R}$ almost surely and $\|\mathbb{E}[Z^{(j)}(Z^{(j)})^T]\| \leq \sigma^2$ for each j , then*

$$\frac{1}{p} \left\| \sum_{j=1}^p Z^{(j)} \right\| \leq \tilde{O} \left(\frac{\mathcal{R}}{p} + \sqrt{\frac{\sigma^2}{p}} \right) \quad (10)$$

holds with probability $1 - n^{-\omega(1)}$.

Since all random variables (or their norms) are not bounded almost surely in our model setting, we make use of a technical lemma that is used in (Arora et al. 2015) to handle the issue.

Lemma 10 ((Arora et al. 2015)). *Suppose a random variable Z satisfies $\mathbb{P}[\|Z\| \geq \mathcal{R}(\log(1/\rho))^C] \leq \rho$ for some constant $C > 0$, then*

- (a) *If $p = n^{O(1)}$, it holds that $\|Z^{(j)}\| \leq \tilde{O}(\mathcal{R})$ for each j with probability $1 - n^{-\omega(1)}$.*
- (b) *$\|\mathbb{E}[Z \mathbf{1}_{\|Z\| \geq \tilde{\Omega}(\mathcal{R})}]\| = n^{-\omega(1)}$.*

This lemma suggests that if $\frac{1}{p} \sum_{i=1}^p Z^{(j)} (1 - \mathbf{1}_{\|Z^{(j)}\| \geq \tilde{\Omega}(\mathcal{R})})$ concentrates around its mean with high probability, then so does $\frac{1}{p} \sum_{i=1}^p Z^{(j)}$ because the part outside the truncation level can be ignored. Since all random variables of our interest are sub-Gaussian or a product of sub-Gaussian that satisfy this lemma, we can apply Lemma 9 to the corresponding truncated random variables with carefully chosen truncation levels. Then the original random variables concentrate likewise.

In the next proofs, we define suitable random variables and identify good bounds of \mathcal{R} and σ^2 for them. Note that in this section, the expectations are taken over y by conditioning on u and v . This aligns with the construction that the estimators of e and $M_{u,v}$ are empirical averages over i.i.d. samples of y , while u and v are kept fixed. Due to the dependency on u and v , these (conditional) expectations inherit randomness from u and v , and we will formulate probabilistic bounds for them.

The application of Bernstein's inequality requires a bound on $\|\mathbb{E}[ZZ^T(1 - \mathbf{1}_{\|Z\| \geq \tilde{\Omega}(\mathcal{R})})]\|$. We achieve that by the following technical lemma, where \tilde{Z} is a standardized version of Z .

Lemma 11. *Suppose a random variable $\tilde{Z}\tilde{Z}^T = aT$ where $a \geq 0$ and T is positive semi-definite. They are both random. Suppose $\mathbb{P}[a \geq \mathcal{A}] = n^{-\omega(1)}$ and $\mathcal{B} > 0$ is a constant. Then,*

$$\|\mathbb{E}[\tilde{Z}\tilde{Z}^T(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\| \leq \mathcal{A}\|\mathbb{E}[T]\| + O(n^{-\omega(1)})$$

Proof. To show this, we make use of the decomposition $\tilde{Z}\tilde{Z}^T = aT$ and a truncation for a . Specifically,

$$\begin{aligned} \|\mathbb{E}[\tilde{Z}\tilde{Z}^T(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\| &= \mathbb{E}[aT(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})] \\ &\leq \|\mathbb{E}[a(1 - \mathbf{1}_{a \geq \mathcal{A}})T(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\| + \|\mathbb{E}[a\mathbf{1}_{a \geq \mathcal{A}}T(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\| \\ &\leq \|\mathbb{E}[a(1 - \mathbf{1}_{a \geq \mathcal{A}})T]\| + \mathbb{E}[a\mathbf{1}_{a \geq \mathcal{A}}\|T\|(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})] \\ &\leq \mathcal{A}\|\mathbb{E}[T]\| + (\mathbb{E}[\|aT\|^2(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\mathbb{E}[\mathbf{1}_{a \geq \mathcal{A}}])^{1/2} \\ &\leq \mathcal{A}\|\mathbb{E}[T]\| + (\mathbb{E}[\|\tilde{Z}\|^4(1 - \mathbf{1}_{\|\tilde{Z}\| \geq \mathcal{B}})]\mathbb{P}[a \geq \mathcal{A}])^{1/2} \\ &\leq \mathcal{A}\|\mathbb{E}[T]\| + \mathcal{B}^2(\mathbb{P}[a \geq \mathcal{A}])^{1/2} \\ &\leq \mathcal{A}\|\mathbb{E}[T]\| + O(n^{-\omega(1)}), \end{aligned}$$

where at the third step we used $T(1 - \mathbf{1}_{\|\tilde{z}\| \geq \mathcal{B}})] \preceq T$ because of the fact that T is the positive semi-definite and $1 - \mathbf{1}_{\|\tilde{z}\| \geq \mathcal{B}} \in \{0, 1\}$. Then, we finish the proof of the lemma. \square

Sample Complexity of Algorithm 1

In Algorithm 1, we empirically compute the ‘‘scores’’ \hat{e} and the reduced weighted covariance matrix $\widehat{M}_{u,v}$ to produce an estimate for each column of A^* . Since the construction of $\widehat{M}_{u,v}$ depends upon the support estimate \widehat{R} given by ranking \hat{e} , we denote it by $\widehat{M}_{u,v}^{\widehat{R}}$. We will show that we only need $p = \widetilde{O}(m)$ samples to be able to recover the support of one particular atom and up to some specified level of column-wise error with high probability.

Lemma 12. *Consider Algorithm 1 in which p is the given number of samples. For any pair u and v , then with high probability a) $\|\hat{e} - e\| \leq O^*(k/m \log^2 n)$ when $p = \widetilde{O}(m)$ and b) $\|\widehat{M}_{u,v}^{\widehat{R}} - M_{u,v}^R\| \leq O^*(k/m \log n)$ when $p = \widetilde{O}(mr)$ where \widehat{R} and R are respectively the estimated and correct support sets of one particular atom.*

Proof of Theorem 2 Using Lemma 12, we are ready to prove the Theorem 2. According to Lemma 1 when $U \cap V = \{i\}$, we can write \hat{e} as

$$\hat{e} = q_i c_i \beta_i \beta_i' A_{R,i}^* \circ A_{R,i}^* + \text{perturbation terms} + (\hat{e} - e),$$

and consider $\hat{e} - e$ as an additional perturbation with the same magnitude $O^*(k/m \log^2 n)$ in the sense of $\|\cdot\|_\infty$ w.h.p. The first part of Lemma 3 suggests that when u and v share exactly one atom i , then the set \widehat{R} including r largest elements of \hat{e} is the same as $\text{supp}(A_i^*)$ with high probability.

Once we have \widehat{R} , we again write $\widehat{M}_{u,v}^{\widehat{R}}$ using Lemma 2 as

$$\widehat{M}_{u,v}^{\widehat{R}} = q_i c_i \beta_i \beta_i' A_{R,i}^* A_{R,i}^{*T} + \text{perturbation terms} + (\widehat{M}_{u,v}^{\widehat{R}} - M_{u,v}^R),$$

and consider $\widehat{M}_{u,v}^{\widehat{R}} - M_{u,v}^R$ as an additional perturbation with the same magnitude $O^*(k/m \log n)$ in the sense of the spectral norm $\|\cdot\|$ w.h.p. Using the second part of Lemma 3, we have the top singular vectors of $\widehat{M}_{u,v}^{\widehat{R}}$ is $O^*(1/\log n)$ -close to $A_{R,i}^*$ with high probability.

Since every vector added to the list L in Algorithm 1 is close to one of the dictionary, then A^0 must be δ -close to A^* . In addition, the nearness of A^0 to A^* is guaranteed via an appropriate projection onto the convex set $\mathcal{B} = \{A | A \text{ close to } A^0 \text{ and } \|A\| \leq 2\|A^*\|\}$. Finally, we finish the proof of Theorem 2. \square

Proof of Lemma 12, Part a For some fixed $l \in [n]$, consider p i.i.d. realizations $Z^{(1)}, Z^{(2)}, \dots, Z^{(p)}$ of the random variable $Z \triangleq \langle y, u \rangle \langle y, v \rangle y_l^2$, then $\hat{e}_l = \frac{1}{p} \sum_{i=1}^p Z^{(i)}$ and $e_l = \mathbb{E}[Z]$. To show that $\|\hat{e} - e\|_\infty \leq O^*(k/m \log^2 n)$ holds with high probability, we first study the concentration for the l -th entry of $\hat{e} - e$ and then take the union bound over all $l = 1, 2, \dots, n$. We derive upper bounds for $|Z|$ and its variance $\mathbb{E}[Z^2]$ in order to apply Bernstein’s inequality in (12) to the truncated version of Z .

Claim 8. $|Z| \leq \widetilde{O}(k)$ and $\mathbb{E}[Z^2] \leq \widetilde{O}(k^2/m)$ with high probability.

Again, the expectation is taken over y by conditioning on u and v , and therefore is still random due to the randomness of u and v . To show Claim 8, we begin with proving the following auxiliary claim.

Claim 9. $\|y\| \leq \widetilde{O}(\sqrt{k})$ and $|\langle y, u \rangle| \leq \widetilde{O}(\sqrt{k})$ with high probability.

Proof. From the generative model, we have

$$\|y\| = \|A_{\bullet,S}^* x_S^* + \varepsilon\| \leq \|A_{\bullet,S}^* x_S^*\| + \|\varepsilon\| \leq \|A_{\bullet,S}^*\| \|x_S^*\| + \|\varepsilon\|,$$

where $S = \text{supp}(x^*)$. From Claim 2, $\|x_S^*\| \leq \widetilde{O}(\sqrt{k})$ and $\|\varepsilon\| \leq \widetilde{O}(\sigma_\varepsilon \sqrt{n})$ w.h.p. In addition, A^* is overcomplete and has bounded spectral norm, then $\|A_{\bullet,S}^*\| \leq \|A^*\| \leq O(1)$. Therefore, $\|y\| \leq \widetilde{O}(\sqrt{k})$ w.h.p., which is the first part of the proof. To bound the second term, we write it as

$$|\langle y, u \rangle| = |\langle A_{\bullet,S}^* x_S^* + \varepsilon, u \rangle| \leq |\langle x_S^*, A_{\bullet,S}^{*T} u \rangle| + |\langle \varepsilon, u \rangle|.$$

Similar to y , we have $\|u\| \leq \widetilde{O}(\sqrt{k})$ w.h.p. and hence $\|A_{\bullet,S}^{*T} u\| \leq \|A_{\bullet,S}^{*T}\| \|u\| \leq O(\sqrt{k})$ with high probability. Since u and x^* are independent sub-Gaussian and $\langle x_S^*, A_{\bullet,S}^{*T} u \rangle$ are sub-exponential with variance at most $O(\sqrt{k})$, $|\langle x_S^*, A_{\bullet,S}^{*T} u \rangle| \leq \widetilde{O}(k)$ w.h.p. Similarly, $|\langle \varepsilon, u \rangle| \leq \widetilde{O}(\sqrt{k})$ w.h.p. Consequently, $|\langle y, u \rangle| \leq \widetilde{O}(\sqrt{k})$ w.h.p., and we conclude the proof of the claim. \square

Proof of Claim 8. We have $Z = \langle y, u \rangle \langle y, v \rangle y_l^2 = \langle y, u \rangle \langle y, v \rangle (\langle A_{i_\bullet}^*, x^* \rangle + \varepsilon_l)^2$ with $\langle y, u \rangle \langle y, v \rangle \leq \tilde{O}(k)$ w.h.p. according to Claim 9. What remains is to bound $y_l^2 = (\langle A_{i_\bullet}^*, x^* \rangle + \varepsilon_l)^2$. Because $\langle A_{i_\bullet}^*, x^* \rangle$ is sub-Gaussian with variance $\mathbb{E}_S(\sum_{i \in S} A_{li}^{*2}) \leq \|A^{*T}\|_{1,2}^2 = O(1)$, then $|\langle A_{i_\bullet}^*, x^* \rangle| \leq O(\log n)$ w.h.p. Similarly for ε_l , $|\varepsilon_l| \leq O(\sigma_\varepsilon \log n)$ w.h.p. Ultimately, $|\langle A_{i_\bullet}^*, x^* \rangle + \varepsilon_l| \leq O(\log n)$, and hence we obtain with high probability the bound $|Z| \leq \tilde{O}(k)$.

To bound the variance term, we write $Z^2 = \langle y, v \rangle^2 y_l^2 \langle y, u \rangle^2 y_l^2$. Note that, from the first part, we get $\langle y, v \rangle^2 y_l^2 \leq \tilde{O}(k)$ and $|Z| \leq \tilde{O}(k)$ w.h.p.. We apply Lemma 11 with some appropriate scaling to both terms, then

$$\mathbb{E}[Z^2(1 - \mathbf{1}_{|Z| \geq \tilde{\Omega}(k)})] \leq \tilde{O}(k) \mathbb{E}[\langle y, u \rangle^2 y_l^2] + O(n^{-\omega(1)}),$$

where $\mathbb{E}[\langle y, u \rangle^2 y_l^2]$ is equal to e_l for pair u, v with $v = u$. From Lemma 1 and its proof in Appendix Section “Analysis of Initialization Algorithm”,

$$\mathbb{E}[\langle y, u \rangle^2 y_l^2] = \sum_{i=1}^m q_i c_i \beta_i^2 A_{li}^{*2} + \text{perturbation terms},$$

in which the perturbation terms are bounded by $O^*(k/m \log^2 n)$ w.h.p. (following Claims 4 and 5). The dominant term $\sum_i q_i c_i \beta_i^2 A_{li}^{*2} \leq (\max_i q_i c_i \beta_i^2) \|A_{i_\bullet}^*\|^2 \leq \tilde{O}(k/m)$ w.h.p. because $|\beta_i| \leq O(\log m)$ (Claim 3). Then we complete the proof of the second part. \square

Proof of Lemma 12, Part a. We are now ready to prove Part a of Lemma 12. We apply Bernstein’s inequality in Lemma 9 for the truncated random variable $Z^{(i)}(1 - \mathbf{1}_{|Z^{(i)}| \geq \tilde{\Omega}(\mathcal{R})})$ with $\mathcal{R} = \tilde{O}(k)$ and variance $\sigma^2 = \tilde{O}(k^2/m)$ from Claim 8, then

$$\left\| \frac{1}{p} \sum_{i=1}^p Z^{(i)}(1 - \mathbf{1}_{|Z^{(i)}| \geq \tilde{\Omega}(\mathcal{R})}) - \mathbb{E}[Z(1 - \mathbf{1}_{|Z| \geq \tilde{\Omega}(\mathcal{R})})] \right\| \leq \frac{\tilde{O}(k)}{p} + \sqrt{\frac{\tilde{O}(k^2/m)}{p}} \leq O^*(k/m \log^2 n), \quad (11)$$

w.h.p. for $p = \tilde{\Omega}(m)$. Then $\hat{e}_l = \frac{1}{p} \sum_{i=1}^p Z^{(i)}$ also concentrates with high probability. Take the union bound over $l = 1, 2, \dots, n$, we get $\|\hat{e} - e\|_\infty \leq O^*(k/m \log^2 n)$ with high probability and complete the proof of 12, Part a. \square

Proof of Lemma 12, Part b Next, we will prove that $\|\widehat{M}_{u,v}^{\widehat{R}} - M_{u,v}^R\| \leq O^*(k/m \log n)$ with high probability. We only need to prove the concentration inequalities for the case when conditioned on the event that \widehat{R} is equivalent to R w.h.p. Again, what we need to derive are an upper norm bound \mathcal{R} of the matrix random variable $Z \triangleq \langle y, u \rangle \langle y, v \rangle y_R y_R^T$ and its variance.

Claim 10. $\|Z\| \leq \tilde{O}(kr)$ and $\|\mathbb{E}[ZZ^T]\| \leq \tilde{O}(k^2 r/m)$ hold with high probability.

Proof. We have $\|Z\| \leq |\langle y, u \rangle \langle y, v \rangle| \|y_R\|^2$ with $|\langle y, u \rangle \langle y, v \rangle| \leq \tilde{O}(k)$ w.h.p. (according to Claim 9) whereas $\|y_R\|^2 = \sum_{i \in R} y_l^2 \leq O(r \log^2 n)$ w.h.p. because $y_l \leq O(\log n)$ w.h.p. (proof of Claim 8). This implies $\|Z\| \leq \tilde{O}(kr)$ w.h.p. The second part is handled similarly as in the proof of Claim 8. We take advantage of the bounds of $\widehat{M}_{u,v}$ in Lemma 2. Specifically, using the first part $\|Z\| \leq \tilde{O}(kr)$ and $\langle y, v \rangle^2 \|y_R\|^2 \leq \tilde{O}(kr)$, and applying Lemma 11, then

$$\|\mathbb{E}[ZZ^T(1 - \mathbf{1}_{\|Z\| \geq \tilde{\Omega}(kr)})]\| \leq \tilde{O}(kr) \mathbb{E}[\langle y, u \rangle^2 y_R y_R^T] + \tilde{O}(kr) O(n^{-\omega(1)}) \leq \tilde{O}(kr) \|M_{u,u}\|,$$

where $M_{u,u}$ arises from the application of Lemma 2. Recall that

$$M_{u,u} = \sum_i q_i c_i \beta_i^2 A_{R,i}^* A_{R,i}^{*T} + \text{perturbation terms},$$

where the perturbation terms are all bounded by $O^*(k/m \log n)$ w.h.p. by Claims 6 and 7. In addition,

$$\left\| \sum_i q_i c_i \beta_i^2 A_{R,i}^* A_{R,i}^{*T} \right\| \leq (\max_i q_i c_i \beta_i^2) \|A_{R,\bullet}^*\|^2 \leq \tilde{O}(k/m) \|A^*\|^2 \leq \tilde{O}(k/m)$$

w.h.p. Finally, the variance bound is $\tilde{O}(k^2 r/m)$ w.h.p. \square

Then, applying Bernstein’s inequality in Lemma 9 to the truncated version of Z with $\mathcal{R} = \tilde{O}(kr)$ and variance $\sigma^2 = \tilde{O}(k^2 r/m)$ and obtain the concentration for the full Z to get

$$\|\widehat{M}_{u,v}^R - M_{u,v}^R\| \leq \frac{\tilde{O}(kr)}{p} + \sqrt{\frac{\tilde{O}(k^2 r/m)}{p}} \leq O^*(k/m \log n)$$

w.h.p. when the number of samples is $p = \tilde{\Omega}(mr)$ under Assumption **A4.1**.

We have proved that $\|\widehat{M}_{u,v}^R - M_{u,v}^R\| \leq O^*(k/m \log n)$ as conditioned on the support consistency event holds w.h.p. $\|\widehat{M}_{u,v}^{\widehat{R}} - M_{u,v}^R\| \leq O^*(k/m \log n)$ is easily followed by the law of total probability through the tail bounds on the conditional and marginal probabilities (i.e. $\mathbb{P}[\|\widehat{M}_{u,v}^R - M_{u,v}^R\| \leq O^*(k/m \log n) | \widehat{R} = R]$) and $\mathbb{P}[\widehat{R} \neq R]$. We finish the proof of Lemma 12, Part b for both cases of the spectral bounds. \square

Proof of Theorem 3 and Sample Complexity of Algorithm 2

In this section, we prove Theorem 3 and identify sample complexity per iteration of Algorithm 2. We divide the proof into two steps: 1) show that when A^s is $(\delta_s, 2)$ -near to A^* for $\delta_s = O^*(1/\log n)$, the approximate gradient estimate \widehat{g}^s is $(\alpha, \beta, \gamma_s)$ -correlated-whp with A^* with $\gamma_s \leq O(k^2/mn) + \alpha o(\delta_s^2)$, and 2) show that the nearness is preserved at each iteration. These correspond to showing the following lemmas:

Lemma 13. *At iteration s of Algorithm 2, suppose that A^s has each column correctly supported and is $(\delta_s, 2)$ -near to A^* and that $\eta = O(m/k)$. Denote $R = \text{supp}(A_{\bullet,i}^s)$, then the update $\widehat{g}_{R,i}^s$ is $(\alpha, \beta, \gamma_s)$ -correlated-whp with $A_{R,i}^*$ where $\alpha = \Omega(k/m)$, $\beta = \Omega(m/k)$ and $\gamma_s \leq O(k^2/mn) + \alpha o(\delta_s^2)$ for $\delta_s = O^*(1/\log n)$.*

Note that this is a finite-sample version of Lemma 6.

Lemma 14. *If A^s is $(\delta_s, 2)$ -near to A^* and number of samples used in step s is $p = \tilde{\Omega}(m)$, then with high probability $\|A^{s+1} - A^*\| \leq 2\|A^*\|$.*

Proof of Theorem 3. The correlation of \widehat{g}_i with A_i^* , described in Lemma 13, implies the descent of column-wise error according to Theorem 1. Along with Lemma 14, the theorem follows directly.

Proof of Lemma 13 We prove Lemma 13 by obtaining a tail bound on the difference between $\widehat{g}_{R,i}^s$ and $g_{R,i}^s$ using the Bernstein's inequality in Lemma 9.

Lemma 15. *At iteration s of Algorithm 2, suppose that A^s has each column correctly supported and is $(\delta_s, 2)$ -near to A^* . For $R = \text{supp}(A_i^s) = \text{supp}(A_i^*)$, then $\|\widehat{g}_{R,i}^s - g_{R,i}^s\| \leq O(k/m) \cdot (o(\delta_s) + O(\epsilon_s))$ with high probability for $\delta_s = O^*(1/\log n)$ and $\epsilon_s = O(\sqrt{k/n})$ when $p = \tilde{\Omega}(m + \sigma_\epsilon^2 \frac{mnr}{k})$.*

To prove this lemma, we study the concentration of $\widehat{g}_{R,i}^s$, which is a sum of random vector of the form $(y - Ax)_{R\text{sgn}(x_i)}$. We consider random variable $Z \triangleq (y - Ax)_{R\text{sgn}(x_i)} | i \in S$, with $S = \text{supp}(x^*)$ and $x = \text{threshold}_{C/2}(A^T y)$. Then, using the following technical lemma to bridge the gap in concentration of the two variables. We adopt this strategy from (Arora et al. 2015) for our purpose.

Claim 11. *Suppose that $Z^{(1)}, Z^{(2)}, \dots, Z^{(N)}$ are i.i.d. samples of the random variable $Z = (y - Ax)_{R\text{sgn}(x_i)} | i \in S$. Then,*

$$\left\| \frac{1}{N} \sum_{j=1}^N Z^{(j)} - \mathbb{E}[Z] \right\| \leq o(\delta_s) + O(\epsilon_s) \quad (12)$$

holds with probability when $N = \tilde{\Omega}(k + \sigma_\epsilon^2 nr)$, $\delta_s = O^*(1/\log n)$ and $\epsilon_s = O(\sqrt{k/n})$.

Proof of Lemma 15. Once we have done the proof of Claim 11, we can easily prove Lemma 15. We recycle the proof of Lemma 43 in (Arora et al. 2015).

Write $W = \{j : i \in \text{supp}(x^{*(j)})\}$ and $N = |W|$, then express $\widehat{g}_{R,i}$ as

$$\widehat{g}_{R,i} = \frac{N}{p} \frac{1}{N} \sum_j (y^{(j)} - Ax^{(j)})_{R\text{sgn}(x_i^{(j)})},$$

where $\frac{1}{|W|} \sum_j (y^{(j)} - Ax^{(j)})_{R\text{sgn}(x_i^{(j)})}$ is distributed as $\frac{1}{N} \sum_{j=1}^N Z^{(j)}$ with $N = |W|$. Note that $\mathbb{E}[(y - Ax)_{R\text{sgn}(x_i)}] = \mathbb{E}[(y - Ax)_{R\text{sgn}(x_i)} \mathbf{1}_{i \in S}] = \mathbb{E}[Z] \mathbb{P}[i \in S] = q_i \mathbb{E}[Z]$ with $q_i = \Theta(k/m)$. Following Claim 11, we have

$$\|\widehat{g}_{R,i}^s - g_{R,i}^s\| \leq O(k/m) \left\| \frac{1}{N} \sum_{j=1}^N Z^{(j)} - \mathbb{E}[Z] \right\| \leq O(k/m) \cdot (o(\delta_s) + O(\epsilon_s)),$$

holds with high probability as $p = \Omega(mN/k)$. Substituting N in Claim 11, we obtain the results in Lemma 15. \square

Proof of Claim 11. We are now ready to prove the claim. What we need are good bounds for $\|Z\|$ and its variance, then we can apply Bernstein's inequality in Lemma 9 for the truncated version of Z , then Z is also concentrates likewise.

Claim 12. $\|Z\| \leq \mathcal{R}$ holds with high probability for $\mathcal{R} = \tilde{O}(\delta_s \sqrt{k} + \mu k / \sqrt{n} + \sigma_\varepsilon \sqrt{r})$ with $\delta_s = O^*(1/\log n)$.

Proof. From the generative model and the support consistency of the encoding step, we have $y = A^* x^* + \varepsilon = A_{\bullet S}^* x_S^* + \varepsilon$ and $x_S = A_{\bullet S}^T y = A_{\bullet S}^T A_{\bullet S}^* x_S^* + A_{\bullet S}^T \varepsilon$. Then,

$$\begin{aligned} (y - Ax)_R &= (A_{R,S}^* x_S^* + \varepsilon_R) - A_{R,S} A_{\bullet S}^T A_{\bullet S}^* x_S^* - A_{R,S} A_{\bullet S}^T \varepsilon \\ &= (A_{R,S}^* - A_{R,S}) x_S^* + A_{R,S} (I_k - A_{\bullet S}^T A_{\bullet S}^*) x_S^* + (I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet} \varepsilon. \end{aligned}$$

Using the fact that x_S^* and ε are sub-Gaussian and that $\|Mw\| \leq \tilde{O}(\sigma_w \|M\|_F)$ holds with high probability for a fixed M and a sub-Gaussian w of variance σ_w^2 , we have

$$\|(y - Ax)_R \text{sgn}(x_i)\| \leq \tilde{O}(\|A_{R,S}^* - A_{R,S}\|_F + \|A_{R,S}(I_k - A_{\bullet S}^T A_{\bullet S}^*)\|_F + \sigma_\varepsilon \|(I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet}\|_F).$$

Now, we need to bound those Frobenius norms. The first quantity is easily bounded as

$$\|A_{R,S}^* - A_{R,S}\|_F \leq \|A_{\bullet S}^* - A_{\bullet S}\|_F \leq \delta_s \sqrt{k}, \quad (13)$$

since A is δ_s -close to A^* . To handle the other two, we use the fact that $\|UV\|_F \leq \|U\| \|V\|_F$. Using this fact for the second term, we have

$$\|A_{R,S}(I_k - A_{\bullet S}^T A_{\bullet S}^*)\|_F \leq \|A_{R,S}\| \|(I_k - A_{\bullet S}^T A_{\bullet S}^*)\|_F,$$

where $\|A_{R,S}\| \leq \|A_{R\bullet}\| \leq O(1)$ due to the nearness. The second part is rearranged to take advantage of the closeness and incoherence properties:

$$\begin{aligned} \|I_k - A_{\bullet S}^T A_{\bullet S}^*\|_F &\leq \|I_k - A_{\bullet S}^{*T} A_{\bullet S}^* - (A_{\bullet S} - A_{\bullet S}^*)^T A_{\bullet S}^*\|_F \\ &\leq \|I_k - A_{\bullet S}^{*T} A_{\bullet S}^*\|_F + \|(A_{\bullet S} - A_{\bullet S}^*)^T A_{\bullet S}^*\|_F \\ &\leq \|I_k - A_{\bullet S}^{*T} A_{\bullet S}^*\|_F + \|A_{\bullet S}^*\| \|A_{\bullet S} - A_{\bullet S}^*\|_F \\ &\leq \mu k / \sqrt{n} + O(\delta_s \sqrt{k}), \end{aligned}$$

where we have used $\|I_k - A_{\bullet S}^{*T} A_{\bullet S}^*\|_F \leq \mu k / \sqrt{n}$ because of the μ -incoherence of A^* , $\|A_{\bullet S} - A_{\bullet S}^*\|_F \leq \delta_s \sqrt{k}$ in (13) and $\|A_{\bullet S}^*\| \leq \|A^*\| \leq O(1)$. Accordingly, the second Frobenius norm is bounded by

$$\|A_{R,S}(I_k - A_{\bullet S}^T A_{\bullet S}^*)\|_F \leq O(\mu k / \sqrt{n} + \delta_s \sqrt{k}). \quad (14)$$

The noise term is handled using the eigen-decomposition $U \Lambda U^T$ of $A_{\bullet S} A_{\bullet S}^T$, then with high probability

$$\|(I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet}\|_F = \|(UU^T - U \Lambda U^T)_{R\bullet}\|_F = \|U_{R\bullet}(I_n - \Lambda)\|_F \leq \|I_n - \Lambda\| \|U_{R\bullet}\|_F \leq O(\sqrt{r}), \quad (15)$$

where the last inequality $\|I_n - \Lambda\| \leq O(1)$ follows by $\|A_{\bullet S}\| \leq \|A\| \leq \|A - A^*\| + \|A^*\| \leq 3\|A^*\| \leq O(1)$ due to the nearness. Putting (13), (14) and (15) together, we obtain the bounds in Claim 12. \square

Next, we determine a bound for the variance of Z .

Claim 13. $\mathbb{E}[\|Z\|^2] = \mathbb{E}[\|(y - Ax)_R \text{sgn}(x_i)\|^2 | i \in S] \leq \sigma^2$ holds with high probability for $\sigma^2 = O(\delta_s^2 k + k^2/n + \sigma_\varepsilon^2 r)$ with $\delta_s = O^*(1/\log n)$.

Proof. We explicitly calculate the variance using the fact that x_S^* is conditionally independent given S , and so is ε . x_S^* and ε are also independent and have zero mean. Then we can decompose the norm into three terms in which the dot product is zero in expectation and the others can be shortened using the fact that $E[x_S^* x_S^{*T}] = I_k$, $E[\varepsilon \varepsilon^T] = \sigma_\varepsilon^2 I_n$.

$$\begin{aligned} \mathbb{E}[\|(y - Ax)_R \text{sgn}(x_i)\|^2 | i \in S] &= \mathbb{E}[\|(A_{R,S}^* - A_{R,S} A_{\bullet S}^T A_{\bullet S}^*) x_S^* + (I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet} \varepsilon\|^2 | i \in S] \\ &= \mathbb{E}[\|A_{R,S}^* - A_{R,S} A_{\bullet S}^T A_{\bullet S}^*\|_F^2 | i \in S] + \sigma_\varepsilon^2 \mathbb{E}[\|(I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet}\|_F^2 | i \in S]. \end{aligned}$$

Then, by re-writing $A_{R,S}^* - A_{R,S} A_{\bullet S}^T A_{\bullet S}^*$ as before, we get the form $(A_{R,S}^* - A_{R,S}) + A_{R,S}(I_k - A_{\bullet S}^T A_{\bullet S}^*)$ in which the first term has norm bounded by $\delta_s \sqrt{k}$. The second is further decomposed as

$$\mathbb{E}[\|A_{R,S}(I_k - A_{\bullet S}^T A_{\bullet S}^*)\|_F^2 | i \in S] \leq \sup_S \|A_{R,S}\|_F^2 \mathbb{E}[\|I_k - A_{\bullet S}^T A_{\bullet S}^*\|_F^2 | i \in S], \quad (16)$$

where $\sup_S \|A_{R,S}\|_F \leq \|A_{R\bullet}\| \leq O(1)$. We will bound $\mathbb{E}[\|I_k - A_{\bullet S}^T A_{\bullet S}^*\|_F^2 | i \in S] \leq O(k \delta_s^2) + O(k^2/n)$ using the proof from (Arora et al. 2015):

$$\begin{aligned} \mathbb{E}[\|I_k - A_{\bullet S}^T A_{\bullet S}^*\|_F^2 | i \in S] &= \mathbb{E}[\sum_{j \in S} (1 - A_{\bullet j}^T A_{\bullet j}^*)^2 + \sum_{j \in S} \|A_{\bullet j}^T A_{\bullet, -j}^*\|^2 | i \in S] \\ &= \mathbb{E}[\sum_{j \in S} \frac{1}{4} \|A_{\bullet j} - A_{\bullet j}^*\|^2] + q_{ij} \sum_{j \neq i} \|A_{\bullet j}^T A_{\bullet, -j}^*\|^2 + q_i \|A_{\bullet i}^T A_{\bullet, -i}^*\|^2 + q_i \|A_{\bullet, -i}^T A_{\bullet i}^*\|^2, \end{aligned}$$

where $A_{\bullet,-i}$ is the matrix A with the i -th column removed, $q_{ij} \leq O(k^2/m^2)$ and $q_i \leq O(k/m)$. For any $j = 1, 2, \dots, m$,

$$\begin{aligned} \|A_{\bullet,j}^T A_{\bullet,-j}^*\|^2 &= \|A_{\bullet,j}^{*T} A_{\bullet,-j}^* + (A_{\bullet,j} - A_{\bullet,j}^*)^T A_{\bullet,-j}^*\|^2 \\ &\leq \sum_{l \neq j} \langle A_{\bullet,j}^*, A_{\bullet,l}^* \rangle^2 + \|(A_{\bullet,j} - A_{\bullet,j}^*)^T A_{\bullet,-j}^*\|^2 \\ &\leq \sum_{l \neq j} \langle A_{\bullet,j}^*, A_{\bullet,l}^* \rangle^2 + \|A_{\bullet,j} - A_{\bullet,j}^*\|^2 \|A_{\bullet,-j}^*\|^2 \leq \mu^2 + \delta_s^2. \end{aligned}$$

The last inequality invokes the μ -incoherence, δ -closeness and the spectral norm of A^* . Similarly, we come up with the same bound for $\|A_{\bullet,i}^T A_{\bullet,-i}^*\|^2$ and $\|A_{\bullet,-i}^T A_{\bullet,i}^*\|^2$. Consequently,

$$\mathbb{E}[\|I_k - A_{\bullet S}^T A_{\bullet S}^*\|_F^2 | i \in S] \leq O(k\delta_s^2) + O(k^2/n). \quad (17)$$

For the last term, we invoke the inequality (15) (Claim 12) to get

$$\mathbb{E}[\|(I_n - A_{\bullet S} A_{\bullet S}^T)_{R\bullet}\|_F^2 | i \in S] \leq r \quad (18)$$

Putting (16), (17) and (18) together and using $\|A_{R\bullet}\| \leq 1$, we obtain the variance bound of Z : $\sigma^2 = O(\delta_s^2 k + k^2/n + \sigma_\epsilon^2 r)$ with $\delta_s = O(1/\log^2 n)$. Finally, we complete the proof. \square

We now apply truncated Bernstein's inequality to the random variable $Z^{(j)}(1 - 1_{\|Z^{(j)}\| \geq \Omega(\mathcal{R})})$ with \mathcal{R} and σ^2 in Claims 12 and 13, which are $\mathcal{R} = \tilde{O}(\delta_s \sqrt{k} + \mu k / \sqrt{n} + \sigma_\epsilon \sqrt{r})$ and $\sigma^2 = O(\delta_s^2 k + k^2/n + \sigma_\epsilon^2 r)$. Then, $(1/N) \sum_{j=1}^N Z^{(j)}$ also concentrates:

$$\left\| \frac{1}{N} \sum_{i=1}^N Z^{(j)} - E[Z] \right\| \leq \tilde{O}\left(\frac{\mathcal{R}}{N}\right) + \tilde{O}\left(\sqrt{\frac{\sigma^2}{N}}\right) = o(\delta_s) + O(\sqrt{k/n})$$

holds with high probability when $N = \tilde{\Omega}(k + \sigma_\epsilon^2 nr)$. Then, we finally finish the proof of Claim 11. \square

Proof of Lemma 13. With Claim 11, we study the concentration of $\hat{g}_{R,i}^s$ around its mean $g_{R,i}^s$. Now, we consider this difference as an error term of the expectation $g_{R,i}^s$ and using Lemma 6 to show the correlation of $\hat{g}_{R,i}^s$. Using the expression in Lemma 5 with high probability, we can write

$$\hat{g}_{R,i}^s = g_{R,i}^s + (g_{R,i}^s - \hat{g}_{R,i}^s) = 2\alpha(A_{R,i} - A_{R,i}^*) + v,$$

where $\|v\| \leq \alpha \|A_{R,i} - A_{R,i}^*\| + O(k/m) \cdot (o(\delta_s) + O(\epsilon_s))$. By Lemma 6, we have $\hat{g}_{R,i}^s$ is $(\alpha, \beta, \gamma_s)$ -correlated-whp with $A_{R,i}^*$ where $\alpha = \Omega(k/m)$, $\beta = \Omega(m/k)$ and $\gamma_s \leq O(k/m) \cdot (o(\delta_s) + O(\sqrt{k/n}))$, then we have done the proof Lemma 13. \square

Proof of Lemma 14 We have shown the correlation of \hat{g}^s with A^* w.h.p. and established the descent property of Algorithm 2. The next step is to show that the nearness is preserved at each iteration. To prove $\|A^{s+1} - A^*\| \leq 2\|A^s\|$ holds with high probability, we recall the update rule

$$A^{s+1} = A^s - \eta \mathcal{P}_H(\hat{g}^s),$$

where $\mathcal{P}_H(\hat{g}^s) = H \circ \hat{g}^s$. Here $H = (h_{ij})$ where $h_{ij} = 1$ if $i \in \text{supp}(A_{\bullet,j})$ and $h_{ij} = 0$ otherwise. Also, note that A^s is $(\delta_s, 2)$ -near to A^* for $\delta_s = O^*(1/\log n)$. We already proved that this holds for the exact expectation g^s in Lemma 8. To prove for \hat{g}^s , we again apply matrix Bernstein's inequality to bound $\|\mathcal{P}_H(g^s) - \mathcal{P}_H(\hat{g}^s)\|$ by $O(k/m)$ because $\eta = \Theta(m/k)$ and $\|A^*\| = O(1)$.

Consider a matrix random variable $Z \triangleq \mathcal{P}_H((y - Ax)\text{sgn}(x)^T)$. Our goal is to bound the spectral norm $\|Z\|$ and, both $\|\mathbb{E}[ZZ^T]\|$ and $\|\mathbb{E}[Z^T Z]\|$ since Z is asymmetric. To simplify our notations, we denote by x_R the vector x by zeroing out the elements not in R . Also, denote $R_i = \text{supp}(h_i)$ and $S = \text{supp}(x)$. Then Z can be written explicitly as

$$Z = [(y - Ax)_{R_1} \text{sgn}(x_1), \dots, (y - Ax)_{R_m} \text{sgn}(x_m)],$$

where many columns are zero since x is k -sparse. The following claims follow from the proof of Claim 42 in (Arora et al. 2015). Here we state and detail some important steps.

Claim 14. $\|Z\| \leq \tilde{O}(k)$ holds with high probability.

Proof. With high probability

$$\|Z\| \leq \sqrt{\sum_{i \in S} \|(y - Ax)_{R_i} \text{sgn}(x_i)\|^2} \leq \sqrt{k} \|(y - Ax)_{R_i}\|$$

where we use Claim 12 with $\|(y - Ax)_R\| \leq \tilde{O}(\delta_s \sqrt{k})$ w.h.p., then $\|Z\| \leq \tilde{O}(k)$ holds w.h.p. \square

Claim 15. $\|\mathbb{E}[ZZ^T]\| \leq O(k^2/n)$ and $\|\mathbb{E}[Z^T Z]\| \leq \tilde{O}(k^2/n)$ with high probability.

Proof. The first term is easily handled. Specifically, with high probability

$$\|\mathbb{E}[ZZ^T]\| \leq \|\mathbb{E}[\sum_{i \in S} (y - Ax)_{R_i} \text{sgn}(x_i)^2 (y - Ax)_{R_i}^T]\| = \|\mathbb{E}[\sum_{i \in S} (y - Ax)_{R_i} (y - Ax)_{R_i}^T]\| \leq O(k^2/n),$$

where the last inequality follows from the proof of Claim 42 in (Arora et al. 2015), which is tedious to be repeated.

To bound $\|\mathbb{E}[Z^T Z]\|$, we use bound of the full matrix $(y - Ax)\text{sgn}(x)^T$. Note that $\|y - Ax\| \leq \tilde{O}(\sqrt{k})$ w.h.p. is similar to what derived in Claim 12. Then with high probability,

$$\|\mathbb{E}[Z^T Z]\| \leq \|\mathbb{E}[\text{sgn}(x)(y - Ax)^T (y - Ax)\text{sgn}(x)^T]\| \leq \tilde{O}(k)\|\mathbb{E}[\text{sgn}(x)\text{sgn}(x)^T]\| \leq \tilde{O}(k^2/m).$$

where $\mathbb{E}[\text{sgn}(x)\text{sgn}(x)^T] = \text{diag}(q_1, q_2, \dots, q_m)$ has norm bounded by $O(k/m)$. We now can apply Bernstein's inequality for the truncated version of Z with $\mathcal{R} = \tilde{O}(k)$ and $\sigma^2 = \tilde{O}(k^2/m)$, then with $p = \tilde{O}(m)$,

$$\|\mathcal{P}_H(g^s) - \mathcal{P}_H(\hat{g}^s)\| \leq \frac{\tilde{O}(k)}{p} + \sqrt{\frac{\tilde{O}(k^2/m)}{p}} \leq O^*(k/m)$$

holds with high probability. Finally, we invoke the bound $\eta = O(m/k)$ and complete the proof. \square

Neural Implementation of Our Approach

We now briefly describe why our algorithm is “neurally plausible”. Basically, similar to the argument in (Arora et al. 2015), we describe at a very high level how our algorithm can be implemented via a neural network architecture. One should note that although both our initialization and descent stages are non-trivial modifications of those in (Arora et al. 2015), both still inherit the nice neural plausibility property.

Neural implementation of Stage 1: Initialization

Recall that the initialization stage includes two main steps: (i) estimate the support of each column of the synthesis matrix, and (ii) compute the top principal component(s) of a certain truncated weighted covariance matrix. Both steps involve simple vector and matrix-vector manipulations that can be implemented plausibly using basic neuronal manipulations.

For the support estimation step, we compute the product $\langle y, u \rangle \langle y, u \rangle y \circ y$, followed by a thresholding. The inner products, $\langle y, u \rangle$ and $\langle y, v \rangle$ can be computed using neurons via an online manner where the samples arrive in sequence; the thresholding can be implemented via a ReLU-type non-linearity.

For the second step, it is well known that the top principal components of a matrix can be computed in a neural (Hebbian) fashion using Oja's Rule (Oja 1992).

Neural implementation of Stage 2: Descent

Our neural implementation of the descent stage (Alg. 2) mimics the architecture of (Arora et al. 2015), which describes a simple two-layer network architecture for computing a single gradient update of A . The only difference in our case is that most of the value in A are set to zero, or in other words, our network is sparse. The network takes values y from the input layer and produce x as the output; there is an intermediate layer in between connecting the middle layer with the output via synapses. The synaptic weights are stored on A . The weights are updated by Hebbian learning. In our case, since A is sparse (with support given by R , as estimated in the first stage), we enforce the condition the corresponding synapses are inactive. In the output layer, as in the initialization stage, the neurons can use a ReLU-type non-linear activation function to enforce the sparsity of x .