

Learning Sparsely Used Overcomplete Dictionaries via Alternating Minimization

Alekh Agarwal, Animashree Anandkumar, Prateek Jain,
Praneeth Netrapalli *

July 30, 2014

Abstract

We consider the problem of sparse coding, where each sample consists of a sparse linear combination of a set of dictionary atoms, and the task is to learn both the dictionary elements and the mixing coefficients. Alternating minimization is a popular heuristic for sparse coding, where the dictionary and the coefficients are estimated in alternate steps, keeping the other fixed. Typically, the coefficients are estimated via ℓ_1 minimization, keeping the dictionary fixed, and the dictionary is estimated through least squares, keeping the coefficients fixed. In this paper, we establish local linear convergence for this variant of alternating minimization and establish that the basin of attraction for the global optimum (corresponding to the true dictionary and the coefficients) is $\mathcal{O}(1/s^2)$, where s is the sparsity level in each sample and the dictionary satisfies RIP. Combined with the recent results of approximate dictionary estimation, this yields provable guarantees for exact recovery of both the dictionary elements and the coefficients, when the dictionary elements are incoherent.

Keywords: Dictionary learning, sparse coding, alternating minimization, RIP, incoherence, lasso.

1 Introduction

A sparse code encodes each sample with a sparse set of elements, termed as dictionary atoms. Specifically, given a set of samples $Y \in \mathbb{R}^{d \times n}$, the generative model is

$$Y = A^* X^*, \quad A^* \in \mathbb{R}^{d \times r}, X^* \in \mathbb{R}^{r \times n},$$

and additionally, each column of X^* has at most s non-zero entries. The columns of A^* correspond to the dictionary atoms, and the columns of X^* correspond to the mixing coefficients of each sample. Each sample is a combination of at most s dictionary atoms. Sparse codes can thus succinctly represent high dimensional observed data.

*A. Agarwal is with Microsoft Research, New York, USA. Email: alekha@microsoft.com. A. Anandkumar is with the Center for Pervasive Communications and Computing, Electrical Engineering and Computer Science Dept., University of California, Irvine, USA 92697. Email: a.anandkumar@uci.edu. P. Jain is with Microsoft Research, Bangalore, India. Email: prajain@microsoft.com. P. Netrapalli is with Dept. of ECE, The University of Texas at Austin. Email: praneethn@utexas.edu. Part of this work was done when A. Anandkumar and P. Netrapalli were visiting Microsoft Research. An extended abstract containing an earlier version of these results appears in COLT 2014.

The problem of sparse coding consists of unsupervised learning of the dictionary and the coefficient matrices. Thus, given only unlabeled data, we aim to learn the set of dictionary atoms or basis functions that provide a good fit to the observed data. Sparse coding is applied in a variety of domains. Sparse coding of natural images has yielded dictionary atoms which resemble the receptive fields of neurons in the visual cortex [26, 27], and has also yielded localized dictionary elements on speech and video data [19, 25].

An important strength of sparse coding is that it can incorporate overcomplete dictionaries, where the number of dictionary atoms r can exceed the observed dimensionality d . It has been argued that having overcomplete representation provides greater flexibility in modeling and more robustness to noise [19], which is crucial for encoding complex signals present in images, speech and video. It has been shown that the performance of most machine learning methods employed downstream is critically dependent on the choice of data representations, and overcomplete representations are the key to obtaining state-of-art prediction results [6].

On the downside, the problem of learning sparse codes is computationally challenging, and is in general, NP-hard [9]. In practice, heuristics are employed based on alternating minimization. At a high level, this consists of alternating steps, where the dictionary is kept fixed and the coefficients are updated and vice versa. Such alternating minimization methods have enjoyed empirical success in a number of settings [18, 10, 2, 20, 35]. In this paper, we carry out a theoretical analysis of the alternating minimization procedure for sparse coding.

1.1 Summary of Results

We consider the alternating minimization procedure where we employ an initial estimate of the dictionary and then use ℓ_1 based minimization for estimating the coefficient matrix, given the dictionary estimate. The dictionary is subsequently re-estimated given the coefficient estimates. We establish local convergence to the true dictionary A^* and coefficient matrix X^* for this procedure whenever A^* satisfies RIP for $2s$ -sparse vectors. In other words, we characterize the “basin of attraction” for the true solution (A^*, X^*) and establish that alternating minimization succeeds in its recovery when a dictionary is initialized with an error of at most $\mathcal{O}(1/s^2)$, where s is the sparsity level. More precisely, the initial dictionary estimate $A(0)$ is required to satisfy

$$\epsilon_0 := \max_{i \in [r]} \min_{z \in \{-1, +1\}} \|zA_i^* - A(0)_i\|_2 = \mathcal{O}\left(\frac{1}{s^2}\right),$$

where A_i^* represents i^{th} column of A^* .

Further when the sparsity level satisfies $s = \mathcal{O}(d^{1/6})$ and the number of samples satisfies $n = \mathcal{O}(r^2)$, we establish a linear rate of convergence for the alternating minimization procedure to the true dictionary even when the dictionary is overcomplete ($r \geq d$), .

For the case of incoherent dictionaries, by combining the above result with recent results on approximate dictionary estimation by Agarwal et. al [1] or Arora et. al [3], we guarantee exact recovery of the true solution (A^*, X^*) when the alternating procedure is initialized with the output of [1] or [3]. If we employ the procedure of Agarwal et. al [1], the overall requirements are as follows: the sparsity level is required to be $s = \mathcal{O}(d^{1/9}, r^{1/8})$, and the number of samples $n = \mathcal{O}(r^2)$ to guarantee exact recovery of the true solution. If we employ the procedure of Arora et. al [3] (in particular their OVERLAPPINGAVERAGE procedure), we can establish exact recovery assuming $s = \mathcal{O}(r^{1/6}, \sqrt{d})$.

1.2 Related Work

Analysis of local optima of non-convex programs for sparse coding: Gribonval and Schnass [13], Geng et al. [12] and Jenatton et al. [15] carry out a theoretical analysis and study the conditions under which the true solution turns out to be a local optimum of a non-convex optimization problem for dictionary recovery. Gribonval and Schnass [13] and Geng et. al [12] both consider the noiseless setting, and analyze the following non-convex program

$$\min \|X\|_1 \quad s.t., Y = AX, \|A_i\|_2 = 1, \forall i \in [r]. \quad (1)$$

Since A and X are both unknown, the constraint $Y = AX$ is non-convex. It is natural to expect the true solution (A^*, X^*) to be a local optimum for (1) under fairly mild conditions, but this turns out to be non-trivial to establish. The difficulties arise from the non-convexity of the problem and the presence of sign-permutation ambiguity which leads to exponentially many equivalent solutions obtained via sign change and permutation. Gribonval and Schnass [13] established that (A^*, X^*) is a local optimum for (1), but limited to the case where the dictionary matrix A is square and hence, did not incorporate the overcomplete setting. Geng et al. [12] extend the analysis to the overcomplete setting, and establish that the true solution is a local optimum of (1) w.h.p. for incoherent dictionaries, when the number of samples n and sparsity level s scale as

$$n = \Omega(\|A\|_2^4 r^3 s), \quad s = \mathcal{O}(\sqrt{d}). \quad (2)$$

In our setting, where the spectral norm is assumed to be $\|A\|_2 < \mu_1 \sqrt{r/d}$, for some constant $\mu_1 > 0$, the sample complexity simplifies as $n = \Omega(r^5 s/d^2)$. Jenatton et al. [15] consider the noisy setting and analyze the modified non-convex program involving ℓ_1 penalty for the coefficient matrix and ℓ_2 penalty for the loss in fitting the samples, and establish that the true solution is in the neighborhood of a local optimum of the modified non-convex program w.h.p. when the number of samples scales as $n = \Omega(\|A\|_2^2 r^3 d s^2)$. In our setting, this reduces to $n = \Omega(r^4 s^2)$. There are significant differences of the above works from ours. While these works establish that (A^*, X^*) is a local optimum of a non-convex program, they do not provide a tractable algorithm to reach this particular solution as opposed to another local optimum. In contrast, we establish guarantees for a simple alternating minimization algorithm and explicitly characterize the ‘‘basin of attraction’’ for the true solution (A^*, X^*) . This provides precise initialization conditions for the alternating minimization to succeed. Moreover, our sample complexity requirements are much weaker and we require only $n = \mathcal{O}(r^2)$ samples for our guarantees to hold.

Alternating minimization for sparse coding: Our analysis in this paper provides a theoretical explanation for the empirical success of alternating minimization, observed in a number of works [18, 10, 2, 20, 35]. These methods are all based on alternating minimization, but differ mostly in how they update the dictionary elements. For instance, Lee et. al. carry out least squares for updating the dictionary [18] similar to the the method of optimal directions [10], while the K-SVD procedure [2], updates the dictionary estimate using a spectral procedure on the residual. However, none of the previous works provide theoretical guarantees on the success of the alternating minimization procedure for sparse coding.

Guaranteed dictionary estimation: Some of the recent works provide theoretical guarantees on the estimation of the true dictionary. Spielman et. al [29] establish exact recovery under ℓ_1

based optimization when the true dictionary A^* is a basis, which rules out the overcomplete setting. Agarwal et. al [1] and Arora et. al [3] propose methods for approximate dictionary estimation in the overcomplete setting. At a high level, both their methods involve a clustering-based approach for finding samples which share a dictionary element, and then using the subset of samples to estimate a dictionary element. Agarwal et. al [1] establish exact recovery of the true solution (A^*, X^*) under a “one-shot” Lasso procedure, when the non-zero coefficients are Bernoulli $\{-1, +1\}$ (or more generally discrete). On the other hand, we assume only mild conditions on the non-zero elements. Arora et. al [3] consider an alternating minimization procedure. However, a key distinction is that their analysis requires *fresh* samples in each iteration, while we consider the same samples for all the iterations. We show *exact* recovery using $n = \Omega(r^2)$ samples, while [3] can only establish that the error is bounded by $\exp[-O(n/r^2)]$. Furthermore, both the above papers [3, 1] assume that the dictionary elements are mutually incoherent. Our local convergence result in this paper assumes only that the dictionary matrix satisfies RIP (which is strictly weaker than incoherence). For the case of incoherent dictionaries, we can employ the procedures of Agarwal et. al [1] or Arora et. al [3] for initializing the alternating procedure and obtain overall guarantees in such scenarios.

Other works on sparse coding: Some of the other recent works are only tangentially related to this paper. For instance, the works [32, 22, 21, 30] provide generalization bounds for predictive sparse coding, without computational considerations, which differs from our generative setting here and algorithmic considerations. Parametric dictionary learning is considered in [34], where the data is fitted to dictionaries with small coherence. Note that we provide guarantees when the underlying dictionary is incoherent, but do not constrain our method to produce an incoherent dictionary. The problem of sparse coding is also closely related to the problem of blind source separation, and we refer the reader to [1] for an extended survey of these works.

Majorization-minimization algorithms for biconvex optimization: Beyond the specific problem of sparse coding, alternating optimization procedures more generally are a natural fit for biconvex optimization problems, where the objective is individually convex in two sets of variables but not jointly convex. Perhaps the most general study of these problems has been carried out in the framework of majorization-minimization schemes [17], or under the name of the EM algorithm in statistics literature. In this generality, the strongest result one can typically provide is a convergence guarantee to a local optimum of the problem. When the bi-convex objective is defined over probability measures, Csiszar presents a fairly general set of conditions on the objective function, under which linear convergence to the global optimum is guaranteed (see, e.g. the recent tutorial [8] for an excellent overview). However, these conditions do not seem to easily hold in the context of dictionary learning. Alternating optimization in related contexts has also been studied in a variety of matrix factorization problems such as low-rank matrix completion and non-negative matrix factorization. Perhaps the most related to our work are similar results for low-rank matrix completion problems by Jain et al. [14].

Notation: Let $[n] := \{1, 2, \dots, n\}$. For a vector v or a matrix W , we will use the shorthand $\text{Supp}(v)$ and $\text{Supp}(W)$ to denote the set of non-zero entries of v and W respectively. $\|w\|_p$ denote the ℓ_p norm of vector w ; by default, $\|w\|$ denotes ℓ_2 norm of w . $\|W\|_2$ denotes the spectral norm (largest singular value) of matrix W . $\|W\|_\infty$ denotes the largest element (in magnitude) of W . For a matrix X , X^i , X_i and X_j^i denote the i^{th} row, i^{th} column and $(i, j)^{\text{th}}$ element of X respectively.

Algorithm 1 AltMinDict($Y, A(0), \epsilon_0$): Alternating minimization for dictionary learning

Input: Samples Y , initial dictionary estimate $A(0)$, accuracy sequence ϵ_t and sparsity level s .

Thresholding function $T_\rho(a) = a$ if $|a| > \rho$ and 0 o.w.

- 1: **for** iterations $t = 0, 1, 2, \dots, T - 1$ **do**
- 2: **for** samples $i = 1, 2, \dots, n$ **do**
- 3: $X(t + 1)_i = \arg \min_{x \in \mathbb{R}^r} \|x\|_1$ such that, $\|Y_i - A(t)x\|_2 \leq \epsilon_t$
- 4: **end for**
- 5: Threshold: $X(t + 1) = X(t + 1) \cdot *(\mathbb{I}[X(t + 1) > 9s\epsilon_t])$
- 6: Estimate $A(t + 1) = YX(t + 1)^+$
- 7: Normalize: $A(t + 1)_i = \frac{A(t+1)_i}{\|A(t+1)_i\|_2}$
- 8: **end for**

Output: $A(T)$

2 Algorithm

Given an initial estimate of the dictionary, we alternate between two procedures, viz., a sparse recovery step for estimating the coefficients given a dictionary, and a least squares step for a dictionary given the estimates of the coefficients. The details of this approach are presented in Algorithm 1.

The sparse recovery step of Algorithm 1 is based on ℓ_1 -regularization, followed by thresholding. The thresholding is required for us to guarantee that the support set of our coefficient estimate $X(t)$ is a *subset* of the true support with high probability. Once we have an estimate of the coefficients, the dictionary is re-estimated through least squares. The overall algorithmic scheme is popular for dictionary learning, and there are a number of variants of the basic method. For instance, the ℓ_1 -regularized problem in step 3 can also be replaced by other robust sparse recovery procedures such as OMP [31] or GraDeS [11]. More generally the exact lasso and least-squares steps may be replaced with other optimization methods for computational efficiency, e.g. [16].

3 Main results and their proofs

In this section, we provide our local convergence result for alternating minimization and also clearly specify all the required assumptions on A^* and X^* . We provide a brief sketch of our proof for each of the steps in Section 3.4.

3.1 Assumptions

We start by formally describing the assumptions needed for the main recovery result of this paper. Without loss of generality, assume that all the elements are normalized: $\|A_i^*\|_2 = 1$, for $i \in [r]$. This is because we can always rescale the dictionary elements and the corresponding coefficients and obtain the same observations.

Assumptions:

- (A1) **Dictionary Matrix satisfying RIP:** The dictionary matrix A^* has a $2s$ -RIP constant of $\delta_{2s} < 0.1$.

- (A2) **Spectral Condition on Dictionary Elements:** The dictionary matrix has bounded spectral norm, for some constant $\mu_1 > 0$, $\|A^*\|_2 < \mu_1 \sqrt{\frac{r}{d}}$.
- (A3) **Non-zero Entries in Coefficient Matrix:** We assume that the non-zero entries of X^* are drawn i.i.d. from a distribution such that $\mathbb{E} \left[(X^{*i}_j)^2 \right] = 1$, and satisfy the following a.s.: $|X^{*i}_j| \leq M, \forall i, j$.
- (A4) **Sparse Coefficient Matrix:** The columns of coefficient matrix have s non-zero entries which are selected uniformly at random from the set of all s -sized subsets of $[r]$, i.e. $|\text{Supp}(X_i^*)| = s, \forall i \in [n]$. We require s to satisfy $s < \frac{d^{1/6}}{c_2 \mu_1^{1/3}}$, for some universal constant c_2 .
- (A5) **Sample Complexity:** For some universal constant $c > 0$ and a given failure parameter $\delta > 0$, the number of samples n needs to satisfy

$$n \geq c_3 \max \left(r^2, rM^2s \right) \log \frac{2r}{\delta}, \quad ,$$

where $c_3 > 0$ is a universal constant.

- (A6) **Initial dictionary with guaranteed error bound:** We assume that we have access to an initial dictionary estimate $A(0)$ such that

$$\hat{\epsilon}_0 := \max_{i \in [r]} \min_{z \in \{-1, +1\}} \|zA_i(0) - A_i^*\|_2 < \frac{1}{2592s^2}.$$

- (A7) **Choice of Parameters for Alternating Minimization:** Algorithm 1 uses a sequence of accuracy parameters $\epsilon_0 = 1/2592s^2$ and

$$\epsilon_{t+1} = \frac{25050\mu_1 s^3}{\sqrt{d}} \epsilon_t. \quad (3)$$

Assumption (A1) regarding the RIP assumption is crucial in establishing our guarantees, since it is critical for analyzing the performance of the compressed sensing subroutine in Algorithm 1 (steps 2-5). It is possible to further weaken this assumption to a Restricted Eigenvalue condition which is often used in the sparse regression literature as well [28, 23]. We will present a more detailed discussion of this condition in the proof sketch. In order to keep the results with cleaner constants, we will continue with the RIP assumption for the rest of the analysis, while mentioning how the result can be extended easily under a more general restricted eigenvalue assumption.

The assumption (A2) provides a bound on the spectral norm of A^* . Note that the RIP and spectral assumptions are satisfied with high probability (w.h.p.) when the dictionary elements are randomly drawn from a mean-zero sub-gaussian distribution.

Assumption (A3) imposes some natural constraints on the non-zero entries of X^* . Assumption(A4) on sparsity in the coefficient matrix is crucial for identifiability of the dictionary learning problem.

Assumption (A5) provides a bound on sample complexity. Assumption (A6) specifies the accuracy of the initial estimate required by Algorithm 1. Recent works [4, 1] provide provable ways of obtaining such an estimate. Please see Section 3.3 for more details.

Assumption (A7) specifies the choice of accuracy parameters used by alternating method in Algorithm 1. Due to Assumption (A4) on sparsity level s , we have that $\frac{25050\mu_1 s^3}{\sqrt{d}} < 1/2$ and the

accuracy parameters in (3) form a decreasing sequence. This implies that in Algorithm 1, the accuracy constraint becomes more stringent with the iterations of the alternating method.

3.2 Guarantees for Alternating Minimization

We now prove a local convergence result for alternating minimization. We assume that we have access to a good initial estimate of the dictionary:

Theorem 3.1 (Local linear convergence). *Under assumptions (A1)-(A7), with probability at least $1 - 2\delta$ the iterate $A(t)$ of Algorithm 1 satisfies the following for all $t \geq 1$:*

$$\min_{z \in \{-1,1\}} \|zA_i(t) - A_i^*\|_2 \leq \sqrt{2}\epsilon_t, 1 \leq i \leq r.$$

Remarks: Note that we have a sign ambiguity in recovery of the dictionary elements, since we can exchange the signs of the dictionary elements and the coefficients to obtain the same observations.

Theorem 3.1 guarantees that we can recover the dictionary A^* to an arbitrary precision ϵ (based on the number of iterations T of Algorithm 1), given $n = \mathcal{O}(r^2)$ samples. We contrast this with the results of [4], who also provide recovery guarantees to an arbitrary accuracy ϵ , but only if the number of samples is allowed to increase as $\mathcal{O}(r^2 \log 1/\epsilon)$.

The consequences of Theorem 3.1 are powerful combined with our Assumption (A4) and the recurrence 3 (since (A4) ensures that ϵ_t forms a decreasing sequence). In particular, it is implied that with high probability we obtain,

$$\min_{z \in \{-1,1\}} \|zA_i(t) - A_i^*\|_2 \leq \hat{\epsilon}_0 2^{-t}.$$

Given the above bound, we need at most $\mathcal{O}\left(\log_2 \frac{\hat{\epsilon}_0}{\epsilon}\right)$ in order to ensure $\|zA_i(T) - A_i^*\|_2 \leq \epsilon$ for all the dictionary elements $i = 1, 2, \dots, r$. In the convex optimization parlance, the result demonstrates a local linear convergence of Algorithm 1 to the globally optimal solution under an initialization condition. Another way of interpreting our result is that the global optimum has a *basin of attraction* of size $\mathcal{O}(1/s^2)$ for our alternating minimization procedure under these assumptions (since we require $\hat{\epsilon}_0 \leq \mathcal{O}(1/s^2)$).

We also recall that the lasso step in Algorithm 1 can be replaced with a different robust sparse recovery procedure, with qualitatively similar results.

3.3 Using Local Convergence for Complete Recovery

In the above section, we showed a local convergence result for Algorithm 1. In particular, Assumption (A6) requires that the initial dictionary estimate be at most $\mathcal{O}\left(\frac{1}{s^2}\right)$ away from A^* . In this section, we use the recent result of [1] to obtain an initialization which satisfies Assumption (A6), and thus, we obtain a full recovery result for the sparsely-used dictionary problem with assumptions only on the model parameters. In order to obtain the initialization from the method of [1], we require the following assumptions:

(B1) **Incoherent Dictionary Elements:** Without loss of generality, assume that all the elements are normalized: $\|A_i^*\|_2 = 1$, for $i \in [r]$. We assume pairwise incoherence condition on the dictionary elements, for some constant $\mu_0 > 0$, $|\langle A_i^*, A_j^* \rangle| < \frac{\mu_0}{\sqrt{d}}$.

(B3) **Non-zero Entries in Coefficient Matrix:** We assume that the non-zero entries of X^* are drawn i.i.d. from a distribution such that $\mathbb{E} \left[(X^{*i}_j)^2 \right] = 1$, and satisfy the following a.s.: $m \leq |X^{*i}_j| \leq M, \forall i, j$.

(B4) **Sparse Coefficient Matrix:** The columns of coefficient matrix have bounded number of non-zero entries s which are selected randomly, i.e.

$$|\text{Supp}(x_i)| = s, \quad \forall i \in [n]. \quad (4)$$

We require s to be

$$s < c_1 \min \left(\frac{m d^{1/4}}{M \sqrt{\mu_0}}, \left(\frac{d m^4}{\mu_1^2 M^4} \right)^{1/9}, r^{1/8} \left(\frac{m}{M} \right)^{1/4} \right),$$

for universal constants $c_1, c_2 > 0$. Constants m, M are as specified above.

(B5) **Sample Complexity:** Given universal constant $c_2 > 0$, choose $\delta > 0$ and the number of samples n such that

$$n := n(d, r, s, \delta) \geq c_2 r^2 \frac{M^2}{m^2} \log \frac{2r}{\delta}.$$

Theorem 3.2 (Specialization of Theorem 2.1 from [1]). *Under assumptions (B1), (A2), (B3) – (B5) and (A7), there exists an algorithm which given Y outputs $A(0)$, such that Assumption (A6) holds with probability greater than $1 - 2n^2\delta$.*

The restatement follows by setting $\alpha = s^{-9/2} \frac{m^2}{M^2}$ in that result which ensures that the error in the initialization is at most $1/s^2$ as required by Assumption (A6). Combining the above theorem with Theorem 3.1 gives the following powerful corollary.

Corollary 3.1 (Exact recovery). *Suppose assumptions (B1), (A2) – (A5), (B3) – (B5) and (A7) hold. If we start Algorithm 1 with the output of Algorithm 1 of [1], then the following holds for all $t \geq 1$:*

$$\min_{z \in \{-1, 1\}} \|z A_i(t) - A_i^*\|_2 \leq \sqrt{2}\epsilon_t, 1 \leq i \leq r.$$

The above result makes use of Lemma A.7 in the appendix which shows that Assumptions (B1) and (B4) imply (A1). Note that the above corollary gives an exact recovery result with the only assumptions being those on the model parameters. We also note that the conclusion of Corollary 3.1 does not crucially rely on initialization specifically by the output of Algorithm 1 of [1], and admits any other initialization satisfying Assumption (A6). As remarked earlier, the recent work of [4] provides an alternative initialization strategy for our alternating minimization procedure. Indeed, under our sample complexity assumption, their OVERLAPPINGAVERAGE method provides a solution with $\hat{\epsilon}_0 = \mathcal{O}(s/\sqrt{r})$ assuming $s = \mathcal{O}(\max(r^{2/5}, \sqrt{d}))$. In particular, if $s = \mathcal{O}(r^{1/6})$, we obtain the desired initial error of $1/s^2$ using that algorithm. The sample complexity of the entire procedure remains identical to that in Assumption (A5).

3.4 Overview of Proof

In this section we outline the key steps in proving Theorem 3.1.

For ease of notation, let us consider just one iteration of Algorithm 1 and denote $X(t+1)$ as X , $A(t+1)$ as A and $A(t)$ as \tilde{A} . Then we have the least-squares update:

$$\begin{aligned} A - A^* &= YX^+ - A^* \\ &= A^*X^*X^+ - A^*XX^+ = A^*\Delta XX^+, \end{aligned}$$

where $\Delta X = X^* - X$. This means that we can understand the error in dictionary recovery by the error in the least squares operator ΔXX^+ . In particular, we can further expand the error in a column p as:

$$A_p - A^*_p = A^*_p(\Delta XX^+)_p^p + A^*_{\setminus p}(\Delta XX^+)_p^{\setminus p},$$

where the notation $\setminus p$ represents the collection of all indices apart from p , i.e., $A^*_{\setminus p}$ denote all the columns of A except the p -th column and $(\Delta XX^+)_p^{\setminus p}$ denotes the off-diagonal elements of the p -th column of (ΔXX^+) .

The above equation indicates that there are two sources of error in our dictionary estimate. The element $(\Delta XX^+)_p^p$ causes the rescaling of A_p relative to A^*_p . However, this is a minor issue since the renormalization would correct it.

More serious is the contribution from the off-diagonal terms $(\Delta XX^+)_p^{\setminus p}$, which corrupt our estimate A_p with other dictionary elements beyond A^*_p . Indeed, a crucial argument in our proof is controlling the contribution of these terms at an appropriately small level. In order to do that, we start by controlling the magnitude of ΔX .

Lemma 3.1 (Error in sparse recovery). *Let $\Delta X := X(t) - X^*$. Assume that $2\mu_0s/\sqrt{d} \leq 0.1$ and $\sqrt{s\epsilon_t} \leq 0.1$. Then, we have $\text{Supp}(\Delta X) \subseteq \text{Supp}(X^*)$ and the error bound $\|\Delta X\|_\infty \leq 9s\epsilon_t$.*

More general Restricted Eigenvalue conditions: The above lemma is the only part of our proof, where we require the RIP assumption. This is in order to invoke the result of Candes [7] regarding the error in compressed sensing with (bounded) deterministic noise. Such results can also be typically established under weaker Restricted Eigenvalue assumptions (henceforth RE). Given a vector $v \in \mathbb{R}^r$, these RE conditions study the norms $\|Av\|_2^2$. A particular form of RE condition for these approximately sparse vectors then posits (see e.g. [28, 23])

$$\|Av\|_2 \geq \gamma\|v\|_2 - \tau\|v\|_1. \quad (5)$$

Under such a condition, it can be readily shown that Lemma 3.1 continues to hold with an error bound which is $\mathcal{O}(s\epsilon_t/(\gamma - s\tau))$. For many random matrix ensembles for A , it is known that $\tau = \mathcal{O}\left(\sqrt{(\log r)/d}\right)$, which means that $\gamma - s\tau$ will be bounded away from zero under our assumptions on the sparsity level s . These RE conditions are the weakest known conditions under which compressed sensing using efficient procedures is possible, and we employ this as a subroutine in our alternating minimization procedure.

The next lemma is very useful in our error analysis, since we establish that any matrix W satisfying $\text{Supp}(W) \subseteq \text{Supp}(X^*)$ has a good bound on its spectral norm (even if the entries depend on A^*, X^*).

Lemma 3.2. *With probability at least $1 - r \exp(-\frac{Cn}{rs})$, for every $r \times n$ matrix W s.t. $\text{Supp}(W) \subseteq \text{Supp}(X^*)$, we have*

$$\|W\|_2 \leq 2\|W\|_\infty \sqrt{\frac{s^2 n}{r}}.$$

A particular consequence of this lemma is that it guarantees the invertibility of the matrix XX^\top , so that the pseudo-inverse X^+ is well-defined for subsequent least squares updates. Next, we present the most crucial step which is controlling the off-diagonal terms $(\Delta XX^+)_p^p$.

Lemma 3.3 (Off-diagonal error bound). *With probability at least $1 - r \exp(-\frac{Cn}{rM^2s}) - r \exp(-Cn/r^2)$, we have uniformly for every $p \in [r]$ and every ΔX such that $\|\Delta X\|_\infty < \frac{1}{288s}$.*

$$\left\| (\Delta XX^+)_p^p \right\|_2 = \left\| (X^* X^+)_p^p \right\|_2 \leq \frac{1968s^2 \|\Delta X\|_\infty}{\sqrt{r}}.$$

The lemma uses the earlier two lemmas along with a few other auxilliary results. Given these lemmas, the proof of the main theorem follows using basic linear algebra arguments. Specifically, for any unit vector w such that $w \perp A_p^*$, we can bound the normalized inner product $\langle w, A_p \rangle / \|A_p\|_2$ which suffices to obtain the result of the theorem.

3.5 Detailed Proof of Theorem 3.1

We now provide a proof of the theorem using the above given lemmas. The proofs of the lemmas are deferred to the appendix. Recall that, we denote $A(t)$ as \tilde{A} and $A(t+1)$ as A . Similarly we denote $X(t)$ and $X(t+1)$ as \tilde{X} and X respectively. Then the goal is to show that A is closer to A^* than \tilde{A} . For the purposes of our analysis, we will find it more convenient to directly work with dot products instead of ℓ_2 -distances (and hence avoid sign ambiguities). With this motivation, we define the following notion of distance between two vectors.

Definition 1. For any two vectors $z, w \in \mathbb{R}^d$, we define the distance between them as follows:

$$\text{dist}(z, w) := \sup_{v \perp w} \frac{\langle v, z \rangle}{\|v\|_2 \|z\|_2} = \sup_{v \perp z} \frac{\langle v, w \rangle}{\|v\|_2 \|w\|_2}.$$

This definition of distance suffices for our purposes due to the following simple lemma.

Lemma 3.4. *For any two unit vectors $u, v \in \mathbb{R}^d$, we have*

$$\text{dist}(u, v) \leq \min_{z \in \{-1, 1\}} \|zu - v\|_2 \leq \sqrt{2} \text{dist}(u, v).$$

Proof: The proof is rather straightforward. Suppose that $\langle u, v \rangle > 0$ so that the minimum happens at $z = 1$. The other case is identical. We can easily rewrite

$$\|u - v\|_2^2 = (2 - 2\langle u, v \rangle) \leq 2(1 - \langle u, v \rangle^2),$$

where the final inequality follows since $0 \leq \langle u, v \rangle \leq 1$. Writing $u = \langle u, v \rangle v + v_\perp$, where $\langle v_\perp, v \rangle = 0$, we see that:

$$1 = \|u\|_2^2 = \langle u, v \rangle^2 + \|v_\perp\|_2^2 = \langle u, v \rangle^2 + \text{dist}(u, v)^2.$$

Substituting this into our earlier bound, we obtain the upper bound. For the lower bound, we note that,

$$\begin{aligned} \text{dist}(u, v)^2 &= 1 - \left(1 - \frac{\|u - v\|^2}{2}\right)^2 \\ &\leq \|u - v\|^2. \end{aligned}$$

□

The distance is naturally extended to matrices for our purposes by applying it columnwise.

Definition 2. For any two $d \times r$ matrices Z and W , we define the distance between them as follows:

$$\text{dist}(Z, W) := \sup_{p \in [r]} \text{dist}(Z_p, W_p).$$

Note that the normalization in the definition of $\text{dist}(z, w)$ ensures that we can apply the distance directly to the result of the least-squares step without worrying about the effects of normalization. This allows us to work with a closed-form expression for A :

$$A = YX^+ = A^*X^*X^+. \quad (6)$$

We are now in a position to prove Theorem 3.1.

Proof of Theorem 3.1: As an induction hypothesis, we have $\text{dist}(\tilde{A}, A^*) < \epsilon_t$, where we recall the definition (3). We will show that for every $p \in [r]$, we will have:

$$\text{dist}(A_p, A^*_p) \leq \epsilon_{t+1} < \frac{23616\mu_1 s^3}{\sqrt{d}} \epsilon_t. \quad (7)$$

This suffices to prove the theorem by appealing to Lemma 3.4.

Now fix any $w \perp A^*_p$ such that $\|w\|_2 = 1$. We first provide a bound on $\langle w, A_p \rangle$. Now, the following holds with high probability:

$$\begin{aligned} \langle w, A_p \rangle &= w^\top A^* X^* X^+_p \stackrel{(\zeta_1)}{\leq} \left\| w^\top A^* \right\|_2 \left\| (X^* X^+)_p^{\setminus p} \right\|_2 \\ &\stackrel{(\zeta_2)}{\leq} \mu_1 \sqrt{\frac{r}{d}} \cdot \frac{1968s^2 \|\Delta X\|_\infty}{\sqrt{r}} \\ &= \frac{17712\mu_1 s^3}{\sqrt{d}} \epsilon_t, \end{aligned} \quad (8)$$

where (ζ_1) follows from the fact that $w^\top A^*_p = 0$ and (ζ_2) follows from Assumption (A2) and Lemma 3.3.

In order to bound $\text{dist}(A, A^*)$, it remains to show a lower bound on $\|A\|_2$. This again follows using basic algebra, given our main lemmas.

$$\begin{aligned} \|A_p\|_2 &= \|A^* X^* X^+_p\|_2 = \|A^* (X - \Delta X) X^+_p\|_2 \\ &\stackrel{(\zeta_1)}{=} \|A^*_p - A^* \Delta X X^+_p\|_2 \\ &\geq \|A^*_p\|_2 - \left\| A^* (\Delta X X^+)_p \right\|_2, \end{aligned}$$

where (ζ_1) follows from the fact that $XX^+ = \mathbb{I}$. We decompose the second term into diagonal and off-diagonal terms of ΔXX^+ , followed by triangle inequality and obtain:

$$\begin{aligned}
\|A_p\|_2 &\geq 1 - \left\| A^*_p (\Delta XX^+)_p^p + A^*_{\setminus p} (\Delta XX^+)_p^{\setminus p} \right\|_2 \\
&\geq 1 - \|A^*_p\|_2 \left\| (\Delta XX^+)_p^p \right\|_2 - \|A^*_{\setminus p}\|_2 \left\| (\Delta XX^+)_p^{\setminus p} \right\|_2 \\
&\geq 1 - 1 \cdot \left\| \Delta XX^T (XX^T)^{-1} \right\|_2 - \|A^*_{\setminus p}\|_2 \left\| (\Delta XX^+)_p^{\setminus p} \right\|_2 \\
&\geq 1 - \underbrace{\| \Delta X \|_2 \| X^T \|_2 \left\| (XX^T)^{-1} \right\|_2}_{\mathcal{T}_1} - \underbrace{\|A^*_{\setminus p}\|_2 \left\| (\Delta XX^+)_p^{\setminus p} \right\|_2}_{\mathcal{T}_2}
\end{aligned}$$

It remains to control \mathcal{T}_1 and \mathcal{T}_2 at an appropriate level. We start from \mathcal{T}_1 . Note that $\|\Delta X\|_2$ is bounded by Lemmas 3.1 and 3.2, while $\|X^T\|_2$ is controlled by Lemma A.3 in the appendix (recall $\|\Delta X\|_\infty \leq 1/(64s)$). Invoking Lemma A.4 to control $\left\| (XX^T)^{-1} \right\|_2$, we obtain the following bound on \mathcal{T}_1 with probability at least $1 - r \exp\left(-\frac{Cn}{rM^2s}\right)$:

$$\mathcal{T}_1 \leq 18\epsilon_t s^2 \sqrt{\frac{n}{r}} \cdot 3s \sqrt{\frac{n}{r}} \cdot \frac{8r}{sn} = 432s^2 \epsilon_t.$$

The second term \mathcal{T}_2 is directly controlled by Lemma 3.3, yielding the following (with probability at least $1 - r \exp\left(-\frac{Cn}{rM^2s}\right) - \exp(-Cn/r^2)$):

$$\mathcal{T}_2 \leq \mu_1 \sqrt{\frac{r}{d}} \frac{1968s^3 \epsilon_t}{\sqrt{r}}.$$

Putting all the terms together, we get:

$$\|A_p\|_2 \geq 1 - 9s^2 \left(48 + \frac{1968s\mu_1}{\sqrt{d}} \right) \epsilon_t \geq \frac{3}{4}, \quad (9)$$

where the inequality follows since $9s^2 \left(48 + \frac{1968s\mu_1}{\sqrt{d}} \right) \epsilon_t \leq 9s^2 \left(48 + \frac{1968s\mu_1}{\sqrt{d}} \right) \epsilon_0 \leq 1/4$ by our assumption (A6) on ϵ_0 . Combining the bounds (8) and (9) yields the desired recursion (7). Appealing to Lemma 3.4 along with our setting of ϵ_t (3) completes the proof of the claim (7). Finally note that the error probability in the theorem is obtained by using the fact that $M \geq 1$, and that the failure probability is purely incurred from the structure of the non-zero entries of X^* , so that it is incurred only once and not at each round. This avoids the need of a union bound over all the rounds, yielding the result. \square

4 Experiments

Alternating minimization/descent approaches have been widely used for dictionary learning and several existing works show effectiveness of these methods on real-world/synthetic datasets [5, 30]. Hence, instead of replicating those results, in this section we focus on illustrating the following three key properties of our algorithms via experiments in a controlled setting: a) advantage of alternating

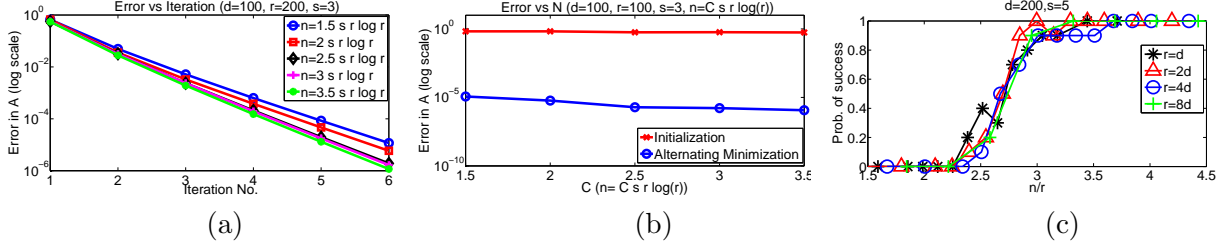


Figure 1: (a): Average error after each alternating minimization step of Algorithm 1 on log-scale. (b): Average error after the initialization procedure (Algorithm 1 of [1]) and after 5 alternating minimization steps of Algorithm 1. (c): Sample complexity requirement of the alternating minimization algorithm. For ease of experiments, we initialize the dictionary using a random perturbation of the true dictionary rather than using Algorithm 1 of [1] which should in fact give better initial point with smaller error.

minimization over one-shot initialization, b) linear convergence of alternating minimization, c) sample complexity of alternating minimization.

Data generation model: Each entry of the dictionary matrix A is chosen i.i.d. from $\mathcal{N}(0, 1/\sqrt{d})$. Note that, random Gaussian matrices are known to satisfy incoherence and the spectral norm bound [33]. The support of each column of X was chosen independently and uniformly from the set of all s -subsets of $[r]$. Similarly, each non-zero element of X was chosen independently from the uniform distribution on $[-2, -1] \cup [1, 2]$. We use the GraDeS algorithm of [11] to solve the sparse recovery step, as it is faster than lasso. We measure error in the recovery of dictionary by $error(A) = \max_i \sqrt{1 - \frac{\langle A_i, A_i^* \rangle^2}{\|A_i\|_2^2 \|A_i^*\|_2^2}}$. The first two plots are for a typical run and the third plot averages over 10 runs. The implementation is in Matlab.

Linear convergence: In the first set of experiments, we fixed $d = 100$, $r = 200$ and measured error after each step of our algorithm for increasing values of n . Figure 1 (a) plots error observed after each iteration of alternating minimization; the first data point refers to the error incurred by the initialization method (Algorithm 1 of [1]). As expected due to Theorem 3.1, we observe a geometric decay in the error.

One-shot vs iterative algorithm: It is conceivable that a good initialization procedure itself is sufficient to obtain an estimate of the dictionary upto reasonable accuracy, without recourse to the alternating minimization procedure of Algorithm 1. Figure 1(b) shows that this is not the case. The figure plots the error in recovery vs the number of samples used for both Algorithm 1 of Agarwal et al. [1] and Algorithm 1. It is clear that the recovery error of the alternating minimization procedure is significantly smaller than that of the initialization procedure. For example, for $n = 2.5sr \log r$ with $s = 3, r = 200, d = 100$, initialization incurs error of .56 while alternating minimization incurs error of 10^{-6} . Note however that the recovery accuracy of the initialization procedure is non-trivial and also crucial to the success of alternating minimization- a random vector in \mathbb{R}^d would give an error of $1 - \frac{1}{d} = 0.99$ (since the inner product is concentrated around $1/\sqrt{d}$), where as the error after initialization procedure is ≈ 0.55 .

Sample complexity: Finally, we study sample complexity requirement of the alternating minimization algorithm which is $n = \mathcal{O}(r^2 \log r)$ according to Theorem 3.1, assuming good enough initialization. Figure 1(c) suggests that in fact only $\mathcal{O}(r)$ samples are sufficient for success of alternating minimization. The figure plots the probability of success with respect to $\frac{n}{r}$ for various

values of r . A trial is said to succeed if at the end of 25 iterations, the error is smaller than 10^{-6} . Since we focus only on the sample complexity of alternating minimization, we use a faster initialization procedure: we initialize the dictionary by randomly perturbing the true dictionary as $A(0) = A^* + Z$, where each element of Z is an $\mathcal{N}(0, 0.5/\sqrt{d})$ random variable. Figure 1 (c) shows that the success probability transitions at nearly the same value for various values of r , suggesting that the sample complexity of the alternating minimization procedure in this regime of $r = \mathcal{O}(d)$ is just $O(r)$.

5 Conclusions

In this paper we provide the first analysis for the local linear convergence of the popular alternating minimization heuristic commonly used for solving dictionary learning problems in practice. Combined with some recent results, this also provides an efficient method for global and exact recovery of the unknown overcomplete dictionary under favorable assumptions. The results are of interest from both theoretical and practical standpoints. From a theoretical standpoint, this is one of the very few results that provides guarantees on a dictionary learned using an efficient algorithm, and one of the first for the overcomplete setting. From a practical standpoint, there is a tremendous interest in the problem, and we believe that an understanding of the theoretical properties of existing methods is critical in designing better methods. Indeed, our work provides some such hints towards designing a better algorithm. For instance, the sparse recovery step in our method decodes the coefficients individually for each sample. We believe that a better method can be designed by jointly decoding all the samples, which allows one to force consistency across samples (for instance, in our random coefficient model, the number of samples per dictionary element is also controlled in addition to the number of dictionary elements per sample).

More interestingly, our work extends a growing body of recent literature on analysis of alternating minimization methods for a variety of non-convex factorization problems [14, 24], where global in addition to local results are being established with appropriate initialization strategies. Of course, results on alternating minimization go much further back, even in non-convex optimization to Csiszar’s seminal works (see, e.g. the recent tutorial [8] for an overview), as well as in convex minimization and projection problems. However, the recent work has been largely motivated by applications of non-convex optimization arising in machine learning. We believe that the emergence of these newer results indicates the possibility of a more general theory of alternating optimization procedures for a broad class of factorization-style non-convex problems, and should be an exciting question for future research

References

- [1] A. Agarwal, A. Anandkumar, and P. Netrapalli. Exact recovery of sparsely used overcomplete dictionaries. *arXiv preprint arXiv:1309.1952*, 2013.
- [2] M. Aharon, M. Elad, and A. Bruckstein. K-svd: An algorithm for designing overcomplete dictionaries for sparse representation. *Signal Processing, IEEE Transactions on*, 54(11):4311–4322, 2006.
- [3] S. Arora, R. Ge, Y. Halpern, D. M. Mimno, A. Moitra, D. Sontag, Y. Wu, and M. Zhu. A practical algorithm for topic modeling with provable guarantees. *ArXiv 1212.4777*, 2012.

- [4] S. Arora, R. Ge, and A. Moitra. New Algorithms for Learning Incoherent and Overcomplete Dictionaries. *ArXiv e-prints*, Aug. 2013.
- [5] K. Balasubramanian, K. Yu, and G. Lebanon. Smooth sparse coding via marginal regression for learning sparse representations. In *ICML*, 2013.
- [6] Y. Bengio, A. Courville, and P. Vincent. Unsupervised feature learning and deep learning: A review and new perspectives. *arXiv preprint arXiv:1206.5538*, 2012.
- [7] E. J. Candes. The restricted isometry property and its implications for compressed sensing. *Comptes Rendus Mathematique*, 346(910):589 – 592, 2008.
- [8] I. Csiszar and P. Shields. Information theory and statistics: A tutorial. *Foundations and Trends in Communications and Information Theory*, 1(4):417–528, 2004.
- [9] G. Davis. *Adaptive nonlinear approximations*. PhD thesis, New York University, 1994.
- [10] K. Engan, S. O. Aase, and J. Hakon Husoy. Method of optimal directions for frame design. In *Acoustics, Speech, and Signal Processing, 1999. Proceedings., 1999 IEEE International Conference on*, volume 5, pages 2443–2446. IEEE, 1999.
- [11] R. Garg and R. Khandekar. Gradient descent with sparsification: an iterative algorithm for sparse recovery with restricted isometry property. In *ICML*, 2009.
- [12] Q. Geng, H. Wang, and J. Wright. On the local correctness of ℓ_1 minimization for dictionary learning. *arXiv preprint arXiv:1101.5672*, 2011. Preprint, URL:<http://arxiv.org/abs/1101.5672>.
- [13] R. Gribonval and K. Schnass. Dictionary identificationsparse matrix-factorization via. *Information Theory, IEEE Transactions on*, 56(7):3523–3539, 2010.
- [14] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*, pages 665–674. ACM, 2013.
- [15] R. Jenatton, R. Gribonval, and F. Bach. Local stability and robustness of sparse dictionary learning in the presence of noise. *arXiv preprint arXiv:1210.0685*, 2012.
- [16] R. Jenatton, J. Mairal, F. R. Bach, and G. R. Obozinski. Proximal methods for sparse hierarchical dictionary learning. In *Proceedings of the 27th International Conference on Machine Learning (ICML-10)*, pages 487–494, 2010.
- [17] K. Lange, D. R. Hunter, and I. Yang. Optimization transfer using surrogate objective functions. *Journal of computational and graphical statistics*, 9(1):1–20, 2000.
- [18] H. Lee, A. Battle, R. Raina, and A. Ng. Efficient sparse coding algorithms. In *Advances in neural information processing systems*, pages 801–808, 2006.
- [19] M. S. Lewicki and T. J. Sejnowski. Learning overcomplete representations. *Neural computation*, 12(2):337–365, 2000.

- [20] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman. Discriminative learned dictionaries for local image analysis. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- [21] A. Maurer, M. Pontil, and B. Romera-Paredes. Sparse coding for multitask and transfer learning. *arXiv preprint arXiv:1209.0738*, 2012.
- [22] N. Mehta and A. G. Gray. Sparsity-based generalization bounds for predictive sparse coding. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 36–44, 2013.
- [23] S. Negahban, P. Ravikumar, M. J. Wainwright, and B. Yu. A unified framework for high-dimensional analysis of M-estimators with decomposable regularizers. *Statistical Science*, 2012. To appear; Original version arxiv:1010.2731v1.
- [24] P. Netrapalli, P. Jain, and S. Sanghavi. Phase retrieval using alternating minimization. In *Advances in Neural Information Processing Systems*, pages 2796–2804, 2013.
- [25] B. A. Olshausen. Sparse coding of time-varying natural images. In *Proc. of the Int. Conf. on Independent Component Analysis and Blind Source Separation*, pages 603–608. Citeseer, 2000.
- [26] B. A. Olshausen et al. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607–609, 1996.
- [27] B. A. Olshausen and D. J. Field. Sparse coding with an overcomplete basis set: A strategy employed by v1? *Vision research*, 37(23):3311–3325, 1997.
- [28] G. Raskutti, M. J. Wainwright, and B. Yu. Restricted eigenvalue properties for correlated gaussian designs. *The Journal of Machine Learning Research*, 11:2241–2259, 2010.
- [29] D. A. Spielman, H. Wang, and J. Wright. Exact recovery of sparsely-used dictionaries. In *Proc. of Conf. on Learning Theory*, 2012.
- [30] J. J. Thiagarajan, K. N. Ramamurthy, and A. Spanias. Learning stable multilevel dictionaries for sparse representation of images. *ArXiv 1303.0448*, 2013.
- [31] J. Tropp and A. Gilbert. Signal recovery from random measurements via orthogonal matching pursuit. *Information Theory, IEEE Transactions on*, 53(12):4655–4666, 2007.
- [32] D. Vainsencher, S. Mannor, and A. M. Bruckstein. The sample complexity of dictionary learning. *The Journal of Machine Learning Research*, 12:3259–3281, 2011.
- [33] R. Vershynin. Introduction to the non-asymptotic analysis of random matrices. *arXiv preprint arXiv:1011.3027*, 2010.
- [34] M. Yaghoobi, L. Daudet, and M. E. Davies. Parametric dictionary design for sparse coding. *Signal Processing, IEEE Transactions on*, 57(12):4800–4810, 2009.
- [35] J. Yang, J. Wright, T. S. Huang, and Y. Ma. Image super-resolution via sparse representation. *Image Processing, IEEE Transactions on*, 19(11):2861–2873, 2010.

A Proofs for alternating minimization

In this section, we will present our proof for the results on alternating minimization. We present the proofs for Theorem 3.1 and the other main lemmas in Section A.1. In Section A.2, we present the auxiliary lemmas and their proofs.

A.1 Proofs of main lemmas

For reader's convenience, we recall Lemmas 3.1, 3.2 and 3.3 from Section 3.4 along with their proofs. The more technical lemmas are deferred to the next section.

We first recall some notation and define additional abbreviations before proving the lemmas. Denote $X_i^{*p} = \chi_i^p M_i^p$, $\forall 1 \leq p \leq r$, $\forall 1 \leq i \leq n$ where $\chi_i^p = 1$ if $p \in \text{Supp}(X_i^*)$ and 0 otherwise and M_i^p are i.i.d. random variables with $\mathbb{E}[M_i^p] = \mu$ and $\mathbb{E}[(M_i^p)^2] = \sigma^2 + \mu^2$. Assumption (A3) gives us:

1. $\mu^2 + \sigma^2 = 1$, and
2. $|M_i^p| \leq M$ a.s.

Lemma 3.1 (Error in sparse recovery). Let $\Delta X := \tilde{X} - X^*$. Assume that $2\mu_0 s/\sqrt{d} \leq 0.1$ and $\sqrt{s\epsilon_t} \leq 0.1$. Then, we have:

1. $\text{Supp}(\Delta X) \subseteq \text{Supp}(X^*)$.
2. $\|\Delta X\|_\infty \leq 9s \cdot \text{dist}(\tilde{A}, A^*) \leq 9s\epsilon_t$.

Proof: In order to establish the lemma, we use a result of Candes regarding the lasso estimator with deterministic noise for the recovery procedure:

$$\hat{x}_i = \arg \min_{x \in \mathbb{R}^r} \|x\|_1 \quad \text{such that,} \quad \|Y_i - Ax\|_2 \leq \epsilon. \quad (10)$$

Theorem A.1 (Theorem 1.2 from [7]). Suppose $Y_i = Ax_i + z_i$, where x_i is s -sparse and $\|z_i\|_2 \leq \epsilon$. Assume further that $\delta_{2s} \leq \sqrt{2} - 1$. Then the solution to Equation (10) obeys the following, for a universal constant C_1 ,

$$\|\hat{x}_i - x_i\|_2 \leq C_1 \epsilon$$

In particular, $C_1 = 8.5$ suffices for $\delta_{2s} \leq 0.2$.

In order to apply the theorem, we need to demonstrate that the RIP condition holds on \tilde{A} . Consider any $2s$ -sparse subset S of $[r]$. We have:

$$\begin{aligned} \sigma_{\min}(\tilde{A}_S) &\geq \sigma_{\min}(A_S^*) - \|A_S^* - \tilde{A}_S\|_2 \stackrel{(\zeta_1)}{\geq} 1 - \delta_{2s} - \|A_S^* - \tilde{A}_S\|_F \quad \text{and,} \\ \sigma_{\max}(\tilde{A}_S) &\leq \sigma_{\max}(A_S^*) + \|A_S^* - \tilde{A}_S\|_2 \stackrel{(\zeta_2)}{\leq} 1 + \delta_{2s} + \|A_S^* - \tilde{A}_S\|_F, \end{aligned}$$

where ζ_1 and ζ_2 follow from Assumption (A1). Recalling the assumption $\sqrt{s\epsilon_t} < 0.1$, and that $\delta_{2s} < 0.1$, we see that the maximum and minimum singular values of \tilde{A}_S are at least $4/5$ and at

most $6/5$ respectively. Appealing to Theorem A.1, we see that this guarantees $\|\Delta X_i\|_2 \leq 9s\epsilon_t$. Since this is also an infinity norm error bound, we obtain the second part of the lemma. The proof of the first part is further implied by the choice of our threshold at a level of $9s\epsilon_t$, which ensures that any non-zero element in X has $|X_p^{*i}| \geq 0$ (since we would have $|X_p^i| \leq 9s\epsilon_t$ by our infinity norm bound otherwise). \square

We now move on to the proof of Lemma 3.2. We point out that the lemma applies uniformly to all matrices W satisfying $\text{Supp}(W) \subseteq \text{Supp}(X^*)$, irrespective of the values of these entries. This might be surprising at first, but is a rather straight forward consequence of random matrix concentration theory.

Lemma 3.2. For every $r \times n$ matrix W s.t. $\text{Supp}(W) \subseteq \text{Supp}(X^*)$, we have (w.p. $\geq 1 - r \exp(-\frac{Cn}{rs})$):

$$\|W\|_2 \leq 2\|W\|_\infty \sqrt{\frac{s^2 n}{r}}.$$

Proof: Since the support of W is a subset of the support of X^* , $W_i^p = \chi_i^p W_i^p$. Now,

$$\begin{aligned} \|W\|_2 &= \max_{u,v \|u\|_2=1, \|v\|_2=1} \sum_{ip} W_i^p u^i v^p = \max_{u,v \|u\|_2=1, \|v\|_2=1} \sum_{ip} \chi_i^p W_i^p u^i v^p \\ &\leq \|W\|_\infty \cdot \max_{u,v \|u\|_2=1, \|v\|_2=1} \sum_{ip} \chi_i^p u^i v^p, \end{aligned}$$

where the inequality holds since the maximum inner product over the all pairs (u, v) from the unit sphere is larger than that over pairs with $u^i v^p \geq 0$ for all i, p . Note that the last expression is equal to $\|W\|_\infty u^\top \chi v$, where we use χ to denote the matrix with the non-zero pattern of the matrix X^* . It suffices to control the operator norm of this matrix for proving the lemma. This can indeed be done by applying Lemmas A.1 and A.2 with $\mu = M = 1$ and $\sigma = 0$. Doing so, yields with probability at least $1 - r \exp(-\frac{Cn}{rs})$

$$\|W\|_2 \leq 2\|W\|_\infty \sqrt{\frac{s^2 n}{r}},$$

which completes the proof. \square

We now finally prove Lemma 3.3, which is our main lemma on the structure of $X^* X^+$. Specifically, the lemma will show how to control the off-diagonal elements of this matrix carefully.

Lemma 3.3 (Off-diagonal error bound). Suppose $\|\Delta X\|_\infty < \frac{1}{288s}$. Then with probability at least $1 - r \exp(-\frac{Cn}{rM^2s}) - r \exp(-Cn/r^2)$, we have uniformly for every $p \in [r]$,

$$\left\| (\Delta X X^+)^{\setminus p} \right\|_2 = \left\| (X^* X^+)^{\setminus p} \right\|_2 \leq \frac{1968s^2 \|\Delta X\|_\infty}{\sqrt{r}}.$$

Proof: For simplicity, we will prove the statement for $p = 1$. We first relate $X^* X^+$ to $\Delta X X^+$.

$$\begin{aligned}
(X^* X^+)_1^{\setminus 1} &= ((X^* - X) X^+)_1^{\setminus 1} \\
&= -(\Delta X X^+)_1^{\setminus 1} \\
&= -\left(\Delta X X^\top (X X^\top)^{-1}\right)_1^{\setminus 1},
\end{aligned}$$

where the first step follows from the fact that $XX^+ = \mathbb{I}$. This proves the first part of the lemma. We now expand the above as follows:

$$\left(\Delta X X^\top (X X^\top)^{-1}\right)_1^{\setminus 1} = \left(\Delta X X^\top\right)_1^{\setminus 1} \left((X X^\top)^{-1}\right)_1^1 + \left(\Delta X X^\top\right)_{\setminus 1}^{\setminus 1} \left((X X^\top)^{-1}\right)_1^{\setminus 1}.$$

Using triangle inequality, we have:

$$\begin{aligned}
\left\| \left(\Delta X X^\top (X X^\top)^{-1}\right)_1^{\setminus 1} \right\|_2 &\leq \underbrace{\left\| \left((X X^\top)^{-1}\right)_1^1 \right\|}_{\mathcal{T}_1} \underbrace{\left\| \left(\Delta X X^\top\right)_1^{\setminus 1} \right\|_2}_{\mathcal{T}_2} \\
&\quad + \underbrace{\left\| \left(\Delta X X^\top\right)_{\setminus 1}^{\setminus 1} \right\|_2}_{\mathcal{T}_3} \underbrace{\left\| \left((X X^\top)^{-1}\right)_1^{\setminus 1} \right\|_2}_{\mathcal{T}_4}. \tag{11}
\end{aligned}$$

We now bound each of the above four quantities. We can easily bound \mathcal{T}_1 via a spectral norm bound on $(X X^\top)^{-1}$. Doing so, we obtain with probability at least $1 - r \exp(-\frac{Cn}{rM^2s})$

$$\mathcal{T}_1 = \left\| \left((X X^\top)^{-1}\right)_1^1 \right\| \leq \left\| (X X^\top)^{-1} \right\|_2 \stackrel{(\zeta_1)}{\leq} \frac{8r}{ns}, \tag{12}$$

where (ζ_1) follows from Lemma A.2. To bound \mathcal{T}_2 , we use Lemma A.6 and obtain with probability at least $1 - r \exp(-\frac{Cn}{r^2}) - r \exp(-\frac{Cn}{rM^2s})$

$$\mathcal{T}_2 = \left\| \left(\Delta X X^\top\right)_1^{\setminus 1} \right\|_2 \leq \frac{6 \|\Delta X\|_\infty s^2 n}{r^{\frac{3}{2}}}, \tag{13}$$

where we recall the assumption $\|\Delta X\|_\infty \leq 1/(64s)$. We now bound \mathcal{T}_3 as follows

$$\begin{aligned}
\mathcal{T}_3 &= \left\| \left(\Delta X X^\top\right)_{\setminus 1}^{\setminus 1} \right\|_2 \leq \left\| (\Delta X)_{\setminus 1}^{\setminus 1} \right\|_2 \left\| (X)_{\setminus 1}^{\setminus 1} \right\|_2 \stackrel{(\zeta_1)}{\leq} 2 \|\Delta X\|_\infty s \sqrt{\frac{n}{r}} \cdot 2(1 + \|\Delta X\|_\infty) s \sqrt{\frac{n}{r}} \\
&< \frac{6 \|\Delta X\|_\infty s^2 n}{r}, \tag{14}
\end{aligned}$$

where (ζ_1) follows from Lemmas 3.2 and A.3 (since $\text{Supp}(\Delta X) \subseteq \text{Supp}(X) \cup \text{Supp}(X^*) = \text{Supp}(X^*)$). Finally, to bound \mathcal{T}_4 , we start by noting the following block decomposition of the matrix XX^\top

$$XX^\top = \begin{bmatrix} X^1(X^1)^\top & X^1(X^{\setminus 1})^\top \\ X^{\setminus 1}X^1{}^\top & X^{\setminus 1}(X^{\setminus 1})^\top \end{bmatrix}.$$

Given this block-structure, we can now invoke Lemma A.8 (Schur complement lemma) to obtain

$$\left((XX^\top)^{-1} \right)_{\setminus 1}^{\setminus 1} = -\frac{1}{X^1(X^1)^\top} BX^{\setminus 1}(X^1)^\top,$$

where,

$$B := \left((XX^\top)^{-1} \right)_{\setminus 1}^{\setminus 1}. \quad (15)$$

Here we recall that B^{-1} is the Schur complement of $X_1X_1^\top$. Using Lemma A.6 and Equation 21 we have with probability at least $1 - r \exp(-\frac{Cn}{rM^2s}) - \exp(-Cn/r^2)$

$$\left\| \left((XX^\top)^{-1} \right)_{\setminus 1}^{\setminus 1} \right\|_2 \leq \frac{1}{\|X^1(X^1)^\top\|} \|B\|_2 \|X^{\setminus 1}(X^1)^\top\|_2 \leq \frac{8r}{sn} \cdot \|B\|_2 \cdot \frac{5s^2n}{r^{\frac{3}{2}}} = \frac{40s}{\sqrt{r}} \|B\|_2. \quad (16)$$

Using the expression (15) and the lower bound on $\sigma_{\min}(X)$ from Lemma A.3, we also have the following bound for $\|B\|_2$ with probability at least $1 - r \exp(-\frac{Cn}{rM^2s})$:

$$\|B\|_2 = \left\| \left((XX^\top)^{-1} \right)_{\setminus 1}^{\setminus 1} \right\|_2 \leq \left\| (XX^\top)^{-1} \right\|_2 \leq \frac{8r}{ns}.$$

Plugging the above into (16), gives us:

$$\left\| \left((XX^\top)^{-1} \right)_{\setminus 1}^{\setminus 1} \right\|_2 \leq \frac{40s}{\sqrt{r}} \cdot \frac{8r}{ns} \leq \frac{320\sqrt{r}}{n}. \quad (17)$$

Combining (12), (13), (14) and (17), we obtain with probability at least $1 - r \exp(-\frac{Cn}{rM^2s}) - \exp(-Cn/r^2)$

$$\begin{aligned} \left\| (XX^{*+})_{\setminus p}^{\setminus p} \right\|_2 &\leq \frac{48 \|\Delta X\|_\infty s}{\sqrt{r}} + \frac{1920 \|\Delta X\|_\infty s^2}{\sqrt{r}} \\ &\leq \frac{1968s^2 \|\Delta X\|_\infty}{\sqrt{r}}. \end{aligned}$$

□

A.2 Main Technical Lemmas

In this section, we state and prove the main technical lemmas used in our results.

Lemma A.1. *We have:*

$$\Sigma := \mathbb{E} \left[X^* X^{*\top} \right] = \left(\frac{s}{r} - \frac{s(s-1)\mu^2}{r(r-1)} \right) \mathbb{I} + \frac{s(s-1)\mu^2}{r(r-1)} \mathbb{1}\mathbb{1}^\top.$$

Proof:

Note that, $\chi_i^p, 1 \leq p \leq r$ all have same distribution. Hence, by symmetry and linearity of expectation, $\mathbb{E}[\chi_i^p] = \frac{1}{r} \mathbb{E} \left[\sum_{q=1}^r \chi_i^q \right] = \frac{s}{r}$. Similarly, $\mathbb{E}[(\chi_i^p)^2] = \frac{1}{r} \mathbb{E} \left[\sum_{q=1}^r (\chi_i^q)^2 \right] = \frac{s}{r}$. Also, $\mathbb{E} \left[\left(\sum_{q=1}^r \chi_i^q \right)^2 \right] = \mathbb{E} \left[\sum_{p,q} \chi_i^p \chi_i^q \right] = r \mathbb{E}[(\chi_i^p)^2] + (r^2 - r) \mathbb{E}[\chi_i^p \chi_i^q]$. Hence, $\mathbb{E}[\chi_i^p \chi_i^q] = \frac{s(s-1)}{r(r-1)}$.

Now, recall that $X_i^{*p} = \chi_i^p M_i^p$. Now, we first consider diagonal terms of Σ :

$$\Sigma_p^p = \mathbb{E}[(X_i^{*p})^2] = \mathbb{E}[(\chi_i^p)^2] \mathbb{E}[(M_i^p)^2] = \frac{s}{r}(\mu^2 + \sigma^2) = \frac{s}{r}. \quad (18)$$

Similarly, using independence of X_i^{*p} and X_i^{*j} , off-diagonal terms of Σ are given by:

$$\Sigma_p^q = \mathbb{E}[\chi_i^p \chi_i^q] \mathbb{E}[M_i^p] \mathbb{E}[M_i^q] = \frac{s(s-1)}{r(r-1)} \mu^2. \quad (19)$$

Lemma now follows by using (18) and (19). \square

In particular, two consequences of the lemma which will be particularly useful are about the extreme singular values of Σ . Recalling that $2s \leq r$ and $\mu^2 \leq 1$ by assumption, we obtain

$$\sigma_{\min}(\Sigma) \geq \frac{s}{2r}, \quad \text{and} \quad \sigma_{\max}(\Sigma) \leq \frac{2s^2}{r}. \quad (20)$$

We next establish some results on the spectrum of the empirical covariance matrix, using a standard result from random matrix theory. For convenience of the reader, we recall the following theorem from [33].

Theorem A.2 (Restatement of Theorem 5.44 from [33]). *Consider a $r \times n$ matrix W where each column w_i of W is an independent random vector with covariance matrix Σ . Suppose further that $\|w_i\|_2 \leq \sqrt{u}$ a.s. for all i . Then for any $t \geq 0$, the following inequality holds with probability at least $1 - r \exp(-ct^2)$:*

$$\left\| \frac{1}{n} W W^T - \Sigma \right\|_2 \leq \max \left(\|\Sigma\|_2^{1/2} \gamma, \gamma^2 \right) \quad \text{where} \quad \gamma = t \sqrt{\frac{u}{n}}.$$

Here $c > 0$ is an absolute numerical constant. In particular, this inequality yields:

$$\|W\|_2 \leq \|\Sigma\|_2^{1/2} \sqrt{n} + t\sqrt{u}.$$

Using the theorem, we can establish the following results on concentration of empirical covariance matrices. Hereafter, C will be a universal constant that can change from line to line.

Lemma A.2. *There exists a universal constant C such that with probability at least $1 - r \exp(-\frac{C\delta^2 ns}{rM^2})$, we have:*

$$\left\| \frac{1}{n} X^* X^{*\top} - \Sigma \right\|_2 \leq \max(\sqrt{2}\delta, \delta^2) \frac{s^2}{r}.$$

In particular, with probability at least $1 - r \exp(-\frac{Cn}{rM^2s})$, we have the bounds

$$\|X^*\|_2 \leq 2\sqrt{\frac{ns^2}{r}} \quad \text{and} \quad \sigma_{\min}(X^*) \geq \sqrt{\frac{ns}{4r}}.$$

Proof:

Note that, $\|X_i^*\|_2 \leq \sqrt{s}M$. Also, $\|\Sigma\|_2 \leq \frac{s}{r} + \frac{s(s-1)\mu^2}{r-1} \leq \frac{2s^2}{r}$. Using Theorem A.2 with $t = \delta\sqrt{\frac{ns}{rM^2}}$, we obtain:

$$\left\| \frac{1}{n} X^* X^{*\top} - \Sigma \right\|_2 \leq \max(\sqrt{2}\delta, \delta^2) \frac{s^2}{r},$$

w.p. greater than $1 - r \exp(-\frac{C\delta^2 ns}{rM^2})$. In order to obtain the second part, we apply the first part of the lemma with $\delta = 1/4\sqrt{2}s$ as well as Lemma A.1 to bound the largest and smallest singular values of XX^\top/n . Taking square roots completes the proof. \square

The next lemma we state is a specialization of Lemma 3.2 to obtain bounds on the spectral norm of our iterates X .

Lemma A.3. *With probability at least $1 - r \exp(-\frac{Cn}{rs}) - r \exp(-\frac{Cn}{rM^2s})$, for every $r \times n$ matrix X s.t. $\text{Supp}(X) \subseteq \text{Supp}(X^*)$, we have:*

$$\|X\|_2 \leq 2 \cdot (1 + \|X - X^*\|_\infty) \cdot s\sqrt{\frac{n}{r}}.$$

Proof:

Let $X = X^* + E_{X^*}$ where $\text{Supp}(E_{X^*}) \subseteq \text{Supp}(X^*)$. Hence, $\|X\|_2 \leq \|X^*\|_2 + \|X - X^*\|_2$. Lemma follows directly using Lemma A.2 and Lemma 3.2. \square

A useful version of the above lemma is when applied to matrices of the form XX^\top . We will need control over the upper and lower singular values of such matrices for our proofs, which we next provide.

Lemma A.4. *With probability at least $1 - r \exp(-\frac{Cn}{rs}) - r \exp(-\frac{Cn}{rM^2s})$, for every $r \times n$ matrix X s.t. $\text{Supp}(X) \subseteq \text{Supp}(X^*)$, we have:*

$$\left\| XX^\top - X^* X^{*\top} \right\|_2 \leq 4 \left(\|X - X^*\|_\infty + \|X - X^*\|_\infty^2 \right) \cdot \frac{s^2 n}{r}.$$

Further assuming $\|X - X^\|_\infty \leq 1/(64s)$, we have with the same probability*

$$\sigma_{\min}(XX^\top) \geq \frac{ns}{8r}.$$

Proof:

Let $X = X^* + E_{X^*}$. Note that $\text{Supp}(E_{X^*}) \subseteq \text{Supp}(X^*)$. Now,

$$\|XX^\top - X^*X^{*\top}\|_2 \leq \|E_{X^*}\|_2(\|E_{X^*}\|_2 + 2\|X^*\|_2).$$

By Lemma 3.2, $\|E_{X^*}\|_2 \leq 2s\sqrt{\frac{n}{r}}\|\mathbb{E}X^*\|_\infty$ with probability at least $\geq 1 - r \exp(-\frac{Cn}{rs})$. Combining this with the bound on $\|X^*\|_2$ from Lemma A.2 completes the proof. The second statement now follows by combining the result with our earlier lower bound on the minimum singular value of X^* in Lemma A.2. \square

A particular consequence of this lemma which will be useful is a lower bound on the diagonal entries of the matrix XX^\top . Indeed, we see that under the assumption $\|X - X^*\|_\infty \leq 1/(64s)$, with probability at least $1 - r \exp(-\frac{Cn}{rs}) - r \exp(-\frac{Cn}{rM^2s})$ we have the lower bound uniformly for all $p = 1, 2, \dots, r$

$$X^p X^{p\top} \geq \frac{ns}{8r}. \quad (21)$$

We finally have the following concentration lemma, which is a simple consequence of the Bernstein concentration bound.

Lemma A.5. *Let χ_i^p be as defined in Section A.1. Then, with probability at least $1 - \exp(-\frac{\delta^2 ns}{3rM^2})$:*

1. $(1 - \delta)\frac{sn}{r} \leq \sum_{i=1}^n \chi_i^p \leq (1 + \delta)\frac{sn}{r}, \forall p \in [r]$, and
2. $(1 - \delta)\frac{sn}{r} \leq \|X^*_{\cdot p}\|_2^2 = \sum_{i=1}^n \chi_i^p (M_i^p)^2 \leq (1 + \delta)\frac{sn}{r} \forall p \in [r]$

Proof:

We start with the proof of the second part, noting that the first part then immediately follows by setting $(M_i^p)^2 \equiv 1$. The second part will follow from a straightforward use of Bernstein's inequality. Note that $|M_i^p| \leq M$ and $\mathbb{E}[(M_i^p)^2] = 1$. As a result, for all $i = 1, 2, \dots, n$ we have $|\chi_i^p (M_i^p)^2| \leq M^2$, and

$$\text{Var}[\chi_i^p (M_i^p)^2] \leq \mathbb{E}[\chi_i^p (M_i^p)^4] \leq M^2 \mathbb{E}[\chi_i^p (M_i^p)^2] = M^2.$$

Also, we have $\mathbb{E}[\chi_i^p (M_i^p)^2] = \mathbb{E}[\chi_i^p] = s/r$. Consequently, we obtain that with probability at least $1 - \exp(-ns\delta^2/(rM^2(1 + \delta/3)))$ we have

$$\left| \sum_{i=1}^n \chi_i^p (M_i^p)^2 - \frac{ns}{r} \right| \leq \frac{\delta ns}{r}.$$

To complete the proof, note that $1 \geq \delta/3$ which yields the stated error probability. Finally, as stated before, we can recover the first part by setting $(M_i^p)^2 \equiv 1$. \square

Lemma A.6. *With probability at least $1 - r \exp(-\frac{Cn}{r^2}) - r \exp(-\frac{Cn}{rM^2})$, for every $r \times n$ matrix X s.t. $\text{Supp}(X) \subseteq \text{Supp}(X^*)$, we have the following bounds uniformly for all $p = 1, 2, \dots, r$*

1. $\left\| (\Delta X X^\top)_p^{\setminus p} \right\|_2 \leq (1 + \sqrt{s} \|\Delta X\|_\infty) \frac{4\sqrt{2}\|\Delta X\|_\infty s^2 n}{r^{\frac{3}{2}}}$, and
2. $\left\| X^{\setminus p} (X^p)^\top \right\|_2 \leq (1 + \sqrt{s} \|\Delta X\|_\infty)^2 \frac{4s^2 n}{r^{\frac{3}{2}}}$,

where $\Delta X := X - X^*$.

Proof: Since X has the same sparsity pattern as X^* , we can rewrite it as $X_i^p = \chi_i^p X_i^p$. We start by proving the first part of the lemma.

Proof of Part 1: Without loss of generality, we will prove the statement for $p = 1$. Let D denote the $n \times n$ diagonal matrix with:

$$D_i^i = \begin{cases} 1, & \text{if } X^{*1}_i \neq 0, \\ 0, & \text{otherwise.} \end{cases}$$

Using this notation, we have $(\Delta X X^\top)_1^{\setminus 1} = (\Delta X D X^\top)_1^{\setminus 1}$. So, we have:

$$\begin{aligned} \left\| (\Delta X X^\top)_1^{\setminus 1} \right\|_2 &= \left\| (\Delta X D X^\top)_1^{\setminus 1} \right\|_2 \\ &\leq \left\| (\Delta X D)_1^{\setminus 1} \right\|_2 \left\| (X^\top)_1 \right\|_2 \\ &\leq \left\| (\Delta X D)_1^{\setminus 1} \right\|_2 \left\| (X^{*\top})_1 \right\|_2 + \left\| (\Delta X^\top)_1 \right\|_2 \\ &\stackrel{(\zeta_1)}{\leq} \left\| (\Delta X D)_1^{\setminus 1} \right\|_2 \cdot \left(\sqrt{\frac{2sn}{r}} + \|\Delta X\|_\infty s \sqrt{\frac{2n}{r}} \right), \end{aligned}$$

with probability at least $1 - r \exp(-Cn/rs)$, where the first term in (ζ_1) follows from the second part of Lemma A.5 (setting $\delta = 1$) and the second is a consequence of Lemma 3.2. In order to control $\left\| (\Delta X D)_1^{\setminus 1} \right\|_2$, we observe that it is a matrix with a random number of columns selected by the matrix D . In particular, conditioned on $\{i : D_i^i = 1\}$, the support of $X^{*\setminus 1}_i$ is independent over $s - 1$ sparse vectors (and the support of ΔX is a subset of the support of X^*). Hence we can easily see that,

$$\begin{aligned} \mathbb{P} \left[\left\| (\Delta X D)_1^{\setminus 1} \right\|_2 > t \right] &\leq \mathbb{P} \left[\left\| (\Delta X D)_1^{\setminus 1} \right\|_2 > t \cap \frac{sn}{2r} < |\{i : D_i^i = 1\}| < \frac{2sn}{r} \right] \\ &\quad + \mathbb{P} \left[|\{i : D_i^i = 1\}| \leq \frac{sn}{2r} \cup |\{i : D_i^i = 1\}| \geq \frac{2sn}{r} \right]. \end{aligned} \quad (22)$$

In order to control the first probability, we note that

$$\begin{aligned} &\mathbb{P} \left[\left\| (\Delta X D)_1^{\setminus 1} \right\|_2 > t \cap \frac{sn}{2r} < |\{i : D_i^i = 1\}| < \frac{2sn}{r} \right] \\ &= \sum_{m=\lfloor sn/2r \rfloor}^{\lceil 2sn/r \rceil} \mathbb{P} \left[\left\| (\Delta X D)_1^{\setminus 1} \right\|_2 > t \cap |\{i : D_i^i = 1\}| = m \right] \\ &= \sum_{m=\lfloor sn/2r \rfloor}^{\lceil 2sn/r \rceil} \mathbb{P} \left[\left\| (\Delta X D)_1^{\setminus 1} \right\|_2 > t \mid |\{i : D_i^i = 1\}| = m \right] \mathbb{P} [|\{i : D_i^i = 1\}| = m]. \end{aligned}$$

Setting $t = 2 \|\Delta X\|_\infty \sqrt{s^2 m/r}$, we obtain as a consequence of Lemma 3.2:

$$\begin{aligned} & \mathbb{P} \left[\left\| (\Delta X D)^{\setminus 1} \right\|_2 > t \cap \frac{sn}{2r} < |\{i : D_i^i = 1\}| < \frac{2sn}{r} \right] \\ &= \sum_{m=\lfloor sn/2r \rfloor}^{\lceil 2sn/r \rceil} r \exp\left(-\frac{Cm}{rs}\right) \mathbb{P} [|\{i : D_i^i = 1\}| = m] \\ &\leq r \exp\left(-\frac{Cn}{2r^2}\right). \end{aligned}$$

The second probability in Equation 22 can be bounded through part 1 of Lemma A.5, since

$$|\{i : D_i^i = 1\}| = \sum_{i=1}^n \chi_i^1.$$

Doing so, we obtain with probability at least $1 - r \exp\left(-\frac{Cn}{2r^2}\right) - r \exp\left(-Cns/(3rM^2)\right)$:

$$\left\| (\Delta X D)^{\setminus 1} \right\|_2 \leq 2 \|\Delta X\|_\infty s \sqrt{\frac{(2sn)}{r}}.$$

This proves part 1.

Proof of Part 2: The proof of this is similar to that of part 1. Wlog, assume $p = 1$. We have:

$$\begin{aligned} \left\| X^{\setminus 1} (X^1)^\top \right\|_2 &= \left\| X^{\setminus 1} D (X^1)^\top \right\|_2 \\ &\leq \left\| (XD)^{\setminus 1} \right\|_2 \left\| (X^\top) \right\|_2 \\ &\leq \left\| (XD)^{\setminus 1} \right\|_2 \cdot 2(1 + \sqrt{s} \|\Delta X\|_\infty) \sqrt{\frac{sn}{r}}. \end{aligned}$$

For the first term above, we have:

$$\left\| (XD)^{\setminus 1} \right\|_2 \leq \left\| (X^* D)^{\setminus 1} \right\|_2 + \left\| (\Delta X D)^{\setminus 1} \right\|_2.$$

The second term in this decomposition was controlled above and the first one can be similarly bounded. Doing so, we obtain with probability at least $1 - r \exp\left(-\frac{Cn}{r^2}\right) - r \exp\left(-\frac{Cn}{rM^2s}\right)$:

$$\left\| (XD)^{\setminus 1} \right\|_2 \leq 2s(1 + \sqrt{s} \|\Delta X\|_\infty) \sqrt{\frac{2sn}{r^3}}.$$

This proves the lemma. □

We begin with an auxiliary result on the RIP constant of an incoherent matrix.

Lemma A.7. *Suppose A^* satisfies Assumption (B1). Then, the $2s$ -RIP constant of A^* , δ_{2s} satisfies $\delta_{2s} < \frac{2\mu_0 s}{\sqrt{d}}$.*

Proof: Consider a $2s$ -sparse unit vector $w \in \mathbb{R}^r$ with $\text{Supp}(w) = S$. We have:

$$\begin{aligned}
\|Aw\|^2 &= \left(\sum_{j \in S} w_j A^*_j \right)^2 = \sum_j w_j^2 \|A^*_j\|^2 + \sum_{j, l \in S, j \neq l} w_j w_l \langle A^*_j, A^*_l \rangle \\
&\geq 1 - \sum_{j, l \in S, j \neq l} |w_j w_l| |\langle A^*_j, A^*_l \rangle| \\
&\geq 1 - \sum_{j, l \in S, j \neq l} |w_j w_l| \frac{\mu_0}{\sqrt{d}} \\
&\geq 1 - \frac{\mu_0}{\sqrt{d}} \|w\|_1^2 \\
&\geq 1 - \frac{\mu_0}{\sqrt{d}} 2s \cdot \|w\|_2^2 = 1 - \frac{2\mu_0 s}{\sqrt{d}}.
\end{aligned}$$

Similarly, we have:

$$\|A^*w\|_2^2 \leq 1 + \frac{2\mu_0 s}{\sqrt{d}}.$$

This proves the lemma. □

Lemma A.8. *We have the following formula for matrix inversion:*

$$\begin{bmatrix} A & B \\ C & D \end{bmatrix}^{-1} = \begin{bmatrix} A^{-1} + A^{-1}BMCA^{-1} & -A^{-1}BM \\ -MCA^{-1} & M \end{bmatrix},$$

where $M^{-1} := (D - CA^{-1}B)$ is the Schur complement of A in the above matrix.