
Complete Dictionary Recovery Using Nonconvex Optimization

Ju Sun
Qing Qu
John Wright

JS4038@COLUMBIA.EDU
QQ2105@COLUMBIA.EDU
JW2966@COLUMBIA.EDU

Department of Electrical Engineering, Columbia University, New York, NY, USA

Abstract

We consider the problem of recovering a complete (i.e., square and invertible) dictionary \mathbf{A}_0 , from $\mathbf{Y} = \mathbf{A}_0\mathbf{X}_0$ with $\mathbf{Y} \in \mathbb{R}^{n \times p}$. This recovery setting is central to the theoretical understanding of dictionary learning. We give the first efficient algorithm that provably recovers \mathbf{A}_0 when \mathbf{X}_0 has $O(n)$ nonzeros per column, under suitable probability model for \mathbf{X}_0 . Prior results provide recovery guarantees when \mathbf{X}_0 has only $O(\sqrt{n})$ nonzeros per column. Our algorithm is based on nonconvex optimization with a spherical constraint, and hence is naturally phrased in the language of manifold optimization. Our proofs give a geometric characterization of the high-dimensional objective landscape, which shows that with high probability there are no spurious local minima. Experiments with synthetic data corroborate our theory. Full version of this paper is available online: <http://arxiv.org/abs/1504.06785>.

1. Introduction

Dictionary learning (DL) is the problem of finding a sparse representation for a collection of input signals. Its applications span classical image processing, visual recognition, compressive signal acquisition, and also recent deep architectures for signal classification. Recent surveys on applications and algorithms of DL include Elad (2010) and Mairal et al (2014).

Formally, given a data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$, DL seeks an approximation $\mathbf{Y} \approx \mathbf{A}\mathbf{X}$, where \mathbf{A} lies in a certain admissible set \mathcal{A} , and \mathbf{X} is as sparse as possible. Typical formulations for DL are nonconvex: the admissible set \mathcal{A} is typically nonconvex, and the observation map $(\mathbf{A}, \mathbf{X}) \mapsto \mathbf{A}\mathbf{X}$ is bilinear. There is also an intrinsic symmetry in the problem

due to sign-permutation ambiguity, which seems to preclude convexification (Gribonval & Schnass (2010), Geng & Wright (2011)). Thus, despite many empirical successes, relatively little is known about the theoretical properties of DL algorithms.

Towards theoretical understanding, it is natural to start with the dictionary recovery (DR) problem: suppose that the data matrix \mathbf{Y} is generated as $\mathbf{Y} = \mathbf{A}_0\mathbf{X}_0$, where $\mathbf{A}_0 \in \mathbb{R}^{n \times m}$ and $\mathbf{X}_0 \in \mathbb{R}^{m \times p}$, and try to recover \mathbf{A}_0 and \mathbf{X}_0 . One might imagine putting favorable structural assumptions on \mathbf{A}_0 and \mathbf{X}_0 to make DR well-posed and amenable to efficient algorithms. However, under natural assumptions for \mathbf{A}_0 and \mathbf{X}_0 , even proving that the target solution is a *local* minimum of certain popular practical DL formulations requires nontrivial analysis (Gribonval & Schnass, 2010; Geng & Wright, 2011; Schnass, 2014a,b; 2015). Obtaining *global* solutions with *efficient* algorithms is a standing challenge.

Existing results on global dictionary recovery pertain only to highly sparse \mathbf{X}_0 . For example, Spielman et al (2012) showed that solving a sequence of certain linear programs can recover a complete dictionary \mathbf{A}_0 from \mathbf{Y} , when \mathbf{X}_0 is a sparse random matrix with $O(\sqrt{n})$ nonzeros per column. Agarwal et al (2013a; 2013b) and Arora et al (2013; 2015) have subsequently given efficient algorithms for the overcomplete setting ($m \geq n$), based on a combination of {clustering or spectral initialization} and local refinement. These algorithms again succeed when \mathbf{X}_0 has $\tilde{O}(\sqrt{n})$ nonzeros per column (The \tilde{O} suppresses some logarithm factors). Barak et al (2014) provides efficient algorithms based on sum-of-square hierarchy that guarantees recovery of complete dictionaries when \mathbf{X}_0 has $O(n^c)$ nonzeros per column for any $c < 1$. Giving *efficient* algorithms which provably succeeds in linear sparsity regime (i.e., $O(n)$ nonzeros per column) is an open problem.¹

¹Recent works, including Arora et al (2014) and Barak et al (2014), contain guarantees for recovery with linear sparsity, but run in super-polynomial (quasipolynomial) time. Aside from efficient recovery, other theoretical work on DL includes results on identifiability (Hillar & Sommer, 2011), generalization bounds (Vainsencher et al., 2011; Mehta & Gray, 2013), and noise

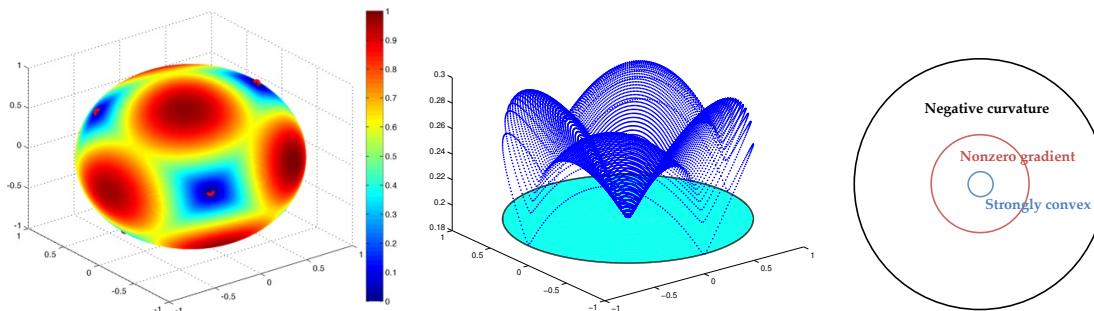


Figure 1. Why is dictionary learning over \mathbb{S}^{n-1} tractable? Assume the target dictionary \mathbf{A}_0 is orthogonal. **Left:** Large sample objective function $\mathbb{E}_{\mathbf{X}_0} [f(\mathbf{q})]$. The only local minima are the columns of \mathbf{A}_0 and their negatives. **Center:** the same function, visualized as a height above the plane \mathbf{a}_1^\perp (\mathbf{a}_1 is the first column of \mathbf{A}_0). **Right:** Around the optimum, the function exhibits a small region of positive curvature, a region of large gradient, and finally a region in which the direction away from \mathbf{a}_1 is a direction of negative curvature.

In this work, we focus on recovering a complete (i.e., invertible) dictionary \mathbf{A}_0 from $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$. We give the first *polynomial-time* algorithm that provably recovers \mathbf{A}_0 when \mathbf{X}_0 has $O(n)$ nonzeros per column, if \mathbf{X}_0 contains i.i.d. Bernoulli-Gaussian entries. We achieve this by formulating complete DR as a nonconvex program with spherical constraint. Under the probability model for \mathbf{X}_0 , we give a geometric characterization of the high-dimensional objective landscape over the sphere, which shows that with high probability (w.h.p.) there are no “spurious” local minima. In particular, the geometric structure allows us to design a Riemannian trust region algorithm over the sphere that provably converges to one local minimum with an *arbitrary initialization*, despite the presence of saddle points.

This paper is organized as follows. Sec. 2 will motivate our nonconvex formulation of DR and formalize the setting. Sec. 3 and 4 will introduce two integral parts of our algorithmic framework: characterization of the high-dimensional function landscape and a Riemannian trust-region algorithm over the sphere. Main theorems are included in Sec. 5, followed by some numerical verification presented in Sec. 6. Sec. 7 will close this paper by discussing implications of this work. Due to space constraint, we only sketch high-level ideas of the proof in this paper. Detailed proofs can be found in the full version (Sun et al., 2015).

2. A Nonconvex Formulation for DR

We assume $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$, where $\mathbf{A}_0 \in \mathbb{R}^{n \times n}$ is a complete matrix and \mathbf{X}_0 follows the Bernoulli-Gaussian (BG) model with rate θ : $[\mathbf{X}_0]_{ij} = \Omega_{ij} G_{ij}$, with $\Omega_{ij} \sim \text{Ber}(\theta)$ and $V_{ij} \sim \mathcal{N}(0, 1)$. We write compactly $\mathbf{X}_0 \sim_{i.i.d.} \text{BG}(\theta)$.

Since $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$ and \mathbf{A}_0 is complete, $\text{row}(\mathbf{Y}) = \text{row}(\mathbf{X}_0)$ ² and rows of \mathbf{X}_0 are sparse vectors in the known subspace $\text{row}(\mathbf{Y})$. Following Spielman et al (2012), we use

stability (Gribonval et al., 2014).

²row (\cdot) denotes the row space.

this fact to first recover the rows of \mathbf{X}_0 , and subsequently recover \mathbf{A}_0 by solving a system of linear equations. In fact, for $\mathbf{X}_0 \sim_{i.i.d.} \text{BG}(\theta)$, rows of \mathbf{X}_0 are the n *sparsest* vectors (directions) in $\text{row}(\mathbf{Y})$ w.h.p. (Spielman et al., 2012). Thus one might try to recover them by solving³

$$\min \|\mathbf{q}^* \mathbf{Y}\|_0 \quad \text{s.t.} \quad \mathbf{q} \neq \mathbf{0}. \quad (2.1)$$

The objective is discontinuous, and the domain is an open set. Known convex relaxations (Spielman et al., 2012; Demanet & Hand, 2014) provably break down beyond the aforementioned $O(\sqrt{n})$ sparsity level. Instead, we work with a *nonconvex* alternative:⁴

$$\min f(\mathbf{q}; \hat{\mathbf{Y}}) \doteq \frac{1}{p} \sum_{k=1}^p h_\mu(\mathbf{q}^* \hat{\mathbf{y}}_k), \quad \text{s.t.} \quad \|\mathbf{q}\|_2 = 1, \quad (2.2)$$

where $\hat{\mathbf{Y}} \in \mathbb{R}^{n \times p}$ is a proxy of \mathbf{Y} and k indexes columns of $\hat{\mathbf{Y}}$. Here $h_\mu(\cdot)$ is chosen to be a convex smooth approximation to the $|\cdot|$ function, namely,

$$h_\mu(z) = \mu \log \cosh\left(\frac{z}{\mu}\right), \quad (2.3)$$

which is infinitely differentiable and μ controls the smoothing level. The spherical constraint is nonconvex. Hence, a-priori, it is unclear whether (2.2) admits efficient algorithms that attain one local optimizer (Murty & Kabadi, 1987). Surprisingly, simple descent algorithms for (2.2) exhibit very striking behavior: on many practical numerical examples⁵, they appear to produce global solutions. Our next section will uncover interesting geometrical structures underlying the phenomenon.

³The notation $*$ denotes matrix transposition.

⁴Similar formulation has been proposed in (Zibulevsky & Pearlmuter, 2001) in the context of blind source separation, see also (Qu et al., 2014).

⁵... not restricted to the model we assume here for \mathbf{A}_0 and \mathbf{X}_0 .

3. High-dimensional Geometry

For the moment, suppose \mathbf{A}_0 is orthogonal, and take $\widehat{\mathbf{Y}} = \mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$ in (2.2). Figure 1 (left) plots $\mathbb{E}_{\mathbf{X}_0} [f(\mathbf{q}; \mathbf{Y})]$ over $\mathbf{q} \in \mathbb{S}^2$ ($n = 3$). Remarkably, $\mathbb{E}_{\mathbf{X}_0} [f(\mathbf{q}; \mathbf{Y})]$ has no spurious local minima. In fact, every local minimum $\widehat{\mathbf{q}}$ produces a row of \mathbf{X}_0 : $\widehat{\mathbf{q}}^* \mathbf{Y} = \alpha e_i^* \mathbf{X}_0$ for some $\alpha \neq 0$.

To better illustrate the point, we take the particular case $\mathbf{A}_0 = \mathbf{I}$ and project the upper hemisphere above the equatorial plane e_3^\perp onto e_3^\perp . The projection is bijective and we equivalently define a reparameterization $g : e_3^\perp \mapsto \mathbb{R}$ of f . Figure 1 (center) plots the graph of g . Obviously the only local minimizers are $\mathbf{0}, \pm e_1, \pm e_2$, and they are also global minimizers. Moreover, the apparent nonconvex landscape has interesting structures around $\mathbf{0}$: when moving away from $\mathbf{0}$, one sees successively a strongly convex region, a nonzero gradient region, and a region where at each point one can always find a direction of negative curvature, as shown schematically in Figure 1 (right). This geometry implies that at any nonoptimal point, there is always at least one direction of descent. Thus, any algorithm that can take advantage of the descent directions will likely converge to one global minimizer, irrespective of initialization.

Two challenges stand out when implementing this idea. For geometry, one has to show similar structure exists for general complete \mathbf{A}_0 , in high dimensions ($n \geq 3$), when the number of observations p is finite (vs. the expectation in the experiment). For algorithms, we need to be able to take advantage of this structure without knowing \mathbf{A}_0 ahead of time. In Sec. 4, we describe a Riemannian trust region method which addresses the later challenge. Now we focus on the first one.

3.1. Geometry for orthogonal \mathbf{A}_0

In this case, we take $\widehat{\mathbf{Y}} = \mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0$. Since $f(\mathbf{q}; \mathbf{A}_0 \mathbf{X}_0) = f(\mathbf{A}_0^* \mathbf{q}; \mathbf{X}_0)$, the landscape of $f(\mathbf{q}; \mathbf{A}_0 \mathbf{X}_0)$ is simply a rotated version of that of $f(\mathbf{q}; \mathbf{X}_0)$, i.e., when $\mathbf{A}_0 = \mathbf{I}$. Hence we will focus on the case when $\mathbf{A}_0 = \mathbf{I}$. Among the $2n$ symmetric sections of \mathbb{S}^{n-1} centered around the signed basis vectors $\pm e_1, \dots, \pm e_n$, we work with the section around e_n as an example. The result will carry over to all sections with the same argument.

We again invoke the projection trick described above, this time onto the equatorial plane e_n^\perp . This can be formally captured by the reparameterization mapping:

$$\mathbf{q}(\mathbf{w}) = \left(\mathbf{w}, \sqrt{1 - \|\mathbf{w}\|^2} \right), \quad \mathbf{w} \in \mathbb{R}^{n-1}, \quad (3.1)$$

where \mathbf{w} is the new variable in e_n^\perp . We study the composition $g(\mathbf{w}; \mathbf{Y}) \doteq f(\mathbf{q}(\mathbf{w}); \mathbf{Y})$ over the set

$$\Gamma \doteq \left\{ \mathbf{w} : \|\mathbf{w}\| < \sqrt{\frac{4n-1}{4n}} \right\}. \quad (3.2)$$

Our next theorem characterizes the properties of $g(\mathbf{w})$. In particular, it shows the favorable structure we observed for $n = 3$ persists in high dimensions, w.h.p.⁶, even when p is large yet finite, for the case \mathbf{A}_0 is orthogonal.

Theorem 3.1. *Suppose $\mathbf{A}_0 = \mathbf{I}$ and hence $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0 = \mathbf{X}_0$. There exist positive constants c_* and C , such that for any $\theta \in (0, 1/2)$ and $\mu < \min \{c_a \theta n^{-1}, c_b n^{-5/4}\}$, whenever $p \geq \frac{C}{\mu^2 \theta^2} n^3 \log \frac{n}{\mu \theta}$, the following hold simultaneously w.h.p.:*

$$\begin{aligned} \nabla^2 g(\mathbf{w}; \mathbf{X}_0) &\succeq \frac{c_* \theta}{\mu} \mathbf{I} & \forall \mathbf{w} \|\mathbf{w}\| &\leq \frac{\mu}{4\sqrt{2}}, \\ \frac{\mathbf{w}^* \nabla g(\mathbf{w}; \mathbf{X}_0)}{\|\mathbf{w}\|} &\geq c_* \theta & \forall \mathbf{w} \frac{\mu}{4\sqrt{2}} &\leq \|\mathbf{w}\| \leq \frac{1}{20\sqrt{5}}, \\ \frac{\mathbf{w}^* \nabla^2 g(\mathbf{w}; \mathbf{X}_0) \mathbf{w}}{\|\mathbf{w}\|^2} &\leq -c_* \theta & \forall \mathbf{w} \frac{1}{20\sqrt{5}} &\leq \|\mathbf{w}\| \leq \sqrt{\frac{4n-1}{4n}}, \end{aligned}$$

and the function $g(\mathbf{w}; \mathbf{X}_0)$ has exactly one local minimizer \mathbf{w}_* over the open set $\Gamma \doteq \left\{ \mathbf{w} : \|\mathbf{w}\| < \sqrt{\frac{4n-1}{4n}} \right\}$, which satisfies

$$\|\mathbf{w}_* - \mathbf{0}\| \leq \min \left\{ \frac{c_c \mu}{\theta} \sqrt{\frac{n \log p}{p}}, \frac{\mu}{16} \right\}. \quad (3.3)$$

In particular, with this choice of p , the probability the claim fails to hold is at most $4np^{-10} + \theta(np)^{-7} + \exp(-0.3\theta np) + c_d \exp(-c_e p \mu^2 \theta^2 / n^2)$. Here c_a to c_e are all positive numerical constants.

In words, when the samples are numerous enough, one sees the strongly convex, nonzero gradient, and negative curvature regions successively when moving away from target solution $\mathbf{0}$, and the local (also global) minimizer of $g(\mathbf{w}; \mathbf{Y})$ is next to $\mathbf{0}$, within a distance of $O(\mu)$.

Note that $\mathbf{q}(\Gamma)$ contains all points $\mathbf{q} \in \mathbb{S}^{n-1}$ such that $\mathbf{q}^* e_n = \max_i |\mathbf{q}^* e_i|$. We can characterize the graph of the function $f(\mathbf{q}; \mathbf{X}_0)$ in the vicinity of some other signed basis vector $\pm e_i$ simply by changing the plane e_n^\perp to e_i^\perp . Doing this $2n$ times (and multiplying the failure probability in Theorem 3.1 by $2n$), we obtain a characterization of $f(\mathbf{q})$ over the entirety of \mathbb{S}^{n-1} .

Corollary 3.2. *Suppose $\mathbf{A}_0 = \mathbf{I}$ and hence $\mathbf{Y} = \mathbf{A}_0 \mathbf{X}_0 = \mathbf{X}_0$. There exist positive constant C , such that for any $\theta \in (0, 1/2)$ and $\mu < \min \{c_a \theta n^{-1}, c_b n^{-5/4}\}$, whenever $p \geq \frac{C}{\mu^2 \theta^2} n^3 \log \frac{n}{\mu \theta}$, with probability at least $1 - 8n^2 p^{-10} - \theta(np)^{-7} - \exp(-0.3\theta np) - c_c \exp(-c_d p \mu^2 \theta^2 / n^2)$, the function $f(\mathbf{q}; \mathbf{X}_0)$ has exactly $2n$ local minimizers over the sphere \mathbb{S}^{n-1} . In particular, there is a bijective map between these minimizers and signed basis vectors $\{\pm e_i\}_i$, such*

⁶In this work, we say some event occurs with high probability when the failure probability is bounded by an inverse polynomial of n and p .

that the corresponding local minimizer \mathbf{q}_* and $\mathbf{b} \in \{\pm \mathbf{e}_i\}_i$ satisfy

$$\|\mathbf{q}_* - \mathbf{b}\| \leq \sqrt{2} \min \left\{ \frac{c_c \mu}{\theta} \sqrt{\frac{n \log p}{p}}, \frac{\mu}{16} \right\}. \quad (3.4)$$

Here c_a to c_d are numerical constants (possibly different from that in the above theorem).

The proof of Theorem 3.1 is conceptually straightforward: one shows that $\mathbb{E}[f(\mathbf{q}; \mathbf{X}_0)]$ has the claimed properties, and then proves that each of the quantities of interest concentrates uniformly about its expectation. The detailed calculations are nontrivial.

3.2. Geometry for complete \mathbf{A}_0

For general complete dictionaries \mathbf{A}_0 , we hope that the function f retains the nice geometric structure discussed above. We can ensure this by “preconditioning” \mathbf{Y} such that the output looks as if being generated from a certain orthogonal matrix, possibly plus a small perturbation. We can then argue that the perturbation does not significantly affect the properties of the graph of the objective function. Write

$$\bar{\mathbf{Y}} = \left(\frac{1}{p\theta} \mathbf{Y} \mathbf{Y}^* \right)^{-1/2} \mathbf{Y}. \quad (3.5)$$

Note that for $\mathbf{X}_0 \sim_{i.i.d.} \text{BG}(\theta)$, $\mathbb{E}[\mathbf{X}_0 \mathbf{X}_0^*] / (p\theta) = \mathbf{I}$. Thus, one expects $\frac{1}{p\theta} \mathbf{Y} \mathbf{Y}^* = \frac{1}{p\theta} \mathbf{A}_0 \mathbf{X}_0 \mathbf{X}_0^* \mathbf{A}_0^*$ to behave roughly like $\mathbf{A}_0 \mathbf{A}_0^*$ and hence $\bar{\mathbf{Y}}$ to behave like

$$(\mathbf{A}_0 \mathbf{A}_0^*)^{-1/2} \mathbf{A}_0 \mathbf{X}_0 = \mathbf{U} \mathbf{V}^* \mathbf{X}_0 \quad (3.6)$$

where we write the SVD of \mathbf{A}_0 as $\mathbf{A}_0 = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$. It is easy to see $\mathbf{U} \mathbf{V}^*$ is an orthogonal matrix. Hence the preconditioning scheme we have introduced is technically sound.

Our analysis shows that $\bar{\mathbf{Y}}$ can be written as

$$\bar{\mathbf{Y}} = \mathbf{U} \mathbf{V}^* \mathbf{X}_0 + \mathbf{\Xi} \mathbf{X}_0, \quad (3.7)$$

where $\mathbf{\Xi}$ is a matrix with small magnitude. Perturbation analysis combines with the the above results for the orthogonal case yields:

Theorem 3.3. *Suppose \mathbf{A}_0 is complete with its condition number $\kappa(\mathbf{A}_0)$. There exist positive constants c_* and C , such that for any $\theta \in (0, 1/2)$ and $\mu < \min\{c_a \theta n^{-1}, c_b n^{-5/4}\}$, when $p \geq \frac{C}{c_*^2 \theta} \max\left\{\frac{n^4}{\mu^4}, \frac{n^5}{\mu^2}\right\} \kappa^8(\mathbf{A}_0) \log^4\left(\frac{\kappa(\mathbf{A}_0)n}{\mu\theta}\right)$ and $\bar{\mathbf{Y}} \doteq \sqrt{p\theta} (\mathbf{Y} \mathbf{Y}^*)^{-1/2} \mathbf{Y}$, $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^* = \text{SVD}(\mathbf{A}_0)$, and let $\tilde{\mathbf{Y}} = \mathbf{V} \mathbf{U}^* \bar{\mathbf{Y}}$, then the following hold simultaneously w.h.p.:*

$$\nabla^2 g(\mathbf{w}; \tilde{\mathbf{Y}}) \succeq \frac{c_* \theta}{2\mu} \mathbf{I} \quad \forall \mathbf{w} \|\mathbf{w}\| \leq \frac{\mu}{4\sqrt{2}},$$

$$\begin{aligned} \frac{\mathbf{w}^* \nabla g(\mathbf{w}; \tilde{\mathbf{Y}})}{\|\mathbf{w}\|} &\geq \frac{1}{2} c_* \theta & \forall \mathbf{w} \frac{\mu}{4\sqrt{2}} \leq \|\mathbf{w}\| \leq \frac{1}{20\sqrt{5}} \\ \frac{\mathbf{w}^* \nabla^2 g(\mathbf{w}; \tilde{\mathbf{Y}}) \mathbf{w}}{\|\mathbf{w}\|^2} &\leq -\frac{1}{2} c_* \theta & \forall \mathbf{w} \frac{1}{20\sqrt{5}} \leq \|\mathbf{w}\| \leq \sqrt{\frac{4n-1}{4n}} \end{aligned}$$

and the function $g(\mathbf{w}; \tilde{\mathbf{Y}})$ has exactly one local minimizer \mathbf{w}_* over the open set $\Gamma \doteq \left\{ \mathbf{w} : \|\mathbf{w}\| < \sqrt{\frac{4n-1}{4n}} \right\}$, which satisfies

$$\|\mathbf{w}_* - \mathbf{0}\| \leq \frac{\mu}{7}. \quad (3.8)$$

In particular, with this choice of p , the probability the claim fails to hold is at most $4np^{-10} + \theta(np)^{-7} + \exp(-0.3\theta np) + p^{-8} + c_d \exp(-c_e p \mu^2 \theta^2 / n^2)$. Here c_a to c_e are all positive numerical constants.

Corollary 3.4. *Suppose \mathbf{A}_0 is complete with its condition number $\kappa(\mathbf{A}_0)$. There exist positive constants c_* and C , such that for any $\theta \in (0, 1/2)$ and $\mu < \min\{c_a \theta n^{-1}, c_b n^{-5/4}\}$, when $p \geq \frac{C}{c_*^2 \theta} \max\left\{\frac{n^4}{\mu^4}, \frac{n^5}{\mu^2}\right\} \kappa^8(\mathbf{A}_0) \log^4\left(\frac{\kappa(\mathbf{A}_0)n}{\mu\theta}\right)$ and $\bar{\mathbf{Y}} \doteq \sqrt{p\theta} (\mathbf{Y} \mathbf{Y}^*)^{-1/2} \mathbf{Y}$, $\mathbf{U} \mathbf{\Sigma} \mathbf{V}^* = \text{SVD}(\mathbf{A}_0)$, with probability at least $1 - 8n^2 p^{-10} - \theta(np)^{-7} - \exp(-0.3\theta np) - p^{-8} - c_d \exp(-c_e p \mu^2 \theta^2 / n^2)$, the function $f(\mathbf{q}; \mathbf{V} \mathbf{U}^* \bar{\mathbf{Y}})$ has exactly $2n$ local minimizers over the sphere \mathbb{S}^{n-1} . In particular, there is a bijective map between these minimizers and signed basis vectors $\{\pm \mathbf{e}_i\}_i$, such that the corresponding local minimizer \mathbf{q}_* and $\mathbf{b} \in \{\pm \mathbf{e}_i\}_i$ satisfy*

$$\|\mathbf{q}_* - \mathbf{b}\| \leq \frac{\sqrt{2}\mu}{7}. \quad (3.9)$$

Here c_a to c_d are numerical constants (possibly different from that in the above theorem).

4. Riemannian Trust Region Algorithm

We do not know \mathbf{A}_0 ahead of time, so our algorithm needs to take advantage of the structure described above without knowledge of \mathbf{A}_0 . Intuitively, this seems possible as the descent direction in the \mathbf{w} space appears to also be a local descent direction for f over the sphere. Another issue is that although the optimization problem has no spurious local minima, it does have many saddle points (Figure. 1). Therefore, certain form of second-order information is needed to help escape the saddle points. Based on these considerations, we describe a Riemannian trust region method (TRM) (Absil et al., 2007; 2009) over the sphere for this purpose.

4.1. Trust region method for Euclidean spaces

For a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ and an unconstrained optimization problem

$$\min_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}), \quad (4.1)$$

typical (second-order) TRM proceeds by successively forming second-order approximations to f at the current iterate,

$$\widehat{f}(\boldsymbol{\delta}; \mathbf{x}_{k-1}) \doteq f(\mathbf{x}_{k-1}) + \nabla^* f(\mathbf{x}_{k-1}) \boldsymbol{\delta} + \frac{1}{2} \boldsymbol{\delta}^* \mathbf{Q}(\mathbf{x}_{k-1}) \boldsymbol{\delta}, \quad (4.2)$$

where $\mathbf{Q}(\mathbf{x}_{k-1})$ is a proxy for the Hessian matrix $\nabla^2 f(\mathbf{x}_{k-1})$, which encodes the second-order geometry. The next iterate is determined by seeking a minimum of $\widehat{f}(\boldsymbol{\delta}; \mathbf{x}_{k-1})$ over a small region, normally an ℓ^2 ball, commonly known as the trust region. Thus, the well known trust region subproblem takes the form

$$\boldsymbol{\delta}_k \doteq \arg \min_{\boldsymbol{\delta} \in \mathbb{R}^n, \|\boldsymbol{\delta}\| \leq \Delta} \widehat{f}(\boldsymbol{\delta}; \mathbf{x}_{k-1}), \quad (4.3)$$

where Δ is called the trust-region radius that controls how far the movement can be made. A ratio

$$\rho_k \doteq \frac{f(\mathbf{x}_{k-1}) - f(\mathbf{x}_{k-1} + \boldsymbol{\delta}_k)}{\widehat{f}(\mathbf{0}) - \widehat{f}(\boldsymbol{\delta}_{k-1})} \quad (4.4)$$

is defined to measure the progress and typically the radius Δ is updated dynamically according to ρ_k to adapt to the local function behavior. If the progress is satisfactory, the next iterate is (perhaps plus some line search improvement)

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \boldsymbol{\delta}_k. \quad (4.5)$$

Detailed introduction to the classical TRM can be found in the texts (Conn et al., 2000; Nocedal & Wright, 2006).

4.2. Trust region method over the sphere

To generalize the idea to smooth manifolds, one natural choice is to form the approximation over the tangent spaces (Absil et al., 2007; 2009). Specific to our spherical manifold, for which the tangent space at an iterate $\mathbf{q}_k \in \mathbb{S}^{n-1}$ is $T_{\mathbf{q}_k} \mathbb{S}^{n-1} \doteq \{\mathbf{v} : \mathbf{v}^* \mathbf{q}_k = 0\}$, we work with the quadratic approximation $\widehat{f} : T_{\mathbf{q}_k} \mathbb{S}^{n-1} \mapsto \mathbb{R}$ defined as

$$\widehat{f}(\mathbf{q}_k, \boldsymbol{\delta}) \doteq f(\mathbf{q}_k) + \langle \nabla f(\mathbf{q}_k), \boldsymbol{\delta} \rangle + \frac{1}{2} \boldsymbol{\delta}^* (\nabla^2 f(\mathbf{q}_k) - \langle \nabla f(\mathbf{q}_k), \mathbf{q}_k \rangle \mathbf{I}) \boldsymbol{\delta}. \quad (4.6)$$

To interpret this approximation, let $\mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} \doteq (\mathbf{I} - \mathbf{q}_k \mathbf{q}_k^*)$ be the orthoprojector onto $T_{\mathbf{q}_k} \mathbb{S}^{n-1}$ and write (4.6) into an equivalent form:

$$\widehat{f}(\mathbf{q}_k, \boldsymbol{\delta}) \doteq f(\mathbf{q}_k) + \left\langle \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} \nabla f(\mathbf{q}_k), \boldsymbol{\delta} \right\rangle + \frac{1}{2} \boldsymbol{\delta}^* \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} (\nabla^2 f(\mathbf{q}_k) - \langle \nabla f(\mathbf{q}_k), \mathbf{q}_k \rangle \mathbf{I}) \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} \boldsymbol{\delta}.$$

The two terms

$$\text{grad} f(\mathbf{q}_k) \doteq \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} \nabla f(\mathbf{q}_k),$$

$$\text{Hess} f(\mathbf{q}_k) \doteq \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}} (\nabla^2 f(\mathbf{q}_k) - \langle \nabla f(\mathbf{q}_k), \mathbf{q}_k \rangle \mathbf{I}) \mathcal{P}_{T_{\mathbf{q}_k} \mathbb{S}^{n-1}}$$

are the Riemannian gradient and Riemannian Hessian of f w.r.t. \mathbb{S}^{n-1} , respectively (Absil et al., 2007; 2009), turning (4.6) into the form of familiar quadratic approximation, as described in (4.2).

Then the Riemannian trust-region subproblem is

$$\min_{\boldsymbol{\delta} \in T_{\mathbf{q}_k} \mathbb{S}^{n-1}, \|\boldsymbol{\delta}\| \leq \Delta} \widehat{f}(\mathbf{q}_k, \boldsymbol{\delta}), \quad (4.7)$$

where $\Delta > 0$ is the familiar trust-region parameter. Taking any orthonormal basis $\mathbf{U}_{\mathbf{q}_k}$ for $T_{\mathbf{q}_k} \mathbb{S}^{n-1}$, we can transform (4.7) into a classical trust-region subproblem:

$$\min_{\|\boldsymbol{\xi}\| \leq \Delta} \widehat{f}(\mathbf{q}_k, \mathbf{U}_{\mathbf{q}_k} \boldsymbol{\xi}), \quad (4.8)$$

for which very efficient numerical algorithms exist (Moré & Sorensen, 1983; Hazan & Koren, 2014). Once we obtain the minimizer $\boldsymbol{\xi}_*$, we set $\boldsymbol{\delta}_* = \mathbf{U}_{\mathbf{q}_k} \boldsymbol{\xi}_*$, which solves (4.7).

One additional issue as compared to the Euclidean setting is that now $\boldsymbol{\delta}_*$ is one vector in the tangent space and additive update leads to a point outside the manifold. To resolve this, we resort to the natural exponential map:

$$\mathbf{q}_{k+1} \doteq \exp_{\mathbf{q}_k} \boldsymbol{\delta}_* = \mathbf{q}_k \cos \|\boldsymbol{\delta}_*\| + \frac{\boldsymbol{\delta}_*}{\|\boldsymbol{\delta}_*\|} \sin \|\boldsymbol{\delta}_*\|, \quad (4.9)$$

which move the sequence to the next iterate ‘‘along the direction’’⁷ of $\boldsymbol{\delta}_*$ while staying over the sphere.

There are many variants of (Riemannian) TRM that allow one to solve the subproblem 4.8 only approximately while still guarantee convergence. For simplicity, we avoid the extra burden caused thereof for analysis by solving the subproblem exactly via SDP relaxation: introduce

$$\tilde{\boldsymbol{\xi}} = [\boldsymbol{\xi}^*, 1]^*, \quad \boldsymbol{\Xi} = \tilde{\boldsymbol{\xi}} \tilde{\boldsymbol{\xi}}^*, \quad \mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{b} \\ \mathbf{b}^* & 0 \end{bmatrix} \quad (4.10)$$

where $\mathbf{A} = \mathbf{U}^* (\nabla^2 f(\mathbf{q}) - \langle \nabla f(\mathbf{q}), \mathbf{q} \rangle \mathbf{I}) \mathbf{U}$ and $\mathbf{b} = \mathbf{U}^* \nabla f(\mathbf{q})$ for any orthobasis \mathbf{U} of $T_{\mathbf{q}_{k-1}} \mathbb{S}^{n-1}$. The subproblem is known to be equivalent to the SDP problem (Fortin & Wolkowicz, 2004):

$$\begin{aligned} & \min_{\boldsymbol{\Xi}} \langle \mathbf{M}, \boldsymbol{\Xi} \rangle, \\ & \text{s.t. } \text{tr}(\boldsymbol{\Xi}) \leq \Delta^2 + 1, \quad \langle \mathbf{E}_{n+1}, \boldsymbol{\Xi} \rangle = 1, \quad \boldsymbol{\Xi} \succeq \mathbf{0}, \end{aligned} \quad (4.11)$$

where $\mathbf{E}_{n+1} = \mathbf{e}_{n+1} \mathbf{e}_{n+1}^*$. The detailed trust region algorithm is presented in Algorithm 1.

4.3. Algorithmic results

Using the geometric characterization in Theorem 3.1 and Theorem 3.3, we prove that when the parameter Δ is suf-

⁷Technically, moving along a curve on the manifold of which $\boldsymbol{\delta}_*$ is the initial tangent vector in certain canonical way.

Algorithm 1 Trust Region Method for Finding a Single Sparse Vector

Input: data matrix $\mathbf{Y} \in \mathbb{R}^{n \times p}$, smoothing parameter μ and parameters $\eta_{vs} = 0.9$, $\eta_s = 0.1$, $\gamma_i = 2$, $\gamma_d = 1/2$, $\Delta_{\max} = 1$, and $\Delta_{\min} = 10^{-16}$.

Output: $\hat{\mathbf{q}} \in \mathbb{S}^{n-1}$

- 1: Initialize $\mathbf{q}^{(0)} \in \mathbb{S}^{n-1}$, $\Delta^{(0)}$ and $k = 1$,
- 2: **while** not converged **do**
- 3: Set $\mathbf{U} \in \mathbb{R}^{n \times (n-1)}$ to be an orthobasis for $\mathbf{q}^{(k-1)\perp}$
- 4: Solve the trust region subproblem

$$\hat{\boldsymbol{\xi}} = \arg \min_{\|\boldsymbol{\xi}\| \leq \Delta^{(k-1)}} \hat{f}(\mathbf{q}^{(k-1)}, \mathbf{U}\boldsymbol{\xi}) \quad (4.12)$$

- 5: Set

$$\begin{aligned} \hat{\boldsymbol{\delta}} &\leftarrow \mathbf{U}\hat{\boldsymbol{\xi}}, \\ \hat{\mathbf{q}} &\leftarrow \mathbf{q}^{(k-1)} \cos \|\hat{\boldsymbol{\delta}}\| + \frac{\hat{\boldsymbol{\delta}}}{\|\hat{\boldsymbol{\delta}}\|} \sin \|\hat{\boldsymbol{\delta}}\|. \end{aligned}$$

- 6: Set

$$\rho_k \leftarrow \frac{f(\mathbf{q}^{(k-1)}) - f(\hat{\mathbf{q}})}{f(\mathbf{q}^{(k-1)}) - \hat{f}(\mathbf{q}^{(k-1)}, \hat{\boldsymbol{\delta}})} \quad (4.13)$$

- 7: **if** $\rho_k \geq \eta_{vs}$ **then**
- 8: Set $\mathbf{q}^{(k)} \leftarrow \hat{\mathbf{q}}$, $\Delta^{(k)} \leftarrow \min(\gamma_i \Delta^{(k-1)}, \Delta_{\max})$.
- 9: **else if** $\rho_k \geq \eta_s$ **then**
- 10: Set $\mathbf{q}^{(k)} \leftarrow \hat{\mathbf{q}}$, $\Delta^{(k)} \leftarrow \Delta^{(k-1)}$.
- 11: **else**
- 12: Set $\mathbf{q}^{(k)} \leftarrow \mathbf{q}^{(k-1)}$, $\Delta^{(k)} \leftarrow \max(\gamma_d \Delta^{(k-1)}, \Delta_{\min})$.
- 13: **end if**
- 14: Set $k = k + 1$.
- 15: **end while**

ficiently small⁸, (1) the trust region step induces at least a fixed amount of decrease to the objective value in the negative curvature and nonzero gradient region; (2) the trust region iterate sequence will eventually move to and stay in the strongly convex region, and converge to the global minima with an asymptotic quadratic rate. In particular, the geometry implies that from *any initialization*, the iterate sequence converges to a close approximation to one local minimizer in a polynomial number of steps.

The following two theorems collect these results, for orthogonal and general complete \mathbf{A}_0 , respectively.

Theorem 4.1 (Orthogonal dictionary). *Suppose the dictionary \mathbf{A}_0 is orthogonal. Then there exists a posi-*

⁸For simplicity of analysis, we have assumed Δ is fixed throughout the analysis. In practice, dynamic updates to Δ tends to lead to faster convergence.

tive constant C , such that for all $\theta \in (0, 1/2)$, and $\mu < \min\{c_a \theta n^{-1}, c_b n^{-5/4}\}$, whenever $\exp(n) \geq p \geq C n^3 \log \frac{n}{\mu \theta} / (\mu^2 \theta^2)$, with probability at least $1 - 8n^2 p^{-10} - \theta(np)^{-7} - \exp(-0.3\theta np) - p^{-10} - c_c \exp(-c_d p \mu^2 \theta^2 / n^2)$, the Riemannian trust-region algorithm with input data matrix $\hat{\mathbf{Y}} = \mathbf{Y}$, any initialization $\mathbf{q}^{(0)}$ on the sphere, and a step size satisfying

$$\Delta \leq \min \left\{ \frac{c_e c_* \theta \mu^2}{n^{5/2} \log^{3/2}(np)}, \frac{c_f c_*^3 \theta^3 \mu}{n^{7/2} \log^{7/2}(np)} \right\}.$$

returns a solution $\hat{\mathbf{q}} \in \mathbb{S}^{n-1}$ which is ε near to one of the local minimizers \mathbf{q}_* (i.e., $\|\hat{\mathbf{q}} - \mathbf{q}_*\| \leq \varepsilon$) in

$$\begin{aligned} \max \left\{ \frac{c_g n^6 \log^3(np)}{c_*^3 \theta^3 \mu^4}, \frac{c_h n}{c_*^2 \theta^2 \Delta^2} \right\} & \left(f(\mathbf{q}^{(0)}) - f(\mathbf{q}_*) \right) \\ & + \log \log \frac{c_i c_* \theta \mu}{\varepsilon n^{3/2} \log^{3/2}(np)} \end{aligned}$$

iterations. Here c_* is as defined in Theorem 3.1, and c_a, c_b are the same numerical constants as defined in Theorem 3.1, c_c to c_i are other positive numerical constants.

Proofs for the complete case basically follows from that the slight perturbation of structure parameters as summarized in Theorem 3.3 (vs. Theorem 3.1) change all algorithm parameter by at most small multiplicative constants.

Theorem 4.2 (Complete dictionary). *Suppose the dictionary \mathbf{A}_0 is complete with condition number $\kappa(\mathbf{A}_0)$. There exists a positive constant C , such that for all $\theta \in (0, 1/2)$, and $\mu < \min\{c_a \theta n^{-1}, c_b n^{-5/4}\}$, whenever $\exp(n) \geq p \geq \frac{C}{c_*^2 \theta} \max\left\{\frac{n^4}{\mu^4}, \frac{n^5}{\mu^2}\right\} \kappa^8(\mathbf{A}_0) \log^4\left(\frac{\kappa(\mathbf{A}_0)n}{\mu \theta}\right)$, with probability at least $1 - 8n^2 p^{-10} - \theta(np)^{-7} - \exp(-0.3\theta np) - 2p^{-8} - c_c \exp(-c_d p \mu^2 \theta^2 / n^2)$, the Riemannian trust-region algorithm with input data matrix $\bar{\mathbf{Y}} \doteq \sqrt{p\theta} (\mathbf{Y}\mathbf{Y}^*)^{-1/2} \mathbf{Y}$ where $\mathbf{U}\Sigma\mathbf{V}^* = \text{SVD}(\mathbf{A}_0)$, any initialization $\mathbf{q}^{(0)}$ on the sphere and a step size satisfying*

$$\Delta \leq \min \left\{ \frac{c_e c_* \theta \mu^2}{n^{5/2} \log^{3/2}(np)}, \frac{c_f c_*^3 \theta^3 \mu}{n^{7/2} \log^{7/2}(np)} \right\}.$$

returns a solution $\hat{\mathbf{q}} \in \mathbb{S}^{n-1}$ which is ε near to one of the local minimizers \mathbf{q}_* (i.e., $\|\hat{\mathbf{q}} - \mathbf{q}_*\| \leq \varepsilon$) in

$$\begin{aligned} \max \left\{ \frac{c_g n^6 \log^3(np)}{c_*^3 \theta^3 \mu^4}, \frac{c_h n}{c_*^2 \theta^2 \Delta^2} \right\} & \left(f(\mathbf{q}^{(0)}) - f(\mathbf{q}_*) \right) \\ & + \log \log \frac{c_i c_* \theta \mu}{\varepsilon n^{3/2} \log^{3/2}(np)} \end{aligned}$$

iterations. Here c_* is as defined in Theorem 3.1, and c_a, c_b are the same numerical constants as defined in Theorem 3.1, c_c to c_i are other positive numerical constants.

5. Main Results

For orthogonal dictionaries, from Theorem 3.1 and its corollary, we know that all the minimizers \hat{q}_* are $O(\mu)$ away from their respective nearest “target” q_* , with $q_*^* \hat{Y} = \alpha e_i^* X_0$ for certain $\alpha \neq 0$ and $i \in [n]$; in Theorem ??, we have shown that w.h.p. the Riemannian TRM algorithm produces a solution $\hat{q} \in \mathbb{S}^{n-1}$ that is ε away to one of the minimizers, say \hat{q}_* . Thus, the \hat{q} returned by the TRM algorithm is $O(\varepsilon + \mu)$ away from q_* . For exact recovery, we use a simple linear programming rounding procedure, which guarantees to exactly produce the optimizer q_* . We then use deflation to sequentially recover other rows of X_0 . Overall, w.h.p. both the dictionary A_0 and sparse coefficient X_0 are exactly recovered up to sign permutation, when $\theta \in \Omega(1)$, for orthogonal dictionaries. The same procedure can be used to recover complete dictionaries, though the analysis is slightly more complicated. Our overall algorithmic pipeline for recovering orthogonal dictionaries is sketched as follows.

- 1. Estimating one row of X_0 by the Riemannian TRM algorithm.** By Theorem 3.1 (resp. Theorem 3.3) and Theorem 4.1 (resp. Theorem 4.2), starting from any, when the relevant parameters are set appropriately (say as μ_* and Δ_*), w.h.p., our Riemannian TRM algorithm finds a local minimizer \hat{q} , with q_* the nearest target that exactly recovers one row of X_0 and $\|\hat{q} - q_*\| \in O(\mu)$ (by setting the target accuracy of the TRM as, say, $\varepsilon = \mu$).
- 2. Recovering one row of X_0 by rounding.** To obtain the target solution q_* and hence recover (up to scale) one row of X_0 , we solve the following linear program:

$$\min_q \|q^* \hat{Y}\|_1, \quad \text{s.t.} \quad \langle r, q \rangle = 1, \quad (5.1)$$

with $r = \hat{q}$. We show that when $\langle \hat{q}, q_* \rangle$ is sufficiently large, implied by μ being sufficiently small, w.h.p. the minimizer of (5.1) is exactly q_* , and hence one row of X_0 is recovered by $q_*^* \hat{Y}$.

- 3. Recovering all rows of X_0 by deflation.** Once ℓ rows of X_0 ($1 \leq \ell \leq n - 2$) have been recovered, say, by unit vectors q_*^1, \dots, q_*^ℓ , one takes an orthonormal basis U for $[\text{span}(q_*^1, \dots, q_*^\ell)]^\perp$, and minimizes the new function $h(z) \doteq f(Uz; \hat{Y})$ on the sphere $\mathbb{S}^{n-\ell-1}$ with the Riemannian TRM algorithm (though conservative, one can again set parameters as μ_* , Δ_* , as in Step 1) to produce a \hat{z} . Another row of X_0 is then recovered via the LP rounding (5.1) with input $r = U\hat{z}$ (to produce $q_*^{\ell+1}$). Finally, by repeating the procedure until depletion, one can recover all the rows of X_0 .
- 4. Reconstructing the dictionary A_0 .** By solving the

linear system $Y = AX_0$, one can obtain the dictionary $A_0 = YX_0^*(X_0X_0^*)^{-1}$.

Formally, we have the following results:

Theorem 5.1 (Orthogonal Dictionary). *Assume the dictionary A_0 is orthogonal and we take $\hat{Y} = Y$. Suppose $\theta \in (0, 1/3)$, $\mu_* < \min\{c_a\theta n^{-1}, c_b n^{-5/4}\}$, and $p \geq Cn^3 \log \frac{n}{\mu_*\theta} / (\mu_*^2\theta^2)$. The above algorithmic pipeline with parameter setting*

$$\Delta_* \leq \min \left\{ \frac{c_c c_* \theta \mu_*^2}{n^{5/2} \log^{5/2}(np)}, \frac{c_d c_*^3 \theta^3 \mu_*}{n^{7/2} \log^{7/2}(np)} \right\},$$

recovers the dictionary A_0 and X_0 in polynomial time, with failure probability bounded by $c_e p^{-6}$. Here c_ is as defined in Theorem 3.1, and c_a through c_e , and C are all positive numerical constants.*

Theorem 5.2 (Complete Dictionary). *Assume the dictionary A_0 is complete with condition number $\kappa(A_0)$ and we take $\hat{Y} = \bar{Y}$. Suppose $\theta \in (0, 1/3)$, $\mu_* < \min\{c_a\theta n^{-1}, c_b n^{-5/4}\}$, and $p \geq \frac{C}{c_f^2 \theta} \max\left\{\frac{n^4}{\mu_*^4}, \frac{n^5}{\mu_*^2}\right\} \kappa^8(A_0) \log^4\left(\frac{\kappa(A_0)n}{\mu\theta}\right)$. The algorithmic pipeline with parameter setting*

$$\Delta_* \leq \min \left\{ \frac{c_c c_* \theta \mu_*^2}{n^{5/2} \log^{5/2}(np)}, \frac{c_d c_*^3 \theta^3 \mu_*}{n^{7/2} \log^{7/2}(np)} \right\},$$

recovers the dictionary A_0 and X_0 in polynomial time, with failure probability bounded by $c_e p^{-6}$. Here c_ is as defined in Theorem 3.1, and c_a through c_f , and C are all positive numerical constants.*

6. Numerical Results

To corroborate our theory, we experiment with dictionary recovery on simulated data⁹. For simplicity, we focus on recovering orthogonal dictionaries and we declare success once a single row of the coefficient matrix is recovered.

Since the problem is invariant to rotations, w.l.o.g. we set the dictionary as $A_0 = I \in \mathbb{R}^{n \times n}$. We fix $p = 5n^2 \log(n)$, and each column of the coefficient matrix $X_0 \in \mathbb{R}^{n \times p}$ has exactly k nonzero entries, chosen uniformly random from $\binom{[n]}{k}$. These nonzero entries are i.i.d. standard normals. This is slightly different from the Bernoulli-Gaussian model we assumed for analysis. For n reasonably large, these two models produce similar behavior. For the sparsity surrogate defined in (2.3), we fix the parameter $\mu = 10^{-2}$. We implement Algorithm 1 with adaptive step size instead of the fixed step size in our analysis.

To see how the allowable sparsity level varies with the dimension, which our theory primarily is about, we vary the

⁹The code is available online: https://github.com/sunju/dl_focm

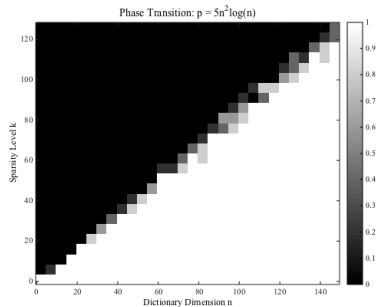


Figure 2. Phase transition for recovering a single sparse vector under the dictionary learning model with $p = 5n^2 \log n$.

dictionary dimension n and the sparsity k both between 1 and 150; for every pair of (k, n) we repeat the simulations independently for $T = 5$ times. Because the optimal solutions are signed coordinate vectors $\{e_i\}_{i=1}^n$, for a solution \hat{q} returned by the TRM algorithm, we define the reconstruction error (RE) to be $\text{RE} = \min_{1 \leq i \leq n} (\|\hat{q} - e_i\|, \|\hat{q} + e_i\|)$. The trial is determined to be a success once $\text{RE} \leq \mu$, with the idea that this indicates \hat{q} is already very near the target and the target can likely be recovered via the LP rounding we described (which we do not implement here). Figure 2 shows the phase transition in the (n, k) plane for the orthogonal case. It is obvious that our TRM algorithm can work well into the linear region whenever $p \in O(n^2 \log n)$. Our analysis is tight up to logarithm factors, and also the polynomial dependency on $1/\mu$, which under the theory is polynomial in n .

7. Discussion

For recovery of complete dictionaries, the LP program approach in (Spielman et al., 2012) that works with $\theta \leq O(1/\sqrt{n})$ only demands $p \geq \Omega(n^2 \log n^2)$, which is recently improved to $p \geq \Omega(n \log^4 n)$ (Luh & Vu, 2015), almost matching the lower bound $\Omega(n \log n)$ (i.e., when $\theta \sim 1/n$). The sample complexity stated in Theorem 5.2 is obviously much higher. It is interesting to see whether such growth in complexity is intrinsic to working in the linear regime. Though our experiments seemed to suggest the necessity of $p \sim O(n^2 \log n)$ even for the orthogonal case, there could be other efficient algorithms that demand much less. Tweaking these three points will likely improve the complexity: (1) The ℓ^1 proxy. The derivative and Hessians of the log cosh function we adopted entail the tanh function, which is not amenable to effective approximation and affects the sample complexity; (2) Geometric characterization and algorithm analysis. It seems working directly on the sphere (i.e., in the q space) could simplify and possibly improve certain parts of the analysis; (3) treating the complete case directly, rather than using (pessimistic) bounds to treat it as a perturbation of the orthogonal case. Particularly, gen-

eral linear transforms may change the space significantly, such that preconditioning and comparing to the orthogonal transforms may not be the most efficient way to proceed.

It is possible to extend the current analysis to other dictionary settings. Our geometric structures and algorithms allow plug-and-play noise analysis. Nevertheless, we believe a more stable way of dealing with noise is to directly extract the whole dictionary, i.e., to consider geometry and optimization (and perturbation) over the orthogonal group. This will require additional nontrivial technical work, but likely feasible thanks to the relatively complete knowledge of the orthogonal group (Edelman et al., 1998; Absil et al., 2009). A substantial leap forward would be to extend the methodology to recovery of *structured* overcomplete dictionaries, such as tight frames. Though there is no natural elimination of one variable, one can consider the marginalization of the objective function wrt the coefficients and work with hidden functions.¹⁰

Under the i.i.d. BG coefficient model, our recovery problem is also an instance of the ICA problem. It is interesting to ask what is vital in making the problem tractable: sparsity or independence. The full version (Sun et al., 2015) includes an experimental study in this direction, which underlines the importance of the sparsity prior. In fact, the preliminary experiments there suggest the independence assumption we made here likely can be removed without losing the favorable geometric structures. In addition, the connection to ICA also suggests the possibility of adapting our geometric characterization and algorithms to the ICA problem. This likely will provide new theoretical insights and computational schemes to ICA.

In the surge of theoretical understanding of nonconvex heuristics (Keshavan et al., 2010; Jain et al., 2013; Hardt, 2014; Hardt & Wootters, 2014; Netrapalli et al., 2014; Jain & Netrapalli, 2014; Netrapalli et al., 2013; Candes et al., 2014; Jain & Oh, 2014; Anandkumar et al., 2014; Yi et al., 2013; Lee et al., 2013; Qu et al., 2014; Lee et al., 2013; Agarwal et al., 2013a;b; Arora et al., 2013; 2015; 2014), the initialization plus local refinement strategy mostly differs from practice, whereby random initializations seem to work well, and the analytic techniques developed are mostly fragmented and highly specialized. The analytic and algorithmic we developed here hold promise to provide a coherent account of these problems. It is interesting to see to what extent we can streamline and generalize the framework.

¹⁰This recent work (Arora et al., 2015) on overcomplete DR has used a similar idea. The marginalization taken there is near to the global optimum of one variable, where the function is well-behaved. Studying the global properties of the marginalization may introduce additional challenges.

Acknowledgments

This work was partially supported by grants ONR N00014-13-1-0492, NSF 1343282, and funding from the Moore and Sloan Foundations and the Wei Family Private Foundation. We thank the area chair and the anonymous reviewers for making painstaking effort to read our long proofs and providing insightful feedback. We also thank Cun Mu and Henry Kuo for discussions related to this project.

References

- Absil, P.-A., Baker, C. G., and Gallivan, K. A. Trust-region methods on riemannian manifolds. *Foundations of Computational Mathematics*, 7(3):303–330, 2007.
- Absil, Pierre-Antoine, Mahoney, Robert, and Sepulchre, Rodolphe. *Optimization Algorithms on Matrix Manifolds*. Princeton University Press, 2009.
- Agarwal, Alekh, Anandkumar, Animashree, Jain, Prateek, Netrapalli, Praneeth, and Tandon, Rashish. Learning sparsely used overcomplete dictionaries via alternating minimization. *arXiv preprint arXiv:1310.7991*, 2013a.
- Agarwal, Alekh, Anandkumar, Animashree, and Netrapalli, Praneeth. Exact recovery of sparsely used overcomplete dictionaries. *arXiv preprint arXiv:1309.1952*, 2013b.
- Anandkumar, Animashree, Ge, Rong, and Janzamin, Majid. Guaranteed non-orthogonal tensor decomposition via alternating rank-1 updates. *arXiv preprint arXiv:1402.5180*, 2014.
- Arora, Sanjeev, Ge, Rong, and Moitra, Ankur. New algorithms for learning incoherent and overcomplete dictionaries. *arXiv preprint arXiv:1308.6273*, 2013.
- Arora, Sanjeev, Bhaskara, Aditya, Ge, Rong, and Ma, Tengyu. More algorithms for provable dictionary learning. *arXiv preprint arXiv:1401.0579*, 2014.
- Arora, Sanjeev, Ge, Rong, Ma, Tengyu, and Moitra, Ankur. Simple, efficient, and neural algorithms for sparse coding. *arXiv preprint arXiv:1503.00778*, 2015.
- Barak, Boaz, Kelner, Jonathan A, and Steurer, David. Dictionary learning and tensor decomposition via the sum-of-squares method. *arXiv preprint arXiv:1407.1543*, 2014.
- Candes, Emmanuel, Li, Xiaodong, and Soltanolkotabi, Mahdi. Phase retrieval via wirtinger flow: Theory and algorithms. *arXiv preprint arXiv:1407.1065*, 2014.
- Conn, Andrew R., Gould, Nicholas I. M., and Toint, Philippe L. *Trust-region Methods*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2000. ISBN 0-89871-460-5.
- Demanet, Laurent and Hand, Paul. Scaling law for recovering the sparsest element in a subspace. *Information and Inference*, 3(4):295–309, 2014.
- Edelman, Alan, Arias, Tomás A, and Smith, Steven T. The geometry of algorithms with orthogonality constraints. *SIAM journal on Matrix Analysis and Applications*, 20(2):303–353, 1998.
- Elad, Michael. *Sparse and redundant representations: from theory to applications in signal and image processing*. Springer, 2010.
- Fortin, Charles and Wolkowicz, Henry. The trust region subproblem and semidefinite programming*. *Optimization methods and software*, 19(1):41–67, 2004.
- Geng, Quan and Wright, John. On the local correctness of ℓ^1 -minimization for dictionary learning. Submitted to *IEEE Transactions on Information Theory*, 2011. Preprint: <http://www.columbia.edu/~jw2966>.
- Gribonval, Rémi and Schnass, Karin. Dictionary identification - sparse matrix-factorization via ℓ^1 -minimization. *IEEE Transactions on Information Theory*, 56(7):3523–3539, 2010.
- Gribonval, Rémi, Jenatton, Rodolphe, and Bach, Francis. Sparse and spurious: dictionary learning with noise and outliers. *arXiv preprint arXiv:1407.5155*, 2014.
- Hardt, Moritz. Understanding alternating minimization for matrix completion. In *Foundations of Computer Science (FOCS), 2014 IEEE 55th Annual Symposium on*, pp. 651–660. IEEE, 2014.
- Hardt, Moritz and Wootters, Mary. Fast matrix completion without the condition number. In *Proceedings of The 27th Conference on Learning Theory*, pp. 638–678, 2014.
- Hazan, Elad and Koren, Tomer. A linear-time algorithm for trust region problems. *arXiv preprint arXiv:1401.6757*, 2014.
- Hillar, Christopher and Sommer, Friedrich T. When can dictionary learning uniquely recover sparse data from subsamples? *arXiv preprint arXiv:1106.3616*, 2011.
- Jain, Prateek and Netrapalli, Praneeth. Fast exact matrix completion with finite samples. *arXiv preprint arXiv:1411.1087*, 2014.
- Jain, Prateek and Oh, Sewoong. Provable tensor factorization with missing data. In *Advances in Neural Information Processing Systems*, pp. 1431–1439, 2014.

- Jain, Prateek, Netrapalli, Praneeth, and Sanghavi, Sujay. Low-rank matrix completion using alternating minimization. In *Proceedings of the forty-fifth annual ACM symposium on Theory of Computing*, pp. 665–674. ACM, 2013.
- Keshavan, Raghunandan H, Montanari, Andrea, and Oh, Sewoong. Matrix completion from a few entries. *Information Theory, IEEE Transactions on*, 56(6):2980–2998, 2010.
- Lee, Kiryung, Wu, Yihong, and Bresler, Yoram. Near optimal compressed sensing of sparse rank-one matrices via sparse power factorization. *arXiv preprint arXiv:1312.0525*, 2013.
- Luh, Kyle and Vu, Van. Dictionary learning with few samples and matrix concentration. *arXiv preprint arXiv:1503.08854*, 2015.
- Mairal, Julien, Bach, Francis, and Ponce, Jean. Sparse modeling for image and vision processing. *Foundations and Trends in Computer Graphics and Vision*, 8(2-3), 2014.
- Mehta, Nishant and Gray, Alexander G. Sparsity-based generalization bounds for predictive sparse coding. *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, 28(1):36–44, 2013.
- Moré, J. J. and Sorensen, D. C. Computing a trust region step. *SIAM J. Scientific and Statistical Computing*, 4: 553–572, 1983.
- Murty, Katta G and Kabadi, Santosh N. Some np-complete problems in quadratic and nonlinear programming. *Mathematical programming*, 39(2):117–129, 1987.
- Netrapalli, Praneeth, Jain, Prateek, and Sanghavi, Sujay. Phase retrieval using alternating minimization. In *Advances in Neural Information Processing Systems*, pp. 2796–2804, 2013.
- Netrapalli, Praneeth, Niranjan, UN, Sanghavi, Sujay, Anandkumar, Animashree, and Jain, Prateek. Non-convex robust pca. In *Advances in Neural Information Processing Systems*, pp. 1107–1115, 2014.
- Nocedal, Jorge and Wright, Stephen. *Numerical Optimization*. Springer, 2006.
- Qu, Qing, Sun, Ju, and Wright, John. Finding a sparse vector in a subspace: Linear sparsity using alternating directions. In *Advances in Neural Information Processing Systems*, pp. 3401–3409, 2014.
- Schnass, Karin. On the identifiability of overcomplete dictionaries via the minimisation principle underlying k-svd. *Applied and Computational Harmonic Analysis*, 37(3): 464–491, 2014a.
- Schnass, Karin. Local identification of overcomplete dictionaries. *arXiv preprint arXiv:1401.6354*, 2014b.
- Schnass, Karin. Convergence radius and sample complexity of itkm algorithms for dictionary learning. *arXiv preprint arXiv:1503.07027*, 2015.
- Spielman, Daniel A, Wang, Huan, and Wright, John. Exact recovery of sparsely-used dictionaries. In *Proceedings of the 25th Annual Conference on Learning Theory*, 2012.
- Sun, Ju, Qu, Qing, and Wright, John. Complete dictionary recovery over the sphere. *arXiv preprint arXiv:1504.06785*, 2015.
- Vainsencher, Daniel, Mannor, Shie, and Bruckstein, Alfred M. The sample complexity of dictionary learning. *J. Mach. Learn. Res.*, 12:3259–3281, November 2011. ISSN 1532-4435.
- Yi, Xinyang, Caramanis, Constantine, and Sanghavi, Sujay. Alternating minimization for mixed linear regression. *arXiv preprint arXiv:1310.3745*, 2013.
- Zibulevsky, Michael and Pearlmutter, Barak. Blind source separation by sparse decomposition in a signal dictionary. *Neural computation*, 13(4):863–882, 2001.