
Complete Dictionary Learning via ℓ_p -norm Maximization

Yifei Shen^{1*}, Ye Xue^{1*}, Jun Zhang², Khaled B. Letaief¹, Vincent Lau¹

¹Hong Kong University of Science and Technology, Hong Kong, China, {yshenaw, ye.xue, eekhaled, eeknlau}@ust.hk

²The Hong Kong Polytechnic University, Hong Kong, China, jun-eie.zhang@polyu.edu.hk

Abstract

Dictionary learning is a classic representation learning method that has been widely applied in signal processing and data analytics. In this paper, we investigate a family of ℓ_p -norm ($p > 2, p \in \mathbb{N}$) maximization approaches for the complete dictionary learning problem from theoretical and algorithmic aspects. Specifically, we prove that the global maximizers of these formulations are very close to the true dictionary with high probability, even when Gaussian noise is present. Based on the generalized power method (GPM), an efficient algorithm is then developed for the ℓ_p -based formulations. We further show the efficacy of the developed algorithm: for the population GPM algorithm over the sphere constraint, it first quickly enters the neighborhood of a global maximizer, and then converges linearly in this region. Extensive experiments demonstrate that the ℓ_p -based approaches enjoy a higher computational efficiency and better robustness than conventional approaches and $p = 3$ performs the best.

1 INTRODUCTION

Dictionary learning is a classic unsupervised representation learning method [16]. Given data \mathbf{Y} , it identifies a representation basis \mathbf{D}_0 and the corresponding coefficients \mathbf{X}_0 such that $\mathbf{Y} \approx \mathbf{D}_0 \mathbf{X}_0$ with \mathbf{X}_0 being sufficiently sparse. Given its powerful capability of exploiting low-dimensional structures in high-dimensional data, dictionary learning has found wide applications in signal and image processing [8].

Solving the dictionary learning problem inevitably involves non-convex optimization [20], and is thus highly challenging. A key ingredient underlying the recent advancement in this area is innovative mathematical formulations. A natural formulation is to minimize the ℓ_0 norm of \mathbf{X}_0 for inducing sparsity. However, ℓ_0 minimization with a non-convex constraint is computationally intractable. Thus, surrogate objective functions are explored, which ideally should come with theoretical guarantees for recovering the dictionary, as well as efficient algorithms. In particular, ℓ_1 -norm minimization based formulations, which promote sparsity in problems such as compressive sensing, have been widely adopted in dictionary learning [3, 10, 21, 24]. While such formulations enjoy theoretical guarantees, they face computational challenges for high-dimensional data. Particularly, existing algorithms can only deal with one row of a dictionary at a time. Given the high computational complexity of classic algorithms for recovering one row, e.g., the Riemannian trust region algorithm or subgradient descent, repeatedly solving the problem for n times to recover the whole dictionary leads to prohibitive complexity. Additionally, ℓ_1 -based methods are known to be sensitive to noise in the observation [24, 25]. Thus, formulations that lead to more efficient and robust methods are needed.

Recently, a novel ℓ_4 -norm maximization formulation was proposed in [25, 26], which is able to recover the entire dictionary at once. It was shown that the global maximizers of the ℓ_4 -based formulation are very close to the true dictionary. Moreover, the concaveness of the formulation enables a fast fixed-point type algorithm, named *matching, stretching, and projection* (MSP), which achieves hundreds of times speedup compared with existing methods. This new formulation is motivated by the fact that maximizing ℓ_{2k+2} -norm promotes spikiness and sparsity [13, 27]. In experiments, ℓ_4 -norm was found to be the best among all the ℓ_{2k+2} -norms in terms of the sample complexity and computa-

* These two authors contributed equally.

tional efficiency.

The essence of the ℓ_{2k+2} -norm approach is to maximize an ℓ_{2k+2} -norm over an ℓ_2 -norm constraint. In principle, maximizing *any* higher-order norm over a lower-order norm constraint leads to sparse and spiky solutions. Thus, we conjecture that the dictionary can be recovered via maximizing *any* ℓ_p -norm ($p > 2$), which is tested numerically in Fig. 1. It is demonstrated that the dictionary is recovered with high probability for $p = 3, 4, 5, 6$. Particularly, the ℓ_3 -based formulation enjoys the lowest sample complexity¹. These observations lead to the following intriguing questions:

- Can all the ℓ_p -norm maximization based formulations provably recover the true dictionary?
- Which p should we pick for practical applications?
- What advantages do they enjoy over the traditional approaches, e.g., ℓ_1 -based approaches?

Unfortunately, the analysis in [26] cannot be extended to ℓ_{2k+1} -norms, and thus could not address the above questions. In this paper, we endeavor to develop more general theoretical results for ℓ_p -norm ($p > 2, p \in \mathbb{N}$) based formulations², which will lead to a better understanding of such methods and also provide guidelines to find more efficient algorithms.

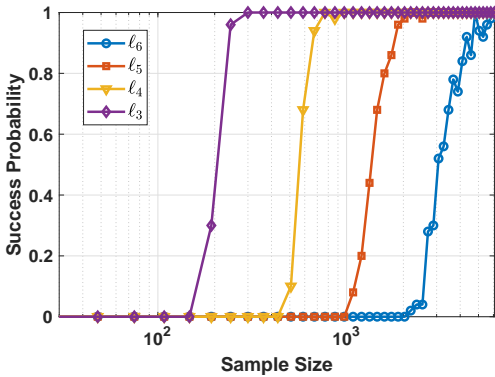


Figure 1: The successful dictionary recover probability with different values of p when $n = 30$ and $\theta = 0.3$.

1.1 Contributions

In this paper, we study the ℓ_p -norm ($p > 2, p \in \mathbb{N}$) maximization based dictionary learning.

¹Sample complexity here means the minimal required number of samples to successfully recover the true dictionary.

²Our analysis can be extended to $2 < p < \infty$ with minor modification. However, due to the high computational complexity of taking fractional power, we will not discuss $p \notin \mathbb{N}$ in this paper.

- We prove that as long as the number of samples is larger than $\Omega\left(nk^{\frac{p}{2}} \log^{\frac{p}{2}+1}(n)\right)$, where n is the size of the dictionary and k is the sparsity level, the global maximizers of all ℓ_p -norm maximization based formulations are very close to the true dictionary with high probability. The dictionary can be recovered even in the presence of the Gaussian noise.
- An efficient algorithm is developed based on the generalized power method [12], which applies to $p > 2$. It is proved that the population generalized power method enjoys a desirable global convergence behavior over the sphere constraint. Specifically, the convergence involves two stages, where the first stage only takes a few iterations and the second stage enjoys a linear convergence rate.
- To guide the practical application, we prove that the ℓ_3 -based approach enjoys the lowest sample complexity and is the most robust among all the ℓ_p -norm maximization based formulations. The experiments will further demonstrate that the ℓ_3 -based approach is also more time-efficient and robust than existing methods, including K-SVD, ℓ_1 and ℓ_4 -based approaches.

1.2 Notations and Terminologies

Asymptotic notations: Throughout the paper, $f(n) = \mathcal{O}(g(n))$ means that there exists a constant $c > 0$ such that $f(n) \leq c|g(n)|$; $f(n) = \Theta(g(n))$ means that there exists constants c_1, c_2 such that $c_1|g(n)| \leq f(n) \leq c_2|g(n)|$; $f(n) = \Omega(g(n))$ means that there exists constants $c_3 > 0$ such that $c_3|g(n)| \leq f(n)$

Norms: $\|\cdot\|_p$ is the element-wise p -th norm of a vector or matrix, i.e., $\|A\|_p = (\sum_i \sum_j |A_{ij}|^p)^{\frac{1}{p}}$. Likewise, $\|\cdot\|_F$ and $\|\cdot\|$ denote the Frobenius norm and operator norm of a matrix, respectively.

Stiefel manifold: The Stiefel manifold $\text{St}(n, m)$ is defined as the subspace of orthonormal N -frames in \mathbb{R}^n , namely,

$$\text{St}(n, m) = \{\Gamma \in \mathbb{R}^{n \times m} : \Gamma^* \Gamma = I_m\} \quad (1)$$

where I_m is the $m \times m$ identity matrix. We further denote the orthogonal group as $\mathbb{O}(n) = \text{St}(n, n)$

Distributions: We denote a Bernoulli distribution with θ non-zero probability as $\text{Ber}(\theta)$. The Bernoulli-Gaussian distribution is denoted by $\mathcal{BG}(\theta)$ and defined as $x = b \cdot g$, where $b \sim \text{Ber}(\theta)$ and $g \sim \mathcal{N}(0, 1)$.

Sign-permutation ambiguity: Note that for a sparse matrix \mathbf{X}_0 and any signed permutation matrix \mathbf{P} , \mathbf{X}_0 and $\mathbf{X}_0\mathbf{P}$ are equally sparse. Thus, we consider that a dictionary \mathbf{D}_0 is successfully recovered if we find any signed permutation of \mathbf{D}_0 .

2 ℓ_p -NORM MAXIMIZATION BASED COMPLETE DICTIONARY LEARNING

In this section, we study the statistical performance of ℓ_p -based formulations of complete dictionary learning and investigate their robustness.

2.1 Problem Formulation

We consider an orthogonal (complete) dictionary learning problem with a Bernoulli-Gaussian model, with the following three justifications. First, it has been demonstrated that the performance of complete bases is competitive to over-complete dictionaries in real applications [4]. Second, the complete dictionary learning problem can be converted into the orthogonal case through a simple preconditioning [21]. Third, the Bernoulli-Gaussian model is a reasonable model for generic sparse coefficients [3, 19, 21, 25, 26].

Specifically, we assume that each sample $\mathbf{y}_i \in \mathbb{R}^n$ is generated from a sparse superposition of an orthogonal dictionary $\mathbf{D}_0 \in \mathbb{O}(n)$, i.e., $\mathbf{y}_i = \mathbf{D}_0\mathbf{x}_i$, where each element in $\mathbf{x}_i \in \mathbb{R}^n$ is i.i.d. Bernoulli-Gaussian, i.e., $x_{i,j} \sim \mathcal{BG}(\theta)$. Denote $\mathbf{Y} = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_r\}$, $\mathbf{X}_0 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r\}$ and thus $\mathbf{Y} = \mathbf{D}_0\mathbf{X}_0$. Our goal is to simultaneously recover the dictionary \mathbf{D}_0 and coefficients \mathbf{X}_0 from the observation \mathbf{Y} . A good estimate of \mathbf{D}_0 , denoted as \mathbf{D} , should maximize the sparsity of the associated coefficients \mathbf{X} . Therefore, a natural ℓ_0 -based formulation of the orthogonal dictionary learning problem is

$$\begin{aligned} & \underset{\mathbf{D}, \mathbf{X}}{\text{minimize}} && \|\mathbf{X}\|_0 \\ & \text{subject to} && \mathbf{Y} = \mathbf{D}\mathbf{X}, \mathbf{D} \in \mathbb{O}(n). \end{aligned} \quad (2)$$

As $\mathbf{D} \in \mathbb{O}(n)$, we can write $\mathbf{X} = \mathbf{D}^*\mathbf{Y}$. Denote $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_n]$ and $\mathbf{A} = \mathbf{D}^*$, and then the formulation (2) can be transformed into

$$\begin{aligned} & \underset{\mathbf{A}}{\text{minimize}} && \|\mathbf{A}\mathbf{Y}\|_0 \\ & \text{subject to} && \mathbf{A} \in \mathbb{O}(n). \end{aligned} \quad (3)$$

Nevertheless, Problem (3) is difficult to solve due to the combinatorial nature of $\|\cdot\|_0$ and the non-convex constraint. To motivate ℓ_p -based formulations, we first consider a simpler case of (3), i.e., solving one column of

the dictionary by

$$\underset{\mathbf{a}}{\text{minimize}} \quad \|\mathbf{a}^*\mathbf{Y}\|_0 \quad \text{subject to} \quad \|\mathbf{a}\|_2 = 1. \quad (4)$$

The essence of the existing heuristic formulations, either the ℓ_1 or ℓ_4 based one, to solve (4) is to promote sparsity over an ℓ_2 constraint. Note that the landscapes of higher-order norms are sharper than lower-order norms. Thus, maximizing *any* higher-order norm over a lower-order norm constraint pushes the variables to extreme values, i.e., 0 or the maximal, which leads to sparse solutions. Fig. 2 illustrates this phenomenon on \mathbb{S}^2 . Therefore, heuristically, we can adopt *any* ℓ_p -based ($p > 2$) formulation for dictionary learning. Specifically, we expect one column of the dictionary to be recovered by solving the following problem

$$\begin{aligned} & \underset{\mathbf{a}}{\text{maximize}} && \|\mathbf{a}^*\mathbf{Y}\|_p^p = \|\mathbf{a}^*\mathbf{D}_0\mathbf{X}_0\|_p^p \\ & \text{subject to} && \|\mathbf{a}\|_2 = 1, \end{aligned} \quad (5)$$

where $p \in \mathbb{N}$ and $p > 2$.

Similarly, m columns of \mathbf{D}_0 can be recovered at once by considering the following optimization problem

$$\begin{aligned} \mathcal{P}: & \underset{\mathbf{A}}{\text{maximize}} && \|\mathbf{A}\mathbf{Y}\|_p^p = \|\mathbf{A}\mathbf{D}_0\mathbf{X}_0\|_p^p \\ & \text{subject to} && \mathbf{A}^* \in \text{St}(n, m). \end{aligned} \quad (6)$$

The whole dictionary can be recovered at once if $m = n$.

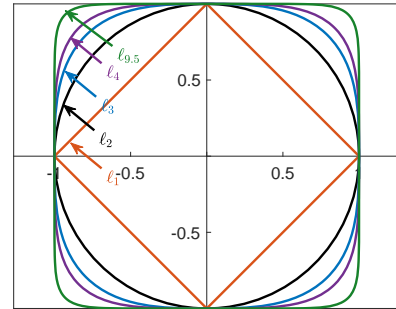


Figure 2: The figure of ℓ_p -norms. Maximizing any ℓ_p -norm ($p > 2$) over the sphere leads to the solutions on the coordinates, i.e., sparse solutions.

2.2 Statistical justification

In this subsection, we offer statistical justification to our ℓ_p -based formulation in (6). We consider $m = n$ for simplicity in this subsection and $m < n$ can be derived in the same way with minor modifications. We prove

that the global maximizers of all ℓ_p -based formulations are very close to the true dictionary, which is formally stated as below.

Theorem 2.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$ with $\theta \in (0, 1)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ be an orthogonal dictionary, and $\mathbf{Y} = \mathbf{D}_0 \mathbf{X}$. Suppose $\hat{\mathbf{A}}$ is a global maximizer to*

$$\underset{\mathbf{A}}{\text{maximize}} \quad \|\mathbf{A}\mathbf{Y}\|_p^p \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n).$$

Provided that the sample size $r = \Omega(\delta^{-2} n \log(n/\delta)(\theta n \log^2 n)^{\frac{p}{2}})$, then for $\delta > 0$, there exists a signed permutation $\mathbf{\Pi}$, such that

$$\frac{1}{n} \|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq C_\theta \delta$$

with probability at least $1 - r^{-1}$ and C_θ is a constant that depends on θ .

Remark. The derivation for the correctness of global maximizers of the ℓ_4 -based formulation in [26] requires a closed-form expression of the expected objective function, which cannot be obtained for ℓ_{2k+1} -based formulations. In order to develop general theoretical results, we first show that if a formulation satisfies *concentration* and *sharpness conditions* defined in Theorem B.1, its global maximizers are very close to the true dictionary. Then we prove that all ℓ_p -based ($p > 2$) formulations satisfy these two conditions. Our result is consistent with the result in [17] when $p = 4$. Please refer to Section B.1 for a detailed proof.

We next present two lemmas to give a better understanding of Theorem 2.1. First, Lemma 2.1 shows that the global maximizers of the population objective are the true dictionary. Second, we figure out how many samples are needed for the concentration of the empirical objective around the population objective in Lemma 2.2.

We first show that the global maximizers are the true dictionary in expectation.

Lemma 2.1. *(Correctness in expectation) Denote \mathcal{D} as the set of global maximizers to*

$$\underset{\mathbf{A}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{Y}} \|\mathbf{A}\mathbf{Y}\|_p^p \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n),$$

then $\mathcal{D} = \{\mathbf{D}_0^ \mathbf{\Pi}^* | \mathbf{\Pi} \in SP(n)\}$, where $SP(n)$ denotes the group of the signed permutation matrices.*

Lemma 2.1 states the correctness of ℓ_p -based formulations, namely, the global maximizers of the population objective are the true dictionary up to some signed permutations. Due to the law of large numbers, the gap between the empirical objective and expectation objective vanishes as the number of samples goes to infinity. Nevertheless, the sample complexity is an important consideration in dictionary learning. We hope to learn the

true dictionary from as few samples as possible. The next proposition states the finite sample concentration, namely, the empirical objective concentrates on the expectation objective as long as the number of samples is $\Omega(nk^{\frac{p}{2}})$, where k is the number of non-zero elements.

Lemma 2.2. *(Concentration bound of the objective) Suppose $\mathbf{X} \in \mathbb{R}^{n \times r}$ follows $\mathcal{BG}(\theta)$. For any given $\theta \in (0, 1)$ and $\delta > 0$, whenever*

$$r \geq C \delta^{-2} n \log(n/\delta) (\theta n \log^2 n)^{\frac{p}{2}},$$

we have

$$\sup_{\mathbf{A} \in \mathbb{O}(n)} \frac{1}{nr} \left| \|\mathbf{A}\mathbf{Y}\|_p^p - \mathbb{E}(\|\mathbf{A}\mathbf{Y}\|_p^p) \right| \leq \delta$$

with probability at least $1 - r^{-1}$.

Remark. We provide experimental validations for the phase transitions of ℓ_p -based formulations in Fig. 5.

From Lemma 2.2, we see that more samples are required as p increases, and thus a smaller p leads to a lower sample complexity. Fig. 1 illustrates the successful recovery probability versus the number of samples for different values of p . It also confirms the exponential increase in the sample complexity as p increases.

2.3 Robustness

In practice, the observations are usually noisy. Intuitively, the regions around the global maximizers of the ℓ_p -norms ($p > 2$) are flatter than that of the ℓ_1 -norm, which leads to a higher tolerance to noise. In this subsection, we investigate the robustness of the ℓ_p -based formulations to Gaussian noise. Other typical noise will be tested in Section 4 via experiments.

We consider noisy measurements $\mathbf{Y}_N = \mathbf{Y} + \mathbf{G}$ where $\mathbf{G} \in \mathbb{R}^{n \times r}$ with $G_{i,j} \sim \mathcal{N}(0, \eta^2)$. We find that when the number of samples goes to infinity, the Gaussian noise will not change the global maximizers. However, the noise makes it harder for the objective value to concentrate and thus more samples are needed compared to the clean objective. The next theorem shows that the global maximizers are very close to the true dictionary under Gaussian noise.

Theorem 2.2. *Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ be an orthogonal dictionary, and $\mathbf{Y}_N = \mathbf{D}_0 \mathbf{X} + \mathbf{G}$, $\mathbf{G} \in \mathbb{R}^{n \times r}$ with $G_{i,j} \sim \mathcal{N}(0, \eta^2)$. Suppose $\hat{\mathbf{A}}$ is a global maximizer to*

$$\underset{\mathbf{A}}{\text{maximize}} \quad \|\mathbf{A}\mathbf{Y}_N\|_p^p \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n)$$

then for $\delta > 0$, there exists a signed permutation $\mathbf{\Pi}$, such that

$$\frac{1}{n} \|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq C_\theta \delta$$

with probability at least $1 - r^{-1}$ as long as $r = \Omega(\delta^{-2}n \log(n/\delta)((1 + \eta^2)n \log n)^{\frac{p}{2}} \xi_\eta^2)$, where $\xi_\eta = (1 + \eta^2)^{p/2} + \eta^p - 2(0.5 + \eta^2)^{p/2}$, and C_θ is a constant that depends on θ .

Remark. Please refer to Section B.2 for a detailed proof.

Theorem 2.2 shows that as the number of samples is sufficiently large, the global maximizers are very close to the true dictionary. It also suggests that the sample complexity increases as p becomes larger, i.e., a smaller p is more robust to Gaussian noise.

3 AN EFFICIENT ALGORITHM

In this section, we develop an efficient algorithm for ℓ_p -based dictionary learning, and investigate its convergence property. Particularly, an interesting two-stage convergence behavior is revealed and explained.

3.1 Algorithm for ℓ_p -based dictionary learning

We develop our algorithm based on the generalized power method (GPM) algorithm [12], which is a general optimization method to deal with concave objective functions. From the GPM algorithm, we can derive efficient algorithms for many applications, e.g., *subspace iteration* for finding k -largest eigenvectors [5], the *matching, stretching, and projection* algorithm for ℓ_4 -based dictionary learning [26], and efficient algorithms for sparse principle component analysis [12]. These algorithms have been shown to enjoy linear or super-linear convergence rates as well as cheap per-iteration cost (the same cost as the gradient method) in experiments, which motivates us to apply GPM to ℓ_p -based dictionary learning.

The GPM algorithm aims at maximizing a convex function $f(\cdot)$ over a compact constraint Q .

$$\max_{x \in Q} \underbrace{f(x)}_{\text{Convex}}$$

In each iteration, the GPM algorithm maximizes a linear surrogate of the objective function. The procedure of the GPM algorithm is shown in Algorithm 1, where $f' \in \partial f$ is any subgradient.

Algorithm 1 Generalized power method [12]

- 1: Initialize $\mathbf{x}^{(0)} \in D$.
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: $\mathbf{x}^{(t+1)} = \operatorname{argmax}_{\mathbf{s} \in D} \langle \mathbf{s}, f'(\mathbf{x}^{(t)}) \rangle$,
 - 4: **end for**
-

We develop an algorithm for ℓ_p -based dictionary learning based on the GPM algorithm. Let $f(\mathbf{A}) = \|\mathbf{A}\mathbf{Y}\|_p^p$ and thus $\nabla f(\mathbf{A}) = (|(\mathbf{A}\mathbf{Y})^{\circ(p-1)}| \circ \operatorname{sign}(\mathbf{A}\mathbf{Y})) \mathbf{Y}^*$. The update for the GPM algorithm is

$$\mathbf{A}^{(t+1)} = \operatorname{argmax}_{\mathbf{s}^* \in \operatorname{St}(n, m)} \langle \mathbf{s}, \nabla f(\mathbf{A}^{(t)}) \rangle. \quad (7)$$

The only thing left is to compute the maximizer over the Stiefel manifold, which is stated in the next lemma.

Lemma 3.1. [12] *Let $\mathbf{C} \in \mathbb{R}^{m \times n}$ with $m \leq n$, and the singular values of \mathbf{C} is denoted by $\sigma_i(\mathbf{C}), i = 1, \dots, m$. Then,*

$$\max_{\mathbf{s} \in \operatorname{St}(n, m)} \langle \mathbf{s}, \mathbf{C} \rangle = \sum_{i=1}^n \sigma_i(\mathbf{C})$$

with maximizer $\mathbf{s} = \operatorname{Polar}(\mathbf{C})$. $\operatorname{Polar}(\mathbf{C})$ denotes the \mathbf{U} factor of the polar decomposition of the matrix \mathbf{C} such that

$$\mathbf{C} = \mathbf{U}\mathbf{P}, \quad \mathbf{U} \in \operatorname{St}(n, m), \quad \mathbf{P} \in S_+^m.$$

where S_+^m denotes positive semi-definite matrices.

The whole algorithm is shown in Algorithm 2, where $\circ^{(p-1)}$ denotes the element-wise $(p-1)$ -th power and \circ denotes the element-wise product.

Algorithm 2 The GPM algorithm for ℓ_p -based dictionary learning

- 1: Initialize $\mathbf{A}^{(0)*} \in \operatorname{St}(n, m)$.
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: $\nabla f(\mathbf{A}^{(t)}) = (|(\mathbf{A}^{(t)}\mathbf{Y})^{\circ(p-1)}| \circ \operatorname{sign}(\mathbf{A}^{(t)}\mathbf{Y})) \mathbf{Y}^*$
 - 4: $\mathbf{A}^{(t+1)} = \operatorname{Polar}(\nabla^* f(\mathbf{A}^{(t)}))^*$
 - 5: **end for**
-

In general, the GPM algorithm only converges to a critical point with a sub-linear rate since \mathcal{P} in (6) is non-convex [12]. Nevertheless, with \mathbf{X}_0 following the Bernoulli-Gaussian model, experiments show that the GPM algorithm converges to the global maximizer in a very fast rate. We explain this phenomenon in the next subsection.

3.2 Global convergence over the sphere constraint

In this subsection, we investigate why and how the GPM algorithm converges to the global maximizer despite the non-convexity of problem \mathcal{P} in (6). Unfortunately, the global optimality analysis over the Stiefel manifold is extremely difficult. In fact, even whether a local maxima of $\|\mathbf{A}\|_4^4$ exists on $\mathbf{A} \in \mathbb{O}(n)$ is still an open problem [14]. We instead analyze a special case of \mathcal{P} when $m = 1$

with population GPM to see how the GPM algorithm converges to the global optimizer for a non-convex problem. The sphere constraint enables a fine characterization of the gradient dynamics, from which we can derive global convergence results.

In the rest of this subsection, we show the global convergence of the population³ GPM algorithm over the sphere constraint. In the population version, we consider to solve the following problem with Algorithm 3.

$$\begin{aligned} & \underset{\mathbf{a}}{\text{maximize}} && \frac{1}{r} \mathbb{E}_{\mathbf{Y}} \|\mathbf{a}^* \mathbf{Y}\|_p^p = \frac{1}{r} \mathbb{E}_{\mathbf{X}_0} \|\mathbf{a}^* \mathbf{D}_0 \mathbf{X}_0\|_p^p \\ & \text{subject to} && \|\mathbf{a}\|_2 = 1. \end{aligned} \quad (8)$$

Algorithm 3 The population GPM algorithm over the sphere constraint

- 1: Initialize $\|\mathbf{a}^{(0)}\|_2 = 1$.
 - 2: **for** $t = 0 \dots T$ **do**
 - 3: $\nabla f(\mathbf{a}^{(t)}) = \mathbf{Y} \left(|(\mathbf{a}^{(t)*} \mathbf{Y})^{\circ(p-1)}| \circ \text{sign}(\mathbf{a}^{(t)*} \mathbf{Y}) \right)^*$
 - 4: $\mathbf{a}^{(t+1)} = \frac{\mathbb{E}_{\mathbf{Y}}[\nabla f(\mathbf{a}^{(t)})]}{\|\mathbb{E}_{\mathbf{Y}}[\nabla f(\mathbf{a}^{(t)})]\|_2}$
 - 5: **end for**
-

We assume $\mathbf{D}_0 = \mathbf{I}$ without loss of generality since the orthogonal transformation has no impact on the convergence [21]. Moreover, according to the statistical analysis, the global maximizer satisfies $\|\mathbf{a}\|_0 = 1$ [19]. Thus, there are $2n$ ground-truth vectors in the dictionary learning problem, i.e., $\mathbf{a} = \pm \mathbf{e}_i$. We can safely assume that the desired global maximizer is \mathbf{e}_n and $a_n \geq |a_i|, \forall i$ at initialization.

Accordingly, given the ground-truth \mathbf{e}_n , a vector \mathbf{a} can be decomposed into two components: a parallel component a_n and a perpendicular components $\mathbf{a}_{-n} := \mathbf{a}_{1:n-1}$. The parallel component is the signal component.

To study the dynamics of the population GPM algorithm, we start by considering the case where the sequences $\{\mathbf{a}^{(t)}\}$ are generated by the population gradient

$$\mathbf{a}^{(t+1)} = \frac{\nabla F(\mathbf{a}^{(t)})}{\|\nabla F(\mathbf{a}^{(t)})\|_2},$$

where $\nabla F(\mathbf{a}^{(t)})$ denotes the population gradient, given by

$$\nabla F(\mathbf{a}^{(t)}) = \frac{1}{r} \mathbb{E}_{\mathbf{Y}} \nabla f(\mathbf{a}^{(t)}) = c_p \mathbb{E}_{\Omega} [\|\mathbf{a}_{\Omega}\|_2^{p-2} \mathbf{a}_{\Omega}],$$

³Due to the concentration, the empirical gradient will be very close to the population gradient when the sample size is large. Hence, we study the population GPM here to reveal the convergence behavior.

where c_p is a constant only related to p , and Ω denotes the support of a random Bernoulli vector $\mathbf{b} \in \mathbb{R}^n$ with $b_i \sim \text{Ber}(\theta)$.

With simple calculations, the dynamics for both the signal and orthogonal components with respect to the global maximizer \mathbf{e}_n are given by

$$\begin{aligned} \text{Signal :} & \quad a_n^{(t+1)} = \frac{c_p \mathbb{E}_{\Omega} [\|\mathbf{a}_{\Omega}^{(t)}\|_2^{p-2} a_{n,\Omega}]}{\|\nabla F(\mathbf{a}^{(t)})\|_2} \\ \text{Orthogonal :} & \quad a_i^{(t+1)} = \frac{c_p \mathbb{E}_{\Omega} [\|\mathbf{a}_{\Omega}^{(t)}\|_2^{p-2} a_{i,\Omega}]}{\|\nabla F(\mathbf{a}^{(t)})\|_2}, i \neq n \end{aligned}$$

To simplify the analysis, we define the signal-to-orthogonal-ratio (SOR) and signal-to-orthogonal at the i -th coordinate ratio (SOR _{i}) as

$$\text{SOR} = \frac{a_n}{\|\mathbf{a}_{-n}\|_2}, \quad \text{and} \quad \text{SOR}_i = \frac{a_n}{a_i}.$$

A unique advantage for studying SOR and SOR _{i} is that they are projection-invariant, i.e., $\frac{(\mathcal{P}_{\mathbf{g}_{n-1}} \mathbf{q})_n}{\|(\mathcal{P}_{\mathbf{g}_{n-1}} \mathbf{q})_{-n}\|_2} = \frac{q_n}{\|\mathbf{q}_{-n}\|_2}$ and $\frac{(\mathcal{P}_{\mathbf{g}_{n-1}} \mathbf{q})_i}{(\mathcal{P}_{\mathbf{g}_{n-1}} \mathbf{q})_i} = \frac{q_n}{q_i}$. This allows us to bypass the study of projection and makes the results easy to interpret.

The SOR and SOR _{i} can also be viewed as an error metric since

$$\|\mathbf{a} - \mathbf{e}_n\|_2^2 = 2 - 2\sqrt{\frac{\text{SOR}^2}{\text{SOR}^2 + 1}},$$

where $\|\mathbf{a} - \mathbf{e}_n\|_2^2$ is the squared ℓ_2 error.

The following proposition shows the evolution of SOR _{i} during the iterations, which will be used to derive the global convergence result.

Proposition 3.1. Denote SOR _{i} ^(t) as the value of SOR _{i} at the t -th iteration and $\mathbf{q} = \mathbf{a}^{(t)}$ as the variable at the t -th iteration. Then SOR _{i} evolves as

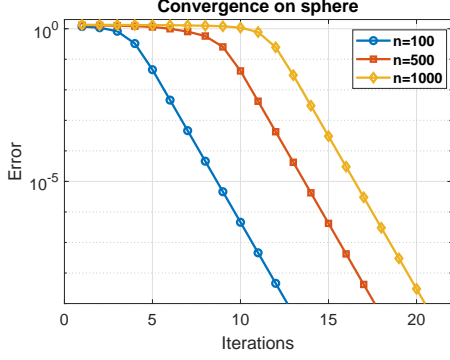
$$\text{SOR}_i^{(t+1)} = \text{SOR}_i^{(t)} (1 + \tau_i(\mathbf{q})),$$

and

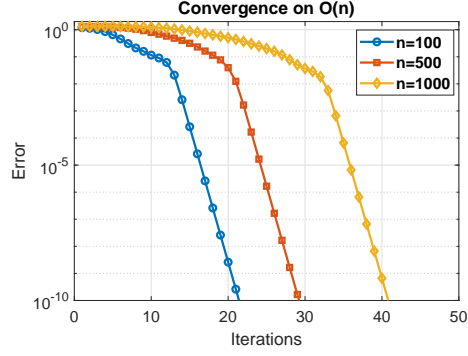
$$\tau_i(\mathbf{q}) = \frac{\mathbb{E}_{\Omega'} \|\mathbf{q}_{\Omega'}, q_n\|_2^k - \mathbb{E}_{\Omega'} \|\mathbf{q}_{\Omega'}, q_i\|_2^k}{\frac{\theta}{1-\theta} \mathbb{E}_{\Omega'} \|\mathbf{q}_{\Omega'}, q_i, q_n\|_2^k + \mathbb{E}_{\Omega'} \|\mathbf{q}_{\Omega'}, q_i\|_2^k},$$

where $\Omega' = \Omega \setminus \{n\} \setminus \{i\}$ and $k = p - 2$. Two properties of $\tau_i(\mathbf{q})$ are listed below

1. $0 \leq \tau_i(\mathbf{q}) \leq \frac{1-\theta}{\theta}$ always holds and $\tau_i(\mathbf{q}) > 0$ if $q_n > q_i$.
2. $\tau_i(\mathbf{q})$ is monotonically increasing in q_n and decreasing in q_i .



(a) Convergence over the sphere. The error metric is $\min_{1 \leq i \leq n} \|\mathbf{a} - \mathbf{D}_0 \mathbf{e}_i\|_2$.



(b) Convergence over the orthogonal group. The error metric is $\min_{\mathbf{\Pi} \in \text{SP}(n)} \frac{\|\mathbf{A}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F}{\|\mathbf{D}_0\|_F}$.

Figure 3: The convergence of the population GPM for ℓ_4 -based formulation. The sparsity level is $\theta = 0.1$. From both figures, we see the first stage is short and the second stage has a linear convergence with rate $\frac{1}{10}$.

Remark. Please refer to Section C.1 for the proof of this proposition.

According to Proposition 3.1, SOR_i grows at an exponential rate. Thus, Algorithm (3) converges to the global maximizer of Problem (8) and the convergence is stated as below.

Theorem 3.1. *Apply Algorithm 3 to solve problem (8) and assume $\mathbf{a}^{(0)}$ follows a uniform distribution over the sphere. Denote $\tau(\mathbf{q}) = \min_{i=1, \dots, n-1} \tau_i(\mathbf{q})$, then there exists $T_\tau \leq \log_{1+\tau(\mathbf{a}^{(0)})}(\sqrt{n})$, $1 \leq i \leq n$, such that*

$$\|\mathbf{a}^{(t)} - \mathbf{D}_0 \mathbf{e}_i\|_2 \leq \left(1 + \tau(\mathbf{a}^{(T_\tau)})\right)^{T_\tau - t}, \forall t \geq T_\tau,$$

almost surely and the convergence rate is given by
 $\lim_{k \rightarrow \infty} \frac{\|\mathbf{a}^{(k+1)} - \mathbf{D}_0 \mathbf{e}_i\|_2}{\|\mathbf{a}^{(k)} - \mathbf{D}_0 \mathbf{e}_i\|_2} = \frac{1}{\theta}$.

Remark. Please refer to Section C.2 for the proof of this theorem.

From Theorem 3.1, we observe a two-stage convergence. Specifically, the first stage only takes T_τ iterations, which is short, and the second stage enjoys a linear convergence rate. Fig. 3(a) illustrates the convergence of the population GPM algorithm for the ℓ_4 -based formulation over the sphere. We clearly see a two-stage convergence in the figure: The first stage only lasts for a few iterations and the second stage has a linear convergence rate of $\frac{1}{10}$. The convergence speed is much faster than Riemannian gradient (RGD), which is the most commonly adopted method on Riemannian manifold [1]. We provide a comparison of the GPM algorithm and RGD in Section A.1. The convergence over $\mathbb{O}(n)$ is shown in Fig. 3(b), from which we also observe the two-stage convergence with a $\frac{1}{10}$ rate.

4 EXPERIMENTS

In this section, we test the performance of the ℓ_p -based approaches under both noiseless and noisy conditions. Specifically, we verify the advantages of the ℓ_3 -based method revealed in the theoretical analysis, especially its robustness. All the experiments are conducted with Matlab 2019a running on Intel Core i7-6700 CPU @ 3.40GHz.

The following three benchmarks are considered:

- **K-SVD [2, 18]:** K-SVD is a classic dictionary learning algorithm. We adopt the K-SVD toolbox [18] for an efficient implementation.
- **Riemannian trust region (RTR) [21]:** This method adopts the smoothed ℓ_1 -based formulation and Riemannian trust region algorithm [1], which enjoys an exact recovery guarantee in the noiseless case.
- **ℓ_4 -based [26]:** It can be viewed as a special case of the ℓ_p -based formulation studied in this paper when $p = 4$. The MSP algorithm proposed in [26] is adopted.

Scalability in the noiseless case: We first test the scalability of different algorithms. Specifically, we compare the accuracy and running time of different methods with different dictionary sizes n and sparsity levels θ in the noiseless case. For the tested problem size, RTR failed to terminate within ten hours. Hence, we do not include it as a baseline. The error and running time are taken average over 10 independent random trials. The results are shown in Table 1.

From Table 1, we see that all ℓ_p -based methods achieve better performance and significant speedup compared to the K-SVD algorithm. As the problem size increases, the error of K-SVD increases while the error of ℓ_p -based approaches remain stable.

We also observe that a smaller p leads to a smaller error, which is consistent with the results in Theorem 2.1. Moreover, the ℓ_3 -based approach enjoys the least running time because of its lower per-iteration complexity. Thus, the ℓ_3 -based approach is the best choice in terms of both accuracy and speed for the noiseless case.

Gaussian noise: A small Gaussian noise usually appears in dictionary learning applications. We consider the noisy observation $\mathbf{Y}_N = \mathbf{Y} + \sigma\mathbf{G}$, where $G_{i,j} \sim \mathcal{N}(0, 1)$. The results of different algorithms are shown in Table 2. While RTR achieves the smallest error on a clean objective (i.e., the cases with $\sigma = 0$), it is very time-consuming. When the noise is large, both RTR and K-SVD fail to recover the dictionary while the ℓ_3 and ℓ_4 -based methods still have a relatively low error. This is because for the ℓ_3 and ℓ_4 -norm, the regions around the global maximizers are very flat, leading to a higher tolerance to noise. Table 2 also shows that the ℓ_3 -based method is the most time-efficient in all cases and achieves the best performance when the noise is present.

Sparse corruptions: Sparse corruptions are another kind of noise usually present in images [6]. We consider the noisy observation $\mathbf{Y}_S = \mathbf{Y} + \sigma\mathbf{B} \circ \mathbf{R}$, where $B_{i,j} \sim \text{Ber}(\vartheta)$ and $R_{i,j}$ has equal probability to be -1 and 1 . The results are shown in Table 3. The performance of different methods are similar to the Gaussian noise case. RTR performs the best when the noise is relatively small but very unstable as the noise increases. The ℓ_3 and ℓ_4 -based methods are more stable as the noise increases and the ℓ_3 -based method is still the most time-efficient and most robust one.

All the above experiments demonstrate that the ℓ_3 -based approach is more time efficient and robust than existing methods, and thus it is a preferred method to use in practice.

5 RELATED WORKS

In this section, we discuss some existing literature related to our work.

Applications of ℓ_p -norm: The ℓ_p -norm plays an important role in machine learning. ℓ_1 -norm can induce sparsity and has been widely applied in compressive sensing [6, 9]. ℓ_2 -norm has strong geometric connections to eigenvectors as well as variance. Thus, it is of-

ten used in different kinds of principle component pursuits [23]. Maximizing ℓ_{2k+2} -norm can thus promote the spikiness, which has found applications in independent component analysis [11], sparse blind deconvolution [13, 27], and dictionary learning [17, 26]. General ℓ_p -norms have also been applied for deconvolution applications in geophysics [7, 15], but with little theory. Our study showed the importance of investigating more general ℓ_p -norm based methods, and developed general frameworks for effective theoretical analysis and algorithm design. In particular, ℓ_3 -norm maximization stands out as a promising method for solving dictionary learning problems.

Independent Component Analysis: ICA factors a data matrix \mathbf{Y} as $\mathbf{Y} = \mathbf{A}\mathbf{X}$ such that \mathbf{A} is square with \mathbf{X} as independent as possible. One popular contrast function for ICA is the kurtosis [11], which has the same formulation with (8) when $p = 4$. It suggests that our general ℓ_p -based formulation and analysis might be able to extend to ICA. The major difference is that ICA is to find statistically independent components while dictionary learning is to find the sparsest representation of the given data. As is revealed in Section 1.5 of [21], it is the sparsity rather than independence shapes the benign landscape of the orthogonal dictionary learning problem.

6 CONCLUSIONS

In this paper, we considered the ℓ_p -based ($p > 2, p \in \mathbb{N}$) formulations for the orthogonal dictionary learning problem. We showed that the global maximizers of these formulations are very close to the true dictionary in both noiseless and noisy cases. In addition, we developed an efficient algorithm based on the generalized power method, and demonstrated its fast global convergence. We further conducted various experiments to show the benefits of adopting the ℓ_3 -based approach.

Acknowledgement

We would like to thank professor Yi Ma of Berkeley EECS Department for his lectures and talks at Tsinghua-Berkeley Shenzhen Institute and Yuexiang Zhai of Berkeley for stimulating discussions during preparation of this manuscript.

Table 1: The performance of different algorithms for noiseless objectives. Since the dictionary recovery is up to some signed permutations, we adopt the error metric $1 - \|AD_0\|_4^4/n$ in [26], which gives 0% error for a perfect recovery.

Settings			ℓ_3 -based		ℓ_4 -based [26]		ℓ_5 -based		K-SVD [18]	
n	θ	$p(\times 10^4)$	Time	Error	Time	Error	Time	Error	Time	Error
100	0.1	4	0.8s	0.056%	1.8s	0.21%	1.7s	0.50%	61s	1.45%
200	0.1	8	4.1s	0.056%	9.3s	0.21%	8.0s	0.51%	131s	3.03%
400	0.1	16	35s	0.056%	50s	0.21%	41s	0.50%	315s	6.45%
100	0.3	4	1.2s	0.094%	3.4s	0.34%	3.1s	0.84%	98s	2.60%
200	0.3	8	10s	0.094%	18s	0.35%	15s	0.85%	215s	6.41%
400	0.3	16	91s	0.096%	122s	0.35%	146s	1.00%	589s	8.25%

Table 2: The performance of different algorithms under Gaussian noise. We set sparsity level $\theta = 0.3$.

Settings			ℓ_3 -based		ℓ_4 -based [26]		RTR [21]		K-SVD [18]	
n	$p(\times 10^4)$	σ	Time	Error	Time	Error	Time	Error	Time	Error
32	1	0	0.05s	0.10%	0.24s	0.4%	100s	0.05%	25s	0.2%
32	1	0.2	0.05s	0.27%	0.24s	0.6%	250s	0.5%	25s	0.37%
32	1	0.4	0.1s	0.79%	0.36s	1.2%	577s	4.27%	25s	2.0%
32	1	0.6	0.2s	2.3%	0.7s	3.4%	823s	57.4%	25s	57.4%
100	4	0	1.2s	0.1%	3.4s	0.35%	863s	0.05%	98s	2.60%
100	4	0.2	2.2s	0.2%	4.2s	0.5%	1643s	0.3%	104s	3.46%
100	4	0.4	3.5s	0.6%	6.1s	1.1%	3796s	5.26%	105s	3.56%
100	4	0.6	8.4s	1.95%	13.5s	2.63%	5412s	50.5%	104s	51.26%

Table 3: The performance of different algorithms under sparse noise. We set sparsity level $\theta = 0.3$. The sparsity of noise is set to $\vartheta = 0.1$.

Settings			ℓ_3 -based		ℓ_4 -based [26]		RTR [21]		K-SVD [18]	
n	$p(\times 10^4)$	σ	Time	Error	Time	Error	Time	Error	Time	Error
32	1	0.5	0.06s	0.20%	0.25s	0.57%	362s	0.10%	25s	0.37%
32	1	1	0.09s	0.50%	0.35s	0.93%	421s	1.4%	25s	2.0%
32	1	1.5	0.14s	1.65%	0.47s	2.26%	420s	13.4%	25s	57.4%
100	1	0.5	1.9s	0.20%	4.0s	0.40%	1649s	0.2%	104s	3.04%
100	4	1	2.6s	0.40%	5.5s	0.80%	2737s	1.3%	105s	3.55%
100	4	1.5	4.6s	1.02%	7.7s	1.49%	5395s	36.8%	104s	5.83%

References

- [1] P-A Absil, Robert Mahony, and Rodolphe Sepulchre. *Optimization algorithms on matrix manifolds*. Princeton University Press, 2009.
- [2] Michal Aharon, Michael Elad, and Alfred Bruckstein. K-SVD: An algorithm for designing overcomplete dictionaries for sparse representation. *IEEE Transactions on Signal Processing*, 54(11):4311–4322, 2006.
- [3] Yu Bai, Qijia Jiang, and Ju Sun. Subgradient descent learns orthogonal dictionaries. In *International Conference on Learning Representations*, 2019.
- [4] Chenglong Bao, Jian-Feng Cai, and Hui Ji. Fast sparsity-based orthogonal dictionary learning for image restoration. In *IEEE International Conference on Computer Vision*, pages 3384–3391, 2013.
- [5] Klaus-Jürgen Bathe. The subspace iteration method—revisited. *Computers & Structures*, 126:177–183, 2013.
- [6] Emmanuel J Candès, Xiaodong Li, Yi Ma, and John Wright. Robust principal component analysis? *Journal of the ACM*, 58(3):11, 2011.
- [7] HWJ Debeye and P Van Riel. ℓ_p -norm deconvolution. *Geophysical Prospecting*, 38(4):381–403, 1990.
- [8] M. Elad. *Sparse and Redundant Representations: From Theory to Applications in Signal and Image Processing*. New York, NY, USA: Springer-Verlag, 2010.
- [9] Simon Foucart and Holger Rauhut. A mathematical introduction to compressive sensing. *Bulletin of the American Mathematical Society*, 54:151–165, 2017.
- [10] Dar Gilboa, Sam Buchanan, and John Wright. Efficient dictionary learning with gradient descent. In *International Conference on Machine Learning*, pages 2252–2259, 2019.
- [11] Aapo Hyvärinen and Erkki Oja. A fast fixed-point algorithm for independent component analysis. *Neural computation*, 9(7):1483–1492, 1997.
- [12] Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power method for sparse principal component analysis. *Journal of Machine Learning Research*, 11:517–553, 2010.
- [13] Yanjun Li and Yoram Bresler. Global geometry of multichannel sparse blind deconvolution on the sphere. In *Advances in Neural Information Processing Systems*, pages 1132–1143, 2018.
- [14] Yi Ma. Complete dictionary learning via ℓ_4 -norm maximization over the orthogonal group. In *International Conference on Computer Vision Workshops*, 2019.
- [15] Kenji Nose-Filho and João MT Romano. On ℓ_p -norm sparse blind deconvolution. In *IEEE International Workshop on Machine Learning for Signal Processing*, pages 1–6. IEEE, 2014.
- [16] Bruno A Olshausen and David J Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607, 1996.
- [17] Qing Qu, Yuexiang Zhai, Xiao Li, Yuqian Zhang, and Zhihui Zhu. Geometric analysis of nonconvex optimization landscapes for overcomplete learning. In *International Conference on Learning Representations*, 2020.
- [18] Ron Rubinfeld, Michael Zibulevsky, and Michael Elad. Efficient implementation of the K-SVD algorithm using batch orthogonal matching pursuit. Technical report, Computer Science Department, Technion, 2008.
- [19] Daniel A Spielman, Huan Wang, and John Wright. Exact recovery of sparsely-used dictionaries. In *Conference on Learning Theory*, pages 1–37, 2012.
- [20] Ju Sun, Qing Qu, and John Wright. When are nonconvex problems not scary? *arXiv preprint arXiv:1510.06096*, 2015.
- [21] Ju Sun, Qing Qu, and John Wright. Complete dictionary recovery over the sphere I: Overview and the geometric picture. *IEEE Transactions on Information Theory*, 63(2):853–884, 2017.
- [22] Roman Vershynin. *High-dimensional probability: An introduction with applications in data science*. Cambridge University Press, 2018.
- [23] René Vidal, Yi Ma, and S Shankar Sastry. *Generalized principal component analysis*. Springer, 2016.
- [24] Yu Wang, Siqi Wu, and Bin Yu. Unique sharp local minimum in ℓ_1 -minimization complete dictionary learning. *Journal of Machine Learning Research*, 21:1–52, 2020.

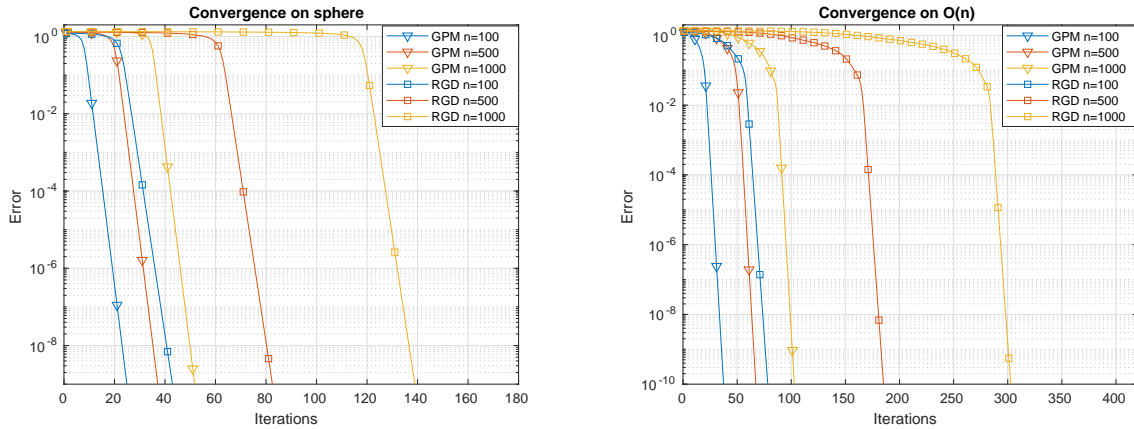
- [25] Yuexiang Zhai, Hermish Mehta, Zhengyuan Zhou, and Ma Yi. Understanding ℓ^4 -based dictionary learning: Interpretation, stability, and robustness. In *International Conference on Learning Representations*, 2020.
- [26] Yuexiang Zhai, Zitong Yang, Zhenyu Liao, John Wright, and Yi Ma. Complete dictionary learning via ℓ_4 -norm maximization over the orthogonal group. *arXiv preprint arXiv:1906.02435*, 2019.
- [27] Yuqian Zhang, Han-Wen Kuo, and John Wright. Structured local optima in sparse blind deconvolution. *IEEE Transactions on Information Theory*, 66(1):419–452, 2020.

Supplementary Materials

A Additional Experiments

A.1 Comparison between GPM and RGD

In this subsection, we present the comparison between the population generalized power method (GPM) and population Riemannian gradient (RGD) for the ℓ_p -based objective. RGD [1] is a popular method for optimization over manifold, which is known for its low computational cost [3]. In the experiment, we set $\theta = 0.3$ and $p = 4$ to compare the convergence speed of the two methods in different scales. For RGD, we fix the step size as $\frac{1}{4}$. The results on the sphere and orthogonal group are shown in Fig. 4(a) and Fig. 4(b) respectively. From both figures, we see that the convergence speed of GPM is much faster than RGD, especially when the problem size is large. In addition, we observe both GPM and RGD have a two stage convergence.



(a) Convergence over the sphere. The error metric is $\min_{1 \leq i \leq n} \|\mathbf{a} - \mathbf{D}_0 \mathbf{e}_i\|_2$.

(b) Convergence over the orthogonal group. The error metric is $\min_{\mathbf{\Pi} \in \text{SP}(n)} \frac{\|\mathbf{A}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F}{\|\mathbf{D}_0\|_F}$.

Figure 4: The convergence of the population GPM for ℓ_4 -based formulation.

A.2 Phase Transition Heatmaps

There are complicated interactions between the number of samples, r , the sparsity level, θ , the order of the norm, p , and the variance, σ of the additive Gaussian noise. For a more comprehensive illustration, we plot the heatmap by fixing two variables among r , θ , p and σ , and varying the other two to show these interactions. The results are shown in Fig. 5, where the color reflects the error.

A.3 Experiments on Real Data

In this subsection, we test the performance on the MNIST dataset to check whether ℓ_p -based methods can learn a sparse representation on real images. We include principal component pursuit (PCA) as a benchmark. For each method, we take the top 5 bases to recover the original images. The results are shown in Fig. 6. With top 5 bases, the recovered images by PCA can hardly be recognized while we can recognize the images recovered by ℓ_3 and ℓ_4 -based method. This suggests that a sparse representation is learned. In addition, the recovered images via ℓ_3 -based method is a bit more identifiable than the images via ℓ_4 -based method from the figure.

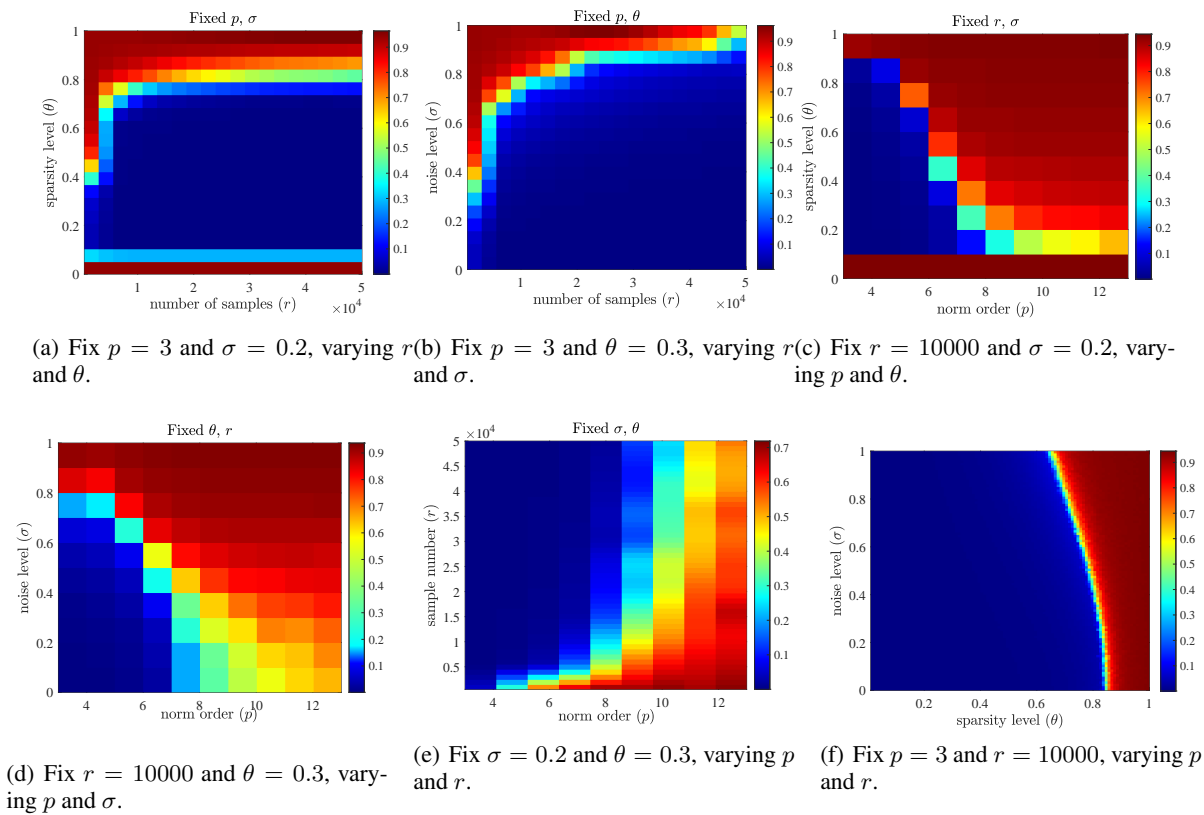


Figure 5: Phase transitions when $n = 50$.

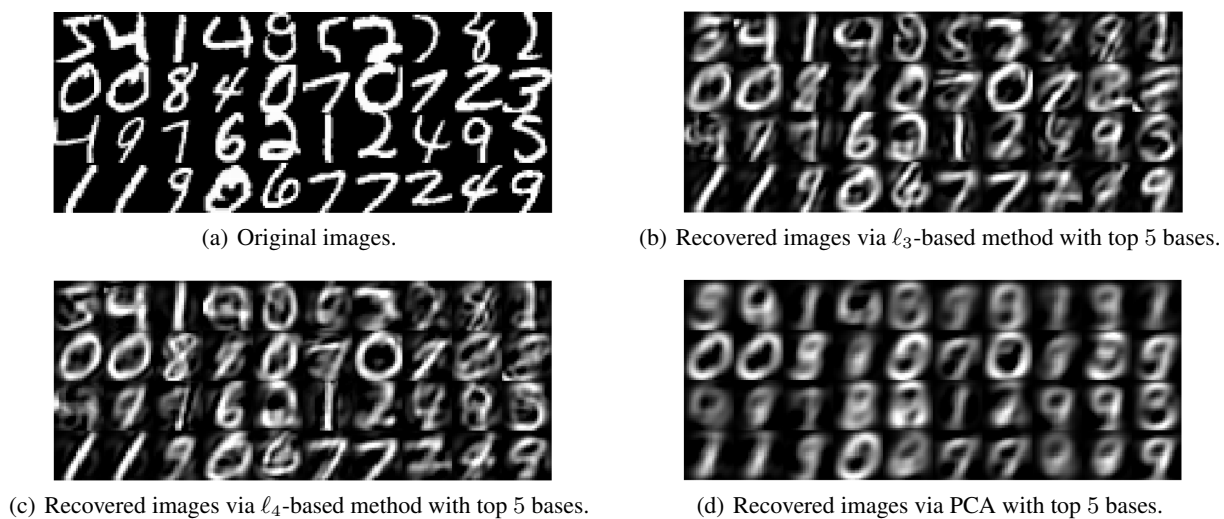


Figure 6: Recovered images with top 5 bases.

B Correctness of the Global Maximizers

In this section, we prove Theorem 2.1 and Theorem 2.2. We first develop a general theory for the correctness of the global maximizers. Consider the following optimization problem to recover the dictionary

$$\underset{\mathbf{A}}{\text{maximize}} \quad f(\mathbf{A}, \mathbf{Y}) \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n).$$

The following result shows that if two conditions are satisfied, the global maximizers to this optimization problem are very close to the true dictionary with high probability.

Theorem B.1. *Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ be an orthogonal dictionary, and $\mathbf{Y} = \mathbf{D}_0 \mathbf{X}$.*

Assume a function $f : \mathbb{R}^{n \times n} \times \mathbb{R}^{n \times r} \rightarrow \mathbb{R}$ that satisfies the following two conditions.

1. (Concentration) Denote $g(\mathbf{A}) = \mathbb{E}_{\mathbf{Y}} f(\mathbf{A}, \mathbf{Y})$, with sample size $r > \Phi(\delta, n)$, and then

$$\sup_{\mathbf{A} \in \mathbb{O}(n)} |f(\mathbf{A}, \mathbf{Y}) - g(\mathbf{A})| \leq \delta \quad (9)$$

holds with high probability.

2. (Sharpness) For all $\mathbf{A} \in \mathbb{O}(n)$, there exists $\alpha > 0$ and $\mathbf{\Pi} \in SP(n)$ such that

$$g(\mathbf{D}_0^* \mathbf{\Pi}^*) - g(\mathbf{A}) \geq \alpha \min_{\mathbf{\Pi} \in SP(n)} \|\mathbf{A} - \mathbf{D}_0^* \mathbf{\Pi}^*\|_F^2. \quad (10)$$

Suppose $\hat{\mathbf{A}}$ is a global maximizer to

$$\underset{\mathbf{A}}{\text{maximize}} \quad f(\mathbf{A}, \mathbf{Y}) \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n). \quad (11)$$

and $r > \Phi(\delta, n)$, then for any $\delta > 0$, there exists a signed permutation $\mathbf{\Pi}$, such that

$$\|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq 2\alpha^{-1} \delta$$

with high probability.

Proof. From concentration, if $r > \Phi(\delta, n)$ we have

$$\left| f(\hat{\mathbf{A}}, \mathbf{Y}) - g(\hat{\mathbf{A}}) \right| \leq \delta, \quad |f(\mathbf{D}_0^*, \mathbf{Y}) - g(\mathbf{D}_0^*)| \leq \delta.$$

with high probability. This leads to

$$g(\mathbf{D}_0^*) - 2\delta \leq g(\hat{\mathbf{A}}) \leq g(\mathbf{D}_0^*). \quad (12)$$

By sharpness, we have

$$\alpha \|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq g(\mathbf{D}_0^*) - g(\hat{\mathbf{A}}) \leq 2\delta.$$

This implies

$$\|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq 2\alpha^{-1} \delta$$

with high probability. □

Remark. The concentration condition determines the sample complexity of the formulation and ensures that the maximal objective values are close to the expected maximal objective value. The sharpness condition measures how close two variables are if their objective values are close, which then leads to the closeness in the variable values.

We check the concentration and sharpness of noiseless and noisy objectives in the following two subsections.

B.1 Proof of Theorem 2.1

In this subsection, we prove Theorem 2.1 via Theorem B.1. Specifically, the concentration and sharpness conditions are established in Lemma B.4 and Lemma B.6 respectively. Then we show the global maximizers of all ℓ_p -based formulation are very close to the true dictionary with high probability via Theorem B.1.

We begin with the correctness of the population objective. We first show the correctness over the sphere and then extend it to $\mathbb{O}(n)$.

Lemma B.1. (*Global optimums over the sphere*) Consider the problem

$$\underset{\mathbf{a} \in \mathbb{S}^{n-1}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{y}} |\mathbf{a}^* \mathbf{y}|^p, \quad (13)$$

where $\mathbf{y} = \mathbf{D}_0(\mathbf{b} \circ \mathbf{g})$ with $\mathbf{b} \sim \text{Ber}(\theta)$ and $\mathbf{g} \sim \mathcal{N}(0, \sigma^2)$, $\gamma_p = \sigma^p 2^{p/2} \frac{\Gamma(\frac{p+1}{2})}{\sqrt{\pi}}$ and $\Gamma(\cdot)$ is the Gamma function. The equality holds when $\|\mathbf{a}^* \mathbf{D}_0\|_0 = 1$ and the maximal objective value is $\gamma_p \theta$.

Proof. We prove this lemma by separating \mathbf{b} and \mathbf{g}

$$\begin{aligned} \mathbb{E}_{\mathbf{y}} |\mathbf{a}^* \mathbf{y}|^p &= \mathbb{E}_{\mathbf{b}, \mathbf{g}} |\mathbf{a}^* \mathbf{D}_0(\mathbf{b} \circ \mathbf{g})|^p = \mathbb{E}_{\mathbf{b}, \mathbf{g}} |(\mathbf{a}^* \mathbf{D}_0 \circ \mathbf{b}) \mathbf{g}|^p = \mathbb{E}_{\mathbf{g}, \Omega} |(\mathbf{a}^* \mathbf{D}_0)_\Omega \mathbf{g}|^p \stackrel{(a)}{=} \gamma_p \mathbb{E}_\Omega \|(\mathbf{a}^* \mathbf{D}_0)_\Omega\|_2^p \\ &\stackrel{(b)}{\leq} \gamma_p \mathbb{E}_\Omega \|(\mathbf{a}^* \mathbf{D}_0)_\Omega\|_2^2 = \gamma_p \theta, \end{aligned}$$

where Ω denotes the support of \mathbf{b} and inequality (a) follows Lemma D.1. The inequality in (b) is because $\|(\mathbf{a}^* \mathbf{D}_0)_\Omega\|_2 \leq 1$ and the equality holds only if $\|\mathbf{a}^* \mathbf{D}_0\|_0 = 1$. \square

Extending Lemma B.1 to $\mathbb{O}(n)$ results in the following lemma, which completes the proof for the correctness of the population objective.

Lemma B.2. (*The only global maximizers to*

$$\underset{\mathbf{A}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{X}_0} \|\mathbf{A} \mathbf{Y}\|_p^p \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n)$$

are $\mathbf{A}^* = \mathbf{D}_0 \mathbf{P}$, where \mathbf{P} is any signed permutation matrix.

Proof. We consider

$$\underset{\mathbf{A}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{Y}} \|\mathbf{A} \mathbf{Y}\|_p^p \quad \text{subject to } \mathbf{A} \in \mathbb{O}(n).$$

Denoting $\mathbf{A} = [\mathbf{a}_1^*; \dots; \mathbf{a}_m^*]$ where $\mathbf{a}_i \in \mathbb{R}^n$ and $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_r]$, we have

$$\mathbb{E}_{\mathbf{Y}} \|\mathbf{A} \mathbf{Y}\|_p^p = \sum_{i=1}^n \sum_{j=1}^r \mathbb{E}_{\mathbf{y}_j} |\mathbf{a}_i^* \mathbf{y}_j|_p^p = r \sum_{i=1}^n \mathbb{E}_{\mathbf{y}_j} |\mathbf{a}_i^* \mathbf{y}_j|_p^p \stackrel{(a)}{\leq} r n \gamma_p \theta. \quad (14)$$

Inequality (a) follows Lemma B.1 and the equality holds only if $\|\mathbf{a}_i \mathbf{D}_0\|_0 = 1, \forall i$ and $\mathbf{A} \in \text{St}(n, m)$. Thus, the global maximum is achieved only if $\mathbf{A}^* = \mathbf{D}_0 \mathbf{P}$, where \mathbf{P} is any signed permutation matrix. \square

We then establish the concentration of the empirical objective. We first show a general heavy-tailed concentration bound over the Stiefel manifold, and then apply it to obtain the required sample complexity for concentration.

Lemma B.3. (*Concentration on Stiefel manifold*) Let $\mathbf{z}_1, \mathbf{z}_2, \dots, \mathbf{z}_r \in \mathbb{R}^{n_1}$ be i.i.d. centered subgaussian random vectors, with $\mathbf{z}_i \equiv_d \mathbf{z}(1 \leq i \leq r)$ such that

$$\mathbb{E}[z_i] = 0, \quad \mathbb{P}(|z_i| > t) \leq 2 \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

For fixed $\mathbf{Q} \in \text{St}(n, m)$, we define a function $f_{\mathbf{Q}} : \mathbb{R}^{n_1} \rightarrow \mathbb{R}^{d_1}$, such that

1. $f_{\mathbf{Q}}(\mathbf{z})$ is a heavy tailed process of \mathbf{z} , in the sense of

$$\mathbb{P}(\|f_{\mathbf{Q}}(\mathbf{z})\|_2 \geq t) \leq 2 \exp(-Ct^{2/p}). \quad (15)$$

2. The expectation $\mathbb{E}[f_{\mathbf{Q}}(\mathbf{z})]$ is bounded and L_f -Lipschitz, i.e.,

$$\|\mathbb{E}[f_{\mathbf{Q}}(\mathbf{z})]\| \leq B_f, \quad \text{and} \quad \|\mathbb{E}[f_{\mathbf{Q}_1}(\mathbf{z})] - \mathbb{E}[f_{\mathbf{Q}_2}(\mathbf{z})]\| \leq L_f \|\mathbf{Q}_1 - \mathbf{Q}_2\|, \forall \mathbf{Q}_1, \mathbf{Q}_2 \in \text{St}(n, m). \quad (16)$$

3. Let $\bar{\mathbf{z}}$ be a truncated vector of \mathbf{z} , such that

$$\mathbf{z} = \bar{\mathbf{z}} + \hat{\mathbf{z}}, \quad \hat{z}_i = \begin{cases} z_i, & \text{if } |z_i| \leq B \\ 0, & \text{otherwise.} \end{cases} \quad (17)$$

with $B = 2\sigma\sqrt{\log(n_1 r)}$. We further assume that

$$\begin{aligned} \|f_{\mathbf{Q}}(\bar{\mathbf{z}})\| &\leq R_1, \quad \mathbb{E}[\|f_{\mathbf{Q}}(\bar{\mathbf{z}})\|^2] \leq R_2 \\ \|f_{\mathbf{Q}_1}(\bar{\mathbf{z}}) - f_{\mathbf{Q}_2}(\bar{\mathbf{z}})\| &\leq \bar{L}_f \|\mathbf{Q}_1 - \mathbf{Q}_2\|, \forall \mathbf{Q}_1, \mathbf{Q}_2 \in \text{St}(n, m). \end{aligned} \quad (18)$$

Then for any $t > 0$, we have

$$\begin{aligned} &\mathbb{P}\left(\sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{mr} \sum_{i=1}^r f_{\mathbf{Q}}(z_i) - \mathbb{E}[f_{\mathbf{Q}}(z_i)] \right\| > t\right) \\ &\leq (n_1 r)^{-1} + \exp\left(-\min\left\{\frac{rm^2 t^2}{64R_2}, \frac{3mrt}{32R_1}\right\} + nm \log\left(\frac{12(L_f + \bar{L}_f)}{t}\right) + \log(d_1)\right). \end{aligned} \quad (19)$$

In other words, for any given δ , whenever

$$r \geq Cn/\delta \log\left(\frac{(L_f + \bar{L}_f)d_1}{\delta}\right) \max\{R_2/(m\delta), R_1\} + \frac{B_f}{\sqrt{n_1}\delta},$$

we have

$$\mathbb{P}\left(\sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{mr} \sum_{i=1}^r f_{\mathbf{Q}}(z_i) - \mathbb{E}[f_{\mathbf{Q}}(z_i)] \right\| > \delta\right) < (n_1 r)^{-1} + (mn)^{-c \log\left(\frac{12(L_f + \bar{L}_f)}{\delta}\right)}.$$

Proof. The proof is heavily based on the concentration of random vectors over the sphere (Theorem F.1 in [17]).

We employ truncation to deal with the heavy-tailed phenomenon and thus the bounds for bounded random variables [22] can be applied here. The truncation level is set as $B = 2\sigma\sqrt{\log(n_1 r)}$. We first separate the concentration into three parts

$$\begin{aligned} &\mathbb{P}\left(\sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}}(z_i) - \mathbb{E}f_{\mathbf{Q}}(\mathbf{z}) \right\| \geq t\right) \\ &\leq \underbrace{\mathbb{P}\left(\sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}}(\bar{\mathbf{z}}) \right\| \geq \frac{t}{2}\right)}_{\mathcal{T}_1} + \underbrace{\mathbb{P}\left(\sup_{\mathbf{Q} \in \text{St}(n, m)} \|\mathbb{E}f_{\mathbf{Q}}(\bar{\mathbf{z}}) - \mathbb{E}f_{\mathbf{Q}}(\mathbf{z})\| \geq \frac{t}{2}\right)}_{\mathcal{T}_2} + \underbrace{\mathbb{P}\left(\max_{1 \leq i \leq r} \|z_i\|_{\infty} \leq B\right)}_{\mathcal{T}_3}. \end{aligned}$$

Denote the ϵ -net on Stiefel manifold as $N(\epsilon)$.

Bound \mathcal{T}_3

$$\mathcal{T}_3 = \mathbb{P}\left(\max_{1 \leq i \leq r} \|z_i\|_{\infty} \geq B\right) \leq n_1 r \mathbb{P}(|z_{i,j}| \geq B) \leq \exp\left(-\frac{B^2}{2\sigma^2} + \log(n_1 r)\right) = (n_1 r)^{-1}.$$

Bound \mathcal{T}_2 Note that

$$\|\mathbb{E}f_{\mathbf{Q}}(z) - \mathbb{E}f_{\mathbf{Q}}(\bar{z})\| \leq \|\mathbb{E}f_{\mathbf{Q}}(z) \circ \mathbf{1}_{z \neq \bar{z}}\|_2 \leq \|\mathbb{E}f_{\mathbf{Q}}(z)\|_2 \|\mathbf{1}_{z \neq \bar{z}}\|_2 \leq B_f r^{-1} n_1^{-\frac{1}{2}}.$$

Thus, we have $\mathcal{T}_2 = 0$ if $r \geq 2B_f t^{-1} n_1^{-\frac{1}{2}}$.

Bound \mathcal{T}_1 To bound \mathcal{T}_1 , we first derive a bound for fixed $\mathbf{Q} \in \text{St}(n, m)$ and then take a union bound over $N(\epsilon)$. By Bernstein inequality for bounded random variables [Theorem 2.8.4 in [22]], we have

$$\mathbb{P}\left(\left\|\mathbb{E}f_{\mathbf{Q}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}}(z)\right\| \geq \frac{t}{2}\right) \leq d_1 \exp\left(-\frac{rt^2}{8R_2 + 8R_1 t/3}\right).$$

Covering over the Stiefel manifold We know that

$$\forall \mathbf{Q} \in \text{St}(n, m), \quad \exists \mathbf{Q}' \in N(\epsilon), \quad \text{such that} \quad \|\mathbf{Q} - \mathbf{Q}'\| \leq \epsilon, \quad \text{and} \quad |N(\epsilon)| \leq \left(\frac{6}{\epsilon}\right)^{mn}.$$

We have

$$\begin{aligned} & \sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}}(z) \right\| = \sup_{\mathbf{Q} \in \text{St}(n, m), \|\mathbf{e}\| \leq \epsilon} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}+\mathbf{e}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}+\mathbf{e}}(z) \right\| \\ & \leq \sup_{\mathbf{Q}' \in N(\epsilon)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}'}(z) \right\| + \sup_{\mathbf{Q}' \in N(\epsilon), \|\mathbf{e}\| \leq \epsilon} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'+\mathbf{e}}(\bar{z}_i) - \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) \right\| \\ & + \sup_{\mathbf{Q}' \in N(\epsilon), \|\mathbf{e}\| \leq \epsilon} \|\mathbb{E}f_{\mathbf{Q}'+\mathbf{e}}(z) - \mathbb{E}f_{\mathbf{Q}'}(z)\|. \end{aligned}$$

Due to the Lipschitz condition in Equation (15) and Equation (18), we have

$$\begin{aligned} & \|\mathbb{E}f_{\mathbf{Q}'+\mathbf{e}}(z) - \mathbb{E}f_{\mathbf{Q}'}(z)\| \leq L_f \|\mathbf{e}\|, \\ & \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'+\mathbf{e}}(\bar{z}_i) - \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) \right\| \leq \|f_{\mathbf{Q}'+\mathbf{e}}(\bar{z}) - f_{\mathbf{Q}'}(\bar{z})\| \leq \bar{L}_f \|\mathbf{e}\|. \end{aligned}$$

This implies

$$\sup_{\mathbf{Q} \in \text{St}(n, m)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}}(z) \right\| \leq \sup_{\mathbf{Q}' \in N(\epsilon)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}'}(z) \right\| + (L_f + \bar{L}_f)\epsilon.$$

Choose $\epsilon \leq \frac{t}{2(L_f + \bar{L}_f)}$, and we have

$$\begin{aligned} \mathcal{T}_1 & \leq \mathbb{P}\left(\sup_{\mathbf{Q}' \in N(\epsilon)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}'}(z) \right\| \geq t - (L_f + \bar{L}_f)\epsilon\right) \\ & \leq \mathbb{P}\left(\sup_{\mathbf{Q}' \in N(\epsilon)} \left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}'}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}'}(z) \right\| \geq t/2\right) \\ & \stackrel{(a)}{\leq} |N(\epsilon)| \mathbb{P}\left(\left\| \frac{1}{r} \sum_{i=1}^r f_{\mathbf{Q}}(\bar{z}_i) - \mathbb{E}f_{\mathbf{Q}}(z) \right\| \geq t/2\right) \\ & \stackrel{(b)}{\leq} \left(\frac{6}{\epsilon}\right)^{mn} d_1 \exp\left(-\frac{pt^2}{32R_2 + 16R_1 t/3}\right) \\ & \leq \exp\left(-\min\left\{\frac{rt^2}{64R_2}, \frac{3rt}{32R_1}\right\}\right) + nm \log\left(\frac{12(L_f + \bar{L}_f)}{t}\right) + \log(d_1). \end{aligned}$$

Inequality (a) can be obtained by taking a union bound over the ϵ -net and inequality (b) follows Lemma D.2.

We finish the proof by taking $t = m\delta$.

□

We then use Lemma B.3 to establish concentration for ℓ_p objectives.

Lemma B.4. (Concentration) Suppose $\mathbf{X} \in \mathbb{R}^{n \times r}$ follows $\mathcal{BG}(\theta)$. For any given $\theta \in (0, 1)$ and $\delta > 0$, whenever

$$r \geq C\delta^{-2}n \log(n/\delta)(\theta n \log^2 n)^{\frac{p}{2}},$$

we have

$$\sup_{\mathbf{A} \in \mathbb{O}(n)} \frac{1}{nr} \left| \|\mathbf{AY}\|_p^p - \mathbb{E}(\|\mathbf{AY}\|_p^p) \right| \leq \delta,$$

with probability at least $1 - r^{-1}$.

Proof. Note that

$$\|\mathbf{AD}_0\mathbf{X}_0\|_p^p = \sum_{i=1}^r \|\mathbf{AD}_0\mathbf{x}_i\|_p^p.$$

To use Lemma B.3, we define $\mathbf{Q} = \mathbf{AD}_0$, $\mathbf{z} = \mathbf{x}$ and $f_{\mathbf{Q}}(\mathbf{z}) = \|\mathbf{Qz}\|_p^p$.

Bound R_2

Denote $\mathbf{Q} \in \mathbb{O}(n)$ and $\mathbf{Q} = [\mathbf{q}_1^*; \dots; \mathbf{q}_n^*]$.

$$\begin{aligned} \mathbb{E}\|f_{\mathbf{Q}}(\bar{\mathbf{z}})\|^2 &\leq \mathbb{E}|f_{\mathbf{Q}}(\mathbf{z})|^2 = \mathbb{E}\|\mathbf{Qz}\|_p^{2p} = \mathbb{E}(\|\mathbf{Qz}\|_p^p)^2 = \mathbb{E}\left(\sum_{i=1}^n |\mathbf{q}_i^* \mathbf{z}|^p\right)^2 \\ &= \sum_{i=1}^n \mathbb{E}|\mathbf{q}_i \mathbf{z}|^{2p} + \sum_{i=1}^n \sum_{j=1, j \neq i}^n \mathbb{E}|\mathbf{q}_i \mathbf{z}|^p \mathbb{E}|\mathbf{q}_j \mathbf{z}|^p \end{aligned} \quad (20)$$

$$\stackrel{(a)}{\leq} \gamma_{2p}\theta n + \theta^2 \gamma_p^2(n^2 - n) = R_2.$$

Bound R_1 By Cauchy inequality, we have

$$\|\mathbf{Q}\bar{\mathbf{z}}\|_p^p \leq \|\mathbf{Q}\bar{\mathbf{z}}\|_2^p \leq \|\mathbf{Q}\|_2^p \|\bar{\mathbf{z}}\|_2^p = (\|\bar{\mathbf{z}}\|_2^2)^{p/2} \leq (B^2 \|\bar{\mathbf{z}}\|_0^2)^{p/2} \stackrel{(a)}{\leq} (B^2 4\theta n \log(r))^{p/2}. \quad (21)$$

with probability at least $1 - \exp(-\theta n)$ and (a) follows Lemma D.3.

Bound L_f, \bar{L}_f Note that sample complexity bound in Lemma B.3 is proportional to $\log(L_f + \bar{L}_f)$ and $\log(L_f + \bar{L}_f) = \Theta(\log n)$ as long as L_f and \bar{L}_f is in the polynomial of n . Thus, it is sufficient to bound them in the polynomials of n . By calculation, we have

$$\bar{L}_f \leq c_1 p n^{p+1} B^p, \quad L_f \leq c_2 n p. \quad (22)$$

Bound B_f

Apply Lemma B.1, we have

$$B_f = \mathbb{E}\|\mathbf{Qz}\|_p^p \leq \sum_{i=1}^n \mathbb{E}|\mathbf{q}_i \mathbf{z}|^p \leq n\theta \gamma_p. \quad (23)$$

We finish the proof by substituting (20-23) into Lemma B.3.

□

In the following two lemmas, we show the sharpness of ℓ_p objectives.

Lemma B.5. Suppose $\mathbf{q} \in \mathbb{S}^{n-1}$, $r = 0.5 \min_{1 \leq i \leq n} \{\|\mathbf{q} - \mathbf{e}_i\|_2^2, \|\mathbf{q} + \mathbf{e}_i\|_2^2\}$, then we have

$$r \leq \frac{C_p}{\theta(1-\theta)} (\theta - \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p), \quad (24)$$

where $C_p = (1 - 2(0.5)^{\frac{p}{2}})^{-1}$.

Proof. Without loss of generality, we assume $r = 0.5 \min\{\|\mathbf{q} - \mathbf{e}_n\|_2^2, \|\mathbf{q} + \mathbf{e}_n\|_2^2\}$ and thus $r = |1 - q_n|^2 = \epsilon^2$.

We first show (24) holds when $\epsilon \leq \frac{1}{\sqrt{2}}$. Let $\Omega' = \Omega \setminus \{n\}$, we have

$$\begin{aligned} \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p &= \mathbb{E}_\Omega (\|\mathbf{q}_\Omega\|_2^2)^{\frac{p}{2}} = \theta \mathbb{E}_{\Omega'} (\|\mathbf{q}_{\Omega'}\|^2 + 1 - \epsilon^2)^{p/2} + (1-\theta) \mathbb{E}_{\Omega'} (\|\mathbf{q}_{\Omega'}\|^2)^{p/2} \\ &\stackrel{(a)}{\leq} \theta(1-\theta)(\epsilon^p + (1-\epsilon^2)^{p/2}) + (1-\theta)^2 \cdot 0 + \theta^2 \cdot 1. \end{aligned}$$

Equality in (a) holds only if $\|\mathbf{q}\|_0 = 2$.

Define $f(\epsilon) = 1 - \epsilon^p - (1 - \epsilon^2)^{p/2}$ and $g(\epsilon) = f(\epsilon) - 2f(\frac{1}{\sqrt{2}})\epsilon^2$. For $0 \leq \epsilon \leq \frac{1}{\sqrt{2}}$, $g(\epsilon) \geq 0$ and the equality holds only if $\epsilon = 0$ or $\epsilon = \frac{1}{\sqrt{2}}$.

Thus, we have

$$\inf_{\mathbf{q} \in \mathbb{S}^{n-1}} \theta - \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p = \theta - \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p = \theta(1-\theta)f(\epsilon) \geq 2\theta(1-\theta)f\left(\frac{1}{\sqrt{2}}\right)\epsilon^2 = 2\theta(1-\theta)f\left(\frac{1}{\sqrt{2}}\right)r,$$

which implies

$$r \leq \frac{C_p}{2\theta(1-\theta)} (\theta - \mathbb{E} \|\mathbf{q}_\Omega\|_2^p),$$

for $\epsilon \leq \frac{1}{\sqrt{2}}$, where $C_p = f^{-1}(\frac{1}{\sqrt{2}})$.

Next, we show that (24) holds when $\epsilon > \frac{1}{\sqrt{2}}$. As $r \leq 1$ always holds, we can take the smallest c such that

$$1 \leq \frac{c}{\theta(1-\theta)} (\theta - \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p), \quad (25)$$

holds for $\epsilon > \frac{1}{\sqrt{2}}$.

Denote $\mathbf{w} = \mathbf{q}_{1:n-1}$ and $g(\mathbf{w}) = \mathbb{E}_\Omega \|\sqrt{1 - \|\mathbf{w}\|_2^2} \mathbf{e}_n\|_2^p$. Then, we have

$$\theta - \mathbb{E}_\Omega \|\mathbf{q}_\Omega\|_2^p = \theta - g(\mathbf{w}) \geq \theta - g\left(\frac{\mathbf{w}}{\sqrt{2\epsilon}}\right) \geq 2\theta(1-\theta)f\left(\frac{1}{\sqrt{2}}\right)\left(\frac{1}{\sqrt{2}}\right)^2 \geq f\left(\frac{1}{\sqrt{2}}\right)\theta(1-\theta).$$

Thus, we take $c = C_p$ in (25) to finish the proof. □

Lemma B.6. (Sharpness) Suppose $\mathbf{Q} \in \mathbb{O}(n)$, and then $\exists \mathbf{P} \in SP(n)$ such that

$$\theta - \frac{1}{nr\gamma_p} \mathbb{E}_Y \|\mathbf{QY}\|_p^p \geq \frac{\theta(1-\theta)}{2nC_p} \|\mathbf{Q} - \mathbf{P}\|_2^2, \quad (26)$$

where $C_p = (1 - 2(0.5)^{\frac{p}{2}})^{-1}$.

Proof. Denote $\mathbf{Q} = [\mathbf{q}_1; \dots; \mathbf{q}_n]$, we have

$$\begin{aligned} \theta - \frac{1}{nr\gamma_p} \mathbb{E}_Y \|\mathbf{QY}\|_p^p &= \theta - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\Omega \|\mathbf{q}_{i,\Omega}\|_2^p = \frac{1}{n} \sum_{i=1}^n (\theta - \mathbb{E}_\Omega \|\mathbf{q}_{i,\Omega}\|_2^p) \stackrel{(a)}{\geq} \frac{\theta(1-\theta)}{2nC_p} \sum_{i=1}^n \min_{1 \leq j \leq n} \{\|\mathbf{q}_i - \mathbf{e}_j\|_2^2, \|\mathbf{q}_i + \mathbf{e}_j\|_2^2\} \\ &\stackrel{(b)}{=} \frac{\theta(1-\theta)}{2nC_p} \|\mathbf{Q} - \mathbf{P}\|_2^2 \end{aligned}$$

where (a) follows Lemma B.5 and (b) follows (A.36) in [26]. □

We give a proof for Theorem 2.1 through Theorem B.1, assisted by Lemma B.4 and Lemma B.6.

Theorem B.2. Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$ with $\theta \in (0, 1)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ is an orthogonal dictionary, and $\mathbf{Y} = \mathbf{D}_0 \mathbf{X}$. Suppose $\hat{\mathbf{A}}$ is a global maximizer to

$$\underset{\mathbf{A}}{\text{maximize}} \quad \|\mathbf{A}\mathbf{Y}\|_p^p \text{ subject to } \mathbf{A} \in \mathbb{O}(n).$$

Provided that the sample size $r = \Omega(\theta\delta^{-2}n \log(n/\delta)(n \log^2 n)^{\frac{p}{2}})$, then for $\delta > 0$, there exists a signed permutation $\mathbf{\Pi}$, such that

$$\frac{1}{n} \|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq C_\theta \delta,$$

with probability at least $1 - r^{-1}$ and C_θ is a constant that depends on θ .

Proof. To use Theorem B.1, we have to check the concentration of empirical objective and sharpness of the population objective.

Concentration From Lemma B.4, for any $\delta > 0$, whenever $r > \Phi(\delta, n) = \Omega(\theta\delta^{-2}n \log(n/\delta)(n \log^2 n)^{\frac{p}{2}})$, we have

$$\sup_{\mathbf{A} \in \mathbb{O}(n)} \frac{1}{nr} \left| \|\mathbf{A}\mathbf{Y}\|_p^p - \mathbb{E}(\|\mathbf{A}\mathbf{Y}\|_p^p) \right| \leq \delta,$$

with probability at least $1 - r^{-1}$.

Sharpness From Lemma B.6, $\forall \mathbf{A} \in \mathbb{O}(n)$, we

$$\frac{1}{nr} (\mathbb{E}(\|\mathbf{D}_0^* \mathbf{Y}\|_p^p) - \mathbb{E}(\|\mathbf{A}\mathbf{Y}\|_p^p)) = \gamma_p \theta - \frac{1}{nr} \mathbb{E}(\|\mathbf{A}\mathbf{Y}\|_p^p) \geq \frac{\gamma_p \theta (1 - \theta)}{2nC_p} \|\mathbf{A} - \mathbf{P}\|_2^2$$

Let $C_\theta = \frac{4C_p}{\theta(1-\theta)\gamma_p}$, by Theorem B.1, we have

$$\frac{1}{n} \|\hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi}\|_F^2 \leq C_\theta \delta,$$

whenever $p \geq \Omega(\theta\delta^{-2}n \log(n/\delta)(n \log^2 n)^{\frac{p}{2}})$. □

B.2 Proof of Theorem 2.2

In this subsection, we prove Theorem 2.2 via Theorem B.1. Specifically, the concentration and sharpness conditions are established in Lemma B.8 and Lemma B.10 respectively. Then we show the global maximizers of all ℓ_p -based formulation are very close to the true dictionary with high probability via Theorem B.1.

Lemma B.7. (*Robustness under Gaussian noise*) Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ is an orthogonal dictionary, $\mathbf{Y}_N = \mathbf{D}_0 \mathbf{X} + \mathbf{G}$, and $\mathbf{G} \in \mathbb{R}^{n \times r}$ with $G_{i,j} \sim \mathcal{N}(0, \eta^2)$. The only global maximizers to

$$\underset{\mathbf{A}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{X}_0, \mathbf{G}} \|\mathbf{A}\mathbf{Y}_N\|_p^p \text{ subject to } \mathbf{A} \in \mathbb{O}(n),$$

are $\mathbf{A}^* = \mathbf{D}_0 \mathbf{\Pi}$, where $\mathbf{\Pi}$ is any signed permutation matrix.

Proof. Consider

$$\underset{\mathbf{A}}{\text{maximize}} \quad \mathbb{E}_{\mathbf{Y}} \|\mathbf{A}\mathbf{Y}\|_p^p \text{ subject to } \mathbf{A} \in \mathbb{O}(n).$$

Denote $\mathbf{A} = [\mathbf{a}_1^*; \dots; \mathbf{a}_m^*]$ where $\mathbf{a}_i \in \mathbb{R}^n$, $\mathbf{Y} = [\mathbf{y}_1, \dots, \mathbf{y}_r]$, and $\mathbf{G} = [\mathbf{g}_1, \dots, \mathbf{g}_m]$.

$$\begin{aligned} \mathbb{E}_{\mathbf{Y}_N} \|\mathbf{A}\mathbf{Y}_N\|_p^p &= \mathbb{E}_{\mathbf{Y}, \mathbf{G}} \|\mathbf{A}(\mathbf{Y} + \mathbf{G})\|_p^p = \mathbb{E}_{\mathbf{Y}, \mathbf{G}} \|\mathbf{A}(\mathbf{Y} + \mathbf{G})\|_p^p \\ &= \sum_{i=1}^n \sum_{j=1}^r \mathbb{E}_{\mathbf{y}_j, \mathbf{g}_j} |\mathbf{a}_i^* \mathbf{y}_j + \mathbf{a}_i^* \mathbf{g}_j|_p^p = r \sum_{i=1}^n \mathbb{E}_{\Omega} (\|(\mathbf{a}_i \mathbf{D}_0)_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}}. \end{aligned} \quad (27)$$

This term achieves its maxima only if $\|\mathbf{a}_i \mathbf{D}_0\|_0 = 1, \forall i$ and $\mathbf{A} \in \mathbb{O}(n)$. Thus, the global maximum is achieved only if $\mathbf{A}^* = \mathbf{D}_0 \mathbf{P}$, where \mathbf{P} is any signed permutation matrix. □

Lemma B.8. (Concentration) Suppose $\mathbf{X} \in \mathbb{R}^{n \times r}$ follows $\mathcal{BG}(\theta)$. For any given $\theta \in (0, 1)$, whenever

$$r \geq C\delta^{-2}n \log(n/\delta)((1 + \eta^2)n \log n)^{\frac{p}{2}},$$

we have

$$\sup_{\mathbf{A} \in \mathbb{O}(n)} \frac{1}{nr} \left| \|\mathbf{A}\mathbf{Y}_N\|_p^p - \mathbb{E}(\|\mathbf{A}\mathbf{Y}_N\|_p^p) \right| \leq \delta,$$

with probability at least $1 - r^{-1}$.

Proof. Similar to the proof for Lemma B.4, we use Lemma B.3 to prove this lemma. Note that

$$\|\mathbf{A}\mathbf{Y}_N\|_p^p = \|\mathbf{A}\mathbf{D}_0(\mathbf{X}_0 + \mathbf{D}_0^*\mathbf{G})\|_p^p = \sum_{i=1}^r \|\mathbf{A}\mathbf{D}_0(\mathbf{x}_i + \mathbf{g}_i)\|_p^p,$$

and we can define $\mathbf{Q} = \mathbf{A}\mathbf{D}_0$ and $f_{\mathbf{Q}}(z) = \|\mathbf{Q}z\|_p^p$ and to use Lemma B.3.

Bound R_2

Denote $\mathbf{Q} \in \mathbb{O}(n)$ and $\mathbf{Q} = [\mathbf{q}_1^*; \dots; \mathbf{q}_n^*]$.

$$\begin{aligned} \mathbb{E}\|f_{\mathbf{Q}}(\bar{z})\|^2 &\leq \mathbb{E}|f_{\mathbf{Q}}(z)|^2 = \mathbb{E}\|\mathbf{Q}z\|_p^{2p} = \mathbb{E}(\|\mathbf{Q}z\|_p^p)^2 = \mathbb{E}\left(\sum_{i=1}^n |\mathbf{q}_i^* z|^p\right)^2 \\ &= \sum_{i=1}^n \mathbb{E}|\mathbf{q}_i z|^{2p} + \sum_{i=1}^n \sum_{j=1, j \neq i}^n \mathbb{E}|\mathbf{q}_i z|^p \mathbb{E}|\mathbf{q}_j z|^p \end{aligned} \quad (28)$$

$$\leq Cn^2(1 + \eta^2)^p = R_2.$$

Bound R_1 By Cauchy inequality, we have

$$\|\mathbf{Q}\bar{z}\|_p^p \leq \|\mathbf{Q}\bar{z}\|_2^p \leq \|\mathbf{Q}\|_2^p \|\bar{z}\|_2^p = (\|\bar{z}\|_2^2)^{p/2} \leq (nB^2)^{p/2} = R_1. \quad (29)$$

Bound L_f, \bar{L}_f Note that the sample complexity bound in Lemma B.3 is propositional to $\log(L_f + \bar{L}_f)$ and $\log(L_f + \bar{L}_f) = \Theta(\log n)$ as long as L_f and \bar{L}_f is in the polynomial of n . Thus, it is sufficient to bound them in the polynomials of n . By calculation, we have

$$\bar{L}_f \leq c_1 pn^{p+1} B^p, \quad L_f \leq c_2 np(1 + \eta)^p. \quad (30)$$

Bound B_f

Applying Lemma B.7, we have

$$B_f = \mathbb{E}\|\mathbf{Q}z\|_p^p \leq \sum_{i=1}^n \mathbb{E}|\mathbf{q}_i z|^p \leq n\gamma_p(1 + \eta^2)^{p/2}. \quad (31)$$

□

Lemma B.9. Suppose $\mathbf{q} \in \mathbb{S}^{n-1}$, $r = 0.5 \min_{1 \leq i \leq n} \{\|\mathbf{q} - \mathbf{e}_i\|_2^2, \|\mathbf{q} + \mathbf{e}_i\|_2^2\}$, and then we have

$$r \leq \frac{C_{\eta,p}}{\theta(1 - \theta)} (\theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \mathbb{E}_{\Omega}(\|\mathbf{q}_{\Omega}\|_2^2 + \eta^2)^{\frac{p}{2}}), \quad (32)$$

where $C_{\eta,p} = ((1 + \eta^2)^{p/2} + \eta^p - 2(0.5 + \eta^2)^{p/2})^{-1}$.

Proof. Without loss of generality, we assume $r = 0.5 \min\{\|\mathbf{q} - \mathbf{e}_n\|_2^2, \|\mathbf{q} + \mathbf{e}_n\|_2^2\}$ and thus $r = |1 - q_n|^2 = \epsilon^2$.

We first show (32) holds when $\epsilon \leq \frac{1}{\sqrt{2}}$. Let $\Omega' = \Omega \setminus \{n\}$, and we have

$$\begin{aligned} \mathbb{E}_\Omega(\|\mathbf{q}_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}} &= \theta \mathbb{E}_{\Omega'}(\|\mathbf{q}_{\Omega'}\|^2 + 1 - \epsilon^2 + \eta^2)^{p/2} + (1 - \theta) \mathbb{E}_{\Omega'}(\|\mathbf{q}_{\Omega'}\|^2 + \eta^2)^{p/2} \\ &\stackrel{(a)}{\leq} \theta(1 - \theta)((\epsilon^2 + \eta^2)^{\frac{p}{2}} + (1 - \epsilon^2 + \eta^2)^{p/2}) + (1 - \theta)^2 \eta^p + \theta^2(1 + \eta^2)^{\frac{p}{2}}. \end{aligned}$$

Equality in (a) holds only if $\|\mathbf{q}\|_0 = 2$.

Define $f(\epsilon) = (1 + \eta^2)^{\frac{p}{2}} + \eta^p - (1 - \epsilon^2 + \eta^2)^{\frac{p}{2}} - (\epsilon^2 + \eta^2)^{\frac{p}{2}}$ and $g(\epsilon) = f(\epsilon) - 2f(\frac{1}{\sqrt{2}})\epsilon^2$. For $0 \leq \epsilon \leq \frac{1}{\sqrt{2}}$, we have $g(\epsilon) \geq 0$ and the equality holds only if $\epsilon = 0$ or $\epsilon = \frac{1}{\sqrt{2}}$.

Thus, we have

$$\begin{aligned} \inf_{\mathbf{q} \in \mathbb{S}^{n-1}} \theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \mathbb{E}_\Omega(\|\mathbf{q}_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}} &= \theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \sup_{\mathbf{q} \in \mathbb{S}^{n-1}} \mathbb{E}_\Omega(\|\mathbf{q}_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}} \\ &= \theta(1 - \theta)f(\epsilon) \geq 2\theta(1 - \theta)f\left(\frac{1}{\sqrt{2}}\right)\epsilon^2 = 2\theta(1 - \theta)f\left(\frac{1}{\sqrt{2}}\right)r, \end{aligned}$$

which implies

$$r \leq \frac{C_{\eta,p}}{2\theta(1 - \theta)}(\theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \mathbb{E}_\Omega(\|\mathbf{q}_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}}),$$

for $\epsilon \leq \frac{1}{\sqrt{2}}$, where $C_{\eta,p} = f^{-1}(\frac{1}{\sqrt{2}})$.

Next, we show that (32) holds when $\epsilon > \frac{1}{\sqrt{2}}$. Considering that we always have $r \leq 1$, we can take the smallest c such that

$$1 \leq \frac{c}{2\theta(1 - \theta)}(\theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \mathbb{E}_\Omega(\|\mathbf{q}_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}}), \quad (33)$$

holds for $\epsilon > \frac{1}{\sqrt{2}}$.

Denote $\mathbf{w} = \mathbf{q}_{1:n-1}$ and define $g(\mathbf{w}) = \mathbb{E}_\Omega(\|[\mathbf{w}, \sqrt{1 - \|\mathbf{w}\|_2^2}]_\Omega\|_2^2 + \eta^2)^{\frac{p}{2}}$. Then we have

$$\begin{aligned} \theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \mathbb{E}_\Omega\|\mathbf{q}_\Omega\|_2^p &= \theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - g(\mathbf{w}) \\ &\geq \theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - g\left(\frac{\mathbf{w}}{\sqrt{2\epsilon}}\right) \geq 2\theta(1 - \theta)f\left(\frac{1}{\sqrt{2}}\right)\left(\frac{1}{\sqrt{2}}\right)^2 = f\left(\frac{1}{\sqrt{2}}\right)\theta(1 - \theta). \end{aligned}$$

Thus, we take $c = C_{\eta,p}$ in (33) to finish the proof. □

Lemma B.10. (Sharpness) Suppose $\mathbf{Q} \in \mathbb{O}(n)$, and then $\exists \mathbf{P} \in SP(n)$ such that

$$\theta(1 + \eta^2)^{\frac{p}{2}} + (1 - \theta)\eta^p - \frac{1}{nr\gamma_p} \mathbb{E}_Y \|\mathbf{QY}\|_p^p \geq \frac{\theta(1 - \theta)}{2nC_{\eta,p}} \|\mathbf{Q} - \mathbf{P}\|_2^2, \quad (34)$$

where $C_{\eta,p} = ((1 + \eta^2)^{p/2} + \eta^p - 2(0.5 + \eta^2)^{p/2})^{-1}$.

Proof. By a similar argument in Lemma B.6. □

Theorem B.3. Let $\mathbf{X} \in \mathbb{R}^{n \times r}$, $x_{i,j} \sim \mathcal{BG}(\theta)$, $\mathbf{D}_0 \in \mathbb{O}(n)$ is orthogonal dictionary, and $\mathbf{Y}_N = \mathbf{D}_0 \mathbf{X} + \mathbf{G}$, and $\mathbf{G} \in \mathbb{R}^{n \times r}$ with $G_{i,j} \sim \mathcal{N}(0, \eta^2)$. Suppose $\hat{\mathbf{A}}$ is a global maximizer to

$$\underset{\mathbf{A}}{\text{maximize}} \quad \|\mathbf{A}\mathbf{Y}_N\|_p^p \text{ subject to } \mathbf{A} \in \mathbb{O}(n),$$

then for $\delta > 0$, there exists a signed permutation $\mathbf{\Pi}$, such that

$$\frac{1}{n} \left\| \hat{\mathbf{A}}^* - \mathbf{D}_0 \mathbf{\Pi} \right\|_F^2 \leq C_\theta \delta,$$

with probability at least $1 - r^{-1}$ as long as $r = \Omega(\delta^{-2} n \log(n/\delta) ((1 + \eta^2) n \log n)^{\frac{p}{2}} \xi_\eta^2)$ where $\xi_\eta = (1 + \eta^2)^{p/2} + \eta^p - 2(0.5 + \eta^2)^{p/2}$ and C_θ is a constant depends on θ .

Proof. By a similar argument in Theorem 2.1 and $C_\theta = \frac{4}{\gamma_p \theta (1 - \theta)}$. \square

C Convergence Result

C.1 Proof of Proposition 3.1

Proposition C.1. Denote $SOR_i^{(t)}$ as the value of SOR_i at the t -th iteration and $\mathbf{q} = \mathbf{a}^{(t)}$ as the variable at the t -th iteration. Then the evolution of SOR_i follows

$$SOR_i^{(t+1)} = SOR_i^{(t)} (1 + \tau_i(\mathbf{q})),$$

and

$$\tau_i(\mathbf{q}) = \frac{\mathbb{E}_{\Omega'} [\|\mathbf{q}_{\Omega'}, q_n\|_2^k] - \mathbb{E}_{\Omega'} [\|\mathbf{q}_{\Omega'}, q_i\|_2^k]}{\frac{\theta}{1-\theta} \mathbb{E}_{\Omega'} [\|\mathbf{q}_{\Omega'}, q_i, q_n\|_2^k] + \mathbb{E}_{\Omega'} [\|\mathbf{q}_{\Omega'}, q_i\|_2^k]},$$

where $\Omega' = \Omega \setminus \{n\} \setminus \{i\}$ and $k = p - 2$. Two properties of $\tau_i(\mathbf{q})$ are listed below

1. $0 \leq \tau_i(\mathbf{q}) \leq \frac{1-\theta}{\theta}$ always holds and $\tau_i(\mathbf{q}) > 0$ if $q_n > q_i$.
2. $\tau_i(\mathbf{q})$ is monotonically increasing in q_n and decreasing in q_i .

Proof. To compute SOR_i , we separate the population gradient into several parts

$$\begin{aligned} \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega] &= \mathbb{P}(n \in \Omega, i \in \Omega) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \in \Omega, i \in \Omega] + \mathbb{P}(n \in \Omega, i \notin \Omega) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \in \Omega, i \notin \Omega] \\ &+ \mathbb{P}(n \notin \Omega, i \in \Omega) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \notin \Omega, i \in \Omega] + \mathbb{P}(n \notin \Omega, i \notin \Omega) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \notin \Omega, i \notin \Omega] \end{aligned}$$

where $k = p - 2$.

The probability for each part is

$$\mathbb{P}(n \in \Omega, i \in \Omega) = \theta^2, \quad \mathbb{P}(n \in \Omega, i \notin \Omega) = \mathbb{P}(n \notin \Omega, i \in \Omega) = \theta(1 - \theta), \quad \mathbb{P}(n \notin \Omega, i \notin \Omega) = (1 - \theta)^2.$$

Thus, denote $\Omega' = \Omega \setminus \{n\} \setminus \{i\}$ we can the i -th

$$\begin{aligned} &\mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega]_n \\ &= \theta^2 \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \in \Omega, i \in \Omega] + \theta(1 - \theta) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \in \Omega, i \notin \Omega] + \theta(1 - \theta) \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \notin \Omega, i \in \Omega] \\ &+ (1 - \theta)^2 \mathbb{E}_\Omega [\|\mathbf{a}_\Omega\|_2^k \mathbf{a}_\Omega | n \notin \Omega, i \notin \Omega] \\ &= \theta^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i, a_n\|_2^k] a_i + \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i\|_2^k] a_i + 0\theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_n\|_2^k] + 0(1 - \theta)^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}\|_2^k] \\ &= a_i (\theta^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i, a_n\|_2^k] + \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i\|_2^k]). \end{aligned}$$

Thus, SOR_i can be computed as

$$\begin{aligned} SOR_i^{(t+1)} &= \frac{a_n (\theta^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i, a_n\|_2^k] + \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_n\|_2^k])}{a_i (\theta^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i, a_n\|_2^k] + \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i\|_2^k])} \\ &= \frac{a_n}{a_i} \left(1 + \frac{\theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_n\|_2^k] - \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i\|_2^k]}{\theta^2 \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i, a_n\|_2^k] + \theta(1 - \theta) \mathbb{E}_{\Omega'} [\|\mathbf{a}_{\Omega'}, a_i\|_2^k]} \right) \\ &= SOR_i^{(t)} (1 + \tau_i(\mathbf{q})). \end{aligned}$$

\square

C.2 Proof of Theorem 3.1

Theorem C.1. Assume we apply Algorithm 3 to solve problem (8) and $\mathbf{a}^{(0)}$ follows a uniform distribution over the sphere, and denote $\tau(\mathbf{q}) = \min_{i=1, \dots, n-1} \tau_i(\mathbf{q})$, then there exists $T_\tau \leq \log_{1+\tau(\mathbf{a}^{(0)})}(\sqrt{n})$, and $1 \leq i \leq n$ such that

$$\|\mathbf{a}^{(t)} - \mathbf{D}_0 \mathbf{e}_i\|_2 \leq \left(1 + \tau(\mathbf{a}^{(T_\tau)})\right)^{T_\tau - t}, \forall t \geq T_\tau,$$

almost surely and the convergence rate $\lim_{k \rightarrow \infty} \frac{\|\mathbf{a}^{(k+1)} - \mathbf{D}_0 \mathbf{e}_i\|_2}{\|\mathbf{a}^{(k)} - \mathbf{D}_0 \mathbf{e}_i\|_2} = \frac{1}{\theta}$.

Proof. As discussed in Section C.2, we assume $\mathbf{D}_0 = \mathbf{I}$ and $i = n$. From Proposition C.1, we know that for $0 < t_1 < t_2$

$$\text{SOR}^{(t_2)} \geq \text{SOR}^{(t_1)} \prod_{i=t_1}^{t_2-1} (1 + \tau(\mathbf{a}^{(i)})) > \text{SOR}^{(t_1)} (1 + \tau(\mathbf{a}^{(t_1)}))^{t_2 - t_1} \quad (35)$$

At initialization, we have

$$\text{SOR}^{(0)} = \frac{a_n^{(0)}}{\|\mathbf{a}_{-n}^{(0)}\|_2} > \frac{1}{\sqrt{n}} \quad (36)$$

almost surely since we assume $a_n \geq a_i, \forall i$.

Thus, combining (35) and (36), there exists $T_\tau \leq \log_{1+\tau(\mathbf{a}^{(0)})}(\sqrt{n})$ such that $\text{SOR}^{(T_\tau)} > 1$. Then, for $t > T_\tau$, we have

$$\begin{aligned} \|\mathbf{a}^{(t)} - \mathbf{e}_n\|_2^2 &= 2 - 2\sqrt{\frac{(\text{SOR}^{(t)})^2}{(\text{SOR}^{(t)})^2 + 1}} \leq \frac{1}{(\text{SOR}^{(t)})^2} \leq \frac{1}{\left(\text{SOR}^{(T_\tau)} (1 + \tau(\mathbf{a}^{(T_\tau)}))^{t - T_\tau}\right)^2} \\ &\leq \left(1 + \tau(\mathbf{a}^{(T_\tau)})\right)^{2(T_\tau - t)}. \end{aligned}$$

□

D Technical Lemmas

Lemma D.1. Suppose $\mathbf{g} \in \mathbb{R}^n$ with $g_i \sim \mathcal{N}(0, \sigma^2)$ and $\mathbf{a} \in \mathbb{R}^n$ is a fixed vector. Then

$$\mathbb{E}(|\mathbf{a}^* \mathbf{g}|^p) = \gamma_p \|\mathbf{a}\|_2^p,$$

where $\gamma_p = \sigma^p 2^{p/2} \frac{\Gamma(\frac{p+1}{2})}{\sqrt{\pi}}$.

Proof. Due to the rotation invariance property of Gaussian, $\mathbf{a}^* \mathbf{g} \sim \mathcal{N}(0, \|\mathbf{a}\|_2)$. Therefore, $\mathbb{E}(|\mathbf{a}^* \mathbf{g}|^p)$ is the p -th moment of an absolute Gaussian random variable with zero mean and variance $\|\mathbf{a}\|_2^2 \sigma^2$. By simple calculation

$$\mathbb{E}(|\mathbf{a}^* \mathbf{g}|^p) = \sigma^p 2^{p/2} \frac{\Gamma(\frac{p+1}{2})}{\sqrt{\pi}} \|\mathbf{a}\|_2^p.$$

□

Lemma D.2. (ϵ -net over Stiefel manifold) [26] There is a covering ϵ -net $N(\epsilon)$ for Stiefel manifold $\mathcal{M} = \{\mathbf{W} \in \mathbb{R}^{n \times m} | \mathbf{W}^* \mathbf{W} = \mathbf{I}, n > m\}$, in operator norm

$$\forall \mathbf{W} \in \mathcal{M}, \exists \mathbf{W}' \in N(\epsilon), \text{ such that } \|\mathbf{W} - \mathbf{W}'\| \leq \epsilon$$

of size $|N(\epsilon)| \leq \left(\frac{6}{\epsilon}\right)^{nm}$.

Proof. See Lemma D.4 in [26].

□

Lemma D.3. ([27]) Let $v \in \mathbb{R}^d$ with each entry following i.i.d. $\text{Ber}(\theta)$, then

$$\mathbb{P}(|\|v\|_0 - \theta d| \geq t\theta d) \leq 2 \exp\left(-\frac{3t^2}{2t+6}\theta d\right)$$

Proof. See Lemma A.4 in [27].

□