Domain Hang Zhang, Afshin Abdi, and Faramarz Fekri School of Electrical and Computer Engineering, Georgia Institute of Technology, Atlanta, GA.

Abstract

Compressive Sensing with a Multiple Convex Sets

In this paper, we study a general framework for compressive sensing assuming the existence of the prior knowledge that x^{\natural} belongs to the union of multiple convex sets, $x^{\natural} \in \bigcup_i C_i$. In fact, by proper choices of these convex sets in the above framework, the problem can be transformed to well known CS problems such as the phase retrieval, quantized compressive sensing, and model-based CS. First we analyze the impact of this prior knowledge on the minimum number of measurements *M* to guarantee the uniqueness of the solution. Then we formulate a universal objective function for signal recovery, which is both computationally inexpensive and flexible. Then, an algorithm based on *multiplicative weight update* and *proximal gradient descent* is proposed and analyzed for signal reconstruction. Finally, we investigate as to how we can improve the signal recovery by introducing regularizers into the objective function.

1 Introduction

In the traditional *compressive sensing* (CS) Eldar & Kutyniok (2012); Foucart & Rauhut (n.d.), sparse signal **x** is reconstructed via

$$\min_{\mathbf{x}} \|\mathbf{x}\|_{1}, \quad \text{s.t.} \quad \mathbf{y} = \mathbf{A}\mathbf{x}, \tag{1.1}$$

where $\mathbf{y} \in \mathbb{R}^M$ denotes the measurement vector and $\mathbf{A} \in \mathbb{R}^{M \times N}$ is the measurement matrix.

In this paper, we assume the existence of extra prior knowledge that **x** lies in the union of some convex sets, $\mathbf{x} \in \bigcup_{i=1}^{L} C_i$, where *L* denotes the number of constraint sets and C_i is the *i*-th convex constraint set. Therefore, we now wish to solve

$$\min_{\mathbf{x}} \|\mathbf{x}\|_{1}, \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \bigcup_{i=1}^{L} \mathcal{C}_{i}$$
(1.2)

This ill-posed inverse problem (i.e., given measurement **y**, solving for **x**) turns out to be a rather general form of CS. For example, setting $\bigcup C_i = \mathbb{R}^n$ simplifies our problem to the traditional CS problem. In the following, we will further show that by appropriate choices of these convex sets, Eq. (1.2) can be transformed to the *phase retrieval* Candes *et al.* (2011); Chen & Candes (2015), *quantized compressive sensing* Dai *et al.* (2009), or *model-based* CS Baraniuk *et al.* (2010) problems.

1.1 Relation with other problems

Phase retrieval Consider the noiseless phase retrieval problem in which the measurements are given by

$$y_i = \left| \left\langle \mathbf{a}_i, \, \mathbf{x} \right\rangle \right|^2, \ 1 \le i \le l, \tag{1.3}$$

where y_i is the *i*-th measurement and \mathbf{a}_i denotes the corresponding coefficients. Considering the first measurement, the constraint $\sqrt{y_1} = |\langle \mathbf{a}_1, \mathbf{x} \rangle|$ can be represented via $\mathbf{x} \in \mathcal{B}_+^{(1)} \cup \mathcal{B}_-^{(1)}$ where $\mathcal{B}_+^{(1)} = \{\mathbf{x} : \langle \mathbf{a}_1, \mathbf{x} \rangle = \sqrt{y_1}\}$ and $\mathcal{B}_-^{(1)} = \{\mathbf{x} : \langle \mathbf{a}_1, \mathbf{x} \rangle = -\sqrt{y_1}\}$. Following these steps, the constraints $\{y_i = \langle \mathbf{a}_i, \mathbf{x} \rangle^2\}_{i=1}^l$ can be transformed to $\mathbf{x} \in \bigcap_i (\mathcal{B}_+^{(i)} \cup \mathcal{B}_-^{(i)}) = \bigcup_{j=1}^{2^l} \mathcal{C}_j$, for some appropriately defined C_j 's given by the intersection of different $\mathcal{B}_{\pm}^{(i)}$. Setting sensing matrix $\mathbf{A} = \mathbf{0}$ will restore the phase retrieval to our setting.

Quantized compressive sensing In this scenario, the measurements are quantized, i.e.,

$$y_i = Q(\langle \mathbf{a}_i, \, \mathbf{x} \rangle), \, 1 \le i \le L, \tag{1.4}$$

where $Q(\cdot)$ is the quantizer. Since $Q^{-1}(\cdot)$ is an interval on real line, C_i would be a convex set and the quantized CS can be easily transformed to Eq. (1.2).

Model-based compressive sensing These lines of works Baraniuk *et al.* (2010); Duarte & Eldar (2011); Silva *et al.* (2011) are the most similar work to our model, where they consider

$$\mathbf{y} = \mathbf{A}\mathbf{x}^{\natural}, \ \mathbf{x}^{\natural} \in \bigcup_{i} \mathcal{L}_{i}.$$
(1.5)

Here, \mathcal{L}_i is assumed to be a linear space whereas the only assumption we make on the models is being a convex set. Hence, their model can be regarded as a special case of our problem.

In Baraniuk *et al.* (2010), the author studied the minimum number of measurements M under different models, i.e., shape of \mathcal{L}_i , and modified CoSaMP algorithms Foucart & Rauhut (n.d.) to reconstruct signal. In Duarte & Eldar (2011), the authors expanded the signal onto different basis and transformed model-based CS to be block-sparse CS. In Silva *et al.* (2011), the author studied model-based CS with incomplete sensing matrix information and reformulated it as a matrix completion problem.

1.2 Our contribution:

Statistical Analysis We analyze the minimum number of measurements to ensure uniqueness of the solution. We first show that the conditions for the uniqueness can be represented as $\min_{u \in E} ||\mathbf{Au}||_2 > 0$, for an appropriate set *E*. Assuming the entries of the sensing matrix **A** are i.i.d. Gaussian, we relate the probability of uniqueness to the number of measurements, *M*. Our results show that depending on the structure of C_i 's, the number of measurements can be reduced significantly.

Optimization Algorithm We propose a novel formulation and the associated optimization algorithm to reconstruct the signal **x**. First, note that existing algorithms on e.g., model-based CS are not applicable to our problem as they rely heavily on the structure of constraint sets. For example, a key idea in model-based CS is to consider expansion of **x** onto the basis of each C_i and then rephrase the constraint as the block sparsity on the representation of **x** on the union of bases. However, such an approach may add complicated constraints on the coefficients of **x** in the new basis, as the sets C_i 's are not necessarily simple subspaces.

Note that although C_i 's are assumed to be convex, their union $\bigcup_i C_i$ is not necessarily a convex set, which makes the optimization problem Eq. (1.2) hard to solve. By introducing an auxiliary variable, **p**, we convert the non-convex optimization problem to a biconvex problem. Using *multiplicative weight update* Arora *et al.* (2012) from online learning theory Shalev-Shwartz *et al.* (2012), we design an algorithm with convergence speed of $\mathcal{O}(T^{-1/2})$ to a local minimum. Further, we investigate improving the performance of the algorithm by incorporating appropriate regularization. Compared to the naive idea of solving *L* simultaneous optimization problems

$$\min_{\mathbf{x}} \|\mathbf{x}\|_{1}, \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}, \quad \mathbf{x} \in \mathcal{C}_{i}, \tag{1.6}$$

and choosing the best solution out of *L* results, our method is computationally less-expensive and more flexible.

2 System Model

Let $\mathbf{A} \in \mathbb{R}^{M \times N}$ be the measurement matrix, and consider the setup

$$\mathbf{y} = \mathbf{A}\mathbf{x}^{\natural}, \text{ and } \mathbf{x}^{\natural} \in \bigcup_{i=1}^{L} \mathcal{C}_{i},$$
 (2.1)

where \mathbf{x}^{\natural} is a *K*-sparse high-dimensional signal, $\mathbf{y} \in \mathbb{R}^{M}$ is the measurement vector, and $C_i \subset \mathbb{R}^N$, i = 1, 2, ..., L, is a convex set.

Due to the sparsity of x^{\natural} , we propose to reconstruct x^{\natural} via

$$\widehat{\mathbf{x}} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{x}\|_{1}, \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}, \ \mathbf{x} \in \bigcup_{i=1}^{L} \mathcal{C}_{i},$$
 (2.2)

where $\hat{\mathbf{x}}$ denotes the reconstructed signal. Let $\mathbf{d} \triangleq \hat{\mathbf{x}} - \mathbf{x}^{\natural}$ be the deviation of the reconstructed signal $\hat{\mathbf{x}}$ from the true signal \mathbf{x}^{\natural} . In the following, we will study the inverse problem in Eq. (2.2) from two perspectives; the statistical and the computational aspects.

3 Statistical Property

In this section, we will find the minimum number of measurements M to $\hat{\mathbf{x}} = \mathbf{x}^{\natural}$, i.e., $\mathbf{d} = \mathbf{0}$.

Definition 1. The tangent cone $\mathcal{T}_{\mathbf{x}}$ for $\|\mathbf{x}\|_1$ is defined as Chandrasekaran et al. (2012)

$$\mathcal{T}_{\mathbf{x}} \triangleq \{ \mathbf{e} : \|\mathbf{x} + t\mathbf{e}\|_1 \le \|\mathbf{x}\|_1, \ \exists \ t \ge 0 \}.$$

$$(3.1)$$

Geometric interpretation of $\mathcal{T}_{\mathbf{x}}$ is that it contains all directions that lead to smaller $\|\cdot\|_1$ originating from \mathbf{x} . In the following analysis, we use \mathcal{T} as a compact notation for $\mathcal{T}_{\mathbf{x}^{\natural}}$. Easily we can prove that $\mathbf{d} \in \mathcal{T}$.

Definition 2. The Gaussian width $\omega(\cdot)$ associated with set U is defined as $\omega(U) \triangleq \mathbb{E} \sup_{\mathbf{x} \in U} \langle \mathbf{g}, \mathbf{x} \rangle$, $\mathbf{g} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, Gordon (1988).

Define cone $\widetilde{C}_{i,j}$ as

$$\widetilde{\mathcal{C}}_{i,j} \triangleq \left\{ \mathbf{z} \mid \mathbf{z} = t(\mathbf{x}_1 - \mathbf{x}_2), \exists t > 0, \mathbf{x}_1 \in \mathcal{C}_i, \mathbf{x}_2 \in \mathcal{C}_j \right\},$$
(3.2)

which denotes the cone consisting of all vectors \mathbf{z} that are parallel with $\mathbf{x}_1 - \mathbf{x}_2$, $\mathbf{x}_1 \in C_i$, and $\mathbf{x}_2 \in C_j$. Then we define event \mathcal{E} as

$$\mathcal{E} \triangleq \left\{ \bigcup_{i,j} \left(\operatorname{null}(\mathbf{A}) \bigcap \mathcal{T} \bigcap \widetilde{\mathcal{C}}_{i,j} \right) = \{\mathbf{0}\} \right\}.$$
(3.3)

Lemma 1. We can guarantee the correct recovery of \mathbf{x} , i.e., $\hat{\mathbf{x}} = \mathbf{x}^{\natural}$, iff we have event \mathcal{E} to be satisfied.

Proof. This proof is fundamentally same as Chandrasekaran *et al.* (2012); Zhang *et al.* (Nov. 2017.). First we prove that \mathcal{E} leads to $\hat{\mathbf{x}} \neq \mathbf{x}^{\natural}$. Provided $\mathbf{d} \triangleq \hat{\mathbf{x}} - \mathbf{x}^{\natural} \neq \mathbf{0}$, we then have a $\hat{\mathbf{x}} \neq \mathbf{x}^{\natural}$ such that $\|\hat{\mathbf{x}}\|_{1} \leq \|\mathbf{x}^{\natural}\|_{1}$. Setting $\mathbf{e} \| \mathbf{d}$, we hav a non-zero $\mathbf{e} \in \text{null}(\mathbf{A}) \cap \mathcal{T} \cap (\bigcup_{i,j} \widetilde{C}_{i,j})$, which violates \mathcal{E} .

Then we prove that $\hat{\mathbf{x}} \neq \mathbf{x}^{\natural}$ implies \mathcal{E} . Assume that there exists non-zero $\mathbf{e} \in \text{null}(\mathbf{A}) \cap \mathcal{T} \cap (\bigcup_{i,j} \vec{C}_{i,j})$. We can show that signal $\mathbf{x}^{\natural} + t\mathbf{e}$, where *t* is some positive constant such that $\|\mathbf{x}^{\natural} + t\mathbf{e}\|_{1} \leq \|\mathbf{x}^{\natural}\|_{1}$, satisfying constraints described by Eq. (2.1). This implies that $\mathbf{d} = t\mathbf{e} \neq \mathbf{0}$ and the wrong recovery of \mathbf{x}^{\natural} .

Since a direct computation of the probability of event \mathcal{E} can be difficult, we analyze the following equivalent event,

$$\min_{\mathbf{x}\in\mathcal{T}\cap\left(\bigcup_{i,j}\widetilde{\mathcal{C}}_{ij}\right)}\|\mathbf{A}\mathbf{x}\|_{2}>0.$$
(3.4)

For the simplicity of analysis, we assume that the entries $A_{i,j}$ of **A** are i.i.d. normal $\mathcal{N}(0,1)$. Using Gordon's escape from mesh theorem Gordon (1988), we obtain the following result that relates $Pr(\mathcal{E})$ with the number of measurements M.

Theorem 1. Let $a_M = \mathbb{E} \|\mathbf{g}\|_2$, where $\mathbf{g} \in \mathcal{N}(\mathbf{0}, \mathbf{I}_{M \times M})$, and $\omega(\cdot)$ denotes the Gaussian width. Provided that $a_M \ge \omega(\mathcal{T})$ and $(1 - 2\epsilon)a_M \ge \omega(\widetilde{C}_{ij})$ for $1 \le i, j \le L$ and $\epsilon > 0$, we have

$$\Pr(\mathcal{E}) \geq 1 - \left(\underbrace{\Pr\left(\min_{\mathbf{u}\in\mathcal{T}^{c}\setminus\{\mathbf{0}\}}\|\mathbf{A}\mathbf{u}\|_{2}>0\right)}_{\mathcal{P}_{1}} \wedge \underbrace{\Pr\left(\min_{\mathbf{u}\in\cap\tilde{\mathcal{C}}_{i}^{c}\setminus\{\mathbf{0}\}}\|\mathbf{A}\mathbf{u}\|_{2}>0\right)}_{\mathcal{P}_{2}}\right), \quad (3.5)$$

where $a \wedge b$ denotes the minimum of a and b, and \mathcal{P}_1 and \mathcal{P}_2 can be bounded as

$$\mathcal{P}_{1} \leq 1 \wedge \exp\left(-\frac{(a_{M} - \omega(\mathcal{T}))^{2}}{2}\right)$$

$$\mathcal{P}_{2} \leq 1 \wedge \frac{3}{2} \exp\left(-\frac{\epsilon^{2} a_{M}^{2}}{2}\right) + \sum_{i \leq j} \exp\left(-\frac{\left((1 - 2\epsilon)a_{M} - \omega(\widetilde{\mathcal{C}}_{ij})\right)^{2}}{2}\right).$$
(3.6)

Thm. 1 links the probability of correct recovery of Eq. (2.2) with the number of measurements M, and the "size" of constraint set. Detailed explanation is given as the following. To ensure high-probability of \mathcal{E} , we would like to $\mathcal{P}_1 \wedge \mathcal{P}_2$ to approach zero, which requires large value of a_M . Meanwhile, a_M is a monotonically increasing function of the sensor number M. Hence, we can obtain the minimum sensor number M requirement by unique recovery via investigating a_M .

Remark 1. Notice that \mathcal{P}_1 is associated with the descent cone \mathcal{T} of the optimization function, namely, $\|\mathbf{x}\|_1$, while \mathcal{P}_2 is associated with the prior knowledge $\mathbf{x} \in \bigcup_i C_i$. Thm. 1 implies that event \mathcal{E} (uniqueness) holds with higher probability than the traditional CS due to the extra constraint $\mathbf{x} \in \bigcup_i C_i$. If we fix $\Pr(\mathcal{E})$, we can separately calculate the corresponding M with and without the constraint $\mathbf{x} \in \bigcup_i C_i$. The difference ΔM would indicate the savings in the number of measurements due to the additional structure $\mathbf{x} \in \bigcup_i C_i$ over the traditional CS.

One simple example is attached below to illustrate the improvement brought by Thm. 1. **Example 2.** *Consider the constraint set*

$$C_i = \{ (0, \cdots, 0, x_i, \cdots, x_{K+i}, 0, \cdots, 0) \},$$
(3.7)

where $1 \le i \le N - K$. We study the asymptotic behavior of Thm. 1 when N is of order $\mathcal{O}(K^c)$, where c > 1 is constant. In the sequel we will show that Thm. 1 gives us the order $M = \mathcal{O}(K)$ to ensure solution uniqueness as K approaches infinity, which gives us the same bound as shown in Baraniuk et al. (2010) and suggests the tightness of our result.

Setting $\epsilon = 1/4$, we can bound \mathcal{P}_2 as

$$\mathcal{P}_2 \leq \frac{3}{2} \exp\left(-\frac{a_M^2}{32}\right) + \frac{N^2}{2} \exp\left(-\frac{(a_M - 2a_{2K})^2}{8}\right),$$
 (3.8)

provided $a_M \ge 2a_K$. With the relation $\frac{M}{\sqrt{M+1}} \le a_M \le \sqrt{M}$ Chandrasekaran et al. (2012); Gordon (1988) and setting M = 3K, we have

$$\mathcal{P}_2 \le c_1 \exp(-c_2 K) + c_3 N^2 \exp(-c_4 K), \tag{3.9}$$

where $c_1, c_2, c_3, c_4 > 0$ are some positive constants. Since $N = O(K^c)$, we can see \mathcal{P}_2 shrinks to zero as K approaches infinity, which implies the solution uniqueness.

Comparing with the traditional CS theory without prior knowledge $x \in \bigcup_i C_i$, our bound reduces the number of measurements from $M = O(K \log N/K) = O(K \log K)$ to M = O(K).

4 Computational Algorithm

Apart from the statistical property, another important aspect of Eq. (2.2) is to design an efficient algorithm. One naive idea is to consider and solve *L* separate optimization problems

$$\widehat{\mathbf{x}}^{(i)} = \operatorname{argmin}_{\mathbf{x}} \|\mathbf{x}\|_{1}, \text{ s.t. } \mathbf{y} = \mathbf{A}\mathbf{x}, \ \mathbf{x} \in \mathcal{C}_{i}, \tag{4.1}$$

and then selecting the best one, i.e., the sparsest reconstructed signal among all $\hat{\mathbf{x}}^{(i)}$'s. However, this method has two drawbacks:

- It requires solving *L* separate optimization problems, which in many applications might be prohibitively large and difficult to handle, but the proposed method is based on one single optimization procedure.
- It is inflexible. For example, some prior knowledge of which C_i the true signal x[↓] is more likely to reside might be available. The above method cannot incorporate such priors.

To overcome the above drawbacks, we (*i*) reformulate Eq. (2.2) to a more tractable objective function, and (*ii*) propose a computationally efficient algorithm to solve it. In the following, we assume that **x** is bounded in the sense that for a constant R, $\|\mathbf{x}\|_2 \leq R$.

4.1 **Reformulation of the objective function**

We introduce an auxiliary variable \mathbf{p} and rewrite the Lagrangian form in Eq. (2.2) as

$$\min_{\mathbf{x}} \min_{\mathbf{p} \in \Delta_L} \sum_{i} p_i \left(\|\mathbf{x}\|_1 + \widetilde{\mathbb{1}}(\mathbf{x} \in \mathcal{C}_i) + \frac{\lambda_1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{\lambda_2 \|\mathbf{x}\|_2^2}{2} \right),$$
(4.2)

where Δ_L is the simplex $\{p_i \ge 0, \sum_i p_i = 1\}$, $\widetilde{\mathbb{1}}(\cdot)$ is the truncated indicator function, which is 0 when its argument is true and is some large finite number *C* otherwise, and $\lambda_1, \lambda_2 > 0$ are the Lagrange multipliers. The term $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2$ is used to penalize for the constraint $\mathbf{y} = \mathbf{A}\mathbf{x}$ while $\|\mathbf{x}\|_2^2$ corresponds to the energy constraint $\|\mathbf{x}\|_2 \le R$. It can be easily shown that solving Eq. (4.2) for large enough *C* ensures $\mathbf{x} \in \bigcup_i^L C_i$.

Algorithm 1 Non-convex Proximal Multiplicative Weighting Algorithm

- Initialization: Initialize all variables with uniform weight $p_i^{(0)} = L^{-1}$ and $\mathbf{x}^{(0)} = \mathbf{0}$.
- For time t = 1 to *T*: We update $p_i^{(t+1)}$ and $\mathbf{x}^{(t)}$ as

$$p_{i}^{(t+1)} \propto p_{i}^{(t)} e^{-\eta_{p}^{(t)} f_{i}(\mathbf{x}^{(t)})}$$

$$\mathbf{x}^{(t+1)} = \operatorname{prox}_{\eta_{w}^{(t)} \parallel \cdot \parallel_{1}} \left[\mathbf{x}^{(t)} - \eta_{x}^{(t)} \sum_{i} p_{i}^{(t)} \left(\nabla_{x} h_{i}(\mathbf{x}^{(t)}) + \lambda_{1} \mathbf{A}^{\top} (\mathbf{A} \mathbf{x}^{(t)} - \mathbf{y}) + \lambda_{2} \mathbf{x}^{(t)} \right) \right],$$
(4.3)

where $p_i^{(t)}$ denotes the *i*th element of $\mathbf{p}^{(t)}$, and the proximal operator $\operatorname{prox}_{\|\cdot\|_1}(\mathbf{x})$ is defined as $\operatorname{argmin}_{\mathbf{z}}\left(\|\mathbf{z}\|_1 + \frac{1}{2}\|\mathbf{z} - \mathbf{x}\|_2^2\right)$ Beck & Teboulle (2009).

• **Output**: Calculate the average value $\bar{\mathbf{p}} = \frac{\sum_t \mathbf{p}^{(t)}}{T}$ and value $\bar{\mathbf{x}} = \frac{\sum_t \mathbf{x}^{(t)}}{T}$. Then output $\hat{\mathbf{x}}$ by projecting $\bar{\mathbf{x}}$ onto the set of $\bigcup_i C_i$.

Apart from the universality, our formulation has the following benefits:

- It is memory efficient. Compared with the naive idea that needs to store *L* different $\hat{\mathbf{x}}^{(i)}$, our method only needs to track one $\hat{\mathbf{x}}$ and one redundant variable **p**. This reduces the storing memory from $\mathcal{O}(NL)$ to $\mathcal{O}(N + L)$.
- It is very flexible. We can easily adjust to the case that x belongs to the intersection, i.e., x ∈ ∩_i C_i via modifying min_{p∈Δi} in Eq. (4.2) to max_{p∈Δi}.

Besides, to the best of our knowledge, this is the first time that such a formulation Eq. (4.2) is proposed. In the following, we will focus on the computational methods. Note that the difficulties in solving Eq. (4.2) are due to two aspects:

- Optimization over **p**: Although classical methods to minimize over **p** with fixed **x**, e.g., *alternative minimization* and ADMM Boyd (n.d.), can calculate local minimum efficiently (due to the biconvexity of Eq. (4.2), they can be easily trapped in the local-minima. This is because some entries in **p** can be set to zero and hence **x** will be kept away from the corresponding set C_i thereafter. To handle this problem, we propose to use *multiplicative weight update* Arora *et al.* (2012) and update **p** with the relation $\mathbf{p}^{(t+1)} \propto \mathbf{p}^{(t)}e^{-\eta_p^{(t)}f_i(\mathbf{x})}$, where $\mathbf{p}^{(t)}$ denotes **p**'s value in the *t*th iteration. This update relation avoids the sudden change of $\mathbf{p}^{(t)}$'s entries from non-zero to zero, which could have forced $\mathbf{x}^{(t)}$ being trapped in a local minimum.
- Optimization over **x**: Due to the non-smoothness of $\tilde{\mathbb{1}}(\mathbf{x} \in C_i)$ and $\|\mathbf{x}\|_1$ in Eq. (4.2) and the difficulties in calculating their sub-gradients, directly minimizing Eq. (4.2) would be computationally prohibitive. We propose to first approximate $\tilde{\mathbb{1}}(\mathbf{x} \in C_i)$ with a smooth function $h_i(\mathbf{x})$ and update $\mathbf{x}^{(t)}$ with the relation Eq. (4.4) used in *proximal gradient descent* Beck (2017).

Definition 3 (L_g -strongly smooth Beck (2017)). Function $g(\cdot) : \mathcal{X} \mapsto \mathbb{R}$ is L_g -strongly smooth iff

$$g(\mathbf{y}) \le g(\mathbf{x}) + \langle \nabla g(\mathbf{x}), \, \mathbf{y} - \mathbf{x} \rangle + \frac{L_g}{2} \|\mathbf{x} - \mathbf{y}\|_2^2, \tag{4.4}$$

for all \mathbf{x} , \mathbf{y} in the domain \mathcal{X} .

4.2 Non-convex Proximal Multiplicative Weighting Algorithm

Here we directly approximate the truncated indicator function $\widetilde{\mathbb{1}}(\mathbf{x} \in C_i)$ by $L_{h,i}$ strongly-smooth convex penalty functions $h_i(\mathbf{x})$, which may be different for different shapes of convex sets. For example, consider the convex set C_i in Example. 2. We may define $h_i(\mathbf{x}) = \sum_{j \notin [i,i+K]}^N x_j^2$, where [a, b] denotes the region from a to b. While for the set $\{\mathbf{x} : \langle \mathbf{a}, \mathbf{x} \rangle \leq b\}$, we may instead adopt the modified log-barrier function with a finite value. Then Eq. (4.2) can be rewritten as

$$\min_{\mathbf{p}} \min_{\mathbf{x}} \mathcal{L}(\mathbf{p}, \mathbf{x}) \triangleq \sum_{i=1}^{L} p_i f_i(\mathbf{x}),$$
(4.5)

where $f_i(\mathbf{x})$ is defined as

$$f_i(\mathbf{x}) \triangleq \|\mathbf{x}\|_1 + h_i(\mathbf{x}) + \frac{\lambda_1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{\lambda_2}{2} \|\mathbf{x}\|_2^2.$$
(4.6)

Hence, the optimization problem in (4.5) can be solved via Alg. 1

Lemma 2. $h(\mathbf{x}) \triangleq \sum_{i} p_{i}h_{i}(\mathbf{x}) + \frac{\lambda_{1} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_{2}^{2}}{2} + \frac{\lambda_{2} \|\mathbf{x}\|_{2}^{2}}{2}$ is strongly-smooth with some positive constant denoted as L_{h} .

Proof. First, we can check that $\frac{\lambda_1 \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2}{2} + \frac{\lambda_2 \|\mathbf{x}\|_2^2}{2}$ is strongly-smooth. Denote the corresponding parameter as $L_{h,0}$. Meanwhile, due to the construction of $h_i(\mathbf{x})$, it is strongly-smooth for every *i*. Since p_i is non-negative for every *i*, we can easily prove the following inequality

$$h(\mathbf{x}_{1}) \ge h(\mathbf{x}_{2}) + \langle \nabla h(\mathbf{x}_{2}), \mathbf{x}_{1} - \mathbf{x}_{2} \rangle + \frac{L_{h}}{2} \|\mathbf{x}_{1} - \mathbf{x}_{2}\|_{2}^{2},$$
(4.7)

*

where L_h is defined as min $(L_{h,i})$, $0 \le i \le L$.

Then we have the following theorem.

Theorem 3. Let $\eta_x^{(t)} = \eta_x \leq L_h^{-1}$, and $\eta_p^{(t)} = R_f^{-1} \sqrt{2 \log L/T}$, where $|f_i(\cdot)| \leq R_f$, $\|\mathbf{x}\|_2 \leq R$. Then we have

$$\left| \frac{\min_{\mathbf{p}} \sum_{t} \mathcal{L}(\mathbf{p}, \mathbf{x}^{(t)})}{T} - \frac{\sum_{t} \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T} \right| + \left| \frac{\sum_{t} \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T} - \frac{\min_{\mathbf{x}} \sum_{t} \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x})}{T} \right|$$

$$\leq \frac{2R^{2}}{\eta_{x}T} + R_{f} \sqrt{\frac{\log L}{T}},$$
(4.8)

where T denotes the number of iterations.

Due to the difficulties in analyzing the global optimum, in Theorem 3 we focus on analyzing the closeness between the average value $\frac{\sum_t \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T}$ to its local minimum. The first term denotes the gap between average value $\frac{\sum_t \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T}$ and the optimal value of $\mathcal{L}(\mathbf{p}, \mathbf{x})$ with $\mathbf{x}^{(t)}$ being fixed. Similarly, the second term represents the gap with $\mathbf{p}^{(t)}$ being fixed. As $T \to \infty$, the sum of these two bounds approaches to zero at the rate of $\mathcal{O}(T^{-1/2})$.

Moreover note that setting $\eta_p^{(t)}$ requires the oracle knowledge of *T*, which is impractical. This artifacts can easily be fixed by the doubling trick (Shalev-Shwartz *et al.*, 2012, §2.3.1). In addition, we have proved the following theorem.

Theorem 4. Let $\eta_w^{(t)} \leq L_h^{-1}$, where $|f_i(\cdot)| \leq R_f$. Then we have

$$\frac{1}{T}\sum_{t} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2\mathcal{L}(\mathbf{p}^{(0)}, \mathbf{x}^{(0)})}{L_{h}T} + \frac{4R_{f}^{2}\sum_{t}\eta_{p}^{(t)}}{L_{h}T}.$$
(4.9)

This theorem discusses the convergence speed with respect to the $\mathbf{x}^{(t)}$ update. Due to the $\mathcal{O}(T^{-1})$ of the first term on the right side of the above inequality, the best convergence rate we can obtain is $\mathcal{O}(T^{-1})$, which is achievable by $\eta_p^{(t)} \propto t^{-2}$. However, using fixed learning rate η_p as in Thm. 3 would result in the convergence rate of $\mathcal{O}(T^{-1/2})$.

4.3 Regularization for p

Another drawback of the naive method is that they cannot exploit the prior knowledge. For example, if we know that the true \mathbf{x}^{\natural} is most likely to reside in set C_1 . With the naive method, we cannot use this information but separately solve Eq. (2.2) for all *L* sets. In the sequel, we will show that our formulation Eq. (4.2) can incorporate such prior knowledge by adding regularizers for \mathbf{p} , and bring certain performance improvement.

Note that we can interpret p_i , the *i*-th element of **p** in Eq. (4.5) as the likelihood of $\mathbf{x}^{\natural} \in C_i$. Without any prior knowledge about which set C_i the true signal \mathbf{x}^{\natural} resides, variable **p** is uniformly distributed among all possible distributions Δ_L . When certain prior information is available, its distribution is skewed towards certain distributions, namely **q**.

In this paper, we adopt $\|\cdot\|_2^2$ to regularize **p** towards **q** and write the modified function $\mathcal{LR}(\mathbf{p}, \mathbf{x})$ as

$$\mathcal{LR}(\mathbf{p}, \mathbf{x}) = \mathcal{L}(\mathbf{p}, \mathbf{x}) + \frac{\lambda_3}{2} \|\mathbf{p} - \mathbf{q}\|_{2}^{2},$$
(4.10)

where $\lambda_3 > 0$ is a constant used to balance $\mathcal{L}(\mathbf{p}, \mathbf{x})$ and $\frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_2^2$. Based on different applications, other norms such as KL-divergence or l_1 norm can be used as the regularizer.

Then we substitute the update equation Eq. (4.3) as

$$\mathbf{p}^{(t+1)} = \mathbb{P}_{\Delta} \left(\mathbf{p}^{(t)} - \eta_p^{(t)} \mathbf{g}^{(t)} \right), \qquad (4.11)$$

where $\mathbf{g}^{(t)} = \nabla_{\mathbf{p}^{(t)}} \mathcal{LR}(\mathbf{p}, \mathbf{x}^{(t)}) = \mathbf{f}(\mathbf{x}^{(t)}) + \lambda_3(\mathbf{p}^{(t)} - \mathbf{q})$, and $\mathbf{f}(\mathbf{x}^{(t)})$ denotes the vector whose *i*th element is $f_i(\mathbf{x}^{(t)})$. Similar as above, we obtain the following theorems.

Theorem 5. Provided that $\|\mathbf{g}^{(t)}\|_2 \leq R_g$, by setting $\eta_x^{(t)} = \eta_x \leq L_h^{-1}$ and $\eta_p^{(t)} = (\lambda t)^{-1}$ we conclude that

$$\left| \frac{\min_{\mathbf{p}} \sum_{t} \mathcal{LR}(\mathbf{p}, \mathbf{x}^{(t)})}{T} - \frac{\sum_{t} \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T} \right| + \left| \frac{\sum_{t} \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}{T} - \frac{\min_{\mathbf{x}} \sum_{t} \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x})}{T} \right|$$

$$\leq \frac{R_{g}^{2} \log T}{2\lambda_{3}T} + \frac{R^{2}}{2\eta_{x}T},$$
(4.12)

Comparing with Thm. 3, Thm. 5 implies that the regularizers improve the optimal rate from $O(T^{-1/2})$ to $O(\log T/T)$. Therefore, our framework can exploit the prior information to improve the recovery performance whereas the naive method of iterative computation fails to achieve as such.

Theorem 6. Provided that $\|\mathbf{g}^{(t)}\|_2 \leq R_g$, by setting $\eta_w^{(t)} = \eta_x \leq L_h^{-1}$ we conclude that

$$\frac{1}{T}\sum_{t} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2\mathcal{LR}(\mathbf{p}^{(0)}, \mathbf{x}^{(0)})}{L_{h}T} + \frac{2R_{g}^{2}}{L_{h}T}\sum_{t} \left(\eta_{p}^{(t)} + \frac{\lambda_{3}\left(\eta_{p}^{(t)}\right)^{2}}{2}\right).$$
(4.13)

In this case, if we set $\eta_p^{(t)}$ as t^{-2} , then $\frac{1}{T}\sum_t \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$ would decrease at the rate of $\mathcal{O}(T^{-1})$, which is the same as Thm. 4.

5 Conclusion

In this paper, we studied the compressive sensing with a multiple convex-set domain. First we analyzed the impact of prior knowledge $\mathbf{x} \in \bigcup_i C_i$ on the minimum number of measurements M to guarantee uniqueness of the solution. We gave an illustrative example and showed that significant savings in M can be achieved. Then we formulated a universal objective function and develop an algorithm for the signal reconstruction. We show that in terms of the speed of convergence to local minimum, our proposed algorithm based on *multiplicative weight update* and *proximal gradient descent* can achieve the optimal rate of $\mathcal{O}(T^{-1/2})$. Further, in terms of $T^{-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$, the optimal speed increases to $\mathcal{O}(T^{-1})$. Moreover, provided that we have a prior knowledge about \mathbf{p} , we show that we can improve the optimal recovery performance by $\|\cdot\|_2^2$ regularizers, and hence increasing the above convergence rate from $\mathcal{O}(T^{-1/2})$ to $\mathcal{O}\left(\frac{\log T}{T}\right)$.

References

- ARORA, SANJEEV, HAZAN, ELAD, & KALE, SATYEN. 2012. The Multiplicative Weights Update Method: a Meta-Algorithm and Applications. *Theory of Computing*, **8**(1), 121–164.
- BARANIUK, RICHARD G, CEVHER, VOLKAN, DUARTE, MARCO F, & HEGDE, CHINMAY. 2010. Model-based compressive sensing. *IEEE Transactions on Information Theory*, **56**(4), 1982–2001.
- BECK, A. 2017. First-Order Methods in Optimization. MOS-SIAM Series on Optimization. Society for Industrial and Applied Mathematics.
- BECK, AMIR, & TEBOULLE, MARC. 2009. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM journal on imaging sciences*, **2**(1), 183–202.
- BOUCHERON, S., LUGOSI, G., & MASSART, P. 2013. Concentration Inequalities: A Nonasymptotic Theory of Independence. Oxford University Press.
- BOYD, STEPHEN. Alternating direction method of multipliers.
- CANDES, EMMANUEL J, ELDAR, YONINA C, STROHMER, THOMAS, & VORONINSKI, VLADISLAV. 2011. Phase Retrieval via Matrix Completion.
- CHANDRASEKARAN, VENKAT, RECHT, BENJAMIN, PARRILO, PABLO A, & WILLSKY, ALAN S. 2012. The convex geometry of linear inverse problems. *Foundations of Computational mathematics*, **12**(6), 805–849.
- CHEN, YUXIN, & CANDES, EMMANUEL. 2015. Solving random quadratic systems of equations is nearly as easy as solving linear systems. *Pages 739–747 of: Advances in Neural Information Processing Systems*.
- DAI, WEI, PHAM, HOA VINH, & MILENKOVIC, OLGICA. 2009. Quantized compressive sensing. *arXiv* preprint arXiv:0901.0749.
- DUARTE, MARCO F, & ELDAR, YONINA C. 2011. Structured compressed sensing: From theory to applications. *IEEE Transactions on Signal Processing*, **59**(9), 4053–4085.
- ELDAR, Y.C., & KUTYNIOK, G. 2012. *Compressed Sensing: Theory and Applications*. Compressed Sensing: Theory and Applications. Cambridge University Press.
- FOUCART, SIMON, & RAUHUT, HOLGER. A mathematical introduction to compressive sensing. Vol. 1.
- GORDON, YEHORAM. 1988. On Milman's inequality and random subspaces which escape through a mesh in \mathbb{R}^n . Pages 84–106 of: Geometric Aspects of Functional Analysis. Springer.
- QIAN, QI, ZHU, SHENGHUO, TANG, JIASHENG, JIN, RONG, SUN, BAIGUI, & LI, HAO. 2018. Robust Optimization over Multiple Domains. *arXiv preprint arXiv:1805.07588*.
- SHALEV-SHWARTZ, SHAI, et al. 2012. Online learning and online convex optimization. Foundations and Trends(R) in Machine Learning, 4(2), 107–194.
- SILVA, JORGE, CHEN, MINHUA, ELDAR, YONINA C, SAPIRO, GUILLERMO, & CARIN, LAWRENCE. 2011. Blind compressed sensing over a structured union of subspaces. *arXiv preprint arXiv:1103.2469*.
- ZHANG, H., ABDI, A., & FEKRI, F. NOV. 2017.. Compressive Sensing with Energy Constraint. In: IEEE Information Theory Workshop (ITW'17).

Appendix A Proof of Theorem 1

Proof. Note that for any-vector non-zero $h \in \text{null}(\mathbf{A}) \cap \mathcal{T} \cap \left(\bigcup_{i,j} \widetilde{\mathcal{C}}_{i,j}\right)$, we can always rescale to make it unit-norm. Hence we can rewrite the event \mathcal{E} as

$$\mathcal{E} = \left\{ \operatorname{null}(\mathbf{A}) \bigcap \mathbb{S}_{2}^{n-1} \bigcap \mathcal{T} \bigcap \left(\bigcup_{i,j} \widetilde{\mathcal{C}}_{i,j} \right) = \emptyset \right\}.$$
(A.1)

For the conciseness of notation, we define \widetilde{C} to be $\widetilde{C} = \bigcup_{i,j} \widetilde{C}_{i,j}$. Then we upper-bound $1 - \Pr(\mathcal{E})$ as

$$1 - \Pr(\mathcal{E}) = \Pr\left(\operatorname{null}(\mathbf{A}) \bigcap S_2^{n-1} \bigcap \mathcal{T} \bigcap \widetilde{\mathcal{C}} \neq \emptyset\right)$$

$$\stackrel{(i)}{\leq} \underbrace{\Pr\left(\operatorname{null}(\mathbf{A}) \bigcap S_2^{n-1} \bigcap \mathcal{T} \neq \emptyset\right)}_{\mathcal{P}_1} \land \underbrace{\Pr\left(\operatorname{null}(\mathbf{A}) \bigcap S_2^{n-1} \bigcap \widetilde{\mathcal{C}} \neq \emptyset\right)}_{\mathcal{P}_2}, \quad (A.2)$$

where (i) is because

$$\left\{ \operatorname{null}(\mathbf{A}) \bigcap \mathbb{S}_{2}^{n-1} \bigcap \mathcal{T} \bigcap \widetilde{\mathcal{C}} \neq \emptyset \right\} \subseteq \left\{ \operatorname{null}(\mathbf{A}) \bigcap \mathbb{S}_{2}^{n-1} \bigcap \mathcal{T} \neq \emptyset \right\},$$

$$\left\{ \operatorname{null}(\mathbf{A}) \bigcap \mathbb{S}_{2}^{n-1} \bigcap \mathcal{T} \bigcap \widetilde{\mathcal{C}} \neq \emptyset \right\} \subseteq \left\{ \operatorname{null}(\mathbf{A}) \bigcap \mathbb{S}_{2}^{n-1} \bigcap \widetilde{\mathcal{C}} \neq \emptyset \right\}.$$

$$(A.3)$$

*

With Lemma 3 and Lemma 4, we can separately bound \mathcal{P}_1 and \mathcal{P}_2 and finish the proof.

Lemma 3. We have
$$\mathcal{P}_1 \leq 1 \wedge \exp\left(-\frac{(a_m - \omega(\mathcal{T}))^2}{2}\right)$$
, if $a_m \geq \omega(\mathcal{T})$

Proof. Note that we have

$$\underbrace{\Pr\left(\operatorname{null}(\mathbf{A})\bigcap \mathbb{S}_{2}^{n-1}\bigcap \mathcal{T}\neq \emptyset\right)}_{\mathcal{P}_{1}} + \underbrace{\Pr\left(\operatorname{null}(\mathbf{A})\bigcap \mathbb{S}_{2}^{n-1}\bigcap \mathcal{T}=\emptyset\right)}_{\mathcal{P}_{1}^{c}} = 1.$$
(A.4)

Then we lower-bound \mathcal{P}_1^c as

$$\mathcal{P}_{1}^{c} = \Pr\left(\min_{\mathbf{u}\in S_{2}^{n-1}\cap\mathcal{T}}\|\mathbf{A}\mathbf{u}\|_{2} > 0\right) \stackrel{(i)}{\geq} 1 - \exp\left(-\frac{\left(a_{m}-\omega\left(\mathcal{T}\right)\right)^{2}}{2}\right),\tag{A.5}$$

provided $a_m \ge \omega(\mathcal{T})$, where (*i*) is because of Corollary 3.3 in Chandrasekaran *et al.* (2012), and $\omega(\cdot)$ denotes the Gaussian width.

Lemma 4. If $\omega(\tilde{\mathcal{C}}_{ij}) \leq 1 - 2\epsilon a_M$, we have

$$\mathcal{P}_{2} \leq 1 \wedge \frac{3}{2} \exp\left(-\frac{\epsilon^{2} a_{M}^{2}}{2}\right) + \sum_{i \leq j} \exp\left(-\frac{\left((1-2\epsilon)a_{M}-\omega(\widetilde{\mathcal{C}}_{ij})\right)^{2}}{2}\right), \tag{A.6}$$

Proof. Note that we have

$$\underbrace{\Pr\left(\operatorname{null}(\mathbf{A})\bigcap S_2^{n-1}\bigcap \widetilde{\mathcal{C}}\neq \emptyset\right)}_{\mathcal{P}_2} + \underbrace{\Pr\left(\operatorname{null}(\mathbf{A})\bigcap S_2^{n-1}\bigcap \widetilde{\mathcal{C}}=\emptyset\right)}_{\mathcal{P}_2^c} = 1.$$
(A.7)

Here we upper-bound \mathcal{P}_2 via lower-bounding \mathcal{P}_2^c . First we define $\mathcal{P}_2^c(d)$ as

$$\mathcal{P}_{2}^{c}(d) \triangleq \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}, \, \mathbf{v} \in \operatorname{null}(\mathbf{A})} \|\mathbf{u} - \mathbf{v}\|_{2} \ge d\right).$$
(A.8)

Then we have $\mathcal{P}_2^c = \lim_{d\to 0} \mathcal{P}_2^c(d)$. The following proof trick is fundamentally the same as that are used in Theorem 4.1 in Gordon (1988) but in a clear format by only keeping the necessary parts for this scenario. We only present it for the self-containing of this paper and do not claim any novelties.

We first define $S_{i,j} = \mathbb{S}_2^{n-1} \cap \widetilde{C}_{i,j}$ and two quantities Q_1 and Q_2 as

$$Q_{1} \triangleq \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}} \|\mathbf{A}\mathbf{u}\|_{2} \ge d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\right),$$

$$Q_{2} \triangleq \Pr\left(\bigcap_{i_{1} \le j_{1}} \bigcap_{\mathbf{u} \in S_{i_{1},j_{1}}} \left(\left(\sum_{i_{2}=1}^{M} g_{i_{2}}^{2}\right)^{1/2} + \sum_{j_{2}} u_{j_{2}}h_{j_{2}} \ge d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} + \epsilon_{2}a_{M}\right)\right),$$
(A.9)

where $a_M = \mathbb{E} ||\mathbf{g}||_2$, $\mathbf{g} \in \mathcal{N}(\mathbf{0}, \mathbf{I}_{M \times M})$, and g_j , h_i are iid standard normal random variables $\mathcal{N}(0, 1)$. The following proof is divided into 3 parts.

Step I. We prove that $\mathcal{P}_2^c(d) + e^{-\epsilon_1^2 a_M^2/2} \ge Q_1$, which is done by

$$Q_{1} = \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}} \|\mathbf{A}\mathbf{u}\|_{2} \ge d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\right)$$

$$= \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}, \mathbf{v} \in \text{null}(\mathbf{A})} \|\mathbf{A}(\mathbf{u}-\mathbf{v})\|_{2} \ge d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\right)$$

$$\stackrel{(i)}{\le} \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}, \mathbf{v} \in \text{null}(\mathbf{A})} \|\mathbf{A}\|_{2}\|\mathbf{u}-\mathbf{v}\|_{2} \ge d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\right)$$

$$\stackrel{(ii)}{\le} \Pr\left(\|\mathbf{A}\|_{2} \ge (1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\right) + \Pr\left(\min_{\mathbf{u} \in \bigcup S_{i,j}, \mathbf{v} \in \text{null}(\mathbf{A})} \|\mathbf{u}-\mathbf{v}\|_{2} \ge d\right)$$

$$\stackrel{(iii)}{\le} \exp\left(-\frac{\epsilon_{1}^{2}(\mathbb{E}\|\mathbf{A}\|_{2})^{2}}{2}\right) + \mathcal{P}_{2}^{c}(d)$$

$$\stackrel{(iv)}{\le} \exp\left(-\frac{\epsilon_{1}^{2}a_{M}^{2}}{2}\right) + \mathcal{P}_{2}^{c}(d),$$
(A.10)

where in (*i*) we use $\|\mathbf{A}\|_2 \|\mathbf{u} - \mathbf{v}\|_2 \ge \|\mathbf{A}(\mathbf{u} - \mathbf{v})\|_2$, in (*ii*) we use the union bound for

$$\begin{cases} \min_{\mathbf{u}\in\cup\mathcal{S}_{i,j},\,\mathbf{v}\in\mathrm{null}(\mathbf{A})} \|\mathbf{A}\|_{2}\|\mathbf{u}-\mathbf{v}\|_{2} \geq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} \\ \\ & \leq \{\|\mathbf{A}\|_{2} \geq (1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2}\} \bigcup \left\{ \min_{\mathbf{u}\in\cup\mathcal{S}_{i,j},\,\,\mathbf{v}\in\mathrm{null}(\mathbf{A})} \|\mathbf{u}-\mathbf{v}\|_{2} \geq d \right\}, \end{cases}$$
(A.11)

, in (*iii*) we use the Gaussian concentration inequality Lipschitz functions (Theorem 5.6 in Boucheron *et al.* (2013)) for $\|\mathbf{A}\|_2$, and in (v) we use $\|\mathbf{A}\|_2 \ge \|\mathbf{A}\mathbf{e}_1\|_2 = \|\sum_{i=1}^M A_{i,1}\|_2$, where \mathbf{e}_1 denotes the canonical basis.

Step II. We prove that $Q_1 + \frac{1}{2}e^{-\epsilon_2^2 a_M^2/2} \ge Q_2$, which is done by

$$Q_{1} + \frac{1}{2}e^{-\epsilon^{2}a_{M}^{2}/2} \stackrel{(i)}{\geq} Q_{1} + \Pr\{g \geq \epsilon_{2}a_{M}\}$$

$$\stackrel{(ii)}{\geq} \Pr\left(\min_{\mathbf{u}\in\cup\mathcal{S}_{i_{1},j_{1}}} \|\mathbf{A}\mathbf{u}\|_{2} + g\|\mathbf{u}\|_{2} \geq d(1+\epsilon)\mathbb{E}\|\mathbf{A}\|_{2} + \epsilon_{2}a_{M}\|\mathbf{u}\|_{2}\right)$$

$$= \Pr\left(\bigcap_{i_{1}\leq j_{1}}\bigcap_{\mathbf{u}\in\mathcal{S}_{i_{1},j_{1}}} \|\mathbf{A}\mathbf{u}\|_{2} + g\|\mathbf{u}\|_{2} \geq d(1+\epsilon)\mathbb{E}\|\mathbf{A}\|_{2} + \epsilon_{2}a_{M}\|\mathbf{u}\|_{2}\right)$$

$$\stackrel{(iii)}{\geq} \underbrace{\Pr\left(\bigcap_{i_{1}\leq j_{1}}\bigcap_{\mathbf{u}\in\mathcal{S}_{i_{1},j_{1}}} \left(\sum_{i_{2}=1}^{M}g_{i_{2}}^{2}\right)^{\frac{1}{2}} + \sum_{j_{2}=1}^{N}u_{j_{2}}h_{j_{2}} \geq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} + \epsilon_{2}a_{M}\right)}_{Q_{2}},$$

$$(A.12)$$

where in (*i*) **g** is a RV satisfying standard normal distribution, in (*ii*) we use the union bound, and (*iii*) comes from Lemma 3.1 in Gordon (1988) and $\|\mathbf{u}\|_2 = 1$.

Step III. We lower bound Q_2 as

$$1 - Q_{2} = \Pr\left(\bigcup_{i_{1} \leq j_{1}} \bigcup_{u \in S_{i_{1},j_{1}}} \left[\left(\sum_{i_{2}=1}^{M} g_{i_{2}}^{2}\right)^{\frac{1}{2}} + \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} + \epsilon_{2}a_{M} \right] \right)$$

$$\leq \Pr\left(\left(\sum_{i_{2}=1}^{M} g_{i_{2}}^{2}\right)^{\frac{1}{2}} \leq (1-\epsilon_{2})a_{M}\right) + \Pr\left(\bigcup_{i_{1} \leq j_{1}} \bigcup_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - (1-2\epsilon_{2})a_{M} \right)$$

$$\leq \Pr\left(\left(\sum_{i_{2}=1}^{M} g_{i_{2}}^{2}\right)^{\frac{1}{2}} - a_{M} \leq -\epsilon_{2}a_{M}\right) + \Pr\left(\bigcup_{i_{1} \leq j_{1}} \bigcup_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - (1-2\epsilon_{2})a_{M} \right)$$

$$\stackrel{(i)}{\leq} \exp\left(-\frac{\epsilon_{2}^{2}a_{M}^{2}}{2}\right) + \Pr\left(\bigcup_{i_{1} \leq j_{1}} \bigcup_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - (1-2\epsilon_{2})a_{M} \right)$$

$$\stackrel{(ii)}{\leq} \exp\left(-\frac{\epsilon_{2}^{2}a_{M}^{2}}{2}\right) + \sum_{i_{1} \leq j_{1}} \Pr\left(\bigcup_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - (1-2\epsilon_{2})a_{M} \right)$$

$$\stackrel{(iii)}{\leq} \exp\left(-\frac{\epsilon_{2}^{2}a_{M}^{2}}{2}\right) + \sum_{i_{1} \leq j_{1}} \Pr\left(\max_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - (1-2\epsilon_{2})a_{M} \right)$$

$$\stackrel{(iii)}{\leq} \exp\left(-\frac{\epsilon_{2}^{2}a_{M}^{2}}{2}\right) + \sum_{i_{1} \leq j_{1}} \Pr\left(\max_{u \in S_{i_{1},j_{1}}} \sum_{j_{2}=1}^{N} u_{j_{2}}h_{j_{2}} \leq (1-2\epsilon_{2})a_{M} - d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} \right)$$

$$\stackrel{(ii)}{\leq} \exp\left(-\frac{\epsilon_{2}^{2}a_{M}^{2}}{2}\right) + \sum_{i_{1} \leq j_{1}} \exp\left(-\frac{\left((1-2\epsilon_{2})a_{M} - d(1+\epsilon_{1})\mathbb{E}\|\mathbf{A}\|_{2} - \omega(\tilde{C}_{ij_{1}})\right)^{2}}{2}\right),$$

$$(A.13)$$

where in (*i*) we use $\mathbb{E}\sqrt{\sum_{i_2=1}^M g_{i_2}^2} = a_M$ and Gaussian concentration inequality in Boucheron *et al.* (2013), in (*ii*) we use union-bound, in (*iii*) we define h' = -h and flip the sign by the symmetry of Gaussian variables, and in (*iv*) we use the definition of $\omega(\tilde{C}_{ij})$. Assuming $(1 - 2\epsilon_2)a_M \ge d(1 + \epsilon_1)\mathbb{E}\|\mathbf{A}\|_2 + \omega(\tilde{C}_{ij})$, we finish the proof via the Gaussian concentration inequality Boucheron *et al.* (2013).

Combining the above together and set $d(1 + \epsilon_1) \rightarrow 0$ while $\epsilon_1 \rightarrow \infty$, we conclude that

$$\mathcal{P}_{2}^{c} \geq 1 - \frac{3}{2} \exp\left(-\frac{\epsilon^{2} a_{M}^{2}}{2}\right) - \sum_{i \leq j} \exp\left(-\frac{\left((1 - 2\epsilon)a_{M} - \omega(\widetilde{\mathcal{C}}_{ij})\right)^{2}}{2}\right), \quad (A.14)$$

#

provided $(1 - 2\epsilon)a_M \ge \omega(\widetilde{C}_{ij})$, and finish the proof.

Appendix B Proof of Theorem 3

Proof. Define \mathbf{p}^* and \mathbf{x}^* as

$$\mathbf{p}^{*} = \operatorname{argmin}_{\mathbf{p}} \sum_{t} \mathcal{L}\left(\mathbf{p}, \mathbf{x}^{(t)}\right), \quad \mathbf{x}^{*} = \operatorname{argmin}_{\mathbf{x}} \sum_{t} \mathcal{L}\left(\mathbf{p}^{(t)}, \mathbf{x}\right). \tag{B.1}$$

respectively First we define \mathcal{T}_1^t and \mathcal{T}_2^t as

$$\mathcal{T}_{1}^{t} = \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^{*}, \mathbf{x}^{(t)});$$

$$\mathcal{T}_{2}^{t} = \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{*}),$$
(B.2)

respectively. Then our goal becomes bounding $|\sum_t T_1^t| + |\sum_t T_2^t|$. With Lemma 5 and Lemma 6, we have finished the proof.

Lemma 5. Define $\mathcal{T}_1^t = \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^*, \mathbf{x}^{(t)})$, where \mathbf{p}^* is defined in Eq. (B.1), then we have

$$0 < \sum_{t} \mathcal{T}_{1}^{t} \le R_{f} \sqrt{T \log L},\tag{B.3}$$

when $\eta_p^{(t)} = R_f \sqrt{2 \log L/T}$.

Proof. Based on the definition of \mathbf{p}^* , we note that $\sum_t \mathcal{T}_1^t$ is non-negative and prove the lower-bound. Then we prove its upper-bound.

Since the function is linear, optimal \mathbf{p}^* must be at the edge of Δ_L and we denote the non-zero entry as i^* . Hence, we could study it via the *multiplicative weight algorithm* analysis Arora *et al.* (2012). First we rewrite the update equation (4.3). Define $\mathbf{w}^{(0)} = \mathbf{1} \in \mathbb{R}^L$ and update $\mathbf{w}^{(t+1)}$ as

$$w_i^{(t+1)} = w_i^{(t)} \exp\left(-\eta_p^{(t)} f_i(\mathbf{x}^{(t)})\right), \ p_i^{(t+1)} = \frac{w_i^{(t+1)}}{\sum_i w_i^{(t+1)}}.$$
 (B.4)

where $(\cdot)_i$ denotes the *i*th element, and $\mathbf{p}^{(t+1)}$ can be regarded as the normalized version of $\mathbf{w}^{(t+1)}$. First we define Ψ_t as

$$\Psi_t = \sum_{i=1}^{L} w_i^{(t)}.$$
(B.5)

Then we have $\Psi_0 = L$ while $\Psi_T \ge w_{i^*}^{(T)} = \exp\left(-\sum_{t=1}^T \eta_p^{(t)} f_{i^*}(\mathbf{x}^{(t)})\right) = \exp\left(-\eta_p \sum_{t=1}^T f_{i^*}(\mathbf{x}^{(t)})\right)$, where $\eta_p^{(t)} = \eta_p = \sqrt{2\log L/T}$.

Then we study the division Ψ_{t+1}/Ψ_t as

$$\begin{split} \Psi_{t+1} &= \sum_{i} w_{i}^{(t+1)} = \sum_{i} w_{i}^{t} \exp\left(-\eta_{p} f_{i}(\mathbf{x}^{(t)})\right) \\ \stackrel{(i)}{\leq} &\sum_{i} w_{i}^{t} \left(1 - \eta_{p} f_{i}(\mathbf{x}^{(t)}) + \frac{\eta_{p}^{2} f_{i}^{2}(\mathbf{x}^{(t)})}{2}\right) \\ &= &\Psi_{t} \left(1 - \eta_{p} \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle + \frac{\eta_{p}^{2} R_{f}^{2}}{2}\right) \end{split}$$
(B.6)
$$\stackrel{(ii)}{\leq} &\Psi_{t} \exp\left(-\eta_{p} \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle + \frac{\eta_{p}^{2} R_{f}^{2}}{2}\right) \end{split}$$

where in (*i*) we use $e^{-x} \le 1 + x + x^2/2$ for $x \ge 0$, and in (*ii*) we use $e^x \ge 1 + x$ for all $x \in \mathbb{R}$. Using the above relation iteratively, we conclude that

$$\frac{\Psi_T}{\Psi_0} \le \exp\left(-\eta_p \sum_t \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle + \frac{\eta_p^2 T R_f^2}{2}\right), \tag{B.7}$$

which gives us

$$\log \Psi_T \leq \log L - \eta_p \sum_t \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle + \frac{T \eta_p^2 R_f^2}{2}.$$
(B.8)

With relation $\log \Psi_T \geq \log w_{i^*}^{(T)}$, we obtain

$$\eta_p \sum_{t} \left(\left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle - \left\langle \mathbf{e}_{i^*}, f(\mathbf{x}^{(t)}) \right\rangle \right) = \eta_p \sum_{t} \left(\left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle - \left\langle \mathbf{p}^*, f(\mathbf{x}^{(t)}) \right\rangle \right) \le \log L + \frac{T \eta_p^2 R_f^2}{2}$$
(B.9)

where \mathbf{e}_{i^*} denotes the canonical basis, namely, has 1 in its i^* th entry and all others to be zero.

#

Lemma 6. Define $\mathcal{T}_2^t = \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^*)$, where \mathbf{x}^* is defined in Eq. (B.1), and set $\eta_x^{(t)} = \eta_x \leq L_h^{-1}$, then we have

$$0 \le \sum_{t} \mathcal{T}_{2}^{t} \le \frac{1}{2\eta_{x}} \| \mathbf{x}^{(0)} - \mathbf{x}^{*} \|_{2}^{2}.$$
(B.10)

Proof. From the definition of \mathbf{x}^* , we can prove the non-negativeness of $\sum_t \mathcal{T}_2^t$. Here we focus on upper-bounding $\sum_t \mathcal{T}_2^t$ by separately analyzing each term \mathcal{T}_2^t . For the conciseness of notation, we drop the time index *t*. Define $h(\mathbf{x})$ as

$$h(\mathbf{x}) = \sum_{i} p_{i} h_{i}(\mathbf{x}) + \frac{\lambda_{1} \|\mathbf{y} - A\mathbf{x}\|_{2}^{2}}{2} + \frac{\lambda_{2} \|\mathbf{x}\|_{2}^{2}}{2}.$$
 (B.11)

First, we rewrite the update equation as

$$\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \underbrace{\left(\mathbf{x}^{(t)} - \mathsf{prox}_{\eta_x \parallel \cdot \parallel_1} \left[\mathbf{x}^{(t)} - \eta_w \nabla h(\mathbf{x}^{(t)})\right]\right)}_{\eta_x H(\mathbf{x}^{(t)})},\tag{B.12}$$

which means that

$$H(\mathbf{x}^{(t)}) = \frac{\mathbf{x}^{(t)} - \operatorname{prox}_{\eta_x \parallel \cdot \parallel_1} \left(\mathbf{x}^{(t)} - \eta_w \nabla h(\mathbf{x}^{(t)}) \right)}{\eta_x},$$
(B.13)

where $\operatorname{prox}_{\eta_{\boldsymbol{x}}\|\cdot\|_1}(\mathbf{x})$ is defined as Beck (2017)

$$\operatorname{prox}_{\eta_{x}\|\cdot\|_{1}}(\mathbf{x}) = \operatorname{argmin}_{\mathbf{z}} \ \eta_{x}\|\mathbf{z}\|_{1} + \frac{1}{2}\|\mathbf{x} - \mathbf{z}\|_{2}^{2}. \tag{B.14}$$

Here we need one important property of $H(\mathbf{x}^{(t)})$, that is widely in the analysis of proximal gradient descent (direct results of Theorem 6.39 in Beck (2017)) and states

$$H(\mathbf{x}^{(t)}) \in \nabla h(\mathbf{x}^{(t)}) + \partial \|\mathbf{x}^{(t+1)}\|_1.$$
(B.15)

Here we consider the update relation $f(\mathbf{x}^{(t+1)}) - f(z)$ as

$$\begin{split} \mathcal{T}_{2}^{t+1} &= \sum_{i} p_{i}^{(t)} \left[f_{i}(\mathbf{x}^{(t+1)}) - f_{i}(\mathbf{x}^{*}) \right] = \|\mathbf{x}^{(t+1)}\|_{1} + h(\mathbf{x}^{(t+1)}) - \|\mathbf{x}^{*}\|_{1} - h(\mathbf{x}^{*}) \\ \stackrel{(i)}{\leq} \left\langle \partial \|\mathbf{x}^{(t+1)}\|_{1}, \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + h(\mathbf{x}^{(t)}) + \left\langle \nabla h(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} - h(\mathbf{x}^{*}) \\ \stackrel{(ii)}{\leq} \left\langle \partial \|\mathbf{x}^{(t+1)}\|_{1}, \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + \left\langle \nabla h(\mathbf{x}^{(t)}), \mathbf{x}^{(t)} - \mathbf{x}^{*} + \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ \stackrel{(iii)}{=} \left\langle \partial \|\mathbf{x}^{(t+1)}\|_{1} + \nabla h(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ \stackrel{(ii)}{\leq} \left\langle H(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ \stackrel{(v)}{=} \frac{1}{\eta_{x}} \left\langle \mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}, \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ = \frac{1}{\eta_{x}} \left\langle \mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}, \mathbf{x}^{(t+1)} - \mathbf{x}^{*} \right\rangle + \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{*}\|_{2}^{2} + \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t)} - \mathbf{x}^{*}\|_{2}^{2} + \frac{1}{\eta_{x}} \left\langle \mathbf{x}^{(t+1)} - \mathbf{x}^{*}, \mathbf{x}^{*} - \mathbf{x}^{(t)} \right\rangle \\ = \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t)} - \mathbf{x}^{*}\|_{2}^{2} + \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{*}\|_{2}^{2} + \frac{1}{\eta_{x}} \left\langle \mathbf{x}^{(t+1)} - \mathbf{x}^{(t+1)} \right\rangle \\ = \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t)} - \mathbf{x}^{*}\|_{2}^{2} - \frac{1}{2\eta_{x}} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{*}\|_{2}^{2}, \end{split}$$
(B.16)

where in (*i*) we use $\|\mathbf{x}^*\|_1 \ge \|\mathbf{x}^{(t+1)}\|_1 + \langle \partial \|\mathbf{x}^{(t+1)}\|_1, \mathbf{x}^* - \mathbf{x}^{(t+1)} \rangle$ based on the definition of subgradients, and $h(\mathbf{x}^{t+1}) \le h(\mathbf{x}^{(t)}) + \langle \nabla h(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} \rangle + L_h \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2/2$ from the L_h smoothness of $h(\cdot)$, in (*ii*) we use $h(\mathbf{x}^*) \ge h(\mathbf{x}^{(t)}) + \langle \nabla h(\mathbf{x}^{(t)}), \mathbf{x}^* - \mathbf{x}^{(t)} \rangle$ since $h(\cdot)$ is convex, in (*iii*) we use $H(\mathbf{x}^{(t)}) \in \nabla h(\mathbf{x}^{(t)}) + \partial \|\mathbf{x}^{(t+1)}\|_1$, and in (*iv*) we use $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta_x H(\mathbf{x}^{(t)})$, and in (*vi*) we use $\eta_x \le L^{-1}$.

Hence, we finishes the proof by

$$\sum_{t} \mathcal{T}_{2}^{t} \leq \sum_{t} \left[\frac{1}{2\eta_{x}} \| \mathbf{x}^{(t-1)} - \mathbf{x}^{*} \|_{2}^{2} - \frac{1}{2\eta_{x}} \| \mathbf{x}^{(t)} - \mathbf{x}^{*} \|_{2}^{2} \right]$$

= $\frac{1}{2\eta_{x}} \| \mathbf{x}^{(0)} - \mathbf{x}^{*} \|_{2}^{2} - \frac{1}{2\eta_{x}} \| \mathbf{x}^{(T)} - \mathbf{x}^{*} \|_{2}^{2} \leq \frac{1}{2\eta_{x}} \| \mathbf{x}^{(0)} - \mathbf{x}^{*} \|_{2}^{2} \leq \frac{4R^{2}}{2\eta_{x}},$ (B.17)

#

where in (*i*) we use $\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2 \le \|\mathbf{x}^*\|_2 + \|\mathbf{x}^{(0)}\| \le 2R$.

Appendix C Proof of Theorem 4

Proof. First we define $h(\mathbf{x}) = \sum_i p_i h_i(\mathbf{x}) + \lambda_1 \|\mathbf{y} - A\mathbf{x}\|_2^2 / 2 + \lambda_2 \|\mathbf{x}\|_2^2 / 2$. Then we consider the term $\mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})$ and have

$$\mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) = \|\mathbf{x}^{(t+1)}\|_{1} - \|\mathbf{x}^{(t)}\|_{1} + \sum_{i} p_{i}^{(t)} \left(h_{i}(\mathbf{x}^{(t+1)}) - h_{i}(\mathbf{x}^{(t)})\right) \\ \stackrel{(i)}{\leq} \left\langle \partial \|\mathbf{x}^{(t+1)}\|_{1}, \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\right\rangle + \sum_{i} p_{i}^{(t)} \left[\left\langle \nabla h_{i}(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2}\right] \\ \stackrel{(ii)}{=} \left\langle \partial \|\mathbf{x}^{(t+1)}\|_{1} + \sum_{i} p_{i}^{(t)} \nabla h_{i}(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ \stackrel{(iii)}{=} \left\langle \frac{\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)}}{\eta_{w}^{(t)}}, \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\right\rangle + \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ = \frac{1}{2} \left(L_{h} - \frac{2}{\eta_{w}^{(t)}}\right) \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \\ \stackrel{(iv)}{\leq} - \frac{L_{h}}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2}, \end{aligned} \tag{C.1}$$

where in (*i*) we have $\|\mathbf{x}^{(t)}\|_1 \ge \|\mathbf{x}^{(t+1)}\|_1 + \langle \partial \|\mathbf{x}^{(t+1)}\|_1$, $\mathbf{x}^{(t)} - \mathbf{x}^{(t+1)} \rangle$ from the definition of subgradient Beck (2017), and $h_i(\mathbf{x}^{(t+1)}) \le h_i(\mathbf{x}^{(t)}) + \langle \nabla h_i(\mathbf{x}^{(t)}), \mathbf{x}^{(t+1)} - \mathbf{x}^{(t)} \rangle + \frac{L_h}{2} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$, in (*ii*) we use the property $\partial \|\mathbf{x}^{(t+1)}\|_1 + \sum_i p_i^{(t)} \nabla h_i(\mathbf{x}^{(t)}) \in H(\mathbf{x}^{(t)})$, in (*iii*) we use $\mathbf{x}^{(t+1)} = \mathbf{x}^{(t)} - \eta_w^{(t)} H(\mathbf{x}^{(t)})$, and in (*iv*) we use $\eta_w^{(t)} \le L^{-1}$. Adopting similar tricks as Qian *et al.* (2018), we could upper-bound $\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_2^2$ as

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2}{L_{h}} \left[\mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) \right]$$

= $\frac{2}{L_{h}} \left[\underbrace{\mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)})}_{\mathcal{T}_{1}^{t}} + \underbrace{\mathcal{L}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathcal{L}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)})}_{\mathcal{T}_{2}^{t}} \right].$ (C.2)

Then we separately discuss bound \mathcal{T}_1^t and \mathcal{T}_2^t . Since most terms of $\sum_t \mathcal{T}_1^t$ will be cancelled after summarization, we focus the analysis on bounding \mathcal{T}_2^t , which is

$$\mathcal{T}_{2}^{t} = \left\langle \mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}, f(\mathbf{x}^{(t+1)}) \right\rangle \leq \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_{1} \underbrace{\|f(\mathbf{x}^{(t+1)})\|_{\infty}}_{\leq R_{f}}, \tag{C.3}$$

where $f(\mathbf{x}^{(t+1)})$ denotes the vector whose *i*th element is $f_i(\mathbf{x}^{(t+1)})$. Notice that we have

$$\|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_{1}^{2} \stackrel{(i)}{\leq} 2D_{KL} \left(\mathbf{p}^{(t+1)} || \mathbf{p}^{(t)}\right) \stackrel{(ii)}{\leq} 2\eta_{p}^{(t)} \left\langle \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)}, f(\mathbf{x}^{(t)}) \right\rangle$$

$$\leq 2\eta_{p}^{(t)} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_{1} \|f(\mathbf{x}^{(t)})\|_{\infty} \leq 2\eta_{p}^{(t)} R_{f} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_{1},$$
(C.4)

which gives us $\|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_1 \le 2\eta_p^{(t)}R_f$, where (*i*) is because of *Pinsker's inequality* (Theorem 4.19 in Boucheron *et al.* (2013)) and (*ii*) is because of Lemma 7. To conclude, we have upper-bound \mathcal{T}_2^t as

$$\mathcal{T}_2^t \le 2\eta_p^{(t)} R_f^2. \tag{C.5}$$

*

Then we finish the proof as

$$\sum_{t} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2\mathcal{L}(\mathbf{p}^{(0)}, \mathbf{x}^{(0)}) - 2\mathcal{L}(\mathbf{p}^{(T+1)}, \mathbf{x}^{(T+1)})}{L_{h}} + \frac{4R_{f}^{2}\sum_{t}\eta_{p}^{(t)}}{L_{h}}$$

$$\stackrel{(i)}{\leq} \frac{2\mathcal{L}(\mathbf{p}^{(0)}, \mathbf{x}^{(0)})}{L_{h}} + \frac{4R_{f}^{2}\sum_{t}\eta_{p}^{(t)}}{L_{h}},$$
(C.6)

where (*i*) is because $\mathcal{L}(\mathbf{p}^{(T+1)}, \mathbf{x}^{(T+1)}) \ge 0$.

Lemma 7. With Alg. 1, we have

$$D_{KL}\left(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}\right) \leq \eta_p^{(t)}\left\langle \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)}, f(\mathbf{x}^{(t)})\right\rangle,$$
(C.7)

where $f(\mathbf{x}^{(t)})$ denotes the vector whose ith element is $f_i(\mathbf{x}^{(t)})$.

Proof. Here we have

$$D_{KL}(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}) = \sum_{i} p_{i}^{(t+1)} \log\left(\frac{p_{i}^{(t+1)}}{p_{i}^{(t)}}\right)$$

= $\sum_{i} p_{i}^{(t+1)} \log\frac{e^{-\eta_{p}^{(t)}f_{i}(\mathbf{x}_{t})}}{Z_{t}} = -\log\left(Z_{t}\right) - \eta_{p}^{(t)}\sum_{i} p_{i}^{(t+1)}f_{i}(\mathbf{x}^{(t)})$
= $-\log\left(Z_{t}\right) - \eta_{p}^{(t)}\left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)})\right\rangle - \eta_{p}^{(t)}\left\langle \mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)})\right\rangle,$ (C.8)

where $Z_t \triangleq \sum_i p_i^{(t)} e^{-\eta_p f_i(\mathbf{x}^{(t)})}$.

Then we have

$$\eta_{p}^{(t)} \left\langle \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)}, f(\mathbf{x}^{(t)}) \right\rangle = D_{KL}(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}) + \log\left(\sum_{i} p_{i}^{(t)} e^{-\eta_{p}^{(t)} f_{i}(\mathbf{x}^{(t)})}\right) + \eta_{p}^{(t)} \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}_{t}) \right\rangle$$

$$\stackrel{(i)}{\geq} D_{KL}(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}) + \log\left[\prod_{i} e^{-\eta_{p}^{(t)} p_{i}^{t} f_{i}(\mathbf{x}^{(t)})}\right] + \eta_{p}^{(t)} \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}^{(t)}) \right\rangle$$

$$= D_{KL}(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}) + \underbrace{\sum_{i} \log\left(e^{-\eta_{p}^{(t)} p_{i}^{(t)} f_{i}(\mathbf{x}^{(t)})}\right) + \eta_{p}^{(t)} \left\langle \mathbf{p}^{(t)}, f(\mathbf{x}_{t}) \right\rangle}_{0} = D_{KL}(\mathbf{p}^{(t+1)}||\mathbf{p}^{(t)}), \tag{C.9}$$

where in (*i*) we use $\sum_{i} p_i x_i \ge \prod_{i} x_i^{p_i}$ such that $\sum_{i} p_i = 1$, $p_i \ge 0$.

Appendix D Proof of Theorem 5

Proof. Define \mathbf{p}^* and \mathbf{x}^* as

$$\mathbf{p}^{*} = \operatorname{argmin}_{\mathbf{p}} \sum_{t} \mathcal{LR}\left(\mathbf{p}, \mathbf{x}^{(t)}\right), \quad \mathbf{x}^{*} = \operatorname{argmin}_{\mathbf{x}} \sum_{t} \mathcal{LR}\left(\mathbf{p}^{(t)}, \mathbf{x}\right). \tag{D.1}$$

respectively First we define \mathcal{T}_1^t and \mathcal{T}_2^t as

$$\mathcal{T}_{1}^{t} = \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{LR}(\mathbf{p}^{*}, \mathbf{x}^{(t)});$$

$$\mathcal{T}_{2}^{t} = \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{*}),$$

(D.2)

*

respectively. Then our goal becomes bounding $|\sum_t \mathcal{T}_1^t| + |\sum_t \mathcal{T}_2^t|$, For term $|\sum_t \mathcal{T}_2^t|$, the analysis stays the same as Lemma 6. Here we focus on bounding $\sum_t \mathcal{T}_1^t$, which proceeds as

$$\sum_{t} \mathcal{T}_{1}^{t} = \sum_{t} \left(\mathcal{L}\mathcal{R}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{L}\mathcal{R}(\mathbf{p}, \mathbf{x}^{(t)}) \right) = \sum_{t} \left\{ -\left\langle \underbrace{\nabla_{t} \mathcal{L}\mathcal{R}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)})}_{\mathbf{g}^{(t)}}, \mathbf{p} - \mathbf{p}^{(t)} \right\rangle - \frac{\lambda_{3}}{2} \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} \right\}$$
$$= \sum_{t} \left\{ \left\langle \mathbf{g}^{(t)}, \mathbf{p}^{(t)} - \mathbf{p} \right\rangle - \frac{\lambda_{3}}{2} \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} \right\}$$
(D.3)

Then we consider the distance $\| \mathbf{p}^{(t+1)} - \mathbf{p} \|_2^2$ which is

$$\|\mathbf{p}^{(t+1)} - \mathbf{p}\|_{2}^{2} = \|\mathbb{P}_{\Delta}(\mathbf{p}^{(t)} - \eta_{p}^{t}\mathbf{g}^{(t)}) - \mathbf{p}\|_{2}^{2} \leq \|\mathbf{p}^{(t)} - \eta_{p}^{(t)}\mathbf{g}^{(t)} - \mathbf{p}\|_{2}^{2}$$

$$= \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} + (\eta_{p}^{(t)})^{2} \|\mathbf{g}^{(t)}\|_{2}^{2} - 2\eta_{p}^{(t)} \langle \mathbf{g}^{(t)}, \mathbf{p}^{(t)} - \mathbf{p} \rangle,$$
(D.4)

where in (i) we use the contraction property for projection, which gives us

$$\left\langle \mathbf{g}^{(t)}, \mathbf{p}^{(t)} - \mathbf{p} \right\rangle \le \frac{\eta_p^{(t)} \|\mathbf{g}^{(t)}\|_2^2}{2} + \frac{\|\mathbf{p}^{(t)} - \mathbf{p}\|_2^2 - \|\mathbf{p}^{(t+1)} - \mathbf{p}\|_2^2}{2\eta_p^{(t)}}$$
 (D.5)

By setting $\eta_p^{(t)} = (\lambda t)^{-1}$, we have

$$\sum_{t} \mathcal{T}_{1}^{t} \leq \frac{R_{g}^{2} \log T}{2\lambda_{3}} + \frac{\lambda_{3}}{2} \sum_{t} t(\|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} - \|\mathbf{p}^{(t+1)} - \mathbf{p}\|_{2}^{2}) - \frac{\lambda_{3}}{2} \sum_{t} \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2}$$

$$= \frac{R_{g}^{2} \log T}{2\lambda_{3}} + \frac{\lambda_{3}}{2} \sum_{t} \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} - \frac{\lambda_{3}(T+1)}{2} \|\mathbf{p}^{(T+1)} - \mathbf{p}\|_{2}^{2} - \frac{\lambda_{3}}{2} \sum_{t} \|\mathbf{p}^{(t)} - \mathbf{p}\|_{2}^{2} \qquad (D.6)$$

$$= \frac{R_{g}^{2} \log T}{2\lambda_{3}} - \frac{\lambda_{3}(T+1)}{2} \|\mathbf{p}^{(T+1)} - \mathbf{p}\|_{2}^{2} \leq \frac{R_{g}^{2} \log T}{2\lambda_{3}}.$$

Hence, we have

$$\frac{\sum_t \mathcal{T}_1^t + \mathcal{T}_2^t}{T} \le \frac{R_g^2 \log T}{2\lambda_3 T} + \frac{R^2}{2\eta_x T'}$$
(D.7)

*

which completes the proof.

Appendix E Proof of Theorem 6

Proof. Following the same procedure in C, we can bound

$$\|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2}{L_{h}} \left[\mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) \right]$$

= $\frac{2}{L_{h}} \left[\underbrace{\mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t)}) - \mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)})}_{\mathcal{T}_{1}^{t}} + \underbrace{\mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)})}_{\mathcal{T}_{2}^{t}} \right].$ (E.1)

Since \mathcal{T}_1^t will cancel themselves after summarization, we focus on bounding \mathcal{T}_2^t . Then we have

$$\mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) = \mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) + \left\langle \nabla_{\mathbf{p}} \mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}), \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)} \right\rangle + \frac{\lambda_3}{2} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_2^2$$

$$\stackrel{(i)}{=} \mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) + \left\langle f(\mathbf{x}^{(t+1)}) + \lambda_3(\mathbf{p}^{(t+1)} - q), \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)} \right\rangle + \frac{\lambda_3}{2} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_2^2,$$
(E.2)

where in (i) we have

$$\nabla_{\mathbf{p}} \mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) = f(\mathbf{x}^{(t+1)}) + \lambda_3 \left(\mathbf{p}^{(t+1)} - \mathbf{q}\right).$$
(E.3)

Then we have

$$\mathcal{LR}(\mathbf{p}^{(t+1)}, \mathbf{x}^{(t+1)}) - \mathcal{LR}(\mathbf{p}^{(t)}, \mathbf{x}^{(t+1)}) = \left\langle \mathbf{g}^{(t+1)}, \mathbf{p}^{(t)} - \mathbf{p}^{(t+1)} \right\rangle + \frac{\lambda_3}{2} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_2^2$$

$$\leq \|\mathbf{g}^{(t+1)}\|_2 \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_2 + \frac{\lambda_3}{2} \|\mathbf{p}^{(t+1)} - \mathbf{p}^{(t)}\|_2^2$$

$$\leq \|\mathbf{g}^{(t+1)}\|_2 \|\eta_p^{(t)}\mathbf{g}^{(t)}\|_2 + \frac{\lambda_3 \left(\eta_p^{(t)}\right)^2}{2} \|\mathbf{g}^{(t)}\|_2^2 \leq R_g^2 \left(\eta_p^{(t)} + \frac{\lambda_3 \left(\eta_p^{(t)}\right)^2}{2}\right).$$
(E.4)

Hence, we conclude that

$$\sum_{t} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_{2}^{2} \leq \frac{2\mathcal{LR}(\mathbf{p}^{(0)}, \mathbf{x}^{(0)})}{L_{h}} + \frac{2R_{g}^{2}}{L_{h}} \sum_{t} \left(\eta_{p}^{(t)} + \frac{\lambda_{3}\left(\eta_{p}^{(t)}\right)^{2}}{2}\right), \quad (E.5)$$

#

where $\|\mathbf{g}^{(t)}\|_{2} \le R_{g}$.