

Exploiting Prior Information in Block-sparse Signals

Sajad Daei, Farzan Haddadi, Arash Amini

Abstract—We study the problem of recovering a block-sparse signal from under-sampled observations. The non-zero values of such signals appear in few blocks, and their recovery is often accomplished using an $\ell_{1,2}$ optimization problem. In applications such as DNA micro-arrays, some extra information about the distribution of non-zero blocks is available; i.e., the number of non-zero blocks in certain subsets of the blocks is known. A typical way to consider the extra information in recovery procedures is to solve a weighted $\ell_{1,2}$ problem. In this paper, we consider a block-sparse model which is accompanied with a partitioning of the blocks; besides the overall block-sparsity level of the signal, we assume to know the block-sparsity of each subset in the partition. Our goal in this work is to minimize the number of required linear measurements for perfect recovery of the signal by tuning the weights of a weighted $\ell_{1,2}$ problem. For this goal, we apply tools from conic integral geometry and derive closed-form expressions for the optimal weights. We show through precise analysis and simulations that the weighted $\ell_{1,2}$ problem with optimal weights significantly outperforms the regular $\ell_{1,2}$ problem. We further show that the optimal weights are robust against the inaccuracies of prior information.

Index Terms—block-sparse, prior information, weighted $\ell_{1,2}$, conic integral geometry.

I. INTRODUCTION

COMPRESSED sensing (CS) aims at recovering a sparse signal $\mathbf{x} \in \mathbb{R}^n$ from a few random linear measurements

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}, \quad (1)$$

where $\mathbf{A} \in \mathbb{R}^{m \times n}$ is called the measurement matrix, \mathbf{e} is the noise vector which is assumed to be i.i.d. Gaussian with variance σ_e^2 , and m is much smaller than n . In this work, the signal is assumed to be block-sparse which occurs in many applications such as DNA micro-arrays [1], direction of arrival (DOA) estimation [2], computational neuroscience [3], and multiple measurement vector (MMV) problem [4]. Consider a block-sparse signal $\mathbf{x} \in \mathbb{R}^n$ which is a concatenation of q blocks $\mathcal{V}_b, b = 1, \dots, q$. The block support is defined as the index set of non-zero blocks, which we denote by \mathcal{B} . Figure 1 shows an example of a block-sparse signal composed of $q = 10$ blocks of length $k = 10$. The block support in this figure is the index set $\{3, \dots, 7\}$; i.e., the signal is 5-block-sparse and contains zero blocks except for $\{\mathcal{V}_3, \dots, \mathcal{V}_7\}$. Ideally, a block-sparse signal is reconstructed using the following optimization problem:

$$\begin{aligned} P_{0,2}^\eta : \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{z}\|_{0,2} &:= \sum_{b=1}^q \mathbf{1}_{\|\mathbf{z}_{\mathcal{V}_b}\|_2 > 0} \\ \text{s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_2 &\leq \eta, \end{aligned} \quad (2)$$

where $\mathbf{1}_{\mathcal{E}}$ denotes the indicator function of the event \mathcal{E} , and η is an upper-bound for $\|\mathbf{e}\|_2$. This problem is computationally

intractable in polynomial time and in general, it is NP-hard. Following Donoho [5], $P_{0,2}^\eta$ can be relaxed as an $\ell_{1,2}$ minimization of the form

$$\begin{aligned} P_{1,2}^\eta : \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{z}\|_{1,2} &:= \sum_{b=1}^q \|\mathbf{z}_{\mathcal{V}_b}\|_2 \\ \text{s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_2 &\leq \eta. \end{aligned} \quad (3)$$

It is known that $P_{1,2}^\eta$ for $\eta = 0$ (noiseless case) finds the original block-sparse signal with high probability if \mathbf{A} comes from a probability law that distributes the null space uniformly with respect to the Haar measure [6]. This includes Gaussian ensembles and partial Fourier matrices. We also assume the same type of random matrices in this paper.

The main challenge in recovering a block-sparse signal is to find the block support. Intuitively, if there exists an additional piece of information about the block-sparsity level of certain subsets of the blocks¹, one can probably solve $P_{1,2}^0$ and $P_{1,2}^\eta$ with fewer measurements or smaller reconstruction error, respectively (see Subsection I-C for more explanations). In some scenarios, such extra information is available. For example, in DNA microarray one might know the occurrence frequency of specific Genes in certain sets of arrays. In computational neuroscience, the behavior of neurons exhibit non-uniform clustered responses [3]. In radar signal processing, an operator might know the range of speed with which an aircraft flies or the range of angular domain it forms in the radar [7]. A typical way to consider such prior information is to apply a weighted $\ell_{1,2}$ minimization formulated as:

$$\begin{aligned} P_{1,2,\mathbf{w}}^\eta : \min_{\mathbf{z} \in \mathbb{R}^n} \|\mathbf{z}\|_{1,2,\mathbf{w}} &:= \sum_{b=1}^q w_b \|\mathbf{z}_{\mathcal{V}_b}\|_2 \\ \text{s.t. } \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_2 &\leq \eta, \end{aligned} \quad (4)$$

where $\mathbf{w} = [w_1, \dots, w_q]^T$. This work explicitly takes prior block information into account by optimally tuning the weights in $P_{1,2,\mathbf{w}}^\eta$. Specifically, the main objective of this paper is to find optimal weights in $P_{1,2,\mathbf{w}}^\eta$ in the sense that they minimize the reconstruction error for $\eta > 0$, and the required number of measurements for $\eta = 0$.

A. Contributions

Although block-sparse signals are widely used in many applications in the literature (see Subsection I-B), the inclusion of prior information about the block-support distribution in improving the performance of the signal recovery is not

¹Our block-sparse model is deterministic in this paper; the normalized block-sparsity of a subset of blocks (the number of nonzero blocks divided by the total number of blocks in the subset) is occasionally referred to as the non-zero probability or non-zero likelihood of the blocks in this subset.

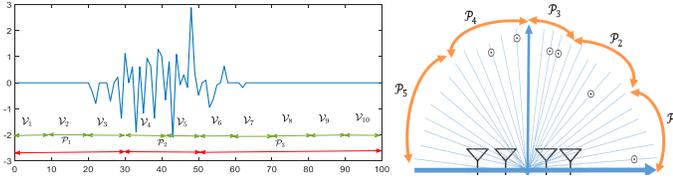


Figure 1: Left image: Illustration of a non-uniform block-sparse signal $\mathbf{x} \in \mathbb{R}^n$ with parameters $n = 100$, $q = 10$, $k = 10$, $L = 3$, $\alpha_1 = \frac{1}{3}$, $\alpha_2 = 1$ and $\alpha_3 = \frac{2}{5}$. Right image: Schematic diagram of DOA estimation of far-field sources. The angular half-space is divided into $q = 30$ angular clusters $\{\mathcal{V}_b\}_{b=1}^q$ with equal length. $L = 5$ block support estimators with accuracies $\alpha_1 = \frac{1}{6}$, $\alpha_2 = \frac{1}{6}$, $\alpha_3 = \frac{1}{2}$, $\alpha_4 = \frac{1}{3}$ and $\alpha_5 = 0$ are considered.

analyzed yet. In this work, we first consider the standard block-sparse model without prior information and derive tight and non-asymptotic bounds for the required sample complexity of $P_{1,2}^0$. Interestingly, the asymptotic form of our bound matches the conjecture in [8]. Then, we consider the weighted $P_{1,2}^\eta$ problem denoted by $P_{1,2,w}^\eta$ in presence of prior information for which we present a new method for obtaining the optimal weights. In particular, we assume to know a partitioning of $[q] := \{1, \dots, q\}$ into L disjoint subsets $\{\mathcal{P}_i\}_{i=1}^L$, about which we know the number of non-zero blocks in each index set \mathcal{P}_i . In simple words, we assume to know the block-sparsity level of each subset in a partitioning of the blocks. These partial block-sparsity levels are commonly expressed via $\alpha_i := \frac{|\mathcal{P}_i \cap \mathcal{B}|}{|\mathcal{P}_i|}$ which are called the accuracy of \mathcal{P}_i s. Obviously, if the partition consists of a single subset $[q]$, we are dealing with a conventional block-sparse model without prior information. Note that the measurement matrix combines all the blocks before generating the measurement vector. Hence, we do not have access to measurements that are solely determined by the blocks of a single partition; otherwise, the recovery problem could be simplified to multiple block-sparse recovery in lower dimensions without prior information. We call signals generated based on this model as *non-uniform* block-sparse signals. The left image of Figure 1 shows a sample signal with this model. The strategy of finding optimal weights is to study the asymptotics of the upper and lower bounds of the statistical dimension of a certain convex cone (these bounds differ from each other in an asymptotically vanishing constant term). In this regard, we prove that our method of obtaining optimal weights is robust to slight inaccuracies of the prior knowledge. With optimal weights in $P_{1,2,w}^0$, we also find that the number of measurements required for exact recovery of a block-sparse vector equals to the sum of the required number of measurements in the simpler (smaller size) problems of recovering each subset via $P_{1,2}^0$ separately.

B. Related works and Key Differences

There are extensive works in CS literature that consider reconstructing low-dimensional signals with prior information; for instance see [9]–[14]. Here, we discuss some of them.

The weighted ℓ_1 minimization for sparse signal recovery is likely initiated by the work [15]. While this work reveals the advantages of weighting, it does not consider any prior information about the distribution of the support. In [16], partial

knowledge about the support is assumed as prior information. In this work, some of the standard sparse recovery techniques are modified to benefit from the available prior information; however, none of the methods rely on weighting. A non-uniform sparse structure similar to the one that is considered in this paper, is proposed in [9] to incorporate prior information about the partial sparsity level of a given partitioning of the elements in sparse signals. The authors apply a weighed ℓ_1 minimization approach and find an upper-bound for the failure probability of the recovery method. The bound is derived based on the internal and external angles of a certain weighted polytope. The weights could then be optimized to minimize the resulting upper-bound. Unfortunately, this approach does not yield explicit expressions for the weights; in addition, since the probability bound is not necessarily tight, the optimality of such weights is questionable. As mentioned by the authors, the extension of this approach to block-sparse signals is very sophisticated.

Using a different approach, [10] investigates the same model with weighted ℓ_1 minimization. Indeed, the approach consists of dividing the signal elements into two subsets; the elements of the first set are penalized with a variable weight, while the weights of the second set are kept fixed at 1. If the first set encompasses at least 50% of the signal support, it is shown in [10] that the weighting technique can improve the recovery performance.

The bipartitioning approach of [10] can be considered as having two estimates of the support. The extension to multiple estimates for the support with varying degrees of accuracy is studied in [11]. It is again confirmed that a multi-level weighting scheme can enhance the recovery performance. The tuning of the weights in [11] is achieved based on empirical results rather than theoretical findings.

A systematic approach for determining the optimal weights is presented in [17]. Indeed, the optimal weights are derived by minimizing the number of required measurements for almost sure sparse recovery of a weighted ℓ_1 minimization. The resulting optimal weights are expressed in terms of the accuracy level (normalized partial sparsity level) of the sets. Our paper could be considered as an extension of [17] to the block-sparse setting. However, our approach to derive theoretical results is substantially different from [17]. For instance, deriving the optimal weights (and proving their uniqueness) for the $P_{1,2,w}^\eta$ problem is much more challenging than the weighted ℓ_1 minimization. Moreover, the construction of a tight upper-bound for the sample complexity that leads to closed-form expressions is nontrivial.

In many practical scenarios, the low-dimensional signal of interest has block-sparse structure rather than simple sparse structure (see Section I-C and [9, Section I] for illustrative examples). In this work, we consider a non-uniform block-sparse model where we have a partitioning of the blocks and their associated accuracy levels. This model extends the non-uniform sparse model considered in [10], [17] to the non-uniform block-sparse signals which is observed in a few works such as [18]–[20]. A distinctive feature of our work, besides the generality of our model, is the closed-form expressions for the optimal weights in the $P_{1,2,w}^\eta$ problem. In contrast to

the result in [10], we show that the minimum 50% accuracy level is not a pre-requisite for $P_{1,2,w}^0$ to outperform the unweighted $P_{1,2}^0$ recovery approach. We study the robustness of the weights under inaccuracies in prior knowledge for the first time, and prove robustness for the introduced optimal weights. The study of robustness has not been accomplished in the past even for the easier sparse models.

C. DOA estimation

DOA estimation is the problem of finding the direction of a few sources from observations on an array of sensors. The location of each source is characterized by the direction of arrival $\theta \in [-\frac{\pi}{2}, \frac{\pi}{2}]$ with respect to the array axis. Let $\theta_i := \pi(\frac{i-1}{q-1} - \frac{1}{2})$, $i = 1, \dots, q$ be a discretization of the interval $[-\frac{\pi}{2}, \frac{\pi}{2}]$ into q grid points representing the potential direction of the sources (see the right block of Figure 1). Assume that $s \ll q$ far-field narrow-band sources with wavelength λ whose directions exactly match the grid angles incident on a q -element uniform linear array (ULA) of sensors with inter-sensor spacing d . Form the so-called array manifold matrix $\mathbf{F} = [\mathbf{a}(\theta_1), \dots, \mathbf{a}(\theta_q)] \in \mathbb{C}^{q \times q}$ where $\mathbf{a}(\theta_l) := e^{j2\pi(d/\lambda)[0, \dots, q-1]^T \sin(\theta_l)}$ is the steering vector caused by the propagation delay from i th potential source to each of the q sensors. The observed signal at ULA elements at time t are given by

$$\mathbf{y}(t) = \mathbf{F} \mathbf{x}(t) + \mathbf{e}(t), \quad (5)$$

where $x_i(t) = \sqrt{p_i} e^{j2\pi i t}$ is the amplitude of the i th potential source, $\mathbf{y}(t) \in \mathbb{C}^q$ is the received vector at time (or snapshot) t at all sensors, and $\mathbf{e}(t)$ is a Gaussian noise term with variance σ_e^2 . In practice, the number of ULA sensors is limited due to cost limitations or physical constraints; hence, it is crucial in DOA estimation to exploit time redundancies (taking several snapshots, *i.e.*, $\mathbf{y}(t_1), \dots, \mathbf{y}(t_k)$). With the assumption that the sources remain fixed in all snapshots, the received signals at all k time instances could be expressed as

$$\mathbf{Y} := [\mathbf{y}(t_1), \dots, \mathbf{y}(t_k)] = \mathbf{F} \mathbf{X} + \mathbf{E}, \quad (6)$$

where $\mathbf{X} = [\mathbf{x}(t_1), \dots, \mathbf{x}(t_k)] \in \mathbb{C}^{q \times k}$ and \mathbf{E} is similarly defined. In practical scenarios, we can not directly observe the received signal at all of the sensors. Instead, we have access to the outputs of $m \ll q$ sensors formed by

$$\mathbf{Y} = \mathbf{A} \mathbf{F} \mathbf{X} + \tilde{\mathbf{E}} \in \mathbb{C}^{m \times k}, \quad (7)$$

where $\mathbf{A} \in \mathbb{C}^{m \times q}$ is the partial identity matrix², $\mathbf{A} \mathbf{F}$ represents the effective measurement matrix and $\tilde{\mathbf{E}} := \mathbf{A} \mathbf{E} \in \mathbb{C}^{m \times k}$. In practical settings, the number of sources are limited. Thus, \mathbf{X} consists mainly of zeros. In addition, if \mathbf{X} is vectorized by concatenating the rows, we achieve a block-sparse vector. Indeed, each row forms a block in the vectorized signal. It is common in radar applications that the operator has prior information (previous measurement of recorded history) about the partial number of sources in a given angular partitioning (see $\{\mathcal{P}_i\}_{i=1}^5$ in the right block of Figure 1). More

²This means that only m random rows of the identity matrix $\mathbf{I}_{q \times q}$ are chosen.

precisely, the radar operator might know the accuracy (α_i s in the right block of Figure 1) of specific angular bands. It is of considerable importance to exploit these information in order to minimize the number of required sensors for finding the sources. We will turn back to this problem in Section VII.

The paper is organized as follows: some concepts from convex geometry are reviewed in Section II. The signal model and our methodology are stated in Section III. The number of measurements required for $P_{1,2,w}^0$ is obtained in Section IV. In Section V, the procedure of finding optimal weights is explained. In Section VI, the robustness of optimal weights with respect to inaccuracies in prior information is discussed. Numerical simulations on synthetic data that support the theory are presented in Section VII. Finally, the paper is concluded in Section VIII.

Notation Throughout the paper, scalars are denoted by lowercase letters, vectors by lowercase boldface letters, and matrices by uppercase boldface letters. The i th element of a vector \mathbf{x} is shown either by $x(i)$ or x_i . $[n]$ refers to $\{1, \dots, n\}$. We reserve calligraphic uppercase letters for sets (e.g. \mathcal{S}). The cardinality of a set \mathcal{S} is shown by $|\mathcal{S}|$. $\bar{\mathcal{S}}$ denotes the complement $[n] \setminus \mathcal{S}$ of a set $\mathcal{S} \subset [n]$. For $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{x}_{\mathcal{S}}$ is the subvector in $\mathbb{R}^{|\mathcal{S}|}$ consisting of the entries indexed by \mathcal{S} , that is, $(\mathbf{x}_{\mathcal{S}})_i = x_{j_i} : \mathcal{S} = \{j_i\}_{i=1}^{|\mathcal{S}|}$. In this paper, $\mathbf{1}_{\mathcal{E}}$ denotes the indicator of a set \mathcal{E} . The null space of linear operators is denoted by $\text{null}(\cdot)$. Given a vector $\mathbf{x} \in \mathbb{R}^n$ and a set $\mathcal{C} \subseteq \mathbb{R}^n$, the set obtained by scaling elements of \mathcal{C} by the elements of \mathbf{x} is denoted by $\mathbf{x} \odot \mathcal{C}$. The polar \mathcal{K}° of a cone $\mathcal{K} \subset \mathbb{R}^n$ is the set of vectors forming non-acute angles with every vector in \mathcal{K} , *i.e.* $\mathcal{K}^\circ = \{\mathbf{v} \in \mathbb{R}^n : \langle \mathbf{v}, \mathbf{z} \rangle \leq 0 \ \forall \mathbf{z} \in \mathcal{K}\}$. For a matrix \mathbf{A} , the operator norm is defined as $\|\mathbf{A}\|_{p \rightarrow q} = \sup_{\|\mathbf{x}\|_p \leq 1} \|\mathbf{A}\mathbf{x}\|_q$. For $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, $\mathbf{x} \leq \mathbf{y}$ stands for component-wise inequality while $\mathbf{x} < \mathbf{y}$ denotes component-wise inequality with strict inequality in at least one component. \mathbb{B}_ϵ^n refers to the ϵ -ball $\mathbb{B}_\epsilon^n = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq \epsilon\}$. Lastly, the notation $(a)_+$ stands for the positive part of a , *i.e.*, $\max\{a, 0\}$.

II. PRELIMINARIES

In this section, basic concepts of conic integral geometry are reviewed.

A. Subdifferential

The subdifferential of a proper³ convex function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$ at $\mathbf{x} \in \mathbb{R}^n$ is given by:

$$\partial f(\mathbf{x}) := \{\mathbf{z} \in \mathbb{R}^n : f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{z}, \mathbf{y} - \mathbf{x} \rangle : \forall \mathbf{y} \in \mathbb{R}^n\}. \quad (8)$$

Proposition 1. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$ be a proper convex function that is 1-homogeneous, *i.e.* $f(\alpha \mathbf{z}) = |\alpha| f(\mathbf{z}) : \forall \alpha \in \mathbb{R}$ and sub-additive, *i.e.* $f(\mathbf{x} + \mathbf{y}) \leq f(\mathbf{x}) + f(\mathbf{y}) : \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$, then we have a simpler form of subdifferential given by:*

$$\partial f(\mathbf{x}) = \left\{ \mathbf{z} \in \mathbb{R}^n : \begin{array}{l} \langle \mathbf{z}, \mathbf{x} \rangle = f(\mathbf{x}), f^*(\mathbf{z}) = 1, \quad \mathbf{x} \neq \mathbf{0} \\ f^*(\mathbf{z}) \leq 1, \quad \mathbf{x} = \mathbf{0} \end{array} \right\} \quad (9)$$

³An everywhere defined function taking values in $(-\infty, \infty]$ with at least one finite value in $(-\infty, \infty)$.

where $f^*(\mathbf{z}) = \sup_{f(\mathbf{y}) \leq 1} \langle \mathbf{z}, \mathbf{y} \rangle$ is the dual function of $f(\mathbf{z})$.

Proof. See Appendix A.

B. Descent Cones

The descent cone of a proper convex function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\pm\infty\}$ at point $\mathbf{x} \in \mathbb{R}^n$ is the conic hull of the directions from \mathbf{x} that f does not increase:

$$\mathcal{D}(f, \mathbf{x}) = \bigcup_{t \geq 0} \{\mathbf{z} \in \mathbb{R}^n : f(\mathbf{x} + t\mathbf{z}) \leq f(\mathbf{x})\}. \quad (10)$$

The descent cone of a convex function is a convex set. There is also a relation between descent cone and subdifferential of a convex function [21, Chapter 23] which is given by:

$$\mathcal{D}^\circ(f, \mathbf{x}) = \text{cone}(\partial f(\mathbf{x})) := \bigcup_{t \geq 0} t \cdot \partial f(\mathbf{x}). \quad (11)$$

C. Statistical Dimension

Definition 1. Statistical Dimension [6]: Let $\mathcal{C} \subseteq \mathbb{R}^n$ be a convex closed cone. Statistical dimension of \mathcal{C} is defined as:

$$\delta(\mathcal{C}) := \mathbb{E}_{\mathbf{g}} \|\mathcal{P}_{\mathcal{C}}(\mathbf{g})\|_2^2 = \mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, \mathcal{C}^\circ), \quad (12)$$

where

$$\mathcal{P}_{\mathcal{C}}(\mathbf{x}) := \arg \min_{\mathbf{z} \in \mathcal{C}} \|\mathbf{z} - \mathbf{x}\|_2, \quad (13)$$

is the projection of $\mathbf{x} \in \mathbb{R}^n$ onto the set \mathcal{C} , and \mathbf{g} is a vector with i.i.d. elements each following a standard normal distribution.

Statistical dimension is a measure of the size of a convex cone and generalizes the concept of dimension for linear subspaces to the class of convex cones.

D. Inverse Problems via Convex Optimization

Convex optimization is a common approach for recovering a structured signal $\mathbf{x}_{n \times 1}$ from m linear measurements of the form (1). Let $f(\cdot)$ be a function that promotes the structure of \mathbf{x} (e.g., sparsity). Then, one might be able to recover the signal \mathbf{x} by solving the problem

$$\mathbf{P}_f^\eta : \min_{\mathbf{z} \in \mathbb{R}^n} f(\mathbf{z}) \quad \text{s.t.} \quad \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2 \leq \eta, \quad (14)$$

where η is a known upper-bound on the norm of the noise \mathbf{e} . For the noiseless case of $\mathbf{e} = \mathbf{0}$ in (1), it is shown in [6] that \mathbf{P}_f^0 exhibits a sharp phase-transition behavior (success/failure) as the number of measurements increases. In addition, the boundary of the transition is determined by the statistical dimension of the descent cone of f at \mathbf{x} , i.e. $\delta(\mathcal{D}(f, \mathbf{x}))$.

In Theorem 1, we present a result regarding the recovery performance of (14) from random measurements, in both noiseless and noisy settings. This theorem is an adaptation of [22, Corollary 3.5] and [6, Theorem 2].

Theorem 1. Let $f(\cdot)$ be a proper convex function that promotes the structure of \mathbf{x} . Let $\mathbf{A}_{m \times n}$ be a random matrix whose

null space is uniformly distributed with respect to the Haar measure. Then, if

$$m \geq \delta(\mathcal{D}(f, \mathbf{x})) + \sqrt{8 \log\left(\frac{4}{\zeta}\right)n}$$

for some $\zeta \in [0, 1]$, then \mathbf{P}_f^0 recovers \mathbf{x} from $\mathbf{y}_{m \times 1} = \mathbf{A}\mathbf{x}$ with probability at least $1 - \zeta$. Alternatively, in the noisy case of $\mathbf{y}_{m \times 1} = \mathbf{A}\mathbf{x} + \mathbf{e}$, where $\|\mathbf{e}\|_2 \leq \eta$, if $\hat{\mathbf{x}}_\eta$ is any solution of \mathbf{P}_f^η , then

$$\|\hat{\mathbf{x}}_\eta - \mathbf{x}\|_2 \leq \frac{2\eta}{(\sqrt{m-1} - \sqrt{\delta(\mathcal{D}(f, \mathbf{x})) - \zeta})_+}, \quad (15)$$

with probability at least $1 - e^{-\frac{\zeta^2}{2}}$.

Also in [6], the following error bound for the statistical dimension is provided.

Theorem 2. [6, Theorem 4.3] For any $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$:

$$0 \leq \inf_{t \geq 0} \mathbb{E} \text{dist}^2(\mathbf{g}, t\partial f(\mathbf{x})) - \delta(\mathcal{D}(f, \mathbf{x})) \leq \frac{2 \sup_{\mathbf{s} \in \partial f(\mathbf{x})} \|\mathbf{s}\|_2}{f\left(\frac{\mathbf{x}}{\|\mathbf{x}\|_2}\right)}. \quad (16)$$

III. MODEL AND METHODOLOGY

Suppose $\mathbf{x}_{n \times 1}$ is a s -block-sparse signal consisting of q blocks $\{\mathcal{V}_b\}_{b=1}^q$ of equal size k , among which only s blocks are non-zero. We are further given a partition of $\{1, \dots, q\}$ into L disjoint subsets $\mathcal{P}_i \subseteq \{1, \dots, q\}$, for which we know

$$\alpha_i = \frac{|\mathcal{P}_i \cap \mathcal{B}|}{|\mathcal{P}_i|}, \quad \rho_i = \frac{|\mathcal{P}_i|}{q}, \quad i = 1, \dots, L, \quad (17)$$

where $\mathcal{B} \subseteq \{1, \dots, q\}$ is the block-support of $\mathbf{x}_{n \times 1}$. Intuitively, α_i (which is called the accuracy of \mathcal{P}_i) stands for the normalized block-sparsity level of \mathbf{x} restricted to \mathcal{P}_i . We call \mathbf{x} a non-uniform block-sparse model with parameters $\{\alpha_i\}_{i=1}^L$ and $\{\rho_i\}_{i=1}^L$; if α_i 's are all equal, the non-uniform model $\mathbf{P}_{1,2,\mathbf{w}}^\eta$ reduces to the uniform model $\mathbf{P}_{1,2}^\eta$. Each set \mathcal{P}_i is associated with a weight $\omega_i \geq 0$ and the resulting weight \mathbf{w} in $\mathbf{P}_{1,2,\mathbf{w}}^\eta$ is

$$\mathbf{w} = \mathbf{D}\boldsymbol{\omega}, \quad (18)$$

where $\mathbf{D} := [\mathbf{1}_{\mathcal{P}_1}, \dots, \mathbf{1}_{\mathcal{P}_L}] \in \mathbb{R}^{q \times L}$. To better distinguish between ω_i and w_i , note that the former penalizes the index set \mathcal{P}_i while the latter penalizes the index set \mathcal{V}_i (see the illustrations in the left image of Figure 1).

In this work, the following questions are answered about a non-uniform block-sparse model:

- 1) How many measurements are required for $\mathbf{P}_{1,2}^0$ and $\mathbf{P}_{1,2,\mathbf{w}}^0$ to successfully recover a s -block-sparse vector from independent linear measurements?
- 2) Given extra prior information, what is the optimal choice of weights in $\mathbf{P}_{1,2,\mathbf{w}}^\eta$?
- 3) How close one can get to the optimal weights if the prior information is slightly inaccurate?

In the remainder of this work, we provide the answer to these questions in three sections.

IV. NUMBER OF MEASUREMENTS FOR SUCCESSFUL RECOVERY

For a fixed tolerance $\zeta \in [0, 1]$, let us denote the normalized number of measurements required for $P_{1,2}^0$ and $P_{1,2,w}^0$ to recover a s -block-sparse vector (with probability $1 - \zeta$) by $m_{q,s}$ and $m_{q,s,w}$, respectively:

$$m_{q,s} := \frac{\delta(\mathcal{D}(\|\cdot\|_{1,2}, \mathbf{x}))}{q}, \quad m_{q,s,w} := \frac{\delta(\mathcal{D}(\|\cdot\|_{1,2,w}, \mathbf{x}))}{q}. \quad (19)$$

Below, we obtain an upper-bound for the number of measurements required for $P_{1,2}^0$ to succeed with probability $1 - \zeta$.

Lemma 1. *Let $\mathbf{x} \in \mathbb{R}^n$ be a non-uniform s -block-sparse vector in \mathbb{R}^n with parameters $\{\rho_i\}_{i=1}^L$ and $\{\alpha_i\}_{i=1}^L$. Then,*

$$m_{q,s,w} \leq \widehat{m}_{q,s,w} \quad (20)$$

for

$$\widehat{m}_{q,s,w} = \inf_{t \geq 0} \Psi_{t,w}(\sigma, \boldsymbol{\rho}, \boldsymbol{\alpha}), \quad (21)$$

where

$$\begin{aligned} \Psi_{t,w}(\sigma, \boldsymbol{\rho}, \boldsymbol{\alpha}) &= \sum_{i=1}^L \rho_i \left(\alpha_i (k + t^2 \omega_i^2) + \frac{(1-\alpha_i)\phi(t\omega_i)}{2^{\frac{k}{2}-1}\Gamma(\frac{k}{2})} \right), \\ \phi(z) &:= \int_z^\infty (u-z)^2 u^{k-1} \exp(-\frac{u^2}{2}) du, \quad k = \frac{n}{q}. \end{aligned} \quad (22)$$

Proof. See Appendix B.

Corollary 1. *By considering $\mathbf{w} = \mathbf{1} \in \mathbb{R}^q$, and using the fact that the normalized block-sparsity level is $\sigma := \frac{\|\mathbf{x}\|_{0,2}}{q} = \sum_{i=1}^L \rho_i \alpha_i$, we reach an upper-bound $\widehat{m}_{q,s}$ for $m_{q,s}$ as*

$$\widehat{m}_{q,s} = \inf_{t \geq 0} \Psi_t(\sigma), \quad (23)$$

where

$$\Psi_t(\sigma) = \sigma(k + t^2) + \frac{(1-\sigma)\phi(t)}{2^{\frac{k}{2}-1}\Gamma(\frac{k}{2})}.$$

Remark. (Prior work) In [8, Lemma 3.2], the same expression as for $\widehat{m}_{q,s}$ is obtained for the normalized minimax MSE of the denoising problem

$$\min_{\mathbf{z} \in \mathbb{R}^n} \tau \|\mathbf{z}\|_{1,2} + \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|_2^2, \quad (24)$$

where $\mathbf{y} = \mathbf{x} + \mathbf{e}$ is the observed noisy vector. It is further conjectured in [8] that this value is equal to the number of measurements required by $P_{1,2}^0$ in the asymptotic regime. We show that this formula describes the required number of measurements in $P_{1,2}^0$ even in the non-asymptotic case (Proposition 2).

In the following Proposition, we demonstrate that the proposed upper-bound in Lemma 1 is asymptotically tight. We use this fact for the optimality of the obtained weights.

Proposition 2. *The normalized number of linear measurements required for $P_{1,2,w}^0$ and $P_{1,2}^0$ to successfully recover a non-uniform s -block-sparse vector in \mathbb{R}^n (i.e. $m_{q,s,w}$ and $m_{q,s}$, respectively) satisfy the following error bounds:*

$$\widehat{m}_{q,s,w} - \frac{2}{\sqrt{qL}} \leq m_{q,s,w} \leq \widehat{m}_{q,s,w}, \quad (25)$$

$$\widehat{m}_{q,s} - \frac{2}{\sqrt{s}q} \leq m_{q,s} \leq \widehat{m}_{q,s}. \quad (26)$$

Proof. See Appendix D.

It is interesting that the error bound in (26) is a special case of the error bound of Proposition (25) where one has s sets of blocks with size $\frac{q}{L}$ that each contributes to the block support with probability $\frac{L}{q}$.

V. OPTIMAL WEIGHTS

Our strategy in finding the optimal weights is to minimize the reconstruction error (15) in the noisy case (under a fixed number of measurements i.e. m), and the required number of measurements in the noiseless case. Based on Theorem 1, both of these objectives lead to the same optimization problem

$$\boldsymbol{\omega}^* = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}_+^L} m_{q,s,D\boldsymbol{\omega}} \in \mathbb{R}_+^L. \quad (27)$$

Instead of the latter minimization, we minimize the upper and lower bounds of statistical dimension (i.e. (25)), simultaneously. So

$$\inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \widehat{m}_{q,s,D\boldsymbol{\omega}} - \frac{2}{\sqrt{qL}} \leq \inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} m_{q,s,D\boldsymbol{\omega}} \leq \inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \widehat{m}_{q,s,D\boldsymbol{\omega}}, \quad (28)$$

where $\mathbf{D} := [\mathbf{1}_{\mathcal{P}_1}, \dots, \mathbf{1}_{\mathcal{P}_L}]_{q \times L}$. In the weighted block sparsity optimization, we call the weight

$$\boldsymbol{\omega}^* = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \widehat{m}_{q,s,D\boldsymbol{\omega}} \in \mathbb{R}_+^L \quad (29)$$

optimal since it asymptotically (as $q \rightarrow \infty$) minimizes simultaneously the number of measurements required for $P_{1,2,w}^0$ to succeed, and the reconstruction error of $P_{1,2,w}^\eta$. In the following lemma, the uniqueness of the optimal weights is shown by proving that $\delta(\mathcal{D}(\|\cdot\|_{1,2,D\boldsymbol{\omega}}, \mathbf{x}))$ is a strictly convex function of $\boldsymbol{\omega} \in \mathbb{R}_+^L$.

Lemma 2. *Assume $\mathcal{C} := \partial \|\cdot\|_{1,2}(\mathbf{x})$ does not contain the origin. We know that \mathcal{C} is compact and $1 \leq \|\mathbf{z}\|_2 \leq \sqrt{q}$ for all $\mathbf{z} \in \mathcal{C}$. Also let $\mathbf{g} \in \mathbb{R}^n$ be a standard normal vector. Consider the function*

$$\begin{aligned} J(\boldsymbol{\nu}) &:= \mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, \boldsymbol{\nu} \odot \mathcal{C}) = \mathbb{E}_{\mathbf{g}} [J_{\mathbf{g}}(\boldsymbol{\nu})] \\ &\text{with } \boldsymbol{\nu} = \mathbf{D}_b \mathbf{D} \boldsymbol{\nu} \in \mathbb{R}^n \text{ for } \boldsymbol{\nu} \in \mathbb{R}_+^L, \end{aligned} \quad (30)$$

where $\mathbf{D}_b := [\mathbf{1}_{\mathcal{V}_1}, \dots, \mathbf{1}_{\mathcal{V}_q}]_{n \times q}$. Then, the function J is strictly convex on $\boldsymbol{\nu} \in \mathbb{R}_{++}^L$ and J has a unique minimizer.

Proof. See Appendix E.

An analytic expression for the optimal weights is given in the following proposition via solving (29).

Proposition 3. *Let \mathbf{x} be a non-uniform s -block-sparse vector in \mathbb{R}^n with parameters $\{\rho_i\}_{i=1}^L$ and $\{\alpha_i\}_{i=1}^L$. Then, there exist unique optimal weights $\boldsymbol{\omega}^* \in \mathbb{R}_+^L$ (up to a positive scaling) that minimize $\widehat{m}_{q,s,D\boldsymbol{\omega}}$. The optimal weights $\boldsymbol{\omega}^* \in \mathbb{R}_+^L$ are obtained via the following integral equations:*

$$\begin{aligned} \alpha_i \omega_i^* &= \frac{1}{2^{\frac{k}{2}-1}\Gamma(\frac{k}{2})} (1 - \alpha_i) \int_{\omega_i^*}^\infty (u - \omega_i^*) u^{k-1} e^{-\frac{u^2}{2}} du \\ i &= 1, \dots, L. \end{aligned} \quad (31)$$

Proof. See Appendix C

The optimal weights in (31) depend only on the accuracy

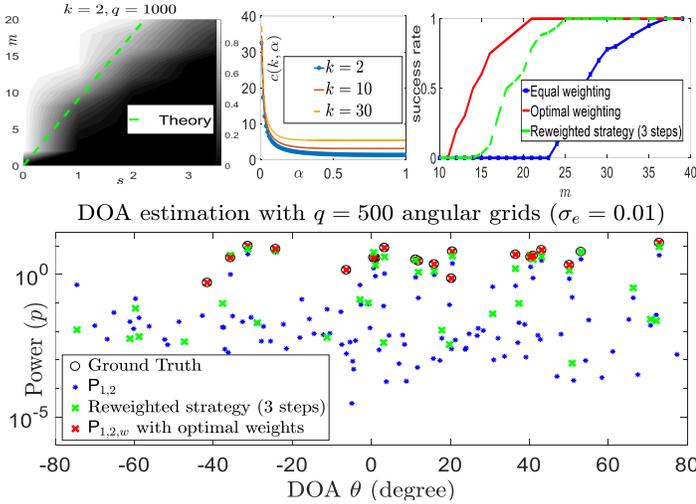


Figure 2: Top left image: This plot shows the empirical probability that $\mathbf{P}_{1,2}^0$ recovers $\mathbf{x} \in \mathbb{R}^{2000}$ with s blocks with nonzero ℓ_2 norm from m Gaussian linear measurements. The green dashed line shows the number of measurements obtained by Corollary 1. Top middle image: This plot shows $c(k, \alpha)$ in (34) versus the accuracy α for $k = 2, 10, 30$. Indeed, it demonstrates the sensitivity of the optimal weight against small perturbations of α . Top right image: This plot shows the probability that $\mathbf{P}_{1,2,w}^0$ succeed to recover $\mathbf{x} \in \mathbb{R}^{1280}$ from Gaussian linear measurements when \mathbf{w} is chosen equally, optimally and via a reweighted strategy (3 steps) proposed in [20]. The parameters we used in this figure, are: $q = 128$, $\sigma = \rho_1 = \frac{10}{128}$, $\alpha_1 = \frac{9}{10}$, $\rho_2 = \frac{118}{128}$, $\alpha_2 = \frac{1}{118}$. The optimal weights obtained via (31) with the aforementioned parameters are $\omega^* = [0.0766, 1]^T$. Bottom image: DOA estimation of $s = 20$ sources in the angular half-space $[-\frac{\pi}{2}, \frac{\pi}{2}]$ with $k = 10$ snapshots, $m = 24$, $q = 500$ and $\frac{d}{\lambda} = 1$. The measurements are contaminated with additive Gaussian noise with $\sigma_e = 0.01$. Also, we consider four angular bands with accuracies $\alpha_1 = \frac{5}{6}$, $\alpha_2 = \frac{13}{14}$, $\alpha_3 = \frac{2}{3}$ and $\alpha_4 = 0$ as prior information.

of \mathcal{P}_i s, i.e., $\{\alpha_i\}_{i=1}^L$ and not their relative size $\{\rho_i\}_{i=1}^L$. Finally, we prove the significant fact that with optimal weights, $\mathbf{P}_{1,2,D\omega^*}^0$ acts as if $\mathbf{x}_{\mathcal{P}_i}$'s are separately recovered via $\mathbf{P}_{1,2}^0$, in the sense of the required number of measurements.

Theorem 3. *Let \mathbf{x} be a non-uniform s -block-sparse vector in \mathbb{R}^n with parameters $\{\rho_i\}_{i=1}^L$ and $\{\alpha_i\}_{i=1}^L$. Then, the number of measurements required for $\mathbf{P}_{1,2,D\omega^*}^0$ is exactly equals the total number of measurements required for $\mathbf{P}_{1,2}^0$ to recover each $\{\mathbf{x}_{\mathcal{P}_i} \in \mathbb{R}^n\}_{i=1}^L$ separately, up to an asymptotically additive vanishing error term i.e.*

$$-\frac{2}{\sqrt{qL}} \leq m_{q,s,D\omega^*} - \sum_{i=1}^L m_{q,\|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}} \leq \frac{2}{\sqrt{q}} \sum_{i=1}^L (\|\mathbf{x}_{\mathcal{P}_i}\|_{0,2})^{-\frac{1}{2}}. \quad (32)$$

VI. ROBUSTNESS ANALYSIS

In this section, we evaluate the robustness of optimal weights if prior information is slightly inaccurate. As stated in the below proposition, under the condition $\alpha \gtrsim \frac{1}{10}$, the optimal weights ω^* in (31) are robust to inaccuracies in the prior information.

Proposition 4. *Let α be the accuracy of a set \mathcal{P} which inaccurately assumed α' in practice. Let ω and ω' be the*

optimal weights, obtained from (31), corresponding to α and α' , respectively. Then, there exists a constant $c(k, \alpha)$ such that

$$|\omega - \omega'| \leq c(k, \alpha) |\alpha - \alpha'| \quad (33)$$

for

$$c(k, \alpha) := \frac{\left(\sqrt{2}h(\alpha) \left(\Gamma(\frac{k}{2}) - \gamma(\frac{k}{2}, \frac{h(\alpha)^2}{2}) \right) + 2\gamma(\frac{k}{2}, \frac{h(\alpha)^2}{2}) \right)^2}{2\sqrt{2}\Gamma(\frac{k}{2})\gamma(\frac{k}{2}, \frac{h(\alpha)^2}{2})}, \quad (34)$$

where $\gamma(a, z) := \int_z^\infty u^{a-1} e^{-u} du$ is incomplete gamma function and $h(\alpha)$ is the nonlinear function in (31) that relates α to the optimal weight ω .

Proof. See Appendix G.

This proposition shows that our method of finding optimal weights in (31) is robust to inaccuracies in prior knowledge.

VII. SIMULATIONS

In this subsection, we numerically verify our theoretical results on the optimal weights. We have employed the CVX MATLAB package [23] to implement optimization problems. First, we investigate the required number of measurements for the scaling of successful recovery of $\mathbf{P}_{1,2}^0$ in terms of the block sparsity. For this purpose, we construct a s -block-sparse $\mathbf{x} \in \mathbb{R}^{2000}$. Then, we form the measurements $\mathbf{y}_{m \times 1} = \mathbf{A}\mathbf{x}$ and obtain an estimate $\hat{\mathbf{x}}$ by solving the problem $\mathbf{P}_{1,2}^0$; we repeat this experiment for 100 realizations of \mathbf{A} . For each experiment, we declare success if $\|\mathbf{x} - \hat{\mathbf{x}}\|_2 \leq 10^{-4}$. The heatmap in the top left image of Figure 2 shows the empirical probability of success for this procedure (black=0%, white=100%) which is consistent with the theory obtained by (23).

In the second experiment, we set $s = 10$ and generate a block-sparse random vector $\mathbf{x} \in \mathbb{R}^{1280}$ with $q = 128$ blocks of equal size $k = 10$. The block support of \mathbf{x} is drawn uniformly at random and the values within each non-zero block are drawn from i.i.d. standard normal distribution. Then, we consider two sets $\mathcal{P}_1, \mathcal{P}_2$ with $\alpha_1 = \frac{9}{10}$, $\alpha_2 = \frac{1}{118}$ that partition the set of blocks $\{1, \dots, 128\}$. There are 100 Monte Carlo trials for each m where in each, we solve $\mathbf{P}_{1,2,D\omega}^0$ with optimal weights ω^* , equal weights and the reweighted strategy (3 steps) proposed in [20] and recover $\mathbf{x} \in \mathbb{R}^{1280}$ from m Gaussian linear measurements. Note that \mathbf{x} is kept fixed in Monte Carlo trials and only \mathbf{A} changes. Optimal weights ω^* are obtained from (31) by MATLAB function `fzero`. The top right image of Figure 2 shows that $\mathbf{P}_{1,2,D\omega}^0$ with optimal weights needs fewer measurements than $\mathbf{P}_{1,2}^0$ and the reweighted method. To better investigate the sensitivity of optimal weights against small perturbations of α , we have plotted the robustness constant $c(k, \alpha)$ (see (34)) in the top middle image of Figure 2. These curves confirm the stability of optimal weights against the inaccuracy of prior information in the range $\alpha \gtrsim \frac{1}{10}$.

In the last experiment, we investigate DOA estimation. We consider $s = 20$ sources. As shown in the right block of Figure 1, these sources are scattered in the angular half-space $[-\frac{\pi}{2}, \frac{\pi}{2}]$ which is divided into q angular grids. In practical scenarios, it is reasonable for an engineer to know the

likelihood of source appearance in some known angular bands (α_i in our non-uniform block-sparse model). This information can be obtained for instance by considering the statistics of previous estimations or the characteristics of the underlying model. In this experiment, we consider four angular bands with accuracies $\alpha_1 = \frac{5}{6}$, $\alpha_2 = \frac{13}{14}$, $\alpha_3 = \frac{2}{3}$ and $\alpha_4 = 0$. We have implemented $P_{1,2,w}^\eta$ using CVX software when w is chosen equally, optimally and via the reweighted strategy proposed in [20]. The optimal weights are calculated using (31). We follow the model (7) with $d = \lambda$, $q = 500$, $\sigma_e = 0.01$, and $m = 24$. As it turns out from bottom image of Figure 2, $P_{1,2,w^*}^\eta$ exactly recovers the direction and the power of the sources in contrast to $P_{1,2}^\eta$ and the reweighted strategy which both fail. This in turn reveals the fact that, with optimal weighting strategy, much less sensors are required for stable recovery.

VIII. CONCLUSION

In this paper, we studied the recovery of block-sparse vectors using a weighted $\ell_{1,2}$ -minimization, when some prior information about the distribution of the block-support is available. We have particularly studied the case where partial block-sparsity level of the vector within a block partitioning is available as prior information. Our goal was to find the optimal weights so as to minimize the number of required measurements for perfect recovery in the noiseless case, and to minimize the reconstruction error in the noisy case. In both cases, we simplified the task to minimizing the statistical dimension of a weighted descent cone. We introduced closed-form expressions that identify the unique optimal weights. We have also shown the robustness of the optimal weights against inaccuracies of the prior information.

APPENDIX

A. Proof of Proposition 1

Assume first that $\mathbf{x} \neq \mathbf{0}$. By setting $\mathbf{y} = \mathbf{0}$ in (8), we have: $\langle \mathbf{z}, \mathbf{x} \rangle \geq f(\mathbf{x}) - f(\mathbf{0})$ and $f(\mathbf{0}) = 0$ due to the homogeneity, $f^*(\mathbf{z}) \geq 1$. Also, setting $\mathbf{y} = \mathbf{v} + \mathbf{x}$ in (8) under condition $f(\mathbf{v}) = 1$ we have:

$$f^*(\mathbf{z}) = \sup_{f(\mathbf{v})=1} \langle \mathbf{z}, \mathbf{v} \rangle \leq \sup_{f(\mathbf{v})=1} (f(\mathbf{v} + \mathbf{x}) - f(\mathbf{x})) \leq 1, \quad (35)$$

where we used sub-additivity of f . Hence, we have $f^*(\mathbf{z}) = 1$ and subsequently $\langle \mathbf{z}, \mathbf{x} \rangle = f(\mathbf{x})$.

On the other hand, if we have \mathbf{z} such that $f^*(\mathbf{z}) = 1$ and $\langle \mathbf{z}, \mathbf{x} \rangle = f(\mathbf{x})$, then for each $\mathbf{y} \in \mathbb{R}^n$:

$$f(\mathbf{x}) + \langle \mathbf{z}, \mathbf{y} - \mathbf{x} \rangle = \langle \mathbf{z}, \mathbf{x} \rangle + \langle \mathbf{z}, \mathbf{y} - \mathbf{x} \rangle \leq f(\mathbf{y}). \quad (36)$$

For $\mathbf{x} = \mathbf{0}$ in (8), $\langle \mathbf{z}, \mathbf{y} \rangle \leq f(\mathbf{y})$ and then $f^*(\mathbf{z}) \leq 1$. Further, $f^*(\mathbf{z}) \leq 1$ implies that $\mathbf{z} \in \partial f(\mathbf{0})$.

B. Proof of Lemma 1

By the definition of statistical dimension for $\mathcal{D}(\|\cdot\|_{1,2,w}, \mathbf{x})$ in (12), the fact that infimum of an affine function is concave

and Jensen's inequality, we can find an upper-bound for $m_{q,s,w}$ as:

$$\overline{m}_{q,s,w} \leq \inf_{t \geq 0} \underbrace{q^{-1} \mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, t\partial\|\cdot\|_{1,2,w}(\mathbf{x}))}_{\Psi_{t,w}(\sigma, \rho, \alpha)} := \widehat{m}_{q,s,w}. \quad (37)$$

The next step is to calculate $\partial\|\cdot\|_{1,2,w}(\mathbf{x})$. From Proposition 1, we have:

$$\partial\|\cdot\|_{1,2,w}(\mathbf{x}) = \{\mathbf{z} \in \mathbb{R}^p : \langle \mathbf{z}, \mathbf{x} \rangle = \|\mathbf{x}\|_{1,2,w}, \|\mathbf{z}\|_{1,2,w}^* = 1\}. \quad (38)$$

It is not hard to show that,

$$\partial\|\cdot\|_{1,2,w}(\mathbf{x}) = \left\{ \mathbf{z} \in \mathbb{R}^n : \begin{array}{ll} \frac{w_b}{\|\mathbf{x}_{\mathcal{V}_b}\|_2} \mathbf{x}_{\mathcal{V}_b}, & b \in \mathcal{B} \\ \|\mathbf{z}_{\mathcal{V}_b}\|_2 \leq w_b, & b \in \overline{\mathcal{B}} \end{array} \right\}. \quad (39)$$

Now to calculate $\Psi_{t,w}(\sigma, \rho, \alpha)$, regarding (39), we compute the distance of the dilated subdifferential of descent cone of the $\ell_{1,2,w}$ norm at $\mathbf{x} \in \mathbb{R}^n$ from a standard Gaussian vector $\mathbf{g} \in \mathbb{R}^n$ which is given by:

$$\begin{aligned} \text{dist}^2(\mathbf{g}, t\partial\|\cdot\|_{1,2,w}(\mathbf{x})) &= \inf_{\mathbf{z} \in \partial\|\cdot\|_{1,2,w}(\mathbf{x})} \|\mathbf{g} - t\mathbf{z}\|_2^2 = \\ &= \sum_{b \in \mathcal{B}} \|\mathbf{g}_{\mathcal{V}_b} - tw_b \frac{\mathbf{x}_{\mathcal{V}_b}}{\|\mathbf{x}_{\mathcal{V}_b}\|_2}\|_2^2 + \sum_{b \in \overline{\mathcal{B}}} \inf_{\|\mathbf{z}_{\mathcal{V}_b}\|_2 \leq w_b} \|\mathbf{g}_{\mathcal{V}_b} - t\mathbf{z}_{\mathcal{V}_b}\|_2^2 = \\ &= \sum_{b \in \mathcal{B}} \|\mathbf{g}_{\mathcal{V}_b} - tw_b \frac{\mathbf{x}_{\mathcal{V}_b}}{\|\mathbf{x}_{\mathcal{V}_b}\|_2}\|_2^2 + \sum_{b \in \overline{\mathcal{B}}} (\|\mathbf{g}_{\mathcal{V}_b}\|_2 - tw_b)_+^2, \end{aligned} \quad (40)$$

where we used triangle inequality in the second part. By taking expectation from both sides, we reach:

$$\begin{aligned} \mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, t\partial\|\cdot\|_{1,2,w}(\mathbf{x})) &= \\ ks + \sum_{b \in \mathcal{B}} (tw_b)^2 + \sum_{b \in \overline{\mathcal{B}}} \mathbb{E}(\underbrace{\|\mathbf{g}_{\mathcal{V}_b}\|_2}_{\zeta} - tw_b)_+^2, \end{aligned} \quad (41)$$

where $k = \frac{n}{q}$. and $\zeta^2 := \|\mathbf{g}_{\mathcal{V}_b}\|_2^2$ is distributed as chi-squared with k degrees of freedom. Moreover,

$$\begin{aligned} \mathbb{E}(\zeta - tw_b)_+^2 &= 2 \int_0^\infty a \mathbb{P}(\zeta^2 \geq (a + tw_b)^2) da = \\ &= \frac{2}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} \int_0^\infty \int_{(tw_b+a)^2}^\infty a u^{\frac{k}{2}-1} e^{-\frac{u}{2}} du da = \\ &= \frac{2}{2^{\frac{k}{2}} \Gamma(\frac{k}{2})} \int_{(tw_b)^2}^\infty \int_0^{\sqrt{u}-tw_b} a u^{\frac{k}{2}-1} e^{-\frac{u}{2}} da du = \\ &= \frac{1}{2^{\frac{k}{2}-1} \Gamma(\frac{k}{2})} \int_{tw_b}^\infty (u - tw_b)^2 u^{k-1} e^{-\frac{u}{2}} du := \frac{\phi(tw_b)}{2^{\frac{k}{2}-1} \Gamma(\frac{k}{2})}, \end{aligned} \quad (42)$$

where in the third line, the order of integration is changed and in the fourth line, a change of variable is used. As a consequence, (41) becomes:

$$\mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, t\partial\|\cdot\|_{1,2,w}(\mathbf{x})) = ks + \sum_{b \in \mathcal{B}} (tw_b)^2 + \frac{\sum_{b \in \overline{\mathcal{B}}} \phi(tw_b)}{2^{\frac{k}{2}-1} \Gamma(\frac{k}{2})}. \quad (43)$$

By normalizing to the number of blocks q and incorporating block prior information using $\mathbf{w} = \mathbf{D}\boldsymbol{\omega} \in \mathbb{R}^q$ we reach

$$\begin{aligned} & \mathbb{E}_{\mathbf{g}} \text{dist}^2(\mathbf{g}, t\partial \|\cdot\|_{1,2,\mathbf{w}}(\mathbf{x})) = \\ & ks + \sum_{i=1}^L |\mathcal{P}_i \cap \mathcal{B}| t^2 \omega_i^2 + |\mathcal{P}_i \cap \bar{\mathcal{B}}| \frac{\phi(t\omega_i)}{2^{2^{-1}\Gamma(\frac{k}{2})}} \\ & = q \left(k\sigma + \sum_{i=1}^L \rho_i \left(\alpha_i t^2 \omega_i^2 + \frac{(1-\alpha_i)\phi(t\omega_i)}{2^{2^{-1}\Gamma(\frac{k}{2})}} \right) \right) \\ & = q \left(\sum_{i=1}^L \rho_i \left(\alpha_i (t^2 \omega_i^2 + k) + (1-\alpha_i) \frac{1}{2^{2^{-1}\Gamma(\frac{k}{2})}} \phi(t\omega_i) \right) \right), \end{aligned} \quad (44)$$

where in the last line above, we benefited the fact that $\sigma = \sum_{i=1}^L \rho_i \alpha_i$.

C. Proof of Proposition 3

Define $\mathcal{C} := \partial \|\cdot\|_{1,2}(\mathbf{x})$ and use Lemma 1 and 2 to obtain:

$$\inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \hat{m}_{q,s,\mathbf{D}\boldsymbol{\omega}} = \inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \inf_{t \in \mathbb{R}_+} \Psi_{t,\mathbf{D}\boldsymbol{\omega}}(\sigma, \boldsymbol{\rho}, \boldsymbol{\alpha}) = \inf_{\boldsymbol{\nu} \in \mathbb{R}_+^L} J_b(\boldsymbol{\nu}),$$

where $\Psi_{t,\mathbf{D}\boldsymbol{\omega}}(\sigma, \boldsymbol{\rho}, \boldsymbol{\alpha})$ is defined in (21). Also, we used a change of variable $\boldsymbol{\nu} = t\boldsymbol{\omega}$ to convert multivariate optimization problem to a single variable optimization problem. Thus, the function $J_b(\boldsymbol{\nu})$ is obtained via the following equation:

$$J_b(\boldsymbol{\nu}) = \sum_{i=1}^L \rho_i \left(\alpha_i (\nu(i)^2 + 1) + \frac{(1-\alpha_i)\phi(\nu(i))}{2^{2^{-1}\Gamma(\frac{k}{2})}} \right). \quad (45)$$

By considering Lemma 2 and $\mathbf{D}_t := [\mathbf{1}_{\mathcal{V}_1}, \dots, \mathbf{1}_{\mathcal{V}_q}]_{n \times q} [\mathbf{1}_{\mathcal{P}_1}, \dots, \mathbf{1}_{\mathcal{P}_L}]_{q \times L}$, the function $J_b(\boldsymbol{\nu})$ is continuous and strictly convex and thus the unique minimizer can be obtained using $\nabla J_b(\boldsymbol{\nu}) = \mathbf{0} \in \mathbb{R}^L$ which leads to

$$2\alpha_i \nu^*(i) + \frac{2(1-\alpha_i)\phi'(\nu^*(i))}{2^{2^{-1}\Gamma(\frac{k}{2})}} = 0 \quad : \quad i = 1, \dots, L. \quad (46)$$

D. Proof of Proposition 2

By the error bound (16) and (39), with $f(\mathbf{x}) = \|\mathbf{x}\|_{1,2,\mathbf{w}}$ and $\mathbf{w} = \sum_{i=1}^L \omega_i \mathbf{1}_{\mathcal{P}_i} \in \mathbb{R}^q$, the numerator of (16) is given by

$$\begin{aligned} 2 \sup_{\mathbf{s} \in \partial \|\cdot\|_{1,2,\mathbf{w}}(\mathbf{x})} \|\mathbf{s}\|_2 & \leq 2 \sqrt{\sum_{i=1}^q w_b^2} = \\ & 2 \sqrt{\sum_{i=1}^L |\mathcal{P}_i| \omega_i^2} = 2 \sqrt{\sum_{i=1}^L q \rho_i \omega_i^2}. \end{aligned}$$

Also, for the denominator we have:

$$\frac{\|\mathbf{x}\|_{1,2,\mathbf{w}}}{\|\mathbf{x}\|_2} \leq \sqrt{\sum_{i \in \mathcal{B}} w_i^2} = \sqrt{\sum_{i=1}^L |\mathcal{P}_i \cap \mathcal{B}| \omega_i^2} = \sqrt{\sum_{i=1}^L q \alpha_i \rho_i \omega_i^2}, \quad (47)$$

where the first inequality in (47) follows from Cauchy-Schwartz inequality. The error bound (16) for $\|\cdot\|_{1,2,\mathbf{w}}$ depends

only on $\mathcal{D}(\|\cdot\|_{1,2,\mathbf{w}}, \mathbf{x})$. Moreover, $\mathcal{D}(\|\cdot\|_{1,2,\mathbf{w}}, \mathbf{x})$ only requires that $\|\mathbf{x}_{\mathcal{V}_b}\|_2 = w_b \quad : \quad \forall b \in \mathcal{B}$. So, a vector

$$\mathbf{z} = \left\{ \begin{array}{l} \|\mathbf{z}_{\mathcal{V}_b}\|_2 = w_b, \quad b \in \mathcal{B} \\ 0, \quad b \in \bar{\mathcal{B}} \end{array} \right\} \in \mathbb{R}^n$$

can be chosen to satisfy equality in (47). Therefore, the error of obtaining the upper-bound of $m_{q,s,\mathbf{w}}$, i.e. $\hat{m}_{q,s,\mathbf{w}}$ is

$$\frac{2\sqrt{\sum_{i=1}^L q \rho_i \omega_i^2}}{q\sqrt{\sum_{i=1}^L q \alpha_i \rho_i \omega_i^2}} \leq \frac{2}{q} \sqrt{\frac{1}{\min_{i \in [q]} \frac{|\mathcal{P}_i \cap \mathcal{B}|}{|\mathcal{P}_i|}}} \leq \frac{2}{\sqrt{qL}}, \quad (48)$$

in which the last inequality follows from the facts that $|\mathcal{P}_i \cap \mathcal{B}| \geq 1$, $|\mathcal{P}_i| \leq \frac{q}{L}$ for at least one $i \in [q]$ and thus $\min_{i \in [q]} \frac{|\mathcal{P}_i \cap \mathcal{B}|}{|\mathcal{P}_i|} \geq \frac{L}{q}$. Further, the error of $\hat{m}_{q,s,\mathbf{w}}$ from $m_{q,s,\mathbf{w}}$ is at most $\frac{2}{\sqrt{qL}}$. For the case of $\ell_{1,2}$, by using the facts $\omega_i = 1 \quad \forall i$ and $s = q \sum_{i=1}^L \rho_i \alpha_i$, it is straightforward to verify that $\frac{2\sqrt{\sum_{i=1}^L q \rho_i \omega_i^2}}{q\sqrt{\sum_{i=1}^L q \alpha_i \rho_i \omega_i^2}} = \frac{2}{\sqrt{qs}}$.

E. Proof of Lemma 2

Convexity. Let $\boldsymbol{\nu}, \tilde{\boldsymbol{\nu}} \in \mathbb{R}_+^L$ and $\theta \in [0, 1]$ with $\mathbf{v} = \mathbf{D}_t \boldsymbol{\nu}$ and $\tilde{\mathbf{v}} = \mathbf{D}_t \tilde{\boldsymbol{\nu}}$. Then we have:

$$\begin{aligned} \forall \epsilon, \tilde{\epsilon} > 0 \exists \mathbf{z}, \tilde{\mathbf{z}} \in \mathcal{C} \quad \text{such that} \\ \|\mathbf{g} - \mathbf{v} \odot \mathbf{z}\|_2 & \leq \text{dist}(\mathbf{g}, \mathbf{v} \odot \mathcal{C}) + \epsilon, \\ \|\mathbf{g} - \tilde{\mathbf{v}} \odot \tilde{\mathbf{z}}\|_2 & \leq \text{dist}(\mathbf{g}, \tilde{\mathbf{v}} \odot \mathcal{C}) + \tilde{\epsilon}. \end{aligned} \quad (49)$$

Since otherwise we have:

$$\begin{aligned} \forall \mathbf{z}, \tilde{\mathbf{z}} \in \mathcal{C} : \|\mathbf{g} - \mathbf{v} \odot \mathbf{z}\|_2 & > \text{dist}(\mathbf{g}, \mathbf{v} \odot \mathcal{C}) + \epsilon, \\ \|\mathbf{g} - \tilde{\mathbf{v}} \odot \tilde{\mathbf{z}}\|_2 & > \text{dist}(\mathbf{g}, \tilde{\mathbf{v}} \odot \mathcal{C}) + \tilde{\epsilon}. \end{aligned} \quad (50)$$

By taking the infimum over $\mathbf{z}, \tilde{\mathbf{z}} \in \mathcal{C}$, we reach a contradiction. We proceed to show convexity of $\text{dist}(\mathbf{g}, (\mathbf{D}_t \boldsymbol{\nu}) \odot \mathcal{C})$ by writing

$$\begin{aligned} \text{dist}(\mathbf{g}, (\theta \mathbf{v} + (1-\theta)\tilde{\mathbf{v}}) \odot \mathcal{C}) & = \\ \inf_{\mathbf{z} \in \mathcal{C}} \|\mathbf{g} - (\theta \mathbf{v} + (1-\theta)\tilde{\mathbf{v}}) \odot \mathbf{z}\|_2 & \\ \leq \inf_{\mathbf{z}_1 \in \mathcal{C}, \mathbf{z}_2 \in \mathcal{C}} \|\mathbf{g} - \theta \mathbf{v} \odot \mathbf{z}_1 - (1-\theta)\tilde{\mathbf{v}} \odot \mathbf{z}_2\|_2 & \leq \\ \theta \|\mathbf{g} - \mathbf{v} \odot \mathbf{z}_1\|_2 + (1-\theta) \|\mathbf{g} - \tilde{\mathbf{v}} \odot \mathbf{z}_2\|_2 & \leq \\ \theta \text{dist}(\mathbf{g}, \mathbf{v} \odot \mathcal{C}) + (1-\theta) \text{dist}(\mathbf{g}, \tilde{\mathbf{v}} \odot \mathcal{C}) + \epsilon + \tilde{\epsilon}. \end{aligned} \quad (51)$$

Since this holds for any ϵ and $\tilde{\epsilon}$, $\text{dist}(\mathbf{g}, (\mathbf{D}_t \boldsymbol{\nu}) \odot \mathcal{C})$ is a convex function. As the square of a non-negative convex function is convex, $J_g(\boldsymbol{\nu})$ is a convex function. At last, the function $J(\boldsymbol{\nu})$ is the average of convex functions, hence is convex. In (51), the first inequality comes from the fact that $\forall \mathbf{z}_1, \mathbf{z}_2 \in \mathcal{C} \quad \exists \mathbf{z} \in \mathcal{C}$:

$$\begin{aligned} \theta \mathbf{v} \odot \mathbf{z}_1 + (1-\theta)\tilde{\mathbf{v}} \odot \mathbf{z}_2 & = \\ \left\{ \begin{array}{l} \mathbf{y}_{\mathcal{V}_b} = (\theta \mathbf{v}_{\mathcal{V}_b} + (1-\theta)\tilde{\mathbf{v}}_{\mathcal{V}_b}) \odot \frac{\mathbf{x}_{\mathcal{V}_b}}{\|\mathbf{x}_{\mathcal{V}_b}\|_2}, \quad b \in \mathcal{B} \\ \mathbf{y}_{n \times 1} : \|\mathbf{y}_{\mathcal{V}_b}\|_2 \leq \theta \|\mathbf{v}\|_{\infty} \|\mathbf{z}_1\|_2 \\ \quad + (1-\theta) \|\tilde{\mathbf{v}}\|_{\infty} \|\mathbf{z}_2\|_2, \quad b \in \bar{\mathcal{B}} \end{array} \right\} \\ \in \left\{ \begin{array}{l} \mathbf{y}_{\mathcal{V}_b} = (\theta \mathbf{v}_{\mathcal{V}_b} + (1-\theta)\tilde{\mathbf{v}}_{\mathcal{V}_b}) \odot \frac{\mathbf{x}_{\mathcal{V}_b}}{\|\mathbf{x}_{\mathcal{V}_b}\|_2}, \quad b \in \mathcal{B} \\ \mathbf{y}_{n \times 1} : \|\mathbf{y}_{\mathcal{V}_b}\|_2 \leq \\ \quad (\theta \|\mathbf{v}\|_{\infty} + (1-\theta) \|\tilde{\mathbf{v}}\|_{\infty}) \|\mathbf{z}_{\mathcal{V}_b}\|_2, \quad b \in \bar{\mathcal{B}} \end{array} \right\} \\ = (\theta \mathbf{v} + (1-\theta)\tilde{\mathbf{v}}) \odot \mathbf{z}. \end{aligned} \quad (52)$$

To verify (52), we argue by contradiction:

$$\forall \mathbf{z} \in \mathcal{C} \exists d \in \overline{\mathcal{B}} : (\theta \|\mathbf{v}\|_\infty + (1-\theta) \|\tilde{\mathbf{v}}\|_\infty) \|\mathbf{z}_{\mathcal{V}_d}\|_2 < \theta \|\mathbf{v}\|_\infty \|\mathbf{z}_{1_{\mathcal{V}_d}}\|_2 + (1-\theta) \|\tilde{\mathbf{v}}\|_\infty \|\mathbf{z}_{2_{\mathcal{V}_d}}\|_2 \leq \theta \|\mathbf{v}\|_\infty + (1-\theta) \|\tilde{\mathbf{v}}\|_\infty. \quad (53)$$

Then, by taking $\mathbf{z}_{\mathcal{V}_d} = \mathbf{e}_i \in \mathbb{R}^k$ for some $i \in [k]$, we reach a contradiction. In the second inequality in (51), we used triangle inequality of norms. The third inequality uses (49).

Strict convexity. We show strict convexity by contradiction. If $J(\boldsymbol{\nu})$ was not strictly convex, there would be vectors $\boldsymbol{\nu}, \tilde{\boldsymbol{\nu}} \in \mathbb{R}_+^L$ with $\mathbf{v} = \mathbf{D}_t \boldsymbol{\nu}, \tilde{\mathbf{v}} = \mathbf{D}_t \tilde{\boldsymbol{\nu}}$ and $\theta \in (0, 1)$ such that

$$\mathbb{E}[J_g(\theta \boldsymbol{\nu} + (1-\theta)\tilde{\boldsymbol{\nu}})] = \mathbb{E}[\theta J_g(\boldsymbol{\nu}) + (1-\theta)J_g(\tilde{\boldsymbol{\nu}})]. \quad (54)$$

For each \mathbf{g} in (54) the left-hand side is smaller than or equal to the right-hand side. Therefore, in (54), $J_g(\theta \boldsymbol{\nu} + (1-\theta)\tilde{\boldsymbol{\nu}})$ and $\theta J_g(\boldsymbol{\nu}) + (1-\theta)J_g(\tilde{\boldsymbol{\nu}})$ are almost surely equal (except at a set of measure zero) with respect to Gaussian measure. Moreover, we have

$$\begin{aligned} J_0(\theta \boldsymbol{\nu} + (1-\theta)\tilde{\boldsymbol{\nu}}) &= \text{dist}^2(\mathbf{0}, (\theta \mathbf{v} + (1-\theta)\tilde{\mathbf{v}}) \odot \mathcal{C}) \leq \\ &\inf_{\mathbf{z}_1, \mathbf{z}_2 \in \mathcal{C}} \|\theta \mathbf{v} \odot \mathbf{z}_1 + (1-\theta)\tilde{\mathbf{v}} \odot \mathbf{z}_2\|_2^2 < \theta \inf_{\mathbf{z}_1 \in \mathcal{C}} \|\mathbf{v} \odot \mathbf{z}_1\|_2^2 + \\ &(1-\theta) \inf_{\mathbf{z}_2 \in \mathcal{C}} \|\tilde{\mathbf{v}} \odot \mathbf{z}_2\|_2^2 = \theta J_0(\boldsymbol{\nu}) + (1-\theta)J_0(\tilde{\boldsymbol{\nu}}), \end{aligned} \quad (55)$$

where the first inequality comes from (52) and the second inequality stems from the strict convexity of $\|\cdot\|_2^2$. From (52), it is easy to verify that the set $\boldsymbol{\nu} \odot \mathcal{C}$ is a convex set. The distance to a convex set e.g. \mathcal{E} i.e. $\text{dist}(\mathbf{g}, \mathcal{E})$ is a 1-lipschitz function (i.e. $|\text{dist}(\mathbf{g}, \mathcal{E}) - \text{dist}(\tilde{\mathbf{g}}, \mathcal{E})| \leq \|\mathbf{g} - \tilde{\mathbf{g}}\|_2 : \forall \mathbf{g}, \tilde{\mathbf{g}} \in \mathbb{R}^n$) and hence continuous with respect to \mathbf{g} . Therefore, $J_g(\boldsymbol{\nu})$ is continuous with respect to \mathbf{g} . So there exist an open ball around $\mathbf{g} = \mathbf{0} \in \mathbb{R}^n$ that similar to (55), we may write the following relation for some $\epsilon > 0$.

$$\exists \mathbf{u} \in \mathbb{B}_\epsilon^n : J_u(\theta \boldsymbol{\nu} + (1-\theta)\tilde{\boldsymbol{\nu}}) < \theta J_u(\boldsymbol{\nu}) + (1-\theta)J_u(\tilde{\boldsymbol{\nu}}). \quad (56)$$

The above statement contradicts with (54) and hence we have strict convexity. Continuity along with convexity of J implies that J is convex on the whole domain $\boldsymbol{\nu} \in \mathbb{R}_+^L$.

Differentiability and continuity. The function $J_g(\boldsymbol{\nu})$ is continuously differentiable and the gradient for $\boldsymbol{\nu} \in \mathbb{R}_+^L$ is

$$\begin{aligned} \nabla_{\boldsymbol{\nu}} J_g(\boldsymbol{\nu}) &= \\ \frac{\partial J_g(\boldsymbol{\nu})}{\partial \boldsymbol{\nu}} &= -2\mathbf{D}_t^T (\mathbf{D}_t \boldsymbol{\nu})^{\odot(-1)} \odot \mathcal{P}_{(\mathbf{D}_t \boldsymbol{\nu}) \odot \mathcal{C}}(\mathbf{g}) \\ &\odot (\mathbf{g} - \mathcal{P}_{\mathbf{D}_t \boldsymbol{\nu} \odot \mathcal{C}}(\mathbf{g})). \end{aligned} \quad (57)$$

Continuity of $\frac{\partial J_g(\boldsymbol{\nu})}{\partial \boldsymbol{\nu}}$ at $\boldsymbol{\nu} \in \mathbb{R}_+^L$ stems from the fact that the projection onto a convex set is continuous. For each compact set $\mathcal{I} \subseteq \mathbb{R}_+^L$ we have:

$$\begin{aligned} \mathbb{E} \sup_{\boldsymbol{\nu} \in \mathcal{I}} \|\nabla_{\boldsymbol{\nu}} J_g(\boldsymbol{\nu})\|_2 &\leq \\ 2\|\mathbf{D}_t\|_{2 \rightarrow 2} \sqrt{q}(\sqrt{n} + 2\sqrt{q}(\sup_{\boldsymbol{\nu} \in \mathcal{I}} \nu_{\max})) &< \infty, \end{aligned} \quad (58)$$

where $\nu_{\max} := \max_{i \in [L]} \nu(i)$. Therefore, we have

$$\nabla_{\boldsymbol{\nu}} J(\boldsymbol{\nu}) = \left(\frac{\partial}{\partial \boldsymbol{\nu}} \right) \mathbb{E} J_g(\boldsymbol{\nu}) = \mathbb{E}[\nabla_{\boldsymbol{\nu}} J_g(\boldsymbol{\nu})] : \forall \boldsymbol{\nu} \in \mathbb{R}_+^L, \quad (59)$$

where in the last equality, we used the Lebesgue's dominated convergence theorem. Also, continuity of $J(\boldsymbol{\nu})$ can be concluded from the continuity of its gradient.

Attainment of the minimum. Suppose that $\boldsymbol{\nu} \geq \|\mathbf{g}\|_2 \mathbf{1}_{L \times 1}$. With this assumption, we may write:

$$\begin{aligned} \text{dist}(\mathbf{g}, (\mathbf{D}_t \boldsymbol{\nu}) \odot \mathcal{C}) &= \inf_{\mathbf{z} \in \mathcal{C}} \|\mathbf{g} - \mathbf{v} \odot \mathbf{z}\|_2 \geq \\ \inf_{\mathbf{z} \in \mathcal{C}} (\|\mathbf{v} \odot \mathbf{z}\|_2 - \|\mathbf{g}\|_2) &\geq \nu_{\min} - \|\mathbf{g}\|_2 \geq 0, \end{aligned} \quad (60)$$

where in (60), $\nu_{\min} := \min_{i \in [L]} \nu(i)$. By squaring (60), we reach

$$J_g(\boldsymbol{\nu}) \geq (\nu_{\min} - \|\mathbf{g}\|_2)^2 : \forall \boldsymbol{\nu} > \|\mathbf{g}\|_2 \mathbf{1}_{L \times 1} \quad (61)$$

Using $\mathbb{E}\|\mathbf{g}\|_2 \geq \frac{n}{\sqrt{n+1}}$ [24, Proposition 8.1] and Marcov's inequality we obtain:

$$\mathbb{P}(\|\mathbf{g}\|_2 \leq \sqrt{n}) \geq 1 - \sqrt{\frac{n}{n+1}}.$$

Then we get:

$$\begin{aligned} J(\boldsymbol{\nu}) &\geq \mathbb{E}[J_g(\boldsymbol{\nu}) \|\mathbf{g}\|_2 \leq \sqrt{n}] \mathbb{P}(\|\mathbf{g}\|_2 \leq \sqrt{n}) \\ &\geq (1 - \sqrt{\frac{n}{n+1}}) \mathbb{E}[(\nu_{\min} - \|\mathbf{g}\|_2)^2 \|\mathbf{g}\|_2 \leq \sqrt{n}] \\ &\geq (1 - \sqrt{\frac{n}{n+1}}) (\nu_{\min} - \sqrt{n})^2, \end{aligned} \quad (62)$$

where the first inequality stems from total probability theorem, the second inequality follows from (61). From (62), we find out that $J(\boldsymbol{\nu}) > J(\mathbf{0})$ when $\boldsymbol{\nu} > (2^{\frac{1}{4}} + 1)\sqrt{n} \mathbf{1}_{L \times 1}$. Therefore, the unique minimizer of the function J must occur in the interval $[\mathbf{0}, (2^{\frac{1}{4}} + 1)\sqrt{n} \mathbf{1}_{L \times 1}]$.

F. Proof of Theorem 3

Using optimal weights, the upper-bound for the normalized number of measurements required for $\mathbb{P}_{1,2, \mathbf{D}\boldsymbol{\omega}^*}^0$ to succeed is:

$$\begin{aligned} \hat{m}_{q,s, \boldsymbol{\omega}^*} &= \inf_{\boldsymbol{\omega} \in \mathbb{R}_+^L} \hat{m}_{q,s, \mathbf{D}\boldsymbol{\omega}} = \\ \sum_{i=1}^L \left[\inf_{\nu_i \in \mathbb{R}_+} \left(\underbrace{\frac{\|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}{q} (\nu_i^2 + k) + (1 - \frac{\|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}{q}) \phi_B(\nu_i)}_{\Psi_{\nu_i, \|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}(\frac{\|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}{q})} \right) \right] \\ &= \sum_{i=1}^L \hat{m}_{q, \|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}. \end{aligned} \quad (63)$$

The expression in the bracket is the upper-bound for normalized number of measurements required for successful recovery of $\mathbf{x}_{\mathcal{P}_i} \in \mathbb{R}^n$ using $\mathbb{P}_{1,2}^0$ i.e. $\hat{m}_{q, \|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}$. Thus, regarding the error bounds obtained in Proposition 2, the relation between $m_{q,s, \mathbf{D}\boldsymbol{\omega}^*}$ and $m_{q, \|\mathbf{x}_{\mathcal{P}_i}\|_{0,2}}$ is given by (32).

G. Proof of Proposition 4

We have that $\lim_{\alpha \rightarrow \alpha'} \frac{|\omega - \omega'|}{|\alpha - \alpha'|} = \left| \frac{\partial \omega}{\partial \alpha} \right|$. By differentiating (31), we reach

$$\frac{\partial \omega}{\partial \alpha} = - \frac{\omega + \frac{1}{2k/2-1\Gamma(k/2)} \int_{\omega}^{\infty} (u-\omega) u^{k-1} e^{-\frac{u^2}{2}} du}{\alpha + \frac{1-\alpha}{2k/2-1\Gamma(k/2)} \int_{\omega}^{\infty} u^{k-1} e^{-\frac{u^2}{2}} du}. \quad (64)$$

Using (31), the above equation reduces to:

$$\frac{\partial \omega}{\partial \alpha} = -\frac{\omega^2 2^{k/2-1} \Gamma(k/2)}{(1-\alpha)^2 \int_{\omega}^{\infty} u^k e^{-\frac{u^2}{2}} du}. \quad (65)$$

By obtaining α from (31) and replacing in (65), we reach:

$$\frac{\partial \omega}{\partial \alpha} = -\frac{\left(\int_{\omega}^{\infty} u^k e^{-\frac{u^2}{2}} du - \omega \int_{\omega}^{\infty} u^{k-1} e^{-\frac{u^2}{2}} + 2 \frac{k}{2} \omega^{-1} \Gamma\left(\frac{k}{2}\right) \omega \right)^2}{2^{\frac{k}{2}-1} \Gamma\left(\frac{k}{2}\right) \int_{\omega}^{\infty} u^k e^{-\frac{u^2}{2}} du} := f(\omega). \quad (66)$$

After some simplification, $f(\omega)$ reduces to

$$f(\omega) = \frac{\left(\sqrt{2} \omega \left(\Gamma\left(\frac{k}{2}\right) - \gamma\left(\frac{k}{2}, \frac{\omega^2}{2}\right) \right) + 2 \gamma\left(\frac{k}{2}, \frac{\omega^2}{2}\right) \right)^2}{2 \sqrt{2} \Gamma\left(\frac{k}{2}\right) \gamma\left(\frac{k}{2}, \frac{\omega^2}{2}\right)}. \quad (67)$$

$f(\omega)$ implicitly depends on the accuracy α since ω is related to α by (31). $f(\omega)$ can be further simplified to a function $c(k, \alpha)$ that only depends on k and α . This is accomplished by obtaining ω corresponding to any $\alpha \in [0, 1]$ by (31) and replacing the result into (67).

REFERENCES

- [1] M. Stojnic, F. Parvaresh, and B. Hassibi, "On the reconstruction of block-sparse signals with an optimal number of measurements," *arXiv preprint arXiv:0804.0041*, 2008.
- [2] M. M. Hyder and K. Mahata, "Direction-of-arrival estimation using a mixed $\ell_{0,2}$ norm approximation," *IEEE Transactions on Signal Processing*, vol. 58, no. 9, pp. 4646–4655, 2010.
- [3] T. Euler and T. Baden, "Computational neuroscience: Species-specific motion detectors," *Nature*, 2016.
- [4] J. Zhu, D. Baron, and F. Krzakala, "Performance limits for noisy multimeasurement vector problems," *IEEE Transactions on Signal Processing*, vol. 65, no. 9, pp. 2444–2454, 2017.
- [5] D. L. Donoho, "For most large underdetermined systems of linear equations the minimal ℓ_1 -norm solution is also the sparsest solution," *Communications on pure and applied mathematics*, vol. 59, no. 6, pp. 797–829, 2006.
- [6] D. Amelunxen, M. Lotz, M. B. McCoy, and J. A. Tropp, "Living on the edge: Phase transitions in convex programs with random data," *Information and Inference: A Journal of the IMA*, vol. 3, no. 3, pp. 224–294, 2014.
- [7] K. V. Mishra, M. Cho, A. Kruger, and W. Xu, "Spectral super-resolution with prior knowledge," *IEEE transactions on signal processing*, vol. 63, no. 20, pp. 5342–5357, 2015.
- [8] D. L. Donoho, I. Johnstone, and A. Montanari, "Accurate prediction of phase transitions in compressed sensing via a connection to minimax denoising," *IEEE transactions on information theory*, vol. 59, no. 6, pp. 3396–3433, 2013.
- [9] M. A. Khajehnejad, W. Xu, A. S. Avestimehr, and B. Hassibi, "Analyzing weighted minimization for sparse recovery with nonuniform sparse models," *IEEE Transactions on Signal Processing*, vol. 59, no. 5, pp. 1985–2001, 2011.
- [10] M. P. Friedlander, H. Mansour, R. Saab, and O. Yilmaz, "Recovering compressively sampled signals using partial support information," *IEEE Transactions on Information Theory*, vol. 58, no. 2, pp. 1122–1134, 2012.
- [11] D. Needell, R. Saab, and T. Woolf, "Weighted-minimization for sparse recovery under arbitrary prior information," *Information and Inference: A Journal of the IMA*, vol. 6, no. 3, pp. 284–309, 2017.
- [12] R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde, "Model-based compressive sensing," *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1982–2001, 2010.
- [13] W. Xu, *Compressive sensing for sparse approximations: constructions, algorithms, and analysis*. PhD thesis, California Institute of Technology, 2010.
- [14] A. Beryehi and R. R. Müller, "Maximum-a-posteriori signal recovery with prior information: Applications to compressive sensing," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4494–4498, IEEE, 2018.
- [15] E. J. Candes, M. B. Wakin, and S. P. Boyd, "Enhancing sparsity by reweighted ℓ_1 minimization," *Journal of Fourier analysis and applications*, vol. 14, no. 5-6, pp. 877–905, 2008.
- [16] R. G. Baraniuk, V. Cevher, M. F. Duarte, and C. Hegde, "Model-based compressive sensing," *IEEE Transactions on Information Theory*, vol. 56, no. 4, pp. 1982–2001, 2010.
- [17] A. Flinthe, "Optimal choice of weights for sparse recovery with prior information," *IEEE Transactions on Information Theory*, vol. 62, no. 7, pp. 4276–4284, 2016.
- [18] S. D. Babacan, S. Nakajima, and M. N. Do, "Bayesian group-sparse modeling and variational inference," *IEEE transactions on signal processing*, vol. 62, no. 11, pp. 2906–2921, 2014.
- [19] S. Som, L. C. Potter, and P. Schniter, "On approximate message passing for reconstruction of non-uniformly sparse signals," in *Proceedings of the IEEE 2010 National Aerospace & Electronics Conference*, pp. 223–229, IEEE, 2010.
- [20] R. Ahmad and P. Schniter, "Iteratively reweighted ℓ_1 approaches to sparse composite regularization," *IEEE transactions on computational imaging*, vol. 1, no. 4, pp. 220–235, 2015.
- [21] R. T. Rockafellar, *Convex analysis*. Princeton university press, 2015.
- [22] J. A. Tropp, "Convex recovery of a structured signal from independent random linear measurements," in *Sampling Theory, a Renaissance*, pp. 67–101, Springer, 2015.
- [23] M. Grant and S. Boyd, "CVX: Matlab software for disciplined convex programming, version 2.1," Mar. 2014.
- [24] S. Foucart and H. Rauhut, *A mathematical introduction to compressive sensing*, vol. 1. Birkhäuser Basel, 2013.



Sajad Daei received the B.Sc. degree in electronic engineering from Guilan University, Guilan, Iran, in 2011, and the M.Sc. degree in communication engineering from Sharif University of Technology, Tehran, Iran, in 2013. He is currently pursuing his Ph.D. at Iran University of Science & Technology. His main research interests include convex optimization, compressed sensing and super resolution.



Farzan Haddadi was born in 1979. He received his B.Sc., M.Sc., and Ph.D. degrees in communication systems in 2001, 2003, and 2010, respectively, from Sharif University of Technology, Tehran, Iran. He joined Iran University of Science & Technology faculty in 2011. His main research interests are array signal processing, statistical signal processing, subspace tracking, and compressed sensing.



Arash Amini received the B.Sc., M.Sc., and Ph.D. degrees in electrical engineering (communications and signal processing) and the B.Sc. degree in petroleum engineering (reservoir) from the Sharif University of Technology, Tehran, Iran, in 2005, 2007, 2011, and 2005, respectively. He was a Researcher with the École Polytechnique fédérale de Lausanne, Lausanne, Switzerland, from 2011 to 2013, working on statistical approaches toward modeling sparsity in continuous-domain. He joined Sharif University of Technology as an assistant professor in 2013, where he is now an associate professor since 2018. He has served as an associate editor of IEEE Signal Processing Letters from 2014 to 2018. His research interests include various topics in statistical signal processing, specially compressed sensing.