

Dense Error Correction via ℓ^1 -Minimization

John Wright, *Student Member*, and Yi Ma, *Senior Member*.

J. Wright and Y. Ma are with the Electrical and Computer Engineering Department, University of Illinois at Urbana-Champaign.
Corresponding author: John Wright, 146 Coordinated Science Laboratory, 1308 West Main Street, Urbana, Illinois 61801, USA.
Email: jnwright@uiuc.edu.

Abstract

This paper studies the problem of recovering a non-negative sparse signal $\mathbf{x} \in \mathbb{R}^n$ from highly corrupted linear measurements $\mathbf{y} = A\mathbf{x} + \mathbf{e} \in \mathbb{R}^m$, where \mathbf{e} is an unknown error vector whose nonzero entries may be unbounded. Motivated by an observation from face recognition in computer vision, this paper proves that for highly correlated (and possibly overcomplete) dictionaries A , any non-negative, sufficiently sparse signal \mathbf{x} can be recovered by solving an ℓ^1 -minimization problem:

$$\min \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1 \quad \text{subject to} \quad \mathbf{y} = A\mathbf{x} + \mathbf{e}.$$

More precisely, if the fraction ρ of errors is bounded away from one and the support of \mathbf{x} grows sublinearly in the dimension m of the observation, then as m goes to infinity, the above ℓ^1 -minimization succeeds for all signals \mathbf{x} and almost all sign-and-support patterns of \mathbf{e} . This result suggests that accurate recovery of sparse signals is possible and computationally feasible even with nearly 100% of the observations corrupted. The proof relies on a careful characterization of the faces of a convex polytope spanned together by the standard crosspolytope and a set of iid Gaussian vectors with nonzero mean and small variance, which we call the “cross-and-bouquet” model. Simulations and experimental results corroborate the findings, and suggest extensions to the result.

Index Terms

Sparse Signal Recovery, Dense Error Correction, ℓ^1 -minimization, Gaussian Matrices, Polytope Neighborliness.

I. INTRODUCTION

Recovery of high-dimensional sparse signals or errors has been one of the fastest growing research areas in signal processing in the past few years. At least two factors have contributed to this explosive growth. On the theoretical side, the progress has been propelled by powerful tools and results from multiple mathematical areas such as measure concentration [1]–[3], statistics [4]–[6], combinatorics [7], and coding theory [8]. On the practical side, a lot of excitement has been generated by remarkable successes in real-world applications in areas such as signal (image or speech) processing [9], communications [10], computer vision and pattern recognition [11]–[13] etc.

A. A Motivating Example

One notable, and somewhat surprising, successful application of sparse representation is automatic face recognition. As described in [11], face recognition can be cast as a sparse representation problem. For

each person, a set of training images are taken under different illuminations. We can view each image as a vector by stacking its columns and put all the training images as column vectors of a matrix, say $A \in \mathbb{R}^{m \times n}$. Then, m is the number of pixels in an image and n is the total number of images for all the subjects of interest. Given a new query image, again we can stack it as a vector $\mathbf{y} \in \mathbb{R}^m$. To identify the image belongs to which subject, we can try to represent \mathbf{y} as a linear combination of all the images, i.e., $\mathbf{y} = A\mathbf{x}$ for some $\mathbf{x} \in \mathbb{R}^n$. Since in practice n can potentially be larger than m , the equations can be underdetermined and the solution \mathbf{x} may not be unique. In this context, it is natural to seek the sparsest solution for \mathbf{x} whose large non-zero coefficients then provide information about the subject's true identity. This can be done by solving the typical ℓ^1 -minimization problem:

$$\min_{\mathbf{x}} \|\mathbf{x}\|_1 \quad \text{subject to} \quad \mathbf{y} = A\mathbf{x}. \quad (1)$$

The problem becomes more interesting if the query image \mathbf{y} is severely occluded or corrupted, as shown in Figure 1 left, column (a). In this case, one needs to solve a corrupted set of linear equations $\mathbf{y} = A\mathbf{x} + \mathbf{e}$, where $\mathbf{e} \in \mathbb{R}^m$ is an unknown vector whose nonzero entries correspond to the corrupted pixels. For sparse errors \mathbf{e} and tall matrices A ($m > n$), Candes and Tao [14] proposed to multiply the equation $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ with a matrix B such that $BA = 0$, and then use ℓ^1 -minimization to recover the error vector \mathbf{e} from the new linear equation $B\mathbf{y} = B\mathbf{e}$.

As we mentioned earlier, in face recognition (and many other applications), n can be larger than m and the matrix A can be full rank. One cannot directly apply the above technique even if the error \mathbf{e} is known to be very sparse. To resolve this difficulty, in [11], the authors proposed to instead seek $[\mathbf{x}, \mathbf{e}]$ together as the sparsest solution to the extended equation $\mathbf{y} = [A \ \mathbf{I}]\mathbf{w}$ with $\mathbf{w} = \begin{bmatrix} \mathbf{x} \\ \mathbf{e} \end{bmatrix} \in \mathbb{R}^{m+n}$, by solving the extended ℓ^1 -minimization problem:

$$\min_{\mathbf{w}} \|\mathbf{w}\|_1 \quad \text{subject to} \quad \mathbf{y} = [A \ \mathbf{I}]\mathbf{w}. \quad (2)$$

This seemingly minor modification to the previous error correction approach has drastic consequences on the performance of robust face recognition. Solving the modified ℓ^1 -minimization enables almost perfect recognition even with more than 60% pixels of the query image are arbitrarily corrupted (see Figure 1 for an example), far beyond the amount of error that can theoretically be corrected by the previous error correction method [14].

Although ℓ^1 -minimization is expected to recover sufficiently sparse solutions with overwhelming probability for general systems of linear equations (see [16]), it is rather surprising that it works for the equation $\mathbf{y} = [A \ \mathbf{I}]\mathbf{w}$ at all. In the application described above, the columns of A are highly correlated.

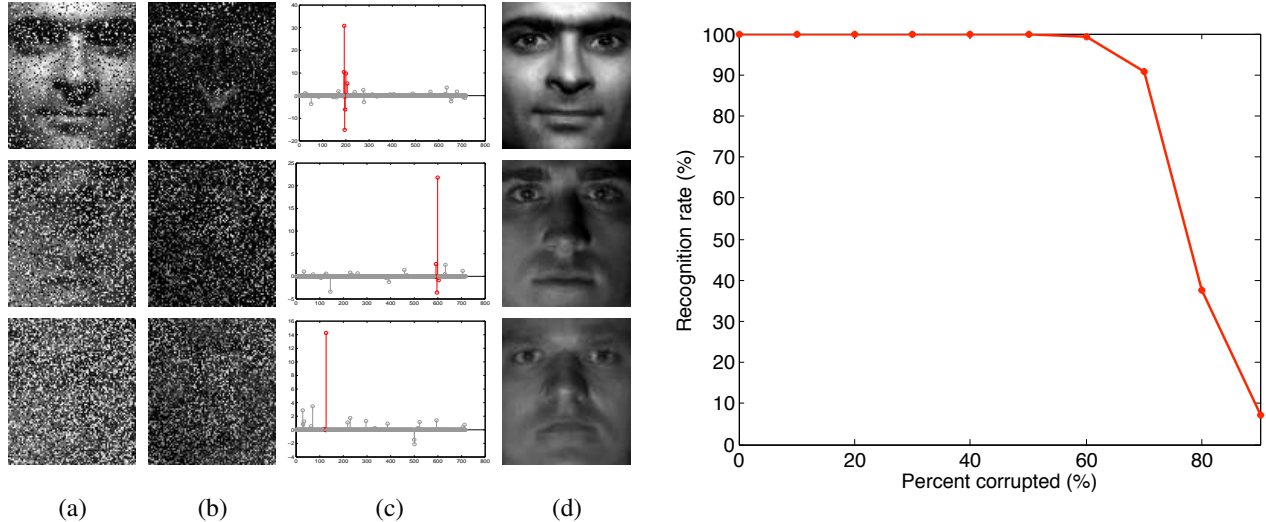


Fig. 1. **Face recognition under random corruption.** Left: (a) Test images y with random corruption from the database presented in [15]. Top row: 30% of pixels are corrupted, Middle row: 50% corrupted, Bottom row: 70% corrupted. (b) Estimated errors \hat{e} . (c) Estimated sparse coefficients \hat{x} . (d) Reconstructed images $y_r = A\hat{x}$. The extended ℓ^1 -minimization (2) correctly recovers and identifies all three corrupted face images. Right: The recognition rate across the entire range of corruption for all the 38 subjects in the database. It performs almost perfectly upto 60% random corruption.

As m becomes large (i.e. the resolution of the image becomes high), the convex hull spanned by all face images of all subjects is only an extremely tiny portion of the unit sphere \mathbb{S}^{m-1} .¹ For example, the images in Figure 1 lie on $\mathbb{S}^{8,063}$. The smallest inner product with their normalized mean is 0.723; they are contained within a spherical cap of volume $\leq 1.47 \times 10^{-229}$. These vectors are tightly bundled together as a “bouquet,” whereas the vectors associated with the identity matrix and its negative $\pm I$ together² form a standard “cross” in \mathbb{R}^m , as illustrated in Figure 2. Notice that such a “cross-and-bouquet” matrix $[A \ I]$ is neither incoherent nor (restrictedly) isometric, at least not uniformly. Also, the density of the desired solution w is not uniform either. The x part of w is usually a very sparse non-negative vector, but the e part can be very dense and have arbitrary signs. Existing results for recovering sparse signals suggest that ℓ^1 -minimization may have difficulty in dealing with such signals, contrary to its empirical success in face recognition.

We have experimented with similar cross-and-bouquet type models where the matrix A is a random matrix with highly correlated column vectors. The simulation results in Section III indicate that what we

¹At first sight, this seems somewhat surprising as faces of different people look so different to human eyes. That is probably because human brain has adapted to distinguish highly correlated visual signals such as faces or voices.

²Here we allow the entries of the error e to assume either positive or negative signs.

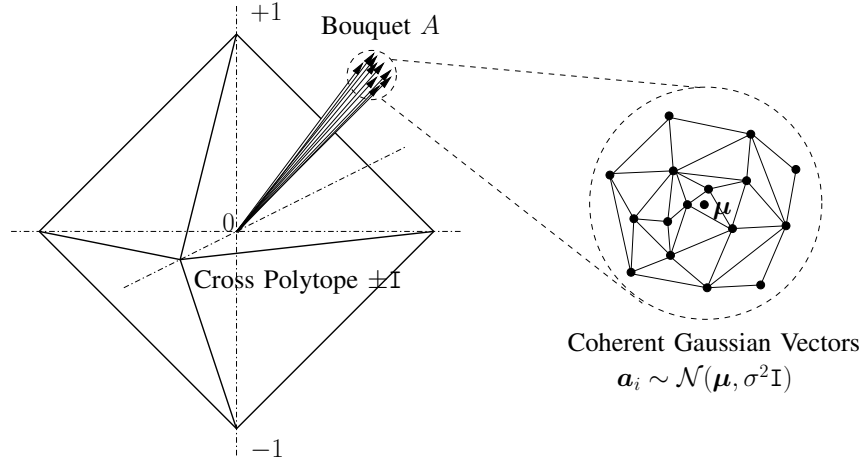


Fig. 2. **The “cross-and-bouquet” model.** Left: the bouquet A and the crosspolytope spanned by the matrix $\pm I$. Right: the tip of the bouquet magnified; it is a collection of iid Gaussian vectors with small variance σ^2 and common mean vector μ . The cross-and-bouquet polytope is spanned by vertices from both the bouquet A and the cross $\pm I$.

have seen in face recognition is not an isolated phenomenon. In fact, the simulations reveal something even more striking and puzzling: As the dimension m increases (and the sample size n grows in proportion), the percentage of errors that the ℓ^1 -minimization (2) can correct seems to approach 100%! This may seem surprising, but this paper explains why this should be expected.

B. The Main Model and Result

Motivated by the above empirical observations, this paper aims to resolve the apparent discrepancy between theory and practice of ℓ^1 -minimization and gives a more careful characterization of its behavior in recovering $[\mathbf{x}, \mathbf{e}]$ from the cross-and-bouquet (CAB) type models:

$$\mathbf{y} = A\mathbf{x} + \mathbf{e} = [A \ I]\mathbf{w}. \quad (3)$$

We model the bouquet, the columns of A , as iid samples from a multivariate Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_m)$, where $\sigma = \nu m^{-1/2}$ with ν sufficiently small, $\|\boldsymbol{\mu}\|_2 = 1$, and $\|\boldsymbol{\mu}\|_\infty \leq C_\mu m^{-1/2}$ for some $C_\mu \in \mathbb{R}_+$. These conditions insure that the bouquet remains tight as the dimension m grows, and that its mean is mostly incoherent with the columns of the cross $\pm I$.

We consider proportional growth for m and n , that is, $n/m \rightarrow \delta \in \mathbb{R}_+$ as $m \rightarrow \infty$. However, the support size of the sparse signal \mathbf{x} is only allowed to grow *sublinearly* in m : $\|\mathbf{x}\|_0 = O(m^{1-\eta})$ for some $\eta > 0$. This condition differs from (and is stronger than) the typical assumption in the sparse representation literature, where the support is often allowed to grow proportionally with the dimension

[16]. In the next subsection, we will explain why the support of the signal \mathbf{x} can only be sublinear if we allow the support of the error \mathbf{e} to be arbitrarily dense. Nevertheless, this sublinear bound of sparsity is more than adequate for signals in many practical problems, including the face recognition problem. There, the support of \mathbf{x} is bounded by a constant – the number of images per subject.

This paper proves that under the above conditions

for any $\rho < 1$, as m goes to infinity, solving the ℓ^1 -minimization problem (2) correctly recovers any non-negative sparse signal \mathbf{x} from almost any error \mathbf{e} with support size $\leq \rho m$.

We leave a more precise statement and the proof of the fact to Section II. In the remainder of this section, we discuss some of the main implications of this result in the broad context of sparse signal recovery, error correction, and some of its potential applications.

C. Relations to Previous Results

a) Restricted isometry and incoherence of the cross-and-bouquet model: As mentioned earlier, typical results in the literature for sparse signal recovery do not apply to equations of the type $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{e}$. The cross-and-bouquet matrix $[\mathbf{A} \ \mathbf{I}]$ is neither highly isometric nor incoherent. As a result, greedy algorithms such as Orthogonal Matching Pursuit [17], [18] succeed only when the error \mathbf{e} is very sparse (see Section III *a*) for the simulation results and comparison with our method). However, this does not mean that the restricted isometry property is irrelevant to the new problem. On the contrary, the proof of our results precisely rely on characterizing a special type of restricted isometry associated with this new problem, see Lemma 5 in Appendix A, which is used in the proof of our main result. Moreover, unlike the typical compressed sensing setting, the solution $[\mathbf{x}, \mathbf{e}]$ sought has very uneven density (or sparsity). This is reminiscent of the block sparsity studied in [19]. However, as we will see, the special block structure of the cross-and-boquet model enables sparse recovery far beyond the breakdown point for general sparse (or block sparse) signals.

b) Error correction: From an error correction viewpoint, the above result seems surprising: One can correctly solve a set of linear equations with *almost all* the equations randomly and arbitrarily corrupted! This is especially surprising considering that the best error-correcting codes (in the binary domain \mathbb{Z}_2), constructed based on expander graphs, normally correct a fixed fraction of errors [20]–[22]. The exact counterpart of our result in the binary domain is not clear.³ While there are superficial

³It is possible that under an analogous growth model (see Section II-A), the LP decoder of [21] could also correct large fractions of binary errors.

similarities between our result and [21], [23] in the use of linear programming for decoding and analysis via polytope geometry, those works do not consider real valued signals. In particular, the negative result of [23] for specific families of binary codes admitting linear programming decoders does not apply here.

We can, however, draw the following comparisons with existing error correction methods in the domain of real numbers:

- When $n < m$, the range of A is a subspace in \mathbb{R}^m . In such an overdetermined case, one could directly apply the method of Candes and Tao [14] mentioned earlier. However, the error vector e needs to be sparse for that approach whereas our result suggests even dense errors (with support far beyond 50%) can be corrected by instead solving the extended ℓ^1 -minimization (2). Thus, even in the overdetermined case, the new method has clear advantages for coherent matrices A . This will be verified by simulations in Section III a).
- The sublinear growth of the support of \mathbf{x} in m is the best one can hope for in the regime of dense errors. In general, we need at least $\|\mathbf{x}\|_0$ uncorrupted linear measurements to recover \mathbf{x} uniquely. If an arbitrary fraction of the m equations can be totally corrupted by e , no fixed fraction of the equations remain good for recovering \mathbf{x} . If, on the other hand, the error e is sparse, then the ℓ^1 -minimization (2) is able to recover \mathbf{x} with linear growth in support, as suggested by the existing theory [14], [16], [24]. Simulation results in Section III d) also confirm this phenomenon. However, in this paper, we are mainly interested in how the ℓ^1 -minimization behaves with dense errors, for $\rho \rightarrow 1$.
- When $n \geq m$, in general the Gaussian matrix A is full rank and the method of Candes and Tao [14] no longer applies.⁴ Our result suggests that as long as A is highly correlated, the ℓ^1 -minimization (2) can still recover the sparse signal \mathbf{x} correctly even if almost all the equations might be corrupted. This is verified by the simulation results in Section III c).

c) *Polytope geometry*: The success of ℓ^1 -minimization in recovering sparse solutions \mathbf{x} from underdetermined systems of linear equations $\mathbf{y} = A\mathbf{x}$ can be viewed as a consequence of a surprising property of high-dimensional polytopes. If the column vectors of A are random samples from a zero-mean Gaussian $\mathcal{N}(0, \mathbf{I})$, and m and n are allowed to grow proportionally, then with overwhelming probability the convex

⁴One could choose to pre-multiply the equation $\mathbf{y} = A\mathbf{x} + e$ with an “approximate orthogonal complement” of A , say the orthogonal complement of the mean vector $\boldsymbol{\mu}$, which is an $(m-1) \times m$ matrix B . Then the equation becomes $B\mathbf{y} = Be + \mathbf{z}$ where $\mathbf{z} = BA\mathbf{x}$. If the norm of \mathbf{x} is bounded, then \mathbf{z} is a signal with small magnitude due to the near-orthogonality of B and A . In this case, one can view \mathbf{z} as a noise term and try to recover e as a sparse signal via ℓ^1 -minimization. However, for e with arbitrary signs, the breakdown point for such ℓ^1 -minimization is less than 50%.

polytope $\text{conv}(A)$ spanned by the columns of A is highly neighborly [24], [25]. Neighborliness provides the necessary and sufficient condition for uniform sparse recovery: the ℓ^1 -minimization (1) correctly recovers \mathbf{x} if and only if the columns associated with the nonzero entries of \mathbf{x} span a face of the polytope $\text{conv}(A)$.

In our case, the columns of the matrix A are iid Gaussian vectors with nonzero mean μ and small variance σ^2 , whereas the vectors of the cross $\pm\mathbf{I}$ are completely fixed. To characterize when the extended ℓ^1 -minimization (2) is able to recover the solution $[\mathbf{x}, \mathbf{e}]$ correctly, we need to examine the geometry of the peculiar convex polytope $\text{conv}(A, \pm\mathbf{I})$ spanned together by the random bouquet A and the fixed cross $\pm\mathbf{I}$. Thus, it comes as no surprise that the proof of our main result relies on a careful study of the geometry of such a “cross-and-bouquet” polytope. As we will show that indeed, the vertices associated with the non-zero entries of \mathbf{x} and \mathbf{e} form a face of the polytope with probability approaching one as the dimension m becomes large. Precisely due to high neighborliness of the cross-and-bouquet polytopes, the extended ℓ^1 -minimization (2) is able to correctly recover the desired solution, even though the part of the solution corresponding to \mathbf{e} might be dense.

D. Implications on Applications

a) Robust reconstruction, classification, and source separation: The new result about the cross-and-bouquet model has strong implications on robust reconstruction, classification, and separation of highly correlated classes of signals such as faces or voices, despite severe corruption. It helps explain the surprising performance of face recognition that we discussed earlier. It further suggests that if the resolution of the image increases in proportion with the size of the database (say, due to the increasing number of subjects), the ℓ^1 -minimization would tolerate even higher level of corruption, far beyond the 60% at the resolution experimented with in [11]. Other applications where this kind of model could be useful and effective include speech recognition/imputation, audio source separation, video segmentation, or activity recognition from motion sensors.

b) Communication through an almost random channel: The result suggests that we can use the cross-and-bouquet model to accurately send information through a highly corrupting channel. Hypothetically, we can imagine a channel through which we can send one real number at a time, say as one packet of binary bits, and each packet has a high probability of being totally corrupted. One can use the sparse vector \mathbf{x} (or its support) to represent useful information, and use a set of highly correlated high-dimensional vectors as the encoding transformation A . The high correlation in A ensures that there is sufficient redundancy built in the encoded message $A\mathbf{x}$ so that the information about \mathbf{x} will not be lost even if many entries

of $A\mathbf{x}$ can be corrupted while being sent through such a channel. Our result suggests that the decoding can be done correctly and efficiently using linear programming.

c) Encryption and information hiding: One can potentially use the cross-and-bouquet model for encryption. For instance, if both the sender and receiver share the same encoding matrix A (say a randomly chosen Gaussian matrix), the sender can deliberately corrupt the message $A\mathbf{x}$ with arbitrary random errors \mathbf{e} before sending it to the receiver. The receiver can use linear programming to decode the information \mathbf{x} , whereas any eavesdropper will not be able to make much sense out of the highly corrupted message $\mathbf{y} = A\mathbf{x} + \mathbf{e}$. Of course, the long-term security of such an encryption scheme relies on the difficulty of learning the encoding matrix A after gathering many instances of corrupted message. It is not even clear whether it is easy to learn A from instances of uncorrupted message $\mathbf{y} = A\mathbf{x}$. Even if the dimensions of the matrix A are given, effectively learning A from a set of observed messages $\mathbf{Y} = [\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k]$ is still a largely open problem, known in the literature as the “dictionary learning” problem. Existing algorithms are iterative or greedy in nature, with no guarantee of global optimality [9]. Although its hardness has not been precisely characterized, we expect dictionary learning from highly corrupted observations to be an even more daunting problem, a challenge for anyone who tries to break this encryption scheme.

II. ROADMAP OF THE PROOF

In this section, we begin with a precise statement of our main result in Section II-A. We then lay out the roadmap for the proof. Section II-B outlines the key geometric picture behind the proof. In Section II-C, we then prove the main result, assuming that two technical conditions in Lemma 2 hold. Section II-D discusses the ideas required to establish these conditions, leaving a number of details to the Appendix.

A. Problem Statement

Motivated by the face recognition example introduced above, we consider the problem of recovering a non-negative⁵ sparse signal $\mathbf{x}_0 \in \mathbb{R}^n$ from highly corrupted observations $\mathbf{y} \in \mathbb{R}^m$:

$$\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}_0,$$

⁵The non-negativity assumption is important: in the highly-coherent systems considered here, ℓ^1 -minimization generally does not recover signals \mathbf{x}_0 with arbitrary signs. Geometrically, this would require vectors from the “bouquet” to “see” through the crosspolytope to vectors that are nearly antipodal to them.

where $\mathbf{e}_0 \in \mathbb{R}^m$ is a sparse vector of errors of arbitrary magnitude. The model for $A \in \mathbb{R}^{m \times n}$ should capture the idea that it consists of small deviations about a mean, hence a ‘‘bouquet.’’ In this paper, we consider the case where the columns of A are iid samples from a Gaussian distribution:

$$A = [\mathbf{a}_1, \dots, \mathbf{a}_n] \in \mathbb{R}^{m \times n}, \quad \mathbf{a}_i \sim_{iid} \mathcal{N}\left(\boldsymbol{\mu}, \frac{\nu^2}{m} \mathbf{I}_m\right), \quad \|\boldsymbol{\mu}\|_2 = 1, \quad \|\boldsymbol{\mu}\|_\infty \leq C_\mu m^{-1/2}. \quad (4)$$

Together, the two assumptions on the mean force it to remain incoherent with the standard basis (or ‘‘cross’’) as the dimension increases.

We study the behavior of the solution to the ℓ^1 -minimization (2) in this model, in the following asymptotic framework, which we term ‘‘weak proportional growth’’:

Assumption 1 (Weak Proportional Growth): A sequence of signal-error problems exhibits weak proportional growth with parameters $\delta > 0, \rho \in (0, 1), C_0 > 0, \eta_0 > 0$, denoted $\text{WPG}_{\delta, \rho, C_0, \eta_0}$ if as $m \rightarrow \infty$,

$$\frac{n}{m} \rightarrow \delta, \quad \frac{\|\mathbf{e}_0\|_0}{m} \rightarrow \rho, \quad \|\mathbf{x}_0\|_0 \leq C_0 m^{1-\eta_0}. \quad (5)$$

This should be contrasted with the ‘‘total proportional growth’’ (TPG) setting of, e.g., [26], in which the number of nonzero entries in the signal \mathbf{x}_0 also grows as a fixed fraction of the dimension. In that setting, one might expect a sharp phase transition in the combined sparsity of $(\mathbf{x}_0, \mathbf{e}_0)$ that can be recovered by ℓ^1 -minimization.⁶ In WPG, on the other hand, we observe a striking phenomenon not seen in TPG: the correction of arbitrary fractions of errors. This comes at the expense of the stronger assumption that $\|\mathbf{x}_0\|_0 = o(m)$, an assumption that is valid in some real applications such as the face recognition example above.

Before stating our main result, we fix some additional notation. For any $n \in \mathbb{Z}_+$, $[n]$ denotes the set $\{1, \dots, n\}$. Let $I = \text{supp}(\mathbf{x}_0) \subset [n]$, $J = \text{supp}(\mathbf{e}_0) \subset [m]$, $\boldsymbol{\sigma} = \text{sgn}(\mathbf{e}_0(J))$, and let $k_1 = |I|$ be the support size of the signal \mathbf{x}_0 and $k_2 = |J|$ the support size of the error \mathbf{e}_0 . For an arbitrary $r_1 \times r_2$ matrix M , if $L_1 \subset [r_1]$ and $L_2 \subset [r_2]$, M_{L_1, L_2} denotes the $|L_1| \times |L_2|$ submatrix of M indexed by these quantities. We use $M_{L_1, \bullet}$ as a shorthand for $M_{L_1, [r_2]}$. M^* denotes the transpose of M . Also, we use $\mathbf{1}_I$ (or $\mathbf{1}_J$) to represent a vector in \mathbb{R}^n (or \mathbb{R}^m) that has ones on the support I (or J) and zeros elsewhere. To reduce confusion between the index set I and the identity matrix, we use \mathbf{I} to denote the latter. Below, where the symbol C occurs with no subscript, it should be read as ‘‘some constant.’’ When used in different sections, it need not refer to the same constant.

⁶Existing results (e.g., [24]) do not prove the existence of phase transitions in inhomogeneous models such as the one considered here. However, simulations suggest that in total proportional growth, such transitions do occur (see Section III d)).

In the following, we say the cross-and-bouquet model is ℓ^1 -recoverable at $(I, J, \boldsymbol{\sigma})$ if for all $\mathbf{x}_0 \geq 0$ with support I and \mathbf{e}_0 with support J and signs $\boldsymbol{\sigma}$, we have

$$(\mathbf{x}_0, \mathbf{e}_0) = \arg \min \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1 \quad \text{subject to} \quad A\mathbf{x} + \mathbf{e} = A\mathbf{x}_0 + \mathbf{e}_0, \quad (6)$$

and the minimizer is uniquely defined. From the geometry of ℓ^1 -minimization, if (6) does not hold for some pair $(\mathbf{x}_0, \mathbf{e}_0)$, then it does not hold for any (\mathbf{x}, \mathbf{e}) with the same signs and support as $(\mathbf{x}_0, \mathbf{e}_0)$ [25]. Understanding ℓ^1 -recoverability at each $(I, J, \boldsymbol{\sigma})$ completely characterizes which solutions to $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ can be correctly recovered. In this language, our main result can be stated more precisely as:

Theorem 1 (Error Correction with the Cross-and-Bouquet Model): For any $\delta > 0$, $\exists \nu_0(\delta) > 0$ such that if $\nu < \nu_0$ and $\rho < 1$, in $\text{WPG}_{\delta, \rho, C_0, \eta_0}$ with A distributed according to (4), error support J chosen uniformly at random from $\binom{[m]}{k_2}$ and error signs $\boldsymbol{\sigma}$ chosen uniformly at random from $\{\pm 1\}^{k_2}$,

$$\lim_{m \rightarrow \infty} P_{A, J, \boldsymbol{\sigma}} \left[\ell^1\text{-recoverability at } (I, J, \boldsymbol{\sigma}) \quad \forall I \in \binom{[n]}{k_1} \right] = 1. \quad (7)$$

In other words, as long as the bouquet is sufficiently tight, asymptotically ℓ^1 -minimization recovers any non-negative sparse signal from almost any error with support size less than 100%.

B. Problem Geometry

We first restate the necessary and sufficient conditions for ℓ^1 -recoverability geometrically, as separation of a higher-dimensional ℓ^1 -ball and an affine subspace (see Figure 3). To witness this separation, we must show the existence of a separating hyperplane, whose normal we will denote by \mathbf{q} .

Lemma 1: Fix $(I, J, \boldsymbol{\sigma})$, and define $\mathbf{w} \doteq A_{J, \bullet}^* \boldsymbol{\sigma} - \mathbf{1}_I \in \mathbb{R}^n$ and

$$G \doteq \begin{bmatrix} A_{J^c, I} & A_{J^c, I^c} \\ 0 & \mathbf{I}_{n-k_1} \end{bmatrix} \in \mathbb{R}^{p \times n}, \quad p = m + n - k_1 - k_2. \quad (8)$$

Suppose G has full column rank n .⁷ The model is ℓ^1 -recoverable at $(I, J, \boldsymbol{\sigma})$ iff

$$\exists \mathbf{q} \in \mathbb{R}^p \quad \text{such that} \quad \|\mathbf{q}\|_\infty < 1 \quad \text{and} \quad G^* \mathbf{q} = \mathbf{w}. \quad (9)$$

Proof: As above, let $\mathbf{y} = A\mathbf{x}_0 + \mathbf{e}_0$. The pair $(\mathbf{x}_0, \mathbf{e}_0)$ is the unique minimum ℓ^1 -norm solution to the equation $\mathbf{y} = A\mathbf{x} + \mathbf{e}$ iff

$$\nexists (\Delta \mathbf{x}, \Delta \mathbf{e}) \neq \mathbf{0} : A\Delta \mathbf{x} = -\Delta \mathbf{e}, \quad \|\mathbf{x} + \Delta \mathbf{x}\|_1 + \|\mathbf{e} + \Delta \mathbf{e}\|_1 \leq \|\mathbf{x}\|_1 + \|\mathbf{e}\|_1. \quad (10)$$

⁷In the model outlined above, this occurs with probability one for m sufficiently large.

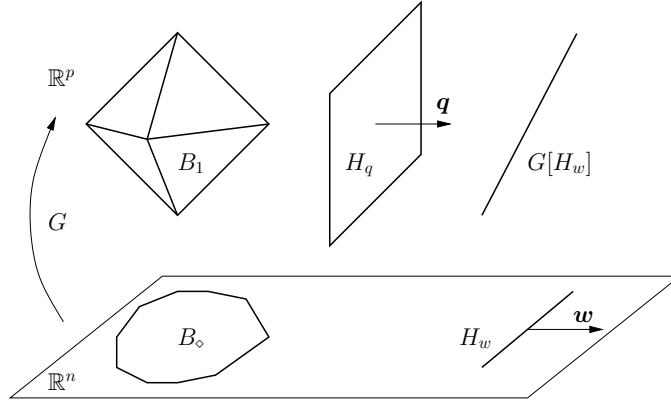


Fig. 3. **Geometry for the proof of Lemma 1.** The unit ball B_\diamond can be separated from H_w in \mathbb{R}^n if and only if in the lifted space \mathbb{R}^p , the ℓ^1 -ball B_1 can be separated from the image of H_w under the injective map G . H_q is the separating hyperplane with a normal vector q . Such an H_q might not be unique in \mathbb{R}^p , and q_0 would be the normal to the special separating hyperplane that contains $G(H_w)$.

Due to the geometry of ℓ^1 -minimization and the convexity of $\|\cdot\|_1$, we lose no generality in assuming that $\mathbf{x} = \mathbf{1}_I$, $\mathbf{e} \in \{-1, 0, 1\}^m$ and $\|\Delta\mathbf{x}\|_\infty < 1$, $\|\Delta\mathbf{e}\|_\infty < 1$. Then,

$$\|\mathbf{x} + \Delta\mathbf{x}\|_1 = \|\mathbf{x}\|_1 + \mathbf{1}_I^* \Delta\mathbf{x} + \|\Delta\mathbf{x}_{I^c}\|_1, \quad \text{and} \quad \|\mathbf{e} + \Delta\mathbf{e}\|_1 = \|\mathbf{e}\|_1 + \mathbf{e}^* \Delta\mathbf{e} + \|\Delta\mathbf{e}_{J^c}\|_1.$$

Substituting into (10) and using $\Delta\mathbf{e} = -A\Delta\mathbf{x}$ yields that (\mathbf{x}, \mathbf{e}) is optimal iff

$$\nexists \Delta\mathbf{x} \neq \mathbf{0} : \|A_{J^c, \bullet} \Delta\mathbf{x}\|_1 + \|\Delta\mathbf{x}_{I^c}\|_1 \leq \langle A^* \mathbf{e} - \mathbf{1}_I, \Delta\mathbf{x} \rangle.$$

Condition (II-B) is satisfied iff

$$\forall \Delta\mathbf{x} \neq \mathbf{0}, \quad \|G\Delta\mathbf{x}\|_1 > \langle \mathbf{w}, \Delta\mathbf{x} \rangle. \quad (11)$$

Let $H_w \subset \mathbb{R}^n$ be the affine subspace $\{\mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle = 1\}$. The function $\|G \cdot\|_1$ defines a norm $\|\cdot\|_\diamond$ on \mathbb{R}^n . Geometrically, (11) is satisfied iff the unit ball B_\diamond of $\|\cdot\|_\diamond$ is contained in the halfspace $H_w^- = \{\mathbf{x} \mid \langle \mathbf{w}, \mathbf{x} \rangle < 1\}$, as illustrated in Figure 3. This unit ball is a convex polytope, given by the inverse image (under the injective map G) of the intersection of $\mathcal{R}(G)$ and the unit ℓ^1 -ball B_1 in \mathbb{R}^p :

$$B_\diamond = G^{-1}[\mathcal{R}(G) \cap B_1(\mathbb{R}^p)]. \quad (12)$$

Now, $B_\diamond \subset H_w^-$ iff $[\mathcal{R}(G) \cap B_1(\mathbb{R}^p)] \subset G[H_w^-]$ iff $B_1(\mathbb{R}^p) \cap G[\text{cl } H_w^+] = \emptyset$. These two closed convex sets are nonintersecting iff there is a hyperplane⁸ $H_q = \{\mathbf{v} \in \mathbb{R}^p \mid \langle \mathbf{q}, \mathbf{v} \rangle = 1\} \subset \mathbb{R}^p$ separating them

⁸Notice H_q cannot contain $0 \in \text{interior}(B_1)$, so the normalization $\langle \mathbf{q}, \mathbf{v} \rangle = 1$ is appropriate.

(see Figure 3 again). We lose no generality in assuming that $B_1 \subset H_q^-$, that $G[\text{cl } H_w^+] \subset \text{cl } H_q^+$, and that H_q meets the relative boundary $\text{rbd } G[\text{cl } H_w^+] = G[H_w]$. The first condition occurs iff $\|\mathbf{q}\|_\infty < 1$, while the second occurs iff $G^*\mathbf{q} = \mathbf{w}$. ■

The most natural candidate for a normal vector \mathbf{q} is the minimum ℓ^2 -norm solution to this equation,

$$\mathbf{q}_0 = (G^\dagger)^*\mathbf{w} = G(G^*G)^{-1}\mathbf{w}. \quad (13)$$

When we use this particular normal \mathbf{q}_0 , we are demanding that the *projection* of B_1 onto $\mathcal{R}(G)$ lie in $G[H_w^-]$. Since the projection contains the intersection, $B_1 \subset \{\langle \mathbf{q}_0, \cdot \rangle < 1\}$ is a sufficient, but not necessary condition. It is not surprising, then, that this condition often does not hold – empirically, $\|\mathbf{q}_0\|_\infty \geq 1$ with high probability. However, as we will see, the set of violations is almost always small, and we can apply a simple iterative scheme to improve \mathbf{q}_0 to a valid separator \mathbf{q} with $\|\mathbf{q}\|_\infty < 1$.

C. Iterative Construction of Separator

Our next lemma argues that if we are given an initial guess at a normal vector $\mathbf{q}_0 \in \mathbb{R}^p$ whose hyperplane $H_{\mathbf{q}_0}$ separates $G[H_w]$ from *most* of the vertices of B_1 , then we can refine \mathbf{q}_0 to a \mathbf{q}_∞ that separates $G[H_w]$ and *all* of the vertices of B_1 . In general, finding such a \mathbf{q}_∞ requires solving a linear programming problem. We will analyze the feasibility of this linear program by considering an iteration similar to the alternating projection method for finding a pair of closest points between two convex sets. In this case, the two convex sets of interest are the hypercube of radius $1 - \varepsilon$ and the affine subspace $\mathbf{q}_0 + \mathcal{R}(G)^\perp$.

In the following lemma, $\mathbf{q}_0 \in \mathbb{R}^p$ is arbitrary (though $\mathbf{q}_0 = G^{\dagger*}\mathbf{w}$ is natural). We will construct a sequence of vectors $(\mathbf{q}_k)_{k=0}^\infty$. Fix a small constant $\varepsilon > 0$, and define the operator θ which takes the part of a vector that protrudes above $1 - \varepsilon$:

$$[\theta \mathbf{x}](i) \doteq \begin{cases} 0, & \text{for } |\mathbf{x}(i)| \leq 1 - \varepsilon, \\ \text{sgn}(\mathbf{x}(i))(|\mathbf{x}(i)| - 1 + \varepsilon), & \text{for } |\mathbf{x}(i)| > 1 - \varepsilon. \end{cases} \quad (14)$$

We iteratively construct \mathbf{q}_∞ by setting

$$\mathbf{q}_{i+1} = \mathbf{q}_i - \pi_{\mathcal{R}(G)^\perp} \theta \mathbf{q}_i = \mathbf{q}_i - \theta \mathbf{q}_i + \pi_{\mathcal{R}(G)} \theta \mathbf{q}_i. \quad (15)$$

Notice that by construction, $G^*\mathbf{q}_k = G^*\mathbf{q}_0 = \mathbf{w}$ for all k . So if $\theta \mathbf{q}_i \rightarrow 0$, then $\|\mathbf{q}_i\|_\infty < 1$ eventually, and \mathbf{q}_∞ is a valid separator.

Before proving that this iteration produces a valid separator with high probability, we first demonstrate its behavior on a simulated example with $m = 3,000$, $\delta = .4$, $\nu = .1$, $\rho = .65$, and $k_1 = 10$. Figure

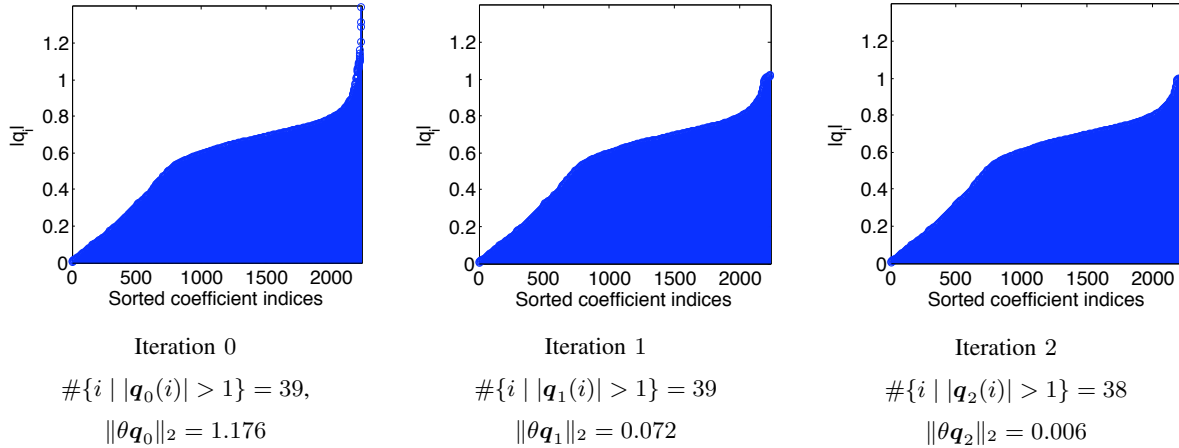


Fig. 4. **Iterative refinement producing a separating hyperplane.** Here, $m = 3000$, $\delta = .4$, $\nu = .1$, $\rho = .65$, $k_1 = 10$. We plot the sorted magnitudes of the entries of \mathbf{q}_i . At left, \mathbf{q}_0 separates $G(H_w)$ from most of the vertices of B_1 : only 39 violations occur. The distinct bimodal characteristic of \mathbf{q}_0 is due to the differences between the statistics of the top (R_1) and bottom (R_2) indices. Applying the iteration decreases $\|\theta \mathbf{q}_i\|_2$ geometrically; after 5 iterations, a valid separator is obtained.

4 plots the sorted absolute values of entries of \mathbf{q}_i . Notice that the sorted coefficients clearly divide into two parts; these correspond to the upper⁹ (R_1) and lower (R_2) indices. The initial separator \mathbf{q}_0 cleanly separates $G[H_w]$ from most of the vertices of B_1 : only 39 entries protrude above $1 - \varepsilon$. These entries are quickly iterated away: $\|\theta \mathbf{q}\|_2$ decreases geometrically until after 5 iterations a valid separator is obtained.

Lemma 2: Suppose $\exists c \in (0, 1)$ such that

$$\xi \doteq \sup_{\|\mathbf{s}\|_0 \leq cp, \mathbf{s} \neq \mathbf{0}} \frac{\|\pi_{\mathcal{R}(G)} \mathbf{s}\|_2}{\|\mathbf{s}\|_2} < 1, \quad (16)$$

and

$$\|\mathbf{q}_0\|_2 + \frac{1}{1 - \xi} \|\theta \mathbf{q}_0\|_2 \leq (1 - \varepsilon) \sqrt{cp}, \quad (17)$$

where G is the matrix defined in (8). Iteratively construct a sequence of vectors $\{\mathbf{q}_i\}$, with $\mathbf{q}_i = \mathbf{q}_{i-1} - \pi_{\mathcal{R}(G)} \theta \mathbf{q}_{i-1}$, where θ threshold-residual operator defined in (14). Then $\lim_{k \rightarrow \infty} \theta \mathbf{q}_k = 0$.

Proof: Let $T_k = \{i \mid |\mathbf{q}_k(i)| > 1 - \varepsilon\} \subset [p]$, and consider the following three statements:

$$\|\mathbf{q}_k\|_2 \leq \|\mathbf{q}_0\|_2 + \|\theta \mathbf{q}_0\|_2 \sum_{i=0}^k \xi^i, \quad \|\theta \mathbf{q}_k\|_2 \leq \|\theta \mathbf{q}_0\|_2 \xi^k, \quad \#T_k \leq cp. \quad (18)$$

⁹Where necessary, we will use $R_1 = \{1, \dots, m - k_2\} \subset [p]$ to index the upper rows of G (corresponding to A), and $R_2 = [p] \setminus R_1$ to index the lower rows.

We will show by induction that these statements hold for all k , establishing the lemma. The first two statements of (18) hold trivially $k = 0$. For $\#T_0$, notice that by (17),

$$\#T_0 \leq \frac{\|\mathbf{q}_0\|_2^2}{(1-\varepsilon)^2} \leq cp.$$

Now, suppose the three statements hold for $0, \dots, k$. Since $\theta\mathbf{q}_k$ has the same signs and smaller magnitude than \mathbf{q}_k , $\|\mathbf{q}_k - \theta\mathbf{q}_k\|_2 \leq \|\mathbf{q}_k\|_2$; combining this with the inductive hypothesis we have

$$\begin{aligned} \|\mathbf{q}_{k+1}\|_2 &= \|\mathbf{q}_k - \theta\mathbf{q}_k + \pi_{\mathcal{R}(G)}\theta\mathbf{q}_k\| \leq \|\mathbf{q}_k - \theta\mathbf{q}_k\| + \|\pi_{\mathcal{R}(G)}\theta\mathbf{q}_k\| \leq \|\mathbf{q}_k\| + \xi^{k+1}\|\theta\mathbf{q}_0\| \\ &\leq \|\mathbf{q}_0\|_2 + \|\theta\mathbf{q}_0\|_2 \sum_{i=0}^{k+1} \xi^i, \end{aligned}$$

Similarly, notice that since $\pi_{\mathcal{R}(G)}\theta\mathbf{q}_k$ dominates $\theta(\mathbf{q}_k - \theta\mathbf{q}_k + \pi_{\mathcal{R}(G)}\theta\mathbf{q}_k)$ elementwise,

$$\|\theta\mathbf{q}_{k+1}\| \leq \|\pi_{\mathcal{R}(G)}\theta\mathbf{q}_k\| \leq \xi\|\theta\mathbf{q}_k\| \leq \xi^{k+1}\|\theta\mathbf{q}_0\|.$$

Finally, for the sparsity result $T_{k+1} \leq cp$, note that

$$\|\mathbf{q}_{k+1}\|_2 \leq \|\mathbf{q}_0\|_2 + \|\theta\mathbf{q}_0\|_2 \sum_{i=0}^{k+1} \xi^i \leq \|\mathbf{q}_0\|_2 + \frac{1}{1-\xi}\|\theta\mathbf{q}_0\|_2 \leq (1-\varepsilon)\sqrt{cp},$$

and so $\theta\mathbf{q}_{k+1}$ must be (cp) -sparse. Since (18) holds for all k , $\|\theta\mathbf{q}_k\|_2 \rightarrow 0$. ■

D. Putting it All Together

By Lemmas 1 and 2, if the two conditions (16) and (17) hold for a given sign and support triplet (I, J, σ) , then (I, J, σ) is ℓ^1 -recoverable.¹⁰ We will show that as $m \rightarrow \infty$, for any sequence of signal supports I , (16) and (17) hold with probability approaching one in the random matrix A and error (J, σ) . The probability that either condition fails for a given I will be small enough to allow a union bound over all I , establishing Theorem 1. We will assume we are in the large error regime, with $\bar{\rho} \doteq 1 - \rho$ lower bounded as specified in the lemmas below. The conclusion still follows for smaller error fractions, since whenever (I, J, σ) is ℓ^1 -recoverable, so is $(I, J', \sigma_{J'})$ for any $J' \subset J$.

In this section, we lay out the main ideas for the rest of the proof, which consists of two parts, one for each of the conditions in Lemma 2. We establish that following two properties hold simultaneously with probability at least $1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$:

- 1) For a small enough constant c , the projection ratio ξ for cm -sparse signals onto $\mathcal{R}(G)$ is bounded below 1 by a polynomial function in ν . More precisely, $\xi < 1 - C\nu^8$ for some constant $C > 0$. As a result, the coefficient $\frac{1}{1-\xi}$ in the second condition (17) is bounded by $C^{-1}\nu^{-8}$.

¹⁰Notice that conditions (16) and (17) depend on (I, J, σ) , through the construction of the matrix G .

- 2) As m goes to infinity, the ℓ^2 -norm of the initial separating normal vector $\|\mathbf{q}_0\|_2$ is bounded above by $\nu O(m^{1/2})$, and $\|\theta\mathbf{q}_0\|_2$ is bounded above by $e^{-\alpha/\nu^2} O(m^{1/2})$ for some constant α .

Putting these results together, the initial separating normal vector \mathbf{q}_0 satisfies:

$$\|\mathbf{q}_0\|_2 + \frac{1}{1-\xi}\|\theta\mathbf{q}_0\|_2 \leq \nu O(m^{1/2}) + C^{-1}\nu^{-8}e^{-\alpha/\nu^2} O(m^{1/2}). \quad (19)$$

If the deviation ν of the bouquet is small enough, the second condition (17) of Lemma 2 will be satisfied, since the right hand side, $(1-\varepsilon)\sqrt{cp} = \Omega(m^{1/2})$ is independent of ν . Hence, by Lemma 2, the initial normal \mathbf{q}_0 will converge to a valid normal vector that separates the ℓ^1 -ball B_1 from the subspace $G[H_w]$, establishing ℓ^1 -recoverability at (I, J, σ) . Comparing the failure probability for the two conditions to the number of subsets $I \subset [n]$ of size $C_0 m^{1-\eta_0}$ then completes the proof of Theorem 1. These arguments are laid out more precisely and quantitatively in Section C of the appendix.

Whereas Lemmas 1 and 2 have simple geometric and algebraic proofs, the above results require more detailed analysis of large Gaussian matrices. We outline the main ideas of their proof in this section, leaving many of the technical details to the appendix. The derivation is based on recent (and now widely-used) results on concentration of Lipschitz functions [3], which state that if \mathbf{x} is a d -dimensional iid $\mathcal{N}(0, 1)$ vector and $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is 1-Lipschitz, then

$$P[|f(\mathbf{x}) - \mathbb{E}f(\mathbf{x})| \geq t] \leq 2 \exp\left(-\frac{2t^2}{\pi^2}\right). \quad (20)$$

Two cases are of particular interest here. First, the norm concentrates as (see, e.g., [27]):

$$P\left[\|\mathbf{x}\|_2 \geq \beta\sqrt{d}\right] \leq \exp\left(-\frac{2(\beta-1)^2}{\pi^2}d\right). \quad (21)$$

Second, as has been widely exploited in the compressed sensing literature (e.g., [14], [16]), the singular values of rectangular Gaussian matrices with aspect ratio α concentrate about the values $1 \pm \sqrt{\alpha}$ predicted by the Marchenko-Pasteur law:

Fact 1 (Concentration of singular values [3]): Let $A \in \mathbb{R}^{m \times n}$, ($m > n$) be a random matrix with entries iid $\mathcal{N}(0, \frac{1}{m})$. Then for any $t > 0$,

$$P\left[\sigma_{\max}(A) > 1 + \sqrt{n/m} + o(1) + t\right] \leq e^{-mt^2/2}, \quad (22)$$

$$P\left[\sigma_{\min}(A) < 1 - \sqrt{n/m} + o(1) - t\right] \leq e^{-mt^2/2}. \quad (23)$$

We will also return to (20) in the proof of Lemma 8 of the appendix.

1) *Projection of Sparse Vectors*: In this subsection, we upper bound the norm of the projection of any sparse vector onto $\mathcal{R}(G)$. Since the lower (R_2) coordinates of

$$G \doteq \begin{bmatrix} A_1 & A_2 \\ 0 & \mathbf{I} \end{bmatrix} = \begin{bmatrix} Z_1 + \boldsymbol{\mu}_{J^c} \mathbf{1}_{k_1}^* & Z_2 + \boldsymbol{\mu}_{J^c} \mathbf{1}_{\delta m - k_1}^* \\ 0 & \mathbf{I} \end{bmatrix}$$

contain an identity matrix, when the variance ν^2/m of the perturbations Z_1, Z_2 is small, we expect that sparse vectors with support on R_2 will be very close to $\mathcal{R}(G)$. The following lemma verifies that this is the case, but argues that distance to $\mathcal{R}(G)$ is at least $\Omega(\nu^8)$. The technical conditions appear complicated, but simply assert that the fraction of nonzeros c is sufficiently small.

Lemma 3 (Projection of Sparse Vectors): Suppose that $\bar{\rho} < \delta$ and $\nu < \min(\frac{1}{9}, (512/\delta)^{1/4})$,

$$c < \min \left\{ \frac{\bar{\rho}}{1024}, \frac{\bar{\rho}}{64(1 + 2C_\mu \bar{\rho}^{-1/2})^2} \right\}, \quad \bar{\rho}H(c/\bar{\rho}) + \delta H(c/\delta) < \frac{\bar{\rho}}{128\pi^2}, \quad (24)$$

where $H(\cdot)$ is the base- e binary entropy function. Then the projection of a sparse vector $\mathbf{s} \in \mathbb{R}^p$ with $\|\mathbf{s}\|_0 \leq cm$ onto the range of G is bounded as

$$\sup_{\|\mathbf{s}\|_0 \leq cm, \mathbf{s} \neq \mathbf{0}} \frac{\|\pi_{\mathcal{R}(G)} \mathbf{s}\|_2}{\|\mathbf{s}\|_2} < 1 - \nu^8 \left(\frac{\sqrt{\bar{\rho}} (\sqrt{\delta} - \sqrt{\bar{\rho}})}{32 + 128 \nu^2 (\sqrt{\delta} + \sqrt{\bar{\rho}})^2} \right)^4 \quad (25)$$

on the complement of a bad event with probability $e^{-Cm(1+o(1))}$.

Proof: The projection of $\mathbf{s} = \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix}$ onto $\mathcal{R}(G)$ solves

$$\min_{\mathbf{r} \in \mathbb{R}^n} \|\begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix} - G\mathbf{r}\|_2^2 = \min_{\mathbf{u}_1, \mathbf{u}_2} \|\begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix} - G \begin{bmatrix} \mathbf{u}_1 \\ \mathbf{s}_2 + \mathbf{u}_2 \end{bmatrix}\|_2^2 = \min_{\mathbf{u}_1, \mathbf{u}_2} \|\mathbf{s}_1 - A_1 \mathbf{u}_1 - A_2(\mathbf{s}_2 + \mathbf{u}_2)\|_2^2 + \mathbf{u}_2^* \mathbf{u}_2.$$

By minimizing the first term, we can write the unique optimal \mathbf{u}_1 in terms of the remaining variables:

$$\mathbf{u}_1 = (A_1^* A_1)^{-1} A_1^* \mathbf{s}_1 - (A_1^* A_1)^{-1} A_1^* A_2 (\mathbf{s}_2 + \mathbf{u}_2)$$

and subsequently, the optimal \mathbf{u}_2 satisfies:

$$-A_2^* \mathbf{s}_1 + A_2^* A_1 \mathbf{u}_1 + A_2^* A_2 (\mathbf{s}_2 + \mathbf{u}_2) + \mathbf{u}_2 = \mathbf{0} \Rightarrow (\mathbf{I} + A_2^* \pi_{A_1^\perp} A_2) \mathbf{u}_2 = A_2^* \pi_{A_1^\perp} \mathbf{s}_1 - A_2^* \pi_{A_1^\perp} A_2 \mathbf{s}_2,$$

where $\pi_{A_1^\perp}$ denotes the projection matrix onto the orthogonal complement of $\mathcal{R}(A_1)$.

Write $A_2^* \pi_{A_1^\perp} = USV^*$ with $U \in \mathbb{R}^{(\delta m - k_1) \times (\bar{\rho} m - k_1)}$ and $V \in \mathbb{R}^{\bar{\rho} m \times (\bar{\rho} m - k_1)}$ orthogonal matrices, and the diagonal of $S \in \mathbb{R}^{(\bar{\rho} m - k_1) \times (\bar{\rho} m - k_1)}$ containing the nonzero singular values of $A_2^* \pi_{A_1^\perp}$. Then if \mathbf{u}_2 is the solution to the above equation

$$\begin{aligned} \|\mathbf{s} - \pi_{\mathcal{R}(G)} \mathbf{s}\|_2 &\geq \|\mathbf{u}_2\|_2 = \|(S^2 + \mathbf{I})^{-1} S V^* [\mathbf{I} \quad -A_2] \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix}\|_2 \\ &= \|(S^2 + \mathbf{I})^{-1} S [V^* \quad -S U^*] \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix}\|_2. \end{aligned} \quad (26)$$

Above is the norm of the product of a diagonal matrix $(S^2 + \mathbf{I})^{-1}S$, a wide matrix $[V^* - SU^*]$, and a sparse vector \mathbf{s} . We will bound it by lower bounding the elements of the diagonal matrix, and then lower bounding the “restricted minimum singular value”

$$\gamma_{cm}([V^* - SU^*]) \doteq \inf_{\|\mathbf{s}\|_0 \leq cp, \mathbf{s} \neq \mathbf{0}} \frac{\|[V^* - SU^*] \mathbf{s}\|_2}{\|\mathbf{s}\|_2}.$$

We first drop the top row of $(S^2 + \mathbf{I})^{-1}S[V^* - SU^*]$. This allows us to uniformly lower bound the diagonal of $(S^2 + \mathbf{I})^{-1}S$. While σ_1 can be quite large due to the inhomogeneous term $(\boldsymbol{\mu}_{J_c} \mathbf{1}^*)$, and hence $\frac{\sigma_1}{\sigma_1^2 + 1}$ can be quite small, for the remaining singular values $\frac{\sigma_i}{\sigma_i^2 + 1}$ is at least on the order of ν . Let $\tilde{S} \in \mathbb{R}^{(\bar{\rho}m - k_1 - 1) \times (\bar{\rho}m - k_1 - 1)}$ be the diagonal matrix obtained by dropping the row and column of S corresponding to the largest singular value; \tilde{V} and \tilde{U} are obtained by dropping the corresponding columns. From (26),

$$\|\mathbf{u}_2\|_2 \geq \left\| (\tilde{S}^2 + \mathbf{I})^{-1} \tilde{S} [\tilde{V}^* - \tilde{S} \tilde{U}^*] \begin{bmatrix} \mathbf{s}_1 \\ \mathbf{s}_2 \end{bmatrix} \right\|_2 \geq \frac{\sigma_{\min}(A_2^* \pi_{A_1^\perp})}{1 + \sigma_2^2(A_2^* \pi_{A_1^\perp})} \gamma_{cm}([\tilde{V}^* - \tilde{S} \tilde{U}^*]) \|\mathbf{s}\|_2, \quad (27)$$

where $\sigma_{\min}(A_2^* \pi_{A_1^\perp})$ is the smallest nonzero singular value and $\sigma_2(A_2^* \pi_{A_1^\perp})$ is the second largest singular value.

a) *Bounding the second largest singular value $\sigma_2(A_2^* \pi_{A_1^\perp})$:* Write $\hat{\boldsymbol{\mu}} \doteq \pi_{A_1^\perp} \boldsymbol{\mu}_{J_c}$, and notice that

$$\begin{aligned} \sigma_2(A_2^* \pi_{A_1^\perp}) &= \inf_{\mathbf{u} \neq \mathbf{0}} \sup_{\mathbf{v} \neq \mathbf{0}} \frac{\|A_2^* \pi_{A_1^\perp} \pi_{\mathbf{u}^\perp} \mathbf{v}\|_2}{\|\mathbf{v}\|_2} = \inf_{\mathbf{u} \neq \mathbf{0}} \sigma_1(A_2^* \pi_{A_1^\perp} \pi_{\mathbf{u}^\perp}) \\ &\leq \sigma_1(A_2^* \pi_{A_1^\perp} \pi_{\hat{\boldsymbol{\mu}}^\perp}) = \sigma_1(Z_2^* \pi_{(\boldsymbol{\mu}_{J_c}, Z_1)^\perp}). \end{aligned}$$

Choose any orthonormal basis for the subspace $\Sigma = (\mathcal{R}(Z_1) + \mathcal{R}(\boldsymbol{\mu}_{J_c}))^\perp$. Since Σ is probabilistically independent of Z_2 , the representation of the projection $Z_2^* \pi_\Sigma$ with respect to the chosen basis is simply distributed as a $(\delta m - k_1) \times (\bar{\rho}m - k_1 - 1)$ random matrix \hat{Z}_2 with entries $\mathcal{N}(0, \nu^2/m)$. Since $\frac{\sqrt{m}}{\nu\sqrt{\delta m - k_1}} \hat{Z}_2$ is $\mathcal{N}(0, \frac{1}{\delta m - k_1})$, by Fact 1,

$$P \left[\sigma_1 \left(\frac{\sqrt{m}}{\nu\sqrt{\delta m - k_1}} \hat{Z}_2 \right) \geq 1 + \sqrt{\frac{\bar{\rho}m - k_1 - 1}{\delta m - k_1}} + t \right] \leq \exp(-(t - o(1))^2 (\delta m - k_1)/2), \quad (28)$$

and so $P \left[\sigma_1(\hat{Z}_2) \geq 2\nu(\sqrt{\delta} + \sqrt{\bar{\rho}}) \right] \leq e^{-Cm(1+o(1))}$. On the complement of this bad event, $\sigma_2^2(A_2^* \pi_{A_1^\perp}) \leq 4\nu^2(\sqrt{\delta} + \sqrt{\bar{\rho}})^2$.

b) *Bounding the smallest nonzero singular value $\sigma_{\min}(A_2^* \pi_{A_1^\perp}) = \inf_{\mathbf{x} \in A_1^\perp} \frac{\|A_2^* \mathbf{x}\|_2}{\|\mathbf{x}\|_2}$:* Let $W \in \mathbb{R}^{\bar{\rho}m \times (\bar{\rho}m - k_1)}$ be a matrix whose columns form an orthonormal basis for A_1^\perp , and let $Q \in \mathbb{R}^{(\delta m - k_1) \times (\delta m - k_1 - 1)}$ be an orthonormal basis for $\mathbf{1}_{\delta m - k_1}^\perp$. Then $\sigma_{\min}(A_2^* \pi_{A_1^\perp}) \geq \sigma_{\min}(Q^* A_2 W)$. Conditioned on A_1 , $Z_2' \doteq Q^* A_2 W \in \mathbb{R}^{(\delta m - k_1 - 1) \times (\bar{\rho}m - k_1)}$ is iid $\mathcal{N}(0, \nu^2/m)$. Applying Fact 1 (with a similar rescaling argument to the one used for $\sigma_{\max}(\hat{Z}_2)$ above) gives that

$$P \left[\sigma_{\min}(Z_2') < \frac{\nu}{2} (\sqrt{\delta} - \sqrt{\bar{\rho}}) \right] \leq e^{-Cm(1+o(1))}. \quad (29)$$

On the complement of this bad event, $\sigma_{\min}(A_2^* \pi_{A_+}) \geq \frac{\nu}{2}(\sqrt{\delta} - \sqrt{\bar{\rho}})$.

Finally, in Lemma 5 of Appendix A, we show that under the stated conditions, the restricted singular value γ_{cm} in (27) satisfies $\gamma_{cm}([\tilde{V}^* \quad -\tilde{S}\tilde{U}^*]) \geq \frac{\nu\sqrt{\bar{\rho}}}{16}$ with probability at least $1 - e^{-Cm(1+o(1))}$. Notice that this bound agrees with (and in fact is looser than) the Marchenko-Pasteur law for a $\bar{\rho}m \times cm$ Gaussian $\mathcal{N}(0, \nu^2/m)$ matrix (i.e., the concentration result of Fact 1). In fact, the proof argues that the two blocks of this matrix are probabilistically independent, and then applies Fact 1 to an equivalent pair of Gaussian matrices. The somewhat technical conditions (24) introduced here are necessary to ensure that a union bound over all subsets of cm columns remains small.

Combining the three results, we have that for all $\mathbf{s} \in \mathbb{R}^P$ with $\|\mathbf{s}\|_0 \leq cm$,

$$\frac{\|\mathbf{s} - \pi_G \mathbf{s}\|_2}{\|\mathbf{s}\|_2} \geq \frac{\nu^2 \sqrt{\bar{\rho}} (\sqrt{\delta} - \sqrt{\bar{\rho}})}{32 + 128 \nu^2 (\sqrt{\bar{\rho}} + \sqrt{\delta})^2} \doteq \beta \quad (30)$$

Notice that $\frac{\|\pi_G \mathbf{s}\|}{\|\mathbf{s}\|} = \sqrt{1 - \left(\frac{\|\mathbf{s} - \pi_G \mathbf{s}\|}{\|\mathbf{s}\|}\right)^2} \leq \sqrt{1 - \beta^2} \leq 1 - \beta^4$, where we have used that $1 - \beta^4 > \sqrt{1 - \beta^2}$ for $\beta < 1/\sqrt{2}$; this is guaranteed for $\nu < (512/\delta)^{1/4}$. Combined with (30), this implies (25). \blacksquare

2) *Initial Separating Hyperplane:* In this section, we analyze the initial separator \mathbf{q}_0 , obtained as the minimum 2-norm solution to the equation $G^* \mathbf{q} = \mathbf{w}$. We upper bound both $\|\mathbf{q}_0\|_2$ and $\|\theta \mathbf{q}_0\|_2$, where the operator θ defined in (14) retains the portion of a vector that protrudes above $1 - \varepsilon$ in absolute value. These bounds provide the second half of the conditions needed in Lemma 2 to show that \mathbf{q}_0 can be refined by alternating projections to give a true separator.

Lemma 4: Suppose $\bar{\rho} < \delta$ and $\nu < \frac{1}{8(\sqrt{\delta}+1)}$. Then for G defined in (8) and $\mathbf{w} = A_{J,\bullet}^* \boldsymbol{\sigma} - \mathbf{1}_I$, \exists constants α_1, α_2 such that $\mathbf{q}_0 = G^{\dagger*} \mathbf{w}$ satisfies

$$\|\mathbf{q}_0\|_2 \leq \alpha_1 \nu m^{1/2} + o(m^{1/2}), \quad (31)$$

$$\|\theta \mathbf{q}_0\|_2 \leq \alpha_2 \exp\left(-\frac{1}{64\nu^2}\right) m^{1/2} + o(m^{1/2}). \quad (32)$$

on the complement of a bad event of probability $\leq e^{-Cm^{1-\eta_0/2}(1+o(1))}$.

Proof: Notice that $G^{\dagger*} = G(G^*G)^{-1} = \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} + \begin{bmatrix} \boldsymbol{\mu}_{J^c} \mathbf{1}^* \\ 0 \end{bmatrix} (G^*G)^{-1}$, where $Z_1 \doteq Z_{J^c, I}$ and $Z_2 \doteq Z_{J^c, J^c}$. Expanding $\mathbf{q}_0 = G^{\dagger*} \mathbf{w}$ gives

$$\begin{aligned} \mathbf{q}_0 &= \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} Z_{J,\bullet}^* \boldsymbol{\sigma} + \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \left(-(G^*G)^{-1} \mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle (G^*G)^{-1} \mathbf{1} \right) \\ &+ \begin{bmatrix} \boldsymbol{\mu}_{J^c} \\ 0 \end{bmatrix} \left(\mathbf{1}^* (G^*G)^{-1} Z_{J,\bullet}^* \boldsymbol{\sigma} - \mathbf{1}^* (G^*G)^{-1} \mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle \mathbf{1}^* (G^*G)^{-1} \mathbf{1} \right). \end{aligned} \quad (33)$$

In this section, we concentrate our efforts on the first term above. In Lemma 7 of Appendix B, we give a more detailed analysis of $(G^*G)^{-1}$, which shows that the remaining terms are all negligible, contributing

$o(m^{1/2})$ to $\|\mathbf{q}_0\|$. This is essentially due to the presence of a large common term $\boldsymbol{\mu}_{J^c}$ in the columns of G : the most significant term in G^*G is $\boldsymbol{\mu}_{J^c}^* \boldsymbol{\mu}_{J^c} \mathbf{1}\mathbf{1}^*$, and $(G^*G)^{-1}$ shrinks $\mathbf{1}$. More precisely, Lemma 7 of Appendix B shows that with probability at least $1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$,

$$\left\| \mathbf{q}_0 - \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} Z_{J^c, \bullet}^* \boldsymbol{\sigma} \right\| \leq Cm^{1/2-\eta_0/4}.$$

This remaining term can be further simplified by splitting out several of the inhomogeneous parts of $(G^*G)^{-1}$. Define $Q \doteq Z_{J^c, \bullet}^* Z_{J^c, \bullet} + \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{I} \end{bmatrix} \in \mathbb{R}^{n \times n}$ and $\boldsymbol{\zeta} \doteq Z_{J^c, \bullet}^* \boldsymbol{\mu}_{J^c} \in \mathbb{R}^n$. In terms of these variables, $G^*G = Q + \boldsymbol{\zeta} \mathbf{1}^* + \mathbf{1} \boldsymbol{\zeta}^* + \alpha \mathbf{1}\mathbf{1}^*$. Applying the matrix inversion lemma,

$$(G^*G)^{-1} = Q^{-1} - Q^{-1/2} M \Xi M^* Q^{-1/2}, \quad (34)$$

where $M = \begin{bmatrix} \frac{Q^{-1/2} \mathbf{1}}{\|Q^{-1/2} \mathbf{1}\|_2} & \frac{Q^{-1/2} \boldsymbol{\zeta}}{\|Q^{-1/2} \boldsymbol{\zeta}\|_2} \end{bmatrix} \in \mathbb{R}^{n \times 2}$, and Ξ is an appropriate 2×2 matrix. Since $\boldsymbol{\vartheta} \doteq Z_{J^c, \bullet}^* \boldsymbol{\sigma} \in \mathbb{R}^n$ is iid $\mathcal{N}(0, \nu^2 \rho)$ independent of G , with high probability it is almost orthogonal to the rank-2 perturbation $\Gamma \doteq Q^{-1/2} M \Xi M^* Q^{-1/2}$: $P[\|\pi_\Gamma \boldsymbol{\vartheta}\| \geq m^{1/2-\eta_0/4}] \asymp e^{-Cm^{1-\eta_0/2}}$.¹¹ Using Fact 1 and block singular value identities, it is not difficult to show¹² that $\|Q^{-1}\| \leq \frac{4}{\nu^2 \bar{\rho}}$ with probability at least $1 - e^{-Cm(1+o(1))}$. Combined with the bound $\|(G^*G)^{-1}\| \leq C_G$ from Lemma 7, we have that $\|\Gamma\| \leq \|(G^*G)^{-1}\| + \|Q^{-1}\| \leq C_G + \frac{4}{\nu^2 \bar{\rho}}$ is bounded by a constant, and

$$\left\| \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \Gamma \boldsymbol{\vartheta} \right\| \leq \left\| \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \right\| \|\Gamma\| \|\pi_\Gamma \boldsymbol{\vartheta}\| \leq \left(1 + 2\nu^2(\sqrt{\bar{\rho}} + \sqrt{\delta})^2\right)^{1/2} \left(C_G + \frac{4}{\nu^2 \bar{\rho}}\right) m^{1/2-\eta_0/4}$$

and the remaining part of \mathbf{q}_0 is

$$\begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} Q^{-1} \boldsymbol{\vartheta} = \begin{bmatrix} Z_1 \\ 0 \end{bmatrix} [Q^{-1}]_{I, \bullet} \boldsymbol{\vartheta} + \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} [Q^{-1}]_{I^c, I} \boldsymbol{\vartheta}_I + \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} [Q^{-1}]_{I^c, I^c} \boldsymbol{\vartheta}_{I^c}.$$

The first two terms involve projections of $\boldsymbol{\vartheta}$ onto k_1 -dimensional subspaces, and hence are of lower order. That is, for $\Sigma \doteq \text{null}([Q^{-1}]_{I, \bullet})^\perp$, we have $P[\|\pi_\Sigma \boldsymbol{\vartheta}\|_2 \geq m^{1/2-\eta_0/4}] \asymp e^{-Cm^{1-\eta_0/2}}$. Since $\|Z_1\|$ and $\|Q^{-1}\|$ are bounded by constants with overwhelming probability, with probability at least $1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$, $\left\| \begin{bmatrix} Z_1 \\ 0 \end{bmatrix} [Q^{-1}]_{I, \bullet} \boldsymbol{\vartheta} \right\| \leq C' m^{1/2-\eta_0/4}$. Identical reasoning shows that on the complement of a bad event of probability $\asymp e^{-Cm^{1-\eta_0/2}}$, $\left\| \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} [Q^{-1}]_{I^c, I} \boldsymbol{\vartheta}_I \right\| \leq C'' m^{1/2-\eta_0/4}$.

This leaves $\begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} [Q^{-1}]_{I^c, I^c} \boldsymbol{\vartheta}_{I^c}$. Expressing Q as $\begin{bmatrix} U & V^* \\ V & W \end{bmatrix}$ and applying the Schur complement formula gives $[Q^{-1}]_{I^c, I^c} = W^{-1} + W^{-1} V (U^{-1} - V^* W^{-1} V)^{-1} V^* W^{-1}$, where $W = Z_2^* Z_2 + \mathbf{I}$, $V = Z_2^* Z_1$, and $U = Z_1^* Z_1$. Because $W \succeq \mathbf{I}$, $\|W^{-1}\| \leq 1$. With probability at least $1 - e^{-Cm(1+o(1))}$, $\|U\| = \|Z_1\|^2 \leq$

¹¹ $\|\pi_\Gamma \boldsymbol{\vartheta}\|$ is distributed as the norm of a 2-dimensional $\mathcal{N}(0, \nu^2 \rho)$ vector. The bound follows from the χ tail bound (21).

¹²Use that $\sigma_{\min}^2 \left(\begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \right) \geq \sigma_{\min}^2(Z_1) - \frac{\|Z_1\|^2 \|Z_2\|^2}{1 - \sigma_{\min}^2(Z_1)}$ and apply Fact 1 to bound each term.

$2\nu^2\bar{\rho}$, $\sigma_{\min}(U) \geq \frac{\nu^2\bar{\rho}}{2}$, and $\|V\| \leq \|Z_1\|\|Z_2\| \leq 2\nu^2(\sqrt{\bar{\rho}\delta} + \bar{\rho})$ and so

$$\|W^{-1}V(U^{-1} - V^*W^{-1}V)^{-1}V^*W^{-1}\| \leq \frac{\|W^{-1}\|^2\|V\|^2}{\sigma_{\min}(U^{-1}) - \|V\|^2\|W^{-1}\|} \leq \frac{8\nu^6(1 + \sqrt{\delta})^2}{1 - 8\nu^6(1 + \sqrt{\delta})^2}$$

is bounded by a constant. Let Σ' denote the k_1 -dimensional range of this matrix. With probability $\geq 1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$, $\|\pi_{\Sigma'}\boldsymbol{\vartheta}\| \leq m^{1/2-\eta_0/4}$, and so

$$\left\| \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} W^{-1}V(U^{-1} - V^*W^{-1}V)^{-1}V^*W^{-1}\boldsymbol{\vartheta} \right\| \leq C''' m^{1/2-\eta_0/4},$$

leaving only $\hat{\mathbf{q}}_0 \doteq \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} (Z_2^*Z_2 + \mathbf{I})^{-1}\boldsymbol{\vartheta}_{I^c}$. With probability at least $1 - e^{-Cm(1+o(1))}$, $\|\boldsymbol{\vartheta}_{I^c}\| \leq \sqrt{2}\nu\sqrt{\bar{\rho}\delta} m^{1/2}$, and so

$$\|\hat{\mathbf{q}}_0\|_2 \leq \left\| \begin{bmatrix} Z_2 \\ \mathbf{I} \end{bmatrix} \right\| \|\boldsymbol{\vartheta}_{I^c}\| \leq \sqrt{1 + \|Z_2\|_2^2} \|\boldsymbol{\vartheta}_{I^c}\| \leq \nu\sqrt{2\delta\rho\left(1 + 2\nu^2\left(\sqrt{\delta} + \sqrt{\bar{\rho}}\right)^2\right)} m^{1/2} \quad (35)$$

establishing the first part of the lemma.

For the second part, we will show that the the upper (R_1) and lower (R_2) parts of $\hat{\mathbf{q}}_0$ can be bounded elementwise by a pair of iid Gaussian vectors. Since for each of these vectors, the Lipschitz function $\|\theta \cdot\|$ is concentrated about its (very small) expectation, the desired result follows. For the upper block, write $Z_2 = QR$, where $Q \in \mathbb{R}^{\bar{\rho}m \times \bar{\rho}m}$ is an orthogonal matrix, and $R \in \mathbb{R}^{\bar{\rho}m \times (\delta m - k_1)}$ is an upper-triangular matrix with non-negative elements on the diagonal. With probability one (as long as $\text{rank}(Z_2) = \bar{\rho}m$), Q and R are uniquely determined by Z_2 . Moreover, Q is a uniform random orthogonal matrix, probabilistically independent of R .¹³ Since $\hat{\mathbf{q}}_0(R_1) = QR(R^*R + \mathbf{I})^{-1}\boldsymbol{\vartheta}_{I^c}$ is the product of a uniform random orthogonal matrix and an independent vector $R(R^*R + \mathbf{I})^{-1}\boldsymbol{\vartheta}_{I^c}$, $\frac{\hat{\mathbf{q}}_0(R_1)}{\|\hat{\mathbf{q}}_0(R_1)\|}$ is uniformly distributed on $\mathbb{S}^{\bar{\rho}m-1}$. With probability $\geq 1 - e^{-Cm(1+o(1))}$, $\|\hat{\mathbf{q}}_0(R_1)\| = \|Z_2(Z_2^*Z_2 + \mathbf{I})\boldsymbol{\vartheta}_{I^c}\| \leq \|Z_2\| \|\pi_{\text{null}(Z_2(Z_2^*Z_2 + \mathbf{I})^{-1})^\perp}\boldsymbol{\vartheta}_{I^c}\| \leq 2\nu^2\sqrt{\bar{\rho}}(\sqrt{\bar{\rho}} + \sqrt{\delta}) m^{1/2}$.¹⁴ Introduce an independent random variable λ_1 distributed as the norm of a $(\bar{\rho}m)$ -dimensional iid $\mathcal{N}(0, \sigma^2)$ vector with $\sigma = 4\nu^2(\sqrt{\bar{\rho}} + \sqrt{\delta})$ (i.e., an appropriately scaled $\chi_{\bar{\rho}m}$ rv), and define

$$\phi_1 \doteq \lambda_1 \frac{\hat{\mathbf{q}}_0(R_1)}{\|\hat{\mathbf{q}}_0(R_1)\|}. \quad (36)$$

Since ϕ_1 is the product of a uniform random unit vector and an appropriate χ random variable, its distribution is iid $\mathcal{N}(0, \sigma^2)$. With probability $1 - e^{-Cm(1+o(1))}$, $\|\phi_1\| \geq \frac{\sigma}{2}\sqrt{\bar{\rho}m} \geq \|\hat{\mathbf{q}}_0(R_1)\|$, so ϕ_1

¹³This follows from the rotational invariance of the Gaussian distribution: left multiplication by an independent orthogonal matrix sampled according to the invariant measure yields an independent pair (Q', R) with $Q'R = Z_2' \doteq_d Z_2$.

¹⁴Here, we have (21) to bound the norm of the projection of $\boldsymbol{\vartheta}$ onto the $(\bar{\rho}m)$ -dimensional subspace $\text{null}(Z_2(Z_2^*Z_2 + \mathbf{I})^{-1})^\perp$.

dominates $\hat{\mathbf{q}}_0(R_1)$ elementwise and $\|\theta\phi_1\| \geq \|\theta\hat{\mathbf{q}}_0(R_1)\|$. Applying Lemma 8 of Appendix B, with probability $1 - e^{-Cm(1+o(1))}$,

$$\|\theta\phi_1\|_2 \leq 4 \exp\left(-\frac{1}{16\sigma^2}\right) \sqrt{\bar{\rho}m} = 4\sqrt{\bar{\rho}} \exp\left(-\frac{1}{256\nu^4(1+\sqrt{\delta})^2}\right) m^{1/2} \leq 4\sqrt{\bar{\rho}} \exp\left(-\frac{1}{64\nu^2}\right). \quad (37)$$

For the lower (R_2) coordinates, write $Z_2^* = [Q_1 \ Q_2] \begin{bmatrix} R_1 \\ 0 \end{bmatrix} \doteq QR$ where $R_1 \in \mathbb{R}^{\bar{\rho}m \times \bar{\rho}m}$ is an upper triangular matrix with nonnegative diagonal elements, Q_1 is an orthogonal matrix, and Q_2 is a random orthobasis for $\mathcal{R}(Q_1)^\perp$ (so that $Q \in \mathbb{R}^{(n-k_1) \times (n-k_1)}$ is an orthogonal matrix). Again from the rotational invariance of the Gaussian distribution, Q is a uniform random orthogonal matrix, independent of R , and

$$\hat{\mathbf{q}}_0(R_2) = (Z_2^*Z_2 + \mathbf{I})^{-1}\boldsymbol{\vartheta}_{I^c} = Q(RR^* + \mathbf{I})^{-1}Q^*\boldsymbol{\vartheta}_{I^c} \doteq Q(RR^* + \mathbf{I})^{-1}\boldsymbol{\gamma}, \quad (38)$$

where $\boldsymbol{\gamma} \doteq Q^*\boldsymbol{\vartheta}_{I^c}$ is an iid $\mathcal{N}(0, \nu^2\rho)$ random vector, *independent of* Q . Hence, $\hat{\mathbf{q}}_0(R_2)$ is the product of a uniform random orthogonal matrix Q , and a probabilistically independent vector $(RR^* + \mathbf{I})^{-1}\boldsymbol{\gamma}$, and its orientation $\frac{\hat{\mathbf{q}}_0(R_2)}{\|\hat{\mathbf{q}}_0(R_2)\|}$ is a uniform random vector on \mathbb{S}^{n-k_1-1} . As above, introduce an independent random variable λ_2 distributed as the norm of an $(n - k_1)$ -dimensional iid $\mathcal{N}(0, 4\nu^2\rho)$ random vector, and define

$$\phi_2 = \lambda_2 \frac{\hat{\mathbf{q}}_0(R_2)}{\|\hat{\mathbf{q}}_0(R_2)\|}. \quad (39)$$

The product of an independent unit vector and (appropriately scaled) χ_{n-k_1} scalar, ϕ_2 is distributed as an iid $\mathcal{N}(0, 4\nu^2\rho)$ vector. With probability at least $1 - e^{-Cm(1+o(1))}$, $\|\phi_2\| \geq \sqrt{2}\nu\sqrt{\rho}\sqrt{n-k_1}$, and $\|\hat{\mathbf{q}}_0(R_2)\| \leq \|\boldsymbol{\vartheta}_{I^c}\| \leq \sqrt{2}\nu\sqrt{\rho}\sqrt{n-k_1}$. Therefore, ϕ_2 dominates $\hat{\mathbf{q}}_0(R_2)$ elementwise, and $\|\theta\phi_2\| \geq \|\theta\hat{\mathbf{q}}_0(R_2)\|$. By Lemma 8,

$$\|\theta\phi_2\|_2 \leq 4\sqrt{\delta} \exp\left(-\frac{1}{64\nu^2\rho}\right) m^{1/2} \leq 4\sqrt{\delta} \exp\left(-\frac{1}{64\nu^2}\right) m^{1/2} \quad (40)$$

Combining the bounds on $\|\theta\phi_1\|$ and $\|\theta\phi_2\|$ gives the second part of the lemma. \blacksquare

III. SIMULATIONS AND EXPERIMENTS

In this section, we perform simulations verifying the conclusions of Theorem 1, and investigating the effect of various model parameters on the error correction capability of the ℓ^1 -minimization (2). In the simulations below we use the publicly available ℓ^1 -magic package [28], except for one (higher-dimensional) face recognition example, which requires a customized interior point method. Since ℓ^1 -recoverability depends only on the signs and support of $(\mathbf{x}_0, \mathbf{e}_0)$, in the simulations below we choose $\mathbf{x}_0(i) \in \{0, 1\}$ and $\mathbf{e}_0(i) \in \{-1, 0, 1\}$. We will judge an output $(\hat{\mathbf{x}}, \hat{\mathbf{e}})$ to be correct if $\max(\|\hat{\mathbf{x}} - \mathbf{x}_0\|_\infty, \|\hat{\mathbf{e}} - \mathbf{e}_0\|_\infty) < 0.01$.

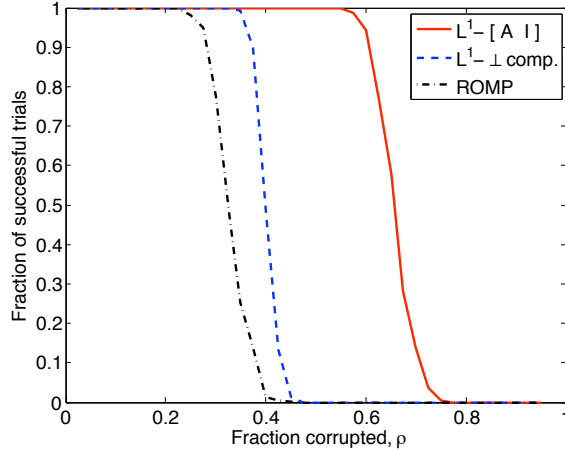


Fig. 5. **Comparison with alternative approaches.** Here, we fix $m = 500$, $\delta = 0.25$, $\nu = 0.05$, and $k_1 = 15$, and compare three approaches to recovering the sparse signal x_0 from error e_0 . The first, denoted “ $L^1 - [A \ I]$ ” solves the extended ℓ^1 minimization advocated in this paper. The second, denoted “ $L^1 - \perp$ comp.” premultiplies by the orthogonal complement of A , and then solves an underdetermined system of linear equations for the sparse error e [14]. The final approach is the greedy Regularized Orthogonal Matching Pursuit (ROMP) [29].

a) *Comparison with alternative approaches:* We first compare the performance of the extended ℓ^1 -minimization

$$\min \|x\|_1 + \|e\|_1 \quad \text{subject to} \quad y = Ax + e$$

to two alternative approaches. The first is the error correction approach of [14], which multiplies by a full rank matrix B such that $BA = 0$,¹⁵ solves

$$\min \|e\|_1 \quad \text{subject to} \quad Be = By,$$

and then subsequently recovers x from the clean system of equations $Ax = y - e$. The second is the Regularized Orthogonal Matching Pursuit (ROMP) algorithm [29], a state-of-the-art greedy method for recovering sparse signals.¹⁶ For this algorithm, we use the implementation from <http://math.ucdavis.edu/~dneedell/>.

For this experiment, the ambient dimension is $m = 500$; the parameters of the CAB model are $\nu = 0.05$ and $\delta = 0.25$. We fix the signal support to be $k_1 = 15$, and vary the fraction of errors from 0 to 0.95.

¹⁵This comparison requires $n \ll m$ although our method is not limited to this case.

¹⁶For the models considered here, less sophisticated greedy methods such as the standard orthogonal matching pursuit fail even for small error fractions.

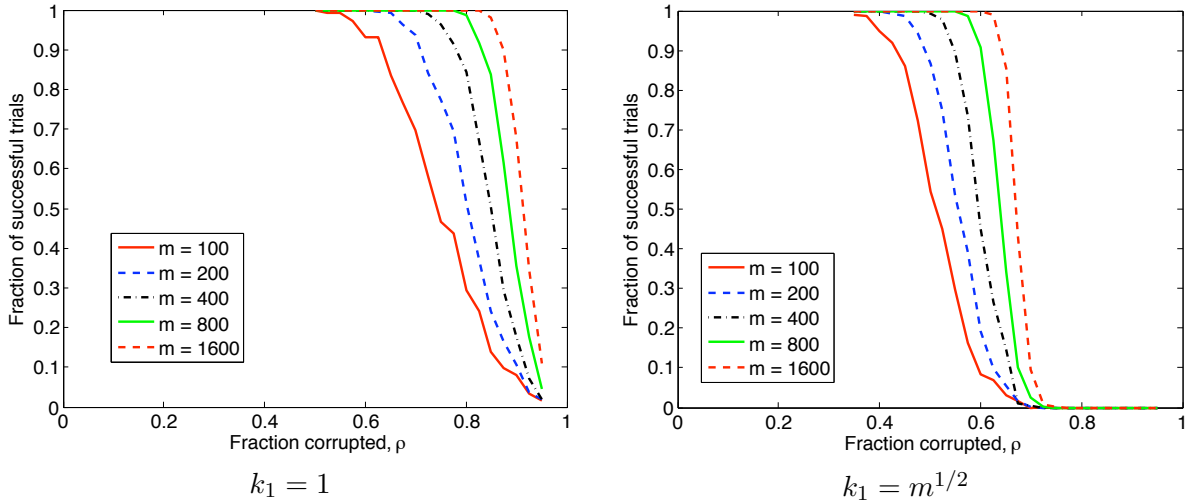


Fig. 6. **Error correction in weak proportional growth.** We fix $\delta = 0.25$, $\nu = 0.05$, and plot the fraction of successful recoveries as a function of the error density ρ , for each $m = 100, 200, 400, 800, 1600$. At left, k_1 is fixed at 1; at right, $k_1 = m^{1/2}$. In both cases, as m increases, the fraction of errors that can be corrected approaches 1.

For each error fraction, we generate 500 independent problems. Figure 5 plots the fraction of successes for each of the three algorithms, as a function of error density ρ . There the extended ℓ^1 -minimization is denoted “ $L^1 - [A \ I]$ ” (red curve), while the alternative approach of [14] is denoted “ $L^1 - \perp \text{comp}$ ” (blue curve). Whereas both ROMP and the ℓ^1 approach of [14] break down around 40% corruption, the extended ℓ^1 -minimization continues to succeed with high probability even beyond 60% corruption.

b) Error correction capacity: While the previous experiment demonstrates the advantages of the extended ℓ^1 -minimization (2) for the CAB model, Theorem 1 suggests that more is true: As the dimension increases, the fraction of errors that the extended ℓ^1 -minimization can correct should approach one. We generate problem instances with $\delta = 0.25$, $\nu = 0.05$, for varying $m = 100, 200, 400, 800, 1600$. For each problem size, and for each error fraction $\rho = 0.05, 0.1, \dots, 0.95$, we generate 500 random problems, and plot the fraction of correct recoveries in Figure 6. At left, we fix $k_1 = 1$, while at right, k_1 grows as $k_1 = m^{1/2}$. In both cases, as m increases, the fraction of errors that can be corrected also increases.

c) Varying model parameters: We next investigate the effect of varying δ (Figure 7 left) and ν (Figure 7 right). We first fix $m = 400$, $\nu = .3$, and consider different bouquet sizes $n = 100, 200, 300, 400, 500$. Figure 7 left plots the fraction of correct trials for varying error densities ρ , for each of these bouquet sizes. For this fixed m , the error correction capability decreases only slightly as n increases.

We next fix $m = 400$, $n = 200$, and consider the effect of varying ν . Figure 7 plots the result for $\nu = .1, .3, .5, .7, .9$. Notice that as ν decreases (i.e., the bouquet becomes tighter), the error correction

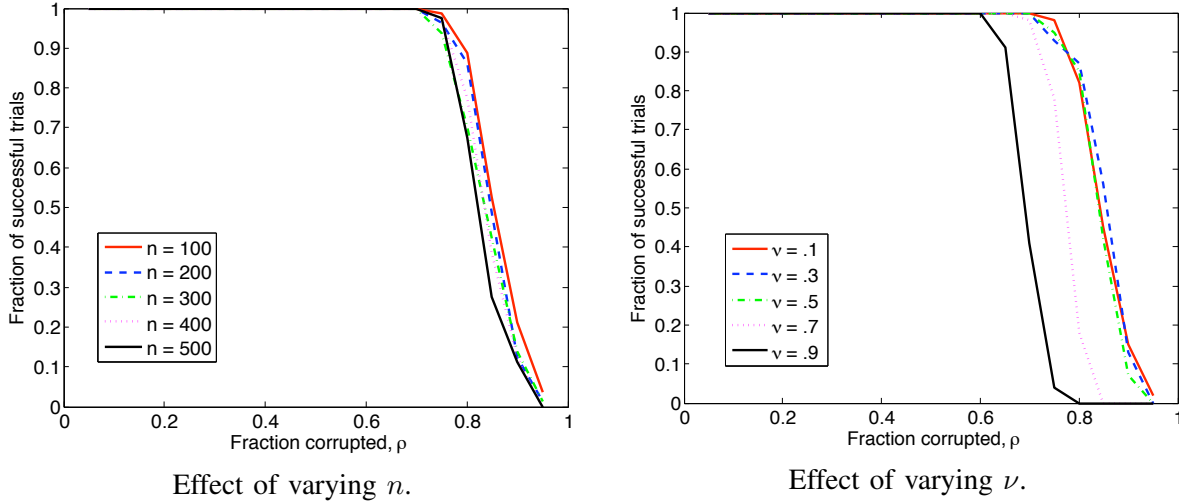


Fig. 7. **Effect of varying n and ν .** At left, we fix $m = 400$, $\nu = .3$, and consider varying $n = 100, 200, \dots, 500$. For each of these model settings, we plot the fraction of correct recoveries as a function of the fraction of errors. Notice that the error correction capacity decreases only slightly as n increases. At right, we fix $m = 400$, $n = 200$, and vary ν from $.1$ to $.9$. Again, we plot the fraction of correct recoveries for each error fraction. As expected from Theorem 1, as ν decreases, the error correction capacity of ℓ^1 increases.

capacity increases: for any fixed fraction of successful trials, the fraction of error that can be corrected increases by approximately 15% as ν decreases from $.9$ to $.5$.

d) Phase transition in total proportional growth: Theorem 1 does not provide any explicit information about the behavior of ℓ^1 -minimization when the signal support k_1 grows proportionally to m : $k_1/m \rightarrow \rho_1 \in (0, 1)$. Based on intuition from more homogeneous polytopes (especially the work of Donoho and Tanner on Gaussian matrices [24]), we might expect that when k_1 also exhibits proportional growth, an asymptotically sharp phase transition between guaranteed recovery and guaranteed failure will occur at some critical error fraction $\rho^* \in (0, 1)$. We investigate this empirically here by again setting $\delta = 0.25$, $\nu = 0.05$, but this time allowing $k_1 = 0.05m$. Figure 8 plots the fraction of correct recovery for varying error fractions ρ , as m grows: $m = 100, 200, 400, 800, 1600$. In this proportional growth setting, we see an increasingly sharp phase transition, near $\rho = 0.6$.

e) Error correction with real face images: Finally, we return to the motivating example of face recognition under varying illumination and random corruption. For this experiment, we use the Extended Yale B face database [15], which tests illumination sensitivity of face recognition algorithms. As in [11], we form the matrix A from images in Subsets 1 and 2, which contain mild-to-moderate illumination variations. Each column of the matrix A is a $w \times h$ face image, stacked as a vector in \mathbb{R}^m ($m = w \times h$).

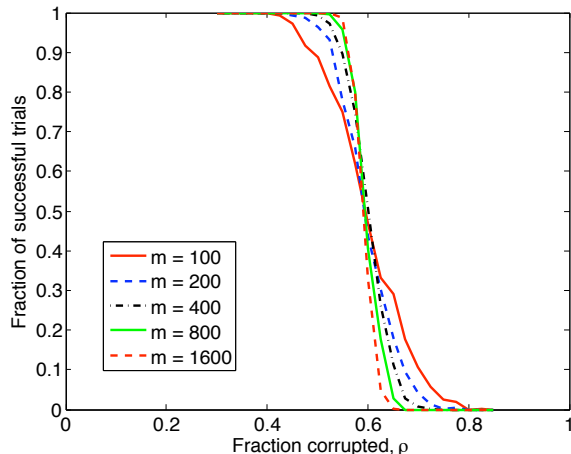


Fig. 8. **Phase transition in total proportional growth.** When the signal support grows in proportion to the dimension ($k_1/m \rightarrow \rho_1 \in (0, 1)$), we observe an asymptotically sharp phase transition in the probability of correct recovery, similar to that investigated in [24]. Here, for $\delta = 0.25$, $\nu = 0.05$, $k_1 = 0.05m$, we indeed see a sharp phase transition at $\rho = 0.6$.

Here, the weak proportional growth setting corresponds to the case when the total number of image pixels grows proportionally to the number n of face images. Since the number of images per subject is fixed, this is the same as the total image resolution growing proportionally to the number of subjects. We vary the image resolutions through the range 34×30 , 48×42 , 68×60 , 96×84 .¹⁷ The matrix A is formed from images of 4, 9, 19, 38 subjects, respectively, corresponding to $\delta \approx 0.09$. Here, $\nu \approx 0.3$. In face recognition, the sublinear growth of $\|\mathbf{x}_0\|_0$ comes from the fact that the observation should ideally be a linear combination of only images of the same subject. Various estimates of the required number of images, k_1 , appear in the literature, ranging from 5 to 9. Here, we fix $k_1 = 7$, and generate the (clean) test image synthetically as a linear combination of k_1 training images from a single subject. The reason for using synthetic linear combinations as opposed to real test images is simply that it allows us to verify whether \mathbf{x}_0 was correctly recovered; in the real data experiments of the introduction of this paper and of [11], success could only be judged in terms of the recognition rate of the entire classification pipeline.

For each resolution considered, and for each error fraction, we generate 75 trials. Figure 9 (left) plots the fraction of successes as a function of the fraction of corruption. Notice that as predicted by Theorem 1, the fraction of errors that can be corrected again approaches 1 as the data size increases. Figure 9 (right) gives a visual demonstration of the algorithm's capability. In the test images in Figure 9 (right,

¹⁷Thus, the total dimension $m = 1020, 2016, 4080, 8064$ grows roughly by a factor of 2 from one curve to the next, similar to the simulations above.

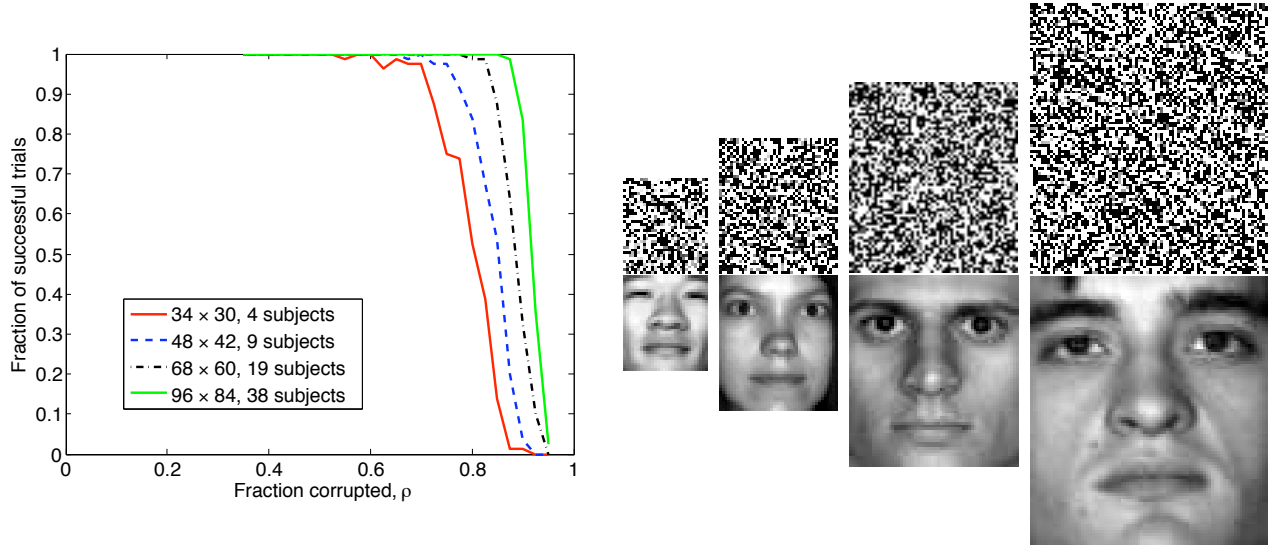


Fig. 9. **Error correction with real face images.** We simulate weak proportional growth in the Extended Yale B face database, with the resolution of the images growing in proportion to the number of subjects. Left: fraction of correct recoveries for varying levels of occlusion. Right: examples of correct recovery for each resolution considered. Top: corrupted test image. The fraction of corruption is chosen so that the probability of correct recovery is 50%. Bottom: clean image, from correctly recovered x_0 .

top), the amount of corruption is chosen to correspond to a 50% probability of success according to the plots in Figure 9 (left). Below each corrupted test image, the “clean” image recovered by our method is shown.

IV. DISCUSSIONS AND FUTURE WORK

a) Compressed sensing for signals with varying sparsity: In the conventional setting for recovering a sparse signal, one often implicitly assumes that each entry of the signal has an equal probability of being nonzero. As a result, one typically requires that the incoherence (or coherence) of the dictionary is somewhat uniform. In this paper, we saw quite a different example. If we view both x and e as the signal that we want to recover, then the sparsity or density of the combined signal is quite uneven – x is very sparse but e can be very dense. Nevertheless, our result suggests that if the incoherence of the dictionary is adaptive to the distribution of the density – more coherent for the sparse part and less for the dense part, then ℓ^1 -minimization will be able to recover such uneven signals even if bounds based on the even sparsity assumption suggest otherwise. Thus, if one has some prior knowledge about which part of the signal is likely to be more sparse or more dense, one can achieve much better performance with ℓ^1 -minimization by using a dictionary with matching incoherence. More generally, for any given distribution of sparsity, one may ask the question whether there exists an optimal dictionary with matching

incoherence such that ℓ^1 -minimization has the highest chance of success.

b) Stability with respect to noise: Although in our model, we do not explicitly consider any noise (say $\mathbf{y} = A\mathbf{x} + \mathbf{e} + \mathbf{z}$, where \mathbf{z} is Gaussian noise), ℓ^1 -minimization is known to be stable under small noise [26]. This is also what we have observed empirically in our simulations and also in experiments with face images: ℓ^1 -minimization for the cross-and-bouquet model is surprisingly stable to measurement or numerical noise. In fact, as the method is able to deal with dense errors regardless of their magnitude, large noisy entries in \mathbf{z} will be treated like errors and be absorbed into \mathbf{e} . However, a more precise characterization of the effect of noise (say Gaussian) on the estimate of the sparse signal \mathbf{x} and the error \mathbf{e} remains an open problem.

c) Neighborliness of polytopes: As we have seen in this paper, a precise characterization of the performance of ℓ^1 -minimization requires us to analyze the geometry of polytopes associated with the specific dictionaries in question. In practice, we often use ℓ^1 -minimization for purposes other than signal reconstruction or error correction. For instance, using machine learning techniques, we can learn from exemplars a dictionary that is optimal for certain tasks such as data classification [13]. The polytope associated with such a dictionary may be very different from those that are normally studied in signal processing or coding theory or error correction, leading to qualitatively different behavior of the ℓ^1 -minimization. Thus, we should expect that in the coming years, many new classes of high-dimensional polytopes with even more interesting properties may arise from other applications and practical problems.

ACKNOWLEDGMENTS

The authors would like to acknowledge helpful conversations with and useful comments from Prof. Robert Fossum (UIUC Math), Prof. Olgica Milenkovic (UIUC ECE), Prof. Sean Meyn (UIUC ECE), and Dr. Gang Hua (Microsoft Live Labs). This work is partially supported by grants NSF CRS-EHS-0509151, NSF CCF-TF-0514955, ONR YIP N00014-05-1-0633, and NSF IIS 07-03756. John Wright is also supported by a Microsoft Fellowship (sponsored by Microsoft Live Labs, Redmond). Finally, Yi Ma would like to thank Microsoft Research Asia, Beijing, China, for its hospitality during his visit there in Summer 2008.

REFERENCES

- [1] T. Figiel, J. Lindenstrauss, and V. D. Milman, "The dimension of almost spherical sections of convex bodies," *Acta Math.*, vol. 139, no. 1-2, pp. 53–94, 2008.
- [2] B. S. Kashin, "The widths of certain finite-dimensional sets and classes of smooth functions," *Izv. Akad. Nauk SSSR Serv. Mat.*, vol. 41, no. 2, pp. 334–351, 2008.
- [3] M. Ledoux, *The Concentration of Measure Phenomenon, Mathematical Surveys and Monographs 89*. American Mathematical Society, 2001.
- [4] R. Tibshirani, "Regression shrinkage and selection via the Lasso," *Journal of the Royal Statistical Society Series B*, vol. 58, pp. 267–288, 1996.
- [5] W. Fu and K. Knight, "Asymptotics for Lasso-type estimators," *Annals of Statistics*, vol. 28, no. 5, pp. 1356–1378, 2000.
- [6] N. Meinshausen and B. Yu, "Lasso-type recovery of sparse representations for high-dimensional data," *to appear in Annals of Statistics*, 2006.
- [7] R. Berinde, A. C. Gilbert, P. Indyk, H. Karloff, and M. J. Strauss, "Combining geometry and combinatorics: A unified approach to sparse signal recovery," *Preprint*, 2008.
- [8] V. Guruswami, J. R. Lee, and A. Razborov, "Almost euclidean subspaces of ℓ_1^n via expander codes," *Electronic Colloquium on Computational Complexity, Report No. 86*, 2007.
- [9] A. Bruckstein, D. L. Donoho, and M. Elad, "From sparse solutions of systems of equations to sparse modeling of signals and images," *To appear in SIAM Review*, 2008.
- [10] W. Bajwa, J. Haupt, G. Raz, and R. Nowak, "Compressed channel sensing," in *Proceedings of Conference on Information Sciences and Systems*, 2008.
- [11] J. Wright, A. Yang, A. Ganesh, S. Sastry, and Y. Ma, "Robust face recognition via sparse representation," *To appear in IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2008.
- [12] J. Yang, J. Wright, T. Huang, and Y. Ma, "Image super-resolution as sparse representation of raw image patches," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [13] J. Mairal, F. Bach, J. Ponce, G. Sapiro, and A. Zisserman, "Discriminative learned dictionaries for local image analysis," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [14] E. Candès and T. Tao, "Decoding by linear programming," *IEEE Trans. Information Theory*, vol. 51, no. 12, 2005.
- [15] K. Lee, J. Ho, and D. Kriegman, "Acquiring linear subspaces for face recognition under variable lighting," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 27, no. 5, pp. 684–698, 2005.
- [16] D. 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.
- [17] J. Tropp, "Greed is good: algorithmic results for sparse approximation," *IEEE Trans. Information Theory*, vol. 50, no. 10, pp. 2231–2242, 2004.
- [18] D. Needell and J. Tropp, "CoSAMP: Iterative signal recovery from incomplete and inaccurate samples," *To appear in Applied and Computational Harmonic Analysis*, 2008.
- [19] Y. Eldar and M. Mishali, "Robust recovery of signals from a union of subspaces," *preprint*, 2008.
- [20] M. Sipser and D. A. Spielman, "Expander codes," *IEEE Transactions on Information Theory*, vol. 42, no. 6, 1996.
- [21] J. Feldman, T. Malkin, R. Servedio, C. Stein, and M. Wainwright, "LP decoding corrects a constant fraction of errors," *IEEE Trans. Information Theory*, vol. 53, no. 1, pp. 82–89, 2007.

- [22] M. R. Capalbo, O. Reingold, S. P. Vadhan, and A. Wigderson, “Randomness conductors and constant-degree lossless expanders,” in *Proceedings of the 34th ACM Symposium on Theory of Computing*, 2002, pp. 659–668.
- [23] N. Kashyap, “A decomposition theorem for binary linear codes,” *IEEE Trans. Information Theory*, vol. 54, no. 7, pp. 3035–3058, 2008.
- [24] D. Donoho and J. Tanner, “Counting faces of randomly projected polytopes when the projection radically lowers dimension,” preprint, <http://www.math.utah.edu/tanner/>, 2007.
- [25] D. Donoho, “Neighborly polytopes and sparse solution of underdetermined linear equations,” preprint, 2005.
- [26] —, “For most large underdetermined systems of linear equations the minimal ℓ^1 -norm near solution approximates the sparsest solution,” preprint, 2004.
- [27] S. Dasgupta, D. Hsu, and N. Verma, “A concentration theorem for projections,” in *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2006.
- [28] E. Candes and J. Romberg, “ ℓ^1 -magic: Recovery of sparse signals via convex programming,” <http://www.acm.caltech.edu/l1magic/>, 2005.
- [29] D. Needell and R. Vershynin, “Signal recovery from inaccurate and incomplete measurements via regularized orthogonalized matching pursuit,” preprint <http://www.math.ucdavis.edu/~dneedell/>, 2007.
- [30] P. Wedin, “Perturbation bounds in connection with singular value decomposition,” *BIT Numerical Mathematics*, vol. 12, pp. 99–111, 1972.
- [31] N. Alon and J. Spencer, *The Probabilistic Method*. Wiley-Interscience, 2001.
- [32] T. Ferguson, *A Course in Large Sample Theory*. Chapman and Hall, 1996.

APPENDIX

TECHNICAL LEMMAS AND RESULTS

A. Restricted Isometry for Sparse Vectors

Here, we give a more precise statement of the restricted isometry property of $[\tilde{V}^* - \tilde{S}\tilde{U}^*]$ used in the proof of Lemma 3. For an arbitrary matrix M , we defined $\gamma_k(M) \doteq \inf_{\|\mathbf{y}\|_0 \leq k, \mathbf{y} \neq \mathbf{0}} \frac{\|M\mathbf{y}\|_2}{\|\mathbf{y}\|_2}$. We are interested in knowing $\gamma_{cm}([\tilde{V}^* - \tilde{S}\tilde{U}^*])$, where \tilde{U} , \tilde{S} , and \tilde{V} come from a (compact) singular value decomposition¹⁸ of $P \doteq A_2^* \pi_{A_1^\perp}$, after dropping the largest singular value. The constants in the following result are less important than the fact that for c sufficiently small, $\gamma_{cm} = \Omega(\nu)$.

Lemma 5 (Restricted Isometry): Suppose that $\bar{\rho} < \delta$, $\nu < 1/9$, and c is sufficiently small:

$$c \leq \min \left\{ \frac{\bar{\rho}}{1024}, \frac{\bar{\rho}}{64(1 + 2C_\mu \bar{\rho}^{-1/2})^2} \right\}, \quad \bar{\rho}H(c/\bar{\rho}) + \delta H(c/\delta) < \frac{\bar{\rho}}{128\pi^2}, \quad (41)$$

where $H(\cdot)$ is the base- e binary entropy function. Let $\mathbf{u}_1, \mathbf{v}_1$ denote the first singular vectors of $P \doteq A_2^* \pi_{A_1^\perp} \in \mathbb{R}^{(\delta m - k_1) \times \bar{\rho} m}$. Then if $\tilde{U}\tilde{S}\tilde{V}^*$ is a compact singular value decomposition of $\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp}$,

$$\gamma_{cm}([\tilde{V}^* - \tilde{S}\tilde{U}^*]) \geq \frac{\nu\sqrt{\bar{\rho}}}{16} \quad (42)$$

on the complement of a bad event of probability $\leq e^{-Cm(1+o(1))}$.

Proof: Notice that the conditional distribution of P given A_1 is Gaussian: $P = Z_2^* \pi_{A_1^\perp} + \mathbf{1}\mu_{J_c}^* \pi_{A_1^\perp} \doteq Z_2^* \pi_{A_1^\perp} + \mathbf{1}\hat{\mu}^*$. We argue that the second term dominates:

a) $\mathbf{1}\hat{\mu}^*$ determines the leading singular vectors: Since the columns of A_1 are k_1 small perturbations of μ_{J_c} , the residual $\|\hat{\mu}\| = \|\pi_{A_1^\perp} \mu_{J_c}\|$ should be small. However, we will see that it is not too small: $\|\pi_{A_1^\perp} \mu_{J_c}\| = \Omega(k_1^{-1/2})$. Choose an orthonormal basis for $\mathbb{R}^{\bar{\rho}m}$, with first basis vector $\frac{\mu_{J_c}}{\|\mu_{J_c}\|}$. The expression of A_1 w.r.t. this basis is $\begin{bmatrix} \mathbf{0} \\ B \end{bmatrix} + \mathbf{e}_1(\mathbf{c}^* + \|\mu_{J_c}\|\mathbf{1}^*) \doteq \begin{bmatrix} \mathbf{0} \\ B \end{bmatrix} + \mathbf{e}_1\mathbf{v}^*$, where B and \mathbf{c} are iid $\mathcal{N}(0, \nu^2/m)$. So, $\left\| \pi_{A_1} \frac{\mu_{J_c}}{\|\mu_{J_c}\|} \right\|_2^2$ can be written as

$$\mathbf{e}_1^* \left(\begin{bmatrix} \mathbf{0} \\ B \end{bmatrix} + \mathbf{e}_1\mathbf{v}^* \right) (\mathbf{v}\mathbf{v}^* + B^*B)^{-1} \left(\begin{bmatrix} \mathbf{0} \\ B^* \end{bmatrix} + \mathbf{v}\mathbf{e}_1^* \right) \mathbf{e}_1 = \frac{\mathbf{v}^*(B^*B)^{-1}\mathbf{v}}{1 + \mathbf{v}^*(B^*B)^{-1}\mathbf{v}}.$$

Applying Fact 1 to the $(\bar{\rho}m - 1) \times k_1$ matrix B , one can easily show that $P \left[\|(B^*B)^{-1}\| > \frac{2}{\nu^2\bar{\rho}} \right] \asymp e^{-Cm}$. By (21) above, the norm of the k_1 -dimensional $\mathcal{N}(0, \nu^2/m)$ vector \mathbf{c} also concentrates: $P \left[\|\mathbf{c}\| > \sqrt{k_1} \right] \asymp e^{-C'mk_1}$. On the complement of these bad events, $\|\mathbf{v}\| \leq \|\mathbf{c}\| + \|\mu_{J_c}\mathbf{1}_{k_1}^*\| = (1 + \|\mu_{J_c}\|)\sqrt{k_1} \leq 2\sqrt{k_1}$,

¹⁸With probability one, the matrices U and V are unique upto multiplication of their columns by a common set of signs. The quantity of interest, γ_k , does not depend on the choice of signs, so there is no ambiguity in writing $\gamma_k([\tilde{V}^* - \tilde{S}\tilde{U}^*])$.

and $\mathbf{v}^*(B^*B)^{-1}\mathbf{v} \leq \frac{8}{\nu^2\bar{\rho}}k_1$. So,

$$\left\| \frac{\boldsymbol{\mu}_{J^c}}{\|\boldsymbol{\mu}_{J^c}\|} - \pi_{A_1} \frac{\boldsymbol{\mu}_{J^c}}{\|\boldsymbol{\mu}_{J^c}\|} \right\|_2^2 = \frac{1}{1 + \mathbf{v}^*(B^*B)^{-1}\mathbf{v}} \geq \frac{1}{1 + \frac{8}{\nu^2\bar{\rho}}k_1}. \quad (43)$$

Lemma 6 below shows that with probability $\geq 1 - e^{-Cm(1+o(1))}$ in the random support of the error \mathbf{e}_0 , $\|\boldsymbol{\mu}_{J^c}\| \geq \bar{\rho}/2$. Together with (43), this implies that $\|\hat{\boldsymbol{\mu}}\| = \|\boldsymbol{\mu}_{J^c} - \pi_{A_1}\boldsymbol{\mu}_{J^c}\|_2 \geq \frac{\bar{\rho}}{2} \sqrt{\frac{1}{1 + \frac{8}{\nu^2\bar{\rho}}k_1}}$. On this good event, $\|\mathbf{1}_{\delta m - k_1} \hat{\boldsymbol{\mu}}^*\|_2 \geq C_1 m^{\eta_0/2}$ for some constant C_1 and m sufficiently large. From Fact 1, $\|Z_2\|$ is bounded by some constant C_2 with probability at least $1 - e^{-Cm(1+o(1))}$. Treating $Z_2^* \pi_{A_1^\perp}$ as a nuisance perturbation of $\mathbf{1} \hat{\boldsymbol{\mu}}^*$ and applying Wedin's perturbation bound for principal subspaces [30] then gives

$$\begin{aligned} \|\pi_{\mathbf{u}_1^\perp} - \pi_{\mathbf{1}^\perp}\| &= \left\| \frac{\mathbf{u}_1 \mathbf{u}_1^*}{\mathbf{u}_1^* \mathbf{u}_1} \left(\mathbf{I} - \frac{\mathbf{1} \mathbf{1}^*}{\mathbf{1}^* \mathbf{1}} + \frac{\mathbf{1} \mathbf{1}^*}{\mathbf{1}^* \mathbf{1}} \right) - \left(\mathbf{I} - \frac{\mathbf{u}_1 \mathbf{u}_1^*}{\mathbf{u}_1^* \mathbf{u}_1} + \frac{\mathbf{u}_1 \mathbf{u}_1^*}{\mathbf{u}_1^* \mathbf{u}_1} \right) \frac{\mathbf{1} \mathbf{1}^*}{\mathbf{1}^* \mathbf{1}} \right\| \\ &\leq 2 \left\| \frac{\mathbf{u}_1 \mathbf{u}_1^*}{\mathbf{u}_1^* \mathbf{u}_1} \left(\mathbf{I} - \frac{\mathbf{1} \mathbf{1}^*}{\mathbf{1}^* \mathbf{1}} \right) \right\| \leq \frac{2 \|Z_2^* \pi_{A_1^\perp}\|}{\|\mathbf{1} \hat{\boldsymbol{\mu}}^*\|} \leq \frac{2 C_2}{C_1 m^{\eta_0/2}}. \end{aligned}$$

Similarly $\|\pi_{\mathbf{v}_1^\perp} - \pi_{\hat{\boldsymbol{\mu}}^\perp}\| \leq \frac{2 C_2}{C_1 m^{\eta_0/2}}$. Write

$$\|\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp} - \pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp}\| \leq \|\pi_{\mathbf{u}_1^\perp} - \pi_{\mathbf{1}^\perp}\| \|P \pi_{\mathbf{v}_1^\perp}\| + \|\pi_{\mathbf{1}^\perp} P\| \|\pi_{\mathbf{v}_1^\perp} - \pi_{\hat{\boldsymbol{\mu}}^\perp}\|.$$

Now, $\|\pi_{\mathbf{1}^\perp} P\| \leq \|Z_2\| \leq C_2$, and $\|P \pi_{\mathbf{v}_1^\perp}\| = \sigma_2(P) \leq \sqrt{2}\nu(\sqrt{\bar{\rho}} + \sqrt{\delta})$ simultaneously with probability $\geq 1 - e^{-Cm(1+o(1))}$ (the second bound was established in part (a) of the proof of Lemma 3). Hence, $\exists C_3$ such that $P[\|\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp} - \pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp}\|_2 > C_3 m^{-\eta_0/2}] \asymp e^{-Cm}$. For an arbitrary matrix W , let $f(W) \doteq \gamma_{cm}([\pi_{\mathcal{R}(W^*)} - W^*])$. We are interested in $f(\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp})$.¹⁹ Using the fact that singular values of submatrices are 1-Lipschitz and applying Wedin's sin Θ theorem [30] to $\pi_{\mathcal{R}(W^*)}$, it is not difficult to show that if $\text{rank}(W + \Delta) = \text{rank}(W)$,

$$|f(W + \Delta) - f(W)| \leq \left(\frac{1}{\sigma_{\min}(W) - \|\Delta\|} + 1 \right) \|\Delta\|, \quad (44)$$

where $\sigma_{\min}(W)$ is the smallest nonzero singular value. Applying this bound with $W = \pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp}$, $\Delta = \pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp} - \pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp}$, and noticing that $\sigma_{\min}(\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp})$ is bounded below by a positive constant with overwhelming probability, we have that $|f(\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp}) - f(\pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp})| < \frac{\nu\sqrt{\bar{\rho}}}{16}$ with probability at least $1 - e^{-Cm(1+o(1))}$. We henceforth restrict our attention to $f(\pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp})$.

b) Analysis via Gaussian measure concentration: Let Σ denote the subspace $(\mathcal{R}(Z_1) + \mathcal{R}(\boldsymbol{\mu}_{J^c}))^\perp$, and let V_0 be some orthonormal basis for this subspace, chosen independently of Z_2 . From the above reasoning, we can restrict our attention to $\pi_{\mathbf{1}^\perp} P \pi_{\hat{\boldsymbol{\mu}}^\perp} = \pi_{\mathbf{1}^\perp} Z_2^* \pi_\Sigma$. Let $\pi_{\mathbf{1}^\perp} Z_2^* \pi_\Sigma = U' S' V'^*$ be a compact singular value decomposition of this matrix. Then,

$$\gamma_{cm} \left(\begin{bmatrix} V'^* & -S' U'^* \end{bmatrix} \right) = \gamma_{cm} \left(V'^* \begin{bmatrix} \mathbf{I} & \pi_\Sigma Z_2 \pi_{\mathbf{1}^\perp} \end{bmatrix} \right) = \gamma_{cm} \left(V_0^* \begin{bmatrix} \mathbf{I} & \pi_\Sigma Z_2 \pi_{\mathbf{1}^\perp} \end{bmatrix} \right).$$

¹⁹Since left multiplication by an orthogonal matrix does not change γ_{cm} , $f(\pi_{\mathbf{u}_1^\perp} P \pi_{\mathbf{v}_1^\perp}) = \gamma_{cm}([\tilde{V}^* - \tilde{S} \tilde{U}^*])$.

Where the final step follows because γ_{cm} is invariant under left multiplication of its argument by an orthogonal matrix. Now, $V_0^* \pi_\Sigma Z_2 = V_0^* Z_2$ is simply distributed as a $(\bar{\rho}m - k_1 - 1) \times (\delta m - k_1)$ iid $\mathcal{N}(0, \nu^2/m)$ random matrix. Finally, introduce an additional uniformly distributed random orthogonal matrix $Q \in \mathbb{R}^{(\bar{\rho}m - k_1 - 1) \times (\bar{\rho}m - k_1 - 1)}$, chosen independently of Z_2 , and define $\Psi \doteq QV_0^* \pi_\Sigma Z_2$. This is again an iid $\mathcal{N}(0, \nu^2/m)$ matrix. Notice then, that $\gamma_{cm} \left(\begin{bmatrix} V'^* & -S'U'^* \end{bmatrix} \right) = \gamma_{cm} \left(\begin{bmatrix} QV_0^* & \Psi\pi_{1^\perp} \end{bmatrix} \right)$. From the rotational invariance of the Gaussian distribution, it is easy to show that Ψ and Q are independent random variables. QV_0^* is the transpose of random orthobasis for Σ ; it can be realized by orthogonalizing the projection of a Gaussian matrix onto Σ . To this end, introduce an iid $\mathcal{N}(0, \nu^2/m)$ matrix $\Phi \in \mathbb{R}^{(\bar{\rho}m - k_1 - 1) \times \bar{\rho}m}$ independent of Σ and Ψ . Then, $\gamma_{cm} \left(\begin{bmatrix} QV_0^* & \Psi\pi_{1^\perp} \end{bmatrix} \right)$ is equal in distribution to $\gamma_{cm} \left(\begin{bmatrix} (\Phi\pi_\Sigma\Phi^*)^{-1/2} \Phi\pi_\Sigma & \Psi\pi_{1^\perp} \end{bmatrix} \right)$. Let $\Lambda \doteq (\Phi\pi_\Sigma\Phi^*)^{-1/2}$, and notice that

$$\begin{aligned} \gamma_{cm} &= \min_{\#L_1 \cup L_2 = cm} \sigma_{\min} \left(\begin{bmatrix} [\Lambda\Phi\pi_\Sigma]_{\bullet, L_1} & [\Psi\pi_{1^\perp}]_{\bullet, L_2} \end{bmatrix} \right) \geq \\ &\min_{\#L_1 = \#L_2 = cm} \min \left\{ \sigma_{\min}([\Lambda\Phi\pi_\Sigma]_{\bullet, L_1}), \sigma_{\min}(\pi_{\Sigma'^\perp} [\Psi\pi_{1^\perp}]_{\bullet, L_2}) \right\} - \max_{\#L_1 = \#L_2 = cm} \left\| \pi_{\Sigma'} [\Psi\pi_{1^\perp}]_{\bullet, L_2} \right\| \end{aligned}$$

where Σ' denotes the subspace $\mathcal{R}([\Lambda\Phi\pi_\Sigma]_{\bullet, L_1})$.

c) *Bounding $\sigma_{\min}[\Lambda\Phi\pi_\Sigma]_{\bullet, L}$* : Applying Fact 1 to $\Phi\pi_\Sigma$ gives that $P[\|\Phi\pi_\Sigma\|_2 \geq 3\nu\sqrt{\bar{\rho}}] \asymp e^{-\bar{\rho}m/2}$. On the complement of this bad event, $\sigma_{\min}(\Lambda) \geq \frac{1}{3\nu\sqrt{\bar{\rho}}}$. Write

$$\begin{aligned} [\Phi\pi_\Sigma]_{\bullet, L} &= \Phi_{\bullet, L} - [\Phi\pi_{\Sigma^\perp}]_{\bullet, L} = \Phi_{\bullet, L}(\mathbf{I} - [\pi_{\Sigma^\perp}]_{L, L}) - \Phi_{\bullet, L^c}[\pi_\Sigma]_{L^c, L} \\ \implies \sigma_{\min}([\Phi\pi_\Sigma]_{\bullet, L}) &\geq \sigma_{\min}(\Phi_{\bullet, L})(1 - \|\pi_{\Sigma^\perp}\|_{L, L}) - \|\pi_{\Phi_{\bullet, L}} \Phi_{\bullet, L^c} \pi_{[\pi_\Sigma]_{L^c, L}}\|. \end{aligned}$$

Straightforward application of Fact 1 shows that $P[\sigma_{\min}(\Phi_{\bullet, L}) \leq \frac{\nu\sqrt{\bar{\rho}}}{2} - \nu\sqrt{c}] \asymp e^{-\bar{\rho}m/8}$, while for any²⁰ $\varepsilon_1 > 0$, $P[\|\pi_{\Phi_{\bullet, L}} \Phi_{\bullet, L^c} \pi_{[\pi_\Sigma]_{L^c, L}}\| \geq 2\nu\sqrt{c} + \nu\sqrt{\bar{\rho}\varepsilon_1}] \asymp e^{-\bar{\rho}\varepsilon_1 m/2}$. Finally, consider the matrix $\Upsilon \doteq \begin{bmatrix} Z_1 & \nu\sqrt{\bar{\rho}} \frac{\mu_{J^c}}{\|\mu_{J^c}\|} \end{bmatrix} \in \mathbb{R}^{\bar{\rho}m \times (k_1 + 1)}$. We are interested in $\|\pi_{\Sigma^\perp}\|_{L, L} = \left\| \Upsilon_{L, \bullet} (\Upsilon^* \Upsilon)^{-1} \Upsilon_{\bullet, L}^* \right\| \leq \frac{\|\Upsilon_{L, \bullet}\|^2}{\sigma_{\min}^2(\Upsilon)}$. It is not difficult to show²¹ that w.p. $\geq 1 - e^{-\frac{\bar{\rho}m}{8}(1 - \varepsilon + o(1))}$, $\sigma_{\min}(\Upsilon) \geq \frac{\nu\sqrt{\bar{\rho}}}{2}$. Meanwhile for any $\varepsilon_2 > 0$, $P[\|[Z_1]_{L, \bullet}\| \geq \nu\sqrt{c} + \nu\sqrt{\bar{\rho}\varepsilon_2}] \asymp e^{-\bar{\rho}\varepsilon_2 m/2}$. On the complement of this bad event (and invoking Lemma 6)

$$\|\Upsilon_{L, \bullet}\| \leq \|[Z_1]_{L, \bullet}\| + \left\| \nu\sqrt{\bar{\rho}} \frac{\mu_{J^c}(L)}{\|\mu_{J^c}\|} \right\| \leq \nu\sqrt{c} + \nu\sqrt{\bar{\rho}\varepsilon_2} + 2\nu C_\mu \sqrt{\frac{c}{\bar{\rho}}} = \nu \left(\sqrt{\bar{\rho}\varepsilon_2} + \sqrt{c} \left(1 + \frac{2C_\mu}{\sqrt{\bar{\rho}}} \right) \right).$$

²⁰Since Φ_{\bullet, L^c} is independent of $\Phi_{\bullet, L}$ and Σ , the norm of $\pi_{\Phi_{\bullet, L}} \Phi_{\bullet, L^c} \pi_{[\pi_\Sigma]_{L^c, L}}$ is simply distributed as the norm of a $cm \times cm$ iid $\mathcal{N}(0, \nu^2/m)$ matrix. By Fact 1, $P[\|\pi_{\Phi_{\bullet, L}} \Phi_{\bullet, L^c} \pi_{[\pi_\Sigma]_{L^c, L}}\| \geq 2\nu\sqrt{c} + t\nu\sqrt{c}] \leq e^{-(t - o(1))^2 cm/2}$. Set $t = \sqrt{\frac{\bar{\rho}\varepsilon_1}{c}}$.

²¹Write $\sigma_{\min}(\Upsilon) \geq \sigma_{\min} \left(\begin{bmatrix} \pi_{\mu_{J^c}^\perp} Z_1 & \nu\sqrt{\bar{\rho}} \frac{\mu_{J^c}}{\|\mu_{J^c}\|} \end{bmatrix} \right) - \|\pi_{\mu_{J^c}} Z_1\| \geq \min \left(\sigma_{\min}(\pi_{\mu_{J^c}^\perp} Z_1), \nu\sqrt{\bar{\rho}} \right) - \|\pi_{\mu_{J^c}} Z_1\|$, apply Fact 1 to the singular value and standard tail bounds to the k_1 -dimensional $\mathcal{N}(0, \nu^2/m)$ vector $\frac{\mu_{J^c}^*}{\|\mu_{J^c}\|} Z_1$.

By the assumptions of the lemma, $\sqrt{c}(1 + \frac{2C_\mu}{\sqrt{\rho}}) \leq \sqrt{\rho}/8$, and $\|[\pi_{\Sigma^\perp}]_{L,L}\| \leq \frac{\|\Upsilon_{L,\bullet}\|^2}{\sigma_{\min}^2(\Upsilon)} \leq 4(\sqrt{\varepsilon_2} + 1/8)^2$.

Setting $\varepsilon_1 = \varepsilon_2 = \frac{1}{64}$, $\|[\pi_{\Sigma^\perp}]_{L,L}\| \leq 1/4$, $\|\pi_{\Phi_{\bullet,L}} \Phi_{\bullet,L^c} \pi_{[\pi_\Sigma]_{L^c,L}}\| \leq 2\nu\sqrt{c} + \nu\sqrt{\rho}/8$, and so

$$\sigma_{\min}([\Phi\pi_\Sigma]_{\bullet,L}) \geq \left(\frac{\nu\sqrt{\rho}}{2} - \nu\sqrt{c}\right) \left(1 - \frac{1}{4}\right) - \left(2\nu\sqrt{c} + \frac{\nu\sqrt{\rho}}{8}\right) = \frac{\nu\sqrt{\rho}}{4} - \frac{11\nu\sqrt{c}}{4}, \quad (45)$$

and $\sigma_{\min}([\Lambda\Phi\pi_\Sigma]_{\bullet,L}) \geq \frac{1}{12} - \frac{11}{12}\sqrt{\frac{c}{\rho}} > \frac{1}{24}$ on the complement of a bad event of probability $e^{-\frac{\bar{\rho}m}{128}(1+o(1))}$.

The number of subsets L of size cm is $e^{\bar{\rho}mH(c/\bar{\rho})(1+o(1))}$. The probability any L is bad is bounded by $e^{\bar{\rho}m(H(c/\bar{\rho}) - \frac{1}{128})(1+o(1))}$, which falls off exponentially when $H(c/\bar{\rho}) < 1/128$. This is guaranteed for $c/\bar{\rho} < 1/1024$.

d) *Bounding $\sigma_{\min}(\pi_{\Sigma'^\perp}[\Psi\pi_{1^\perp}]_{\bullet,L})$* : Recall that Σ' denotes the cm -dimensional range of $[\Lambda\Phi\pi_\Sigma]_{\bullet,L_1}$. Choose any orthonormal basis for the $[(\bar{\rho} - c)m - k_1 - 1]$ -dimensional subspace Σ'^\perp . The expression of the columns of $\pi_{\Sigma'^\perp}\Psi$ with respect to this basis is a $((\bar{\rho} - c)m - k_1 - 1) \times (\delta m - k_1)$ matrix $\tilde{\Psi}$ with entries $\mathcal{N}(0, \nu^2/m)$. Split $\tilde{\Psi}\pi_{1^\perp}$ as

$$[\tilde{\Psi}\pi_{1^\perp}]_{\bullet,L} = \tilde{\Psi}_{\bullet,L} - \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L^c} \mathbf{1}_{L^c} \mathbf{1}_L^* - \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L} \mathbf{1}_L \mathbf{1}_L^*.$$

Using the independence of $\frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L^c} \mathbf{1}$ and $\tilde{\Psi}_{\bullet,L}$ and applying Fact 1, it is not difficult to show²² that

$$P \left[\sigma_{\min} \left(\tilde{\Psi}_{\bullet,L} - \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L^c} \mathbf{1} \mathbf{1}^* \right) \leq \frac{\nu\sqrt{\rho - c}}{2} - \nu\sqrt{c} \right] \leq e^{-\frac{(\bar{\rho} - c)m}{8}(1+o(1))}. \quad (46)$$

For the other term, $\left\| \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L} \mathbf{1} \mathbf{1}^* \right\| \leq \|\tilde{\Psi}_{\bullet,L}\| \frac{c}{\delta} (1 + o(1))$. From Fact 1, $P \left[\|\tilde{\Psi}_{\bullet,L}\| \geq 3\nu\sqrt{\rho} \right] \leq e^{-\frac{(\bar{\rho} - c)m}{2}(1+o(1))}$. On the complement of this event, $\left\| \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L} \mathbf{1} \mathbf{1}^* \right\| \leq \frac{3\nu c \sqrt{\rho}}{\delta} (1 + o(1)) \leq 4\nu\sqrt{c}\sqrt{\frac{c}{\delta}}$ eventually. Since $\sqrt{\frac{c}{\delta}} < \sqrt{\frac{c}{\rho}} \leq \frac{1}{32}$, $\left\| \frac{1}{\delta m - k_1} \tilde{\Psi}_{\bullet,L} \mathbf{1} \mathbf{1}^* \right\| \leq \frac{\nu\sqrt{c}}{8}$. All together, with probability at least $1 - e^{-(\bar{\rho} - c)m/8(1+o(1))}$,

$$\sigma_{\min}([\tilde{\Psi}\pi_{1^\perp}]_{\bullet,L}) \geq \sigma_{\min} \left(\tilde{\Psi}_{\bullet,L} - \frac{1}{m} \tilde{\Psi}_{\bullet,L^c} \mathbf{1} \mathbf{1}^* \right) - \left\| \frac{\tilde{\Psi}_{\bullet,L} \mathbf{1} \mathbf{1}^*}{\delta m - k_1} \right\| \geq \frac{7\nu\sqrt{\rho}}{16} - \frac{9}{8}\nu\sqrt{c} \geq \frac{3}{8}\nu\sqrt{\rho}$$

There are $\asymp e^{\bar{\rho}mH(c/\bar{\rho})}$ subsets L_1 of size cm and $\asymp e^{\delta mH(c/\delta)}$ subsets L_2 of size cm . The total number of choices of L_1, L_2 is asymptotic to $e^{(\bar{\rho}H(\frac{c}{\bar{\rho}}) + \delta H(\frac{c}{\delta}))m}$, and the probability that any pair is bad is bounded by a function asymptotic to $\exp\left(\left(\bar{\rho}H(c/\bar{\rho}) + \delta H(c/\delta) - \frac{\bar{\rho} - c}{8}\right)m(1 + o(1))\right)$. Under the assumptions of the lemma, the exponent is negative.

²²Translation does not substantially affect the bound on σ_{\min} in Fact 1: for an $m \times n$ iid $\mathcal{N}(0, 1/m)$ matrix M and an independent translation \mathbf{x} , $\sigma_{\min}(M + \mathbf{x}\mathbf{1}^*) \geq \sigma_{\min}(\pi_{\mathbf{x}^\perp} M)$, which obeys the same concentration result, now applied to an $(m - 1) \times n$ matrix. Appropriate rescaling of the $((\bar{\rho} - c)m - k_1 - 1) \times cm$ $\mathcal{N}(0, \nu^2/m)$ matrix $\tilde{\Psi}_{\bullet,L}$ yields the desired expression.

e) *Bounding the cross-coherence* $\left\| \pi_{\Sigma''} [\Psi \pi_{1^\perp}]_{\bullet, L_2} \right\|$: Let Σ'' denote the subspace $\mathcal{R}([\Lambda \Phi \pi_\Sigma]_{\bullet, L_1})$. Notice that Σ'' and Ψ are probabilistically independent. Now,

$$\left\| \pi_{\Sigma''} [\Psi \pi_{1^\perp}]_{\bullet, L_2} \right\| \leq \left\| \pi_{\Sigma''} \Psi_{\bullet, L_2} \right\| + \left\| \pi_{\Sigma''} \Psi \frac{\mathbf{1}_{\delta m - k_1} \mathbf{1}_{cm}^*}{\delta m - k_1} \right\| \leq \left\| \pi_{\Sigma''} \Psi_{\bullet, L_2} \right\| + \left\| \frac{\pi_{\Sigma''} \Psi \mathbf{1}}{\sqrt{\delta m - k_1}} \right\|$$

eventually, since $\frac{\|\mathbf{1}_{cm}\|}{\sqrt{\delta m - k_1}} = \sqrt{\frac{cm}{\delta m - k_1}} < 1$ eventually. Now, $\|\pi_{\Sigma''} \Psi_{\bullet, L_2}\|$ is distributed as the norm of a $cm \times cm$ iid $\mathcal{N}(0, \nu^2/m)$ matrix, and so for any $\varepsilon_1 > 0$,

$$P \left[\|\pi_{\Sigma''} \Psi_{\bullet, L_2}\| \geq 2\nu\sqrt{c} + \varepsilon_1\nu\sqrt{\bar{\rho}} \right] \asymp e^{-\varepsilon_1^2 \bar{\rho} m / 2}. \quad (47)$$

Similarly, $\frac{1}{\sqrt{\delta m - k_1}} \pi_{\Sigma''} \Psi \mathbf{1}$ is has the same norm as a cm -dimensional iid $\mathcal{N}(0, \nu^2/m)$ vector, so

$$P \left[\left\| \frac{1}{\sqrt{\delta m - k_1}} \pi_{\Sigma''} \Psi \mathbf{1} \right\| \geq \nu\sqrt{c} + \varepsilon_2\nu\sqrt{\bar{\rho}} \right] \leq e^{-2\varepsilon_2^2 \bar{\rho} m / \pi^2}. \quad (48)$$

On the complement of these two bad events, $\left\| \pi_{\Sigma''} [\Psi \pi_{1^\perp}]_{\bullet, L_2} \right\| \leq (\varepsilon_1 + \varepsilon_2)\nu\sqrt{\bar{\rho}} + 3\nu\sqrt{c}$. Set $\varepsilon_1 = \varepsilon_2 = 1/16$. Then w.p. $\geq 1 - e^{-\frac{\bar{\rho} m}{128\pi^2}(1+o(1))}$, $\left\| \pi_{\Sigma''} [\Psi \pi_{1^\perp}]_{\bullet, L_2} \right\| \leq \frac{\nu\sqrt{\bar{\rho}}}{8} + 3\nu\sqrt{c} \leq \frac{\nu\sqrt{\bar{\rho}}}{4}$. We again union bound over L_1, L_2 . The number of such pairs is asymptotic to $e^{(\bar{\rho}H(\frac{c}{\bar{\rho}}) + \delta H(\frac{c}{\delta}))m}$, and the probability of some bad pair is bounded by a function asymptotic to $\exp\left(\left(\bar{\rho}H\left(\frac{c}{\bar{\rho}}\right) + \delta H\left(\frac{c}{\delta}\right) - \frac{\bar{\rho}}{128\pi^2}\right)m\right)$. Under the hypotheses of the lemma, the coefficient of this exponent is negative.

f) *Pulling the bounds together*: For $\nu < 1/9$, $\frac{3\nu\sqrt{\bar{\rho}}}{8} < \frac{1}{24} \leq \min_{L_1} \sigma_{\min}([\Lambda \Phi \pi_\Sigma]_{\bullet, L_1})$, and so this quantity lower bounds $\min_{L_1, L_2} \left\{ \sigma_{\min}([\Lambda \Phi \pi_\Sigma]_{\bullet, L_1}), \sigma_{\min}(\pi_{\Sigma^\perp} [\Psi \pi_{1^\perp}]_{\bullet, L_2}) \right\}$. So, w.p. $\geq 1 - e^{-Cm(1+o(1))}$, $\gamma_{cm}([\Phi \pi_\Sigma \Phi^*]^{-1/2} \Phi \pi_\Sigma \Psi \pi_{1^\perp}) \geq \frac{3}{8}\nu\sqrt{\bar{\rho}} - \frac{1}{4}\nu\sqrt{\bar{\rho}} = \frac{\nu\sqrt{\bar{\rho}}}{8}$. Since

$$\left| \gamma_{cm}([\tilde{V}^* - \tilde{S}\tilde{U}^*]) - \gamma_{cm}([\Phi \pi_\Sigma \Phi^*]^{-1/2} \Phi \pi_\Sigma \Psi \pi_{1^\perp}) \right| \leq \frac{\nu\sqrt{\bar{\rho}}}{16},$$

the desired bound follows. \blacksquare

Lemma 6: Let J^c be chosen uniformly at random from $\binom{[m]}{\bar{\rho}m}$, and let $\boldsymbol{\mu} \in \mathbb{R}^m$ with $\|\boldsymbol{\mu}\|_2 = 1$ and $\|\boldsymbol{\mu}\|_\infty \leq C_\mu m^{-1/2}$. Then $\|\boldsymbol{\mu}_{J^c}\|_2 \geq \bar{\rho}/2$ on the complement of a bad event of probability $\leq e^{-Cm(1+o(1))}$.

Proof: Form the subset J^c by choosing $\bar{\rho}m$ indices $j_1 \dots j_{\bar{\rho}m}$, with j_i chosen uniformly at random from $[m] \setminus \{j_1 \dots j_{i-1}\}$. Let $Y_0, Y_1, \dots, Y_{\bar{\rho}m}$ denote the Doob process associated with $\|\boldsymbol{\mu}_{J^c}\|_2^2$: $Y_0 \doteq \mathbb{E}[\|\boldsymbol{\mu}_{J^c}\|_2^2] = \bar{\rho}$ and $Y_k \doteq \mathbb{E}[\|\boldsymbol{\mu}_{J^c}\|_2^2 \mid j_1 \dots j_k]$. Then, letting $X_k \doteq \sum_{i=1}^k \boldsymbol{\mu}_{j_i}^2$, $Y_k = X_k + \frac{1-X_k}{m-k}(\bar{\rho}m - k) = \frac{\rho m X_k + 1}{m-k}$, and

$$|Y_{k+1} - Y_k| = \left| \frac{\rho m (X_k + \boldsymbol{\mu}_{j_{k+1}}^2) + 1}{m-k-1} - \frac{\rho m X_k + 1}{m-k} \right| \leq \frac{\rho m X_k + \rho^2 m^2 \boldsymbol{\mu}_{j_{k+1}}^2 + 1}{\rho^2 m^2} \leq \frac{1}{\rho m} + \frac{C_\mu}{m} + \frac{1}{\rho^2 m^2}.$$

The above is $\leq C' m^{-1}$ for appropriate constant C' . By Azuma's inequality (Theorem 7.2.1 of [31]),

$$P[|Y_{\bar{\rho}m} - \bar{\rho}| \geq t] \leq 2 \exp\left(-\frac{t^2}{2\bar{\rho}m(C'/m)^2}\right) \asymp \exp(-Cm). \quad (49)$$

\blacksquare

B. Technical Lemmas for Initial Separating Hyperplane

This section contains two results used above for controlling the initial separator \mathbf{q}_0 . We first justify the assertion that $\begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} Z_{J^*,\bullet}^* \boldsymbol{\sigma}$ is the only term that contributes $O(m^{1/2})$ to $\|\mathbf{q}_0\|$, and then close with a measure concentration result for $\|\theta \cdot\|$, also used in the proof of Lemma 4.

Lemma 7 (Lower order terms in \mathbf{q}_0): Suppose that $\bar{\rho} < \delta$ and $\nu < \frac{1}{8(\sqrt{\delta}+1)}$. There exist constants (wrt m) C_G and C_q such that

$$\|(G^*G)^{-1}\| \leq C_G \quad \text{and} \quad \|\mathbf{q}_0 - \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} Z_{J^*,\bullet}^* \boldsymbol{\sigma}\| \leq C_q m^{1/2-\eta_0/4} \quad (50)$$

simultaneously on the complement of a bad event of probability $\leq e^{-Cm^{1-\eta_0/2}(1+o(1))}$.

Proof: Write $Q = \begin{bmatrix} Z_1^* Z_1 & Z_1^* Z_2 \\ Z_2^* Z_1 & Z_2^* Z_2 + \mathbf{I} \end{bmatrix} \in \mathbb{R}^{n \times n}$, and $\boldsymbol{\zeta} = Z_{J^c, \bullet}^* \boldsymbol{\mu}_{J^c} \in \mathbb{R}^n$. Then $G^*G = Q + \boldsymbol{\zeta} \mathbf{1}^* + \mathbf{1} \boldsymbol{\zeta}^* + \alpha \mathbf{1} \mathbf{1}^*$, where $\alpha = \boldsymbol{\mu}_{J^c}^* \boldsymbol{\mu}_{J^c}$. So,

$$(G^*G)^{-1} = Q^{-1} - Q^{-1} \begin{bmatrix} \mathbf{1} & \boldsymbol{\zeta} \end{bmatrix} \begin{bmatrix} \mathbf{1}^* Q^{-1} \mathbf{1} & \mathbf{1}^* Q^{-1} \boldsymbol{\zeta} + 1 \\ \mathbf{1}^* Q^{-1} \boldsymbol{\zeta} + 1 & \boldsymbol{\zeta}^* Q^{-1} \boldsymbol{\zeta} - \alpha \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}^* \\ \boldsymbol{\zeta}^* \end{bmatrix} Q^{-1}. \quad (51)$$

Set $b \doteq \mathbf{1}^* Q^{-1} \mathbf{1}$, $c \doteq \mathbf{1}^* Q^{-1} \boldsymbol{\zeta}$, $d \doteq \boldsymbol{\zeta}^* Q^{-1} \boldsymbol{\zeta}$, and write $(G^*G)^{-1} = Q^{-1} - Q^{-1/2} M \Xi M^* Q^{-1/2}$ with

$$M = \begin{bmatrix} \frac{Q^{-1/2} \mathbf{1}}{\|Q^{-1/2} \mathbf{1}\|_2} & \frac{Q^{-1/2} \boldsymbol{\zeta}}{\|Q^{-1/2} \boldsymbol{\zeta}\|_2} \end{bmatrix} \quad \text{and} \quad \Xi = \frac{\begin{bmatrix} b(\alpha-d) & -\sqrt{bd}(c+1) \\ -\sqrt{bd}(c+1) & bd \end{bmatrix}}{b(\alpha-d) + (c+1)^2}. \quad (52)$$

We next bound the quadratic terms b , c , and d . Applying Fact 1 to the $\delta m \times \bar{\rho} m$ iid $\mathcal{N}(0, \nu^2/m)$ matrix $Z_{J^c, \bullet} = [Z_1 \ Z_2]$ gives that $\|Z_{J^c, \bullet}\|_2 \leq \sqrt{2}\nu \left(\sqrt{\delta} + \sqrt{\bar{\rho}} \right)$ w.p. $\geq 1 - e^{-Cm(1+o(1))}$. On the complement of that bad event,

$$b = \mathbf{1}^* Q^{-1} \mathbf{1} \geq \frac{\|\mathbf{1}\|_2^2}{\|Q\|} \geq \frac{\delta m}{1 + \|Z_{J^c, \bullet}\|_2^2} \geq \frac{\delta m}{1 + 2\nu^2(\sqrt{\delta} + \sqrt{\bar{\rho}})^2} \doteq C_b m. \quad (53)$$

Similarly, $b \leq \delta m / \sigma_{\min}(Q)$. It is not difficult to show²³ that for any block matrix $M = \begin{bmatrix} A & B \\ 0 & \mathbf{I} \end{bmatrix}$ with $\sigma_{\min}(A) < 1$,

$$\sigma_{\min}^2(M) \geq \sigma_{\min}^2(A) - \frac{\|A\|^2 \|B\|^2}{1 - \sigma_{\min}^2(A)}.$$

By Fact 1, on the complement of an event of probability $\asymp e^{-Cm}$,

$$\sigma_{\min}^2(Z_1) \geq \frac{\nu^2 \bar{\rho}}{2}, \quad \|Z_1\|^2 \leq \|Z_1\|_2^2 \leq 2\nu^2 \bar{\rho}, \quad \|Z_2\|^2 \leq \|Z_2\|_2^2 \leq 2\nu^2 \left(\sqrt{\delta} + \sqrt{\bar{\rho}} \right)^2.$$

On the good event above, for $\nu < \frac{1}{\sqrt{2}}$, $\sigma_{\min}^2(Z_1) \leq \|Z_1\|^2 < 1$. Plugging in, $\sigma_{\min}(Q) = \sigma_{\min} \left(\begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \right) \geq \sigma_{\min}^2(Z_1) - \frac{\|Z_1\|^2 \|Z_2\|^2}{1 - \sigma_{\min}^2(Z_1)} \geq \frac{\nu^2 \bar{\rho}}{2} - \frac{4\nu^4 \bar{\rho} (\sqrt{\delta} + \sqrt{\bar{\rho}})^2}{1 - 2\nu^2 \bar{\rho}} \geq \frac{\nu^2 \bar{\rho}}{4}$ for ν sufficiently small (e.g., $\nu < \frac{1}{8(\sqrt{\delta}+1)}$ suffices), and so $b \leq \frac{4\delta}{\nu^2 \bar{\rho}} m$ w.p. $\geq 1 - e^{-Cm(1+o(1))}$.

²³Write $\sigma_{\min}^2(M) \geq \min_{\|\mathbf{x}_1\|_2^2 + \|\mathbf{x}_2\|_2^2 = 1} (\|A\mathbf{x}_1\|_2 - \|B\mathbf{x}_2\|_2)^2 + \|\mathbf{x}_2\|_2^2$. Setting $\lambda = \|\mathbf{x}_1\|_2^2$, the previous is $\geq \min_{\lambda \in [0,1]} \sigma_{\min}^2(A) + (1 - \sigma_{\min}^2(A))(1 - \lambda) - 2\|A\| \|B\| \sqrt{1 - \lambda}$, which is minimized at $\sqrt{1 - \lambda} = \frac{\|A\| \|B\|}{1 - \sigma_{\min}^2(A)}$.

For $c = \mathbf{1}^* Q^{-1} \zeta$, notice that $\zeta = Z_{J^c, \bullet}^* \mu_{J^c}$ is iid $\mathcal{N}(0, \nu^2 \alpha / m)$. Write $Q = Z_{J^c, \bullet}^* \pi_{\mu_{J^c}^\perp} Z_{J^c, \bullet} + \begin{bmatrix} 0 & 0 \\ 0 & \mathbf{I} \end{bmatrix} + \frac{1}{\alpha} \zeta \zeta^* \doteq L + \frac{1}{\alpha} \zeta \zeta^*$, then $Q^{-1} = L^{-1} - L^{-1} \zeta \frac{1}{\alpha + \zeta^* L^{-1} \zeta} \zeta^* L^{-1}$, and $|\mathbf{1}^* Q^{-1} \zeta| = \left| \mathbf{1}^* L^{-1} \zeta \left(\frac{\alpha}{\alpha + \zeta^* L^{-1} \zeta} \right) \right| \leq |\mathbf{1}^* L^{-1} \zeta|$. An identical argument²⁴ to the one given above for Q shows that on the complement of an event of probability $\asymp e^{-Cm}$, $\sigma_{\min}(L) \geq \frac{\nu^2 \bar{\rho}}{4}$, and so $\|L^{-1} \mathbf{1}\|_2 \leq \frac{4\sqrt{\delta}}{\nu^2 \bar{\rho}} m^{1/2}$. Since ζ is independent of L , $\left\langle \frac{L^{-1} \mathbf{1}}{\|L^{-1} \mathbf{1}\|_2}, \zeta \right\rangle$ is simply an $\mathcal{N}(0, \nu^2 \alpha / m)$ random variable, and so for any $\varepsilon > 0$

$$P \left[|\mathbf{1}^* L^{-1} \zeta| > \varepsilon m^{1/2} \right] \leq P \left[\|L^{-1} \mathbf{1}\|_2 > \frac{4\sqrt{\delta}}{\nu^2 \bar{\rho}} m^{1/2} \right] + P \left[\left| \left\langle \frac{L^{-1} \mathbf{1}}{\|L^{-1} \mathbf{1}\|_2}, \zeta \right\rangle \right| > \varepsilon \frac{\nu^2 \bar{\rho}}{4\sqrt{\delta}} \right] \asymp e^{-C_\varepsilon m}$$

for some constant C_ε (where we have controlled the second part via standard Gaussian tail bounds²⁵).

So, with overwhelming probability, $|c| = |\mathbf{1}^* Q^{-1} \zeta| \leq \varepsilon m^{1/2}$.

The final quadratic term is $d = \zeta^* Q^{-1} \zeta = \zeta^* L^{-1} \zeta \frac{\zeta^* L^{-1} \zeta}{\alpha + \zeta^* L^{-1} \zeta} \leq \zeta^* L^{-1} \zeta$. The norm of the δm -dimensional $\mathcal{N}(0, \nu^2 \alpha / m)$ vector ζ concentrates: by (21), $\|\zeta\|_2 \leq \sqrt{2} \nu \sqrt{\alpha \delta}$ with probability at least $1 - e^{-Cm(1+o(1))}$. We exploit the fact that although $\|L^{-1}\| = O(\nu^{-2})$, for most vectors L is well-conditioned (due to the presence of the identity matrix in $\begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix}$). Consider the subspace $\Sigma = \{x \mid x_I = 0\} \subset \mathbb{R}^n$. Since for all $x \in \Sigma$, $\|Lx\|_2 \geq \|x\|_2$, $\|L^{-1}|_{L\Sigma}\| \leq 1$, and

$$\begin{aligned} \zeta^* L^{-1} \zeta &= \zeta^* (L^{-1}|_{L\Sigma} \pi_{L\Sigma} \zeta + L^{-1} \pi_{(L\Sigma)^\perp} \zeta) \\ &\leq \|\zeta\|_2^2 \|L^{-1}|_{L\Sigma}\|_2 + \|L^{-1}\|_2 \|\zeta\|_2 \|\pi_{(L\Sigma)^\perp} \zeta\|_2 \leq 2\nu^2 \alpha \delta + \frac{4\sqrt{2\alpha\delta}}{\nu\bar{\rho}} \|\pi_{(L\Sigma)^\perp} \zeta\|_2. \end{aligned}$$

The norm $\|\pi_{(L\Sigma)^\perp} \zeta\|$ of the projection of ζ onto an independent k_1 -dimensional subspace is distributed as the norm of a k_1 -dimensional $\mathcal{N}(0, \nu^2 \alpha / m)$ vector: $P[\|\pi_{(L\Sigma)^\perp} \zeta\| \geq \varepsilon' \nu \sqrt{\alpha}] \asymp e^{-2\varepsilon'^2 m / \pi^2}$. For appropriate ε , with overwhelming probability, $d \leq \zeta^* L^{-1} \zeta \leq 4\nu^2 \alpha \delta$.

The denominator of Ξ in (52) is $b(\alpha - d) + (c + 1)^2 \geq C_b \alpha (1 - 4\nu^2 \delta) m$. By Lemma 6, $\alpha = \|\mu_{J^c}\|_2^2 \geq \bar{\rho} / 2$ w.p. $\geq 1 - e^{-Cm(1+o(1))}$, and so the denominator is $\geq C_{\text{denom}} m$ with overwhelming probability. Since each of the terms in the numerator is $\leq Cm$ with overwhelming probability, $\|\Xi\| \leq C_\Xi$ for appropriate constant C_Ξ . Since the columns of M have unit norm, $\|M\| \leq 2$, and

$$\|(G^* G)^{-1}\| \leq \|Q^{-1}\| + \|Q^{-1}\| \|M\|^2 \|\Xi\| \leq \frac{4}{\nu^2 \bar{\rho}} + \frac{4}{\nu^2 \bar{\rho}} 4 C_\Xi \doteq C_G,$$

a constant, establishing the first assertion of the lemma.

²⁴Consider instead $\sigma_{\min}^2 \left(\begin{bmatrix} \pi_{\mu_{J^c}^\perp} Z_1 & \pi_{\mu_{J^c}^\perp} Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \right)$. The singular values of $\pi_{\mu_{J^c}^\perp} Z_2$ are distributed as those of a $(\bar{\rho} m - 1) \times (\delta m - k_1)$ iid $\mathcal{N}(0, \nu^2 / m)$ matrix. The bounds given by Fact 1 are essentially the same as those for Z_2 .

²⁵For example, if X is $\mathcal{N}(0, \sigma^2)$, $P[|X| \geq \sigma t] \leq t^{-1} e^{-t^2/2}$.

We next extend the above reasoning to bound $(G^*G)^{-1}\mathbf{1}$ and $\mathbf{1}^*(G^*G)^{-1}\mathbf{1}$. Notice that

$$(G^*G)^{-1}\mathbf{1} = \frac{c+1}{b(\alpha-d)+(c+1)^2} Q^{-1}\mathbf{1} - \frac{1}{\alpha-d+(c+1)^2/b} Q^{-1}\boldsymbol{\zeta} \doteq \lambda_1 Q^{-1}\mathbf{1} + \lambda_2 Q^{-1}\boldsymbol{\zeta}.$$

For any $\varepsilon > 0$, $|\lambda_1| \leq \frac{|c+1|}{b(\alpha-d)} \leq \frac{\varepsilon m^{1/2+1}}{C_b m^{\frac{\rho}{2}(1-4\nu^2\delta)}}$ with overwhelming probability. Hence for any $\varepsilon'' > 0$, $|\lambda_1| \leq \varepsilon'' m^{-1/2}$ for m sufficiently large, on the complement of a bad event of probability $\asymp e^{-Cm}$. Similarly, $|\lambda_2| \leq \frac{1}{\alpha-d} \leq \frac{2}{\bar{\rho}(1-4\nu^2\delta)}$, and so

$$\|(G^*G)^{-1}\mathbf{1}\|_2 \leq |\lambda_1| \|Q^{-1}\mathbf{1}\| + |\lambda_2| \|Q^{-1}\boldsymbol{\zeta}\| \leq \frac{4\varepsilon''\sqrt{\delta}}{\nu^2\bar{\rho}} + \frac{8\sqrt{2\delta}}{\nu^2\bar{\rho}^2(1-4\nu^2\delta)} \doteq C_1.$$

Similarly, $\mathbf{1}^*(G^*G)^{-1}\mathbf{1} = \frac{b}{b(\alpha-d)+(c+1)^2} \leq \frac{2}{\bar{\rho}(1-4\nu^2\delta)} \doteq C_2$.

We need one more bound, for $|\langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle|$. Consider the Martingale $(X_i)_{i=0}^{\rho m}$ given by $X_0 = 0$, $X_i = \sum_{j=1}^i \boldsymbol{\mu}_J(j) \boldsymbol{\sigma}(j)$. We are interested in $X_{\rho m} = \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle$. Since $|X_i - X_{i-1}| \leq |\boldsymbol{\mu}_J(i)|$, by Hoeffding's inequality [31],

$$P[|X_{\rho m}| \geq t] \leq 2 \exp\left(-\frac{t^2}{2 \sum_{j=1}^{\rho m} \boldsymbol{\mu}_J^2(j)}\right) \leq 2e^{-\frac{t^2}{2}}, \quad (54)$$

and so with probability $\geq 1 - e^{-Cm^{1-\eta_0/2}}$, $|\langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle| \leq m^{1/2-\eta_0/4}$.

With these results in hand, recall that

$$\begin{aligned} \mathbf{q}_0 &= \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} (G^*G)^{-1} Z_{J,\bullet}^* \boldsymbol{\sigma} + \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \left(-(G^*G)^{-1}\mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle (G^*G)^{-1}\mathbf{1} \right) \\ &+ \begin{bmatrix} \boldsymbol{\mu}_{J^c} \\ \mathbf{0} \end{bmatrix} \left(-\mathbf{1}^*(G^*G)^{-1}\mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle \mathbf{1}^*(G^*G)^{-1}\mathbf{1} \right) + \begin{bmatrix} \boldsymbol{\mu}_{J^c} \\ \mathbf{0} \end{bmatrix} \mathbf{1}^*(G^*G)^{-1} Z_{J,\bullet}^* \boldsymbol{\sigma}. \end{aligned} \quad (55)$$

The second term of (55), $\left\| \begin{bmatrix} Z_1 & Z_2 \\ 0 & \mathbf{I} \end{bmatrix} \left(-(G^*G)^{-1}\mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle (G^*G)^{-1}\mathbf{1} \right) \right\|$ is bounded above by

$$\left(1 + \sqrt{2}\nu(\sqrt{\delta} + \sqrt{\bar{\rho}})\right) \left(C_G \sqrt{C_0} m^{1/2-\eta_0/2} + m^{1/2-\eta_0/4} C_1\right)$$

w.p. $\geq 1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$. Similarly, for the third term of (55)

$$\left\| \begin{bmatrix} \boldsymbol{\mu}_{J^c} \\ \mathbf{0} \end{bmatrix} \left(-\mathbf{1}^*(G^*G)^{-1}\mathbf{1}_I + \langle \boldsymbol{\mu}_J, \boldsymbol{\sigma} \rangle \mathbf{1}^*(G^*G)^{-1}\mathbf{1} \right) \right\| \leq C_1 C_0 m^{1/2-\eta_0/4} + C_2 m^{1/2-\eta_0/4}.$$

For the final term of (55), $\boldsymbol{\vartheta} \doteq Z_{J,\bullet}^* \boldsymbol{\sigma}$ is distributed as an iid $\mathcal{N}(0, \nu^2 \rho)$ vector, independent of G , and so

$$P\left[\left|\left\langle \frac{(G^*G)^{-1}\mathbf{1}}{\|(G^*G)^{-1}\mathbf{1}\|}, \boldsymbol{\vartheta} \right\rangle\right| \geq m^{1/2-\eta_0/4}\right] \asymp e^{-Cm^{1-\eta_0/2}}. \quad (56)$$

On the complement of this bad event,

$$\|\boldsymbol{\mu}_{J^c} \mathbf{1}^*(G^*G)^{-1}\boldsymbol{\vartheta}\| \leq \|(G^*G)^{-1}\mathbf{1}\| \cdot \left|\left\langle \frac{(G^*G)^{-1}\mathbf{1}}{\|(G^*G)^{-1}\mathbf{1}\|}, \boldsymbol{\vartheta} \right\rangle\right| \leq C_1 m^{1/2-\eta_0/4}. \quad (57)$$

■

Lemma 8 (Concentration for Gaussian tops): Fix $\sigma \leq 1$, $\varepsilon \leq 1/2$. Let \mathbf{x} be a d -dimensional random vector with entries iid $\mathcal{N}(0, \sigma^2)$, and let θ be the operator that takes the part of \mathbf{x} above $1 - \varepsilon$:

$$\theta : \mathbb{R}^d \rightarrow \mathbb{R}^d \text{ such that } [\theta \mathbf{x}](i) = \begin{cases} \text{sgn}(\mathbf{x}(i))(|\mathbf{x}(i)| - 1 + \varepsilon), & |\mathbf{x}(i)| > 1 - \varepsilon, \\ 0, & \text{else.} \end{cases} \quad (58)$$

Then $P \left[\|\theta \mathbf{x}\|_2 \geq 4e^{-\frac{1}{16\sigma^2}} d^{1/2} \right] \asymp e^{-C_\sigma d}$, where C_σ is a constant (w.r.t. d) depending only on σ .

Proof: Let $\mathbf{y} \in \mathbb{R}^d$ be iid $\mathcal{N}(0, 1)$, then $\|\theta \mathbf{x}\|_2$ is equal in distribution to $\|\theta \sigma \mathbf{y}\|_2$. Now, $\mathbb{E} \|\theta \sigma \mathbf{y}\|_2^2 = d \cdot \mathbb{E}(\theta \mathbf{x}(i))^2 = \frac{d}{\sigma} \sqrt{\frac{2}{\pi}} \int_{1-\varepsilon}^{\infty} t^2 e^{-t^2/2\sigma^2} dt$. Integrating by parts²⁶ yields

$$d^{-1} \mathbb{E} \|\theta \sigma \mathbf{y}\|_2^2 = \frac{(1-\varepsilon)\sigma}{\sqrt{\pi/2}} e^{-\frac{(1-\varepsilon)^2}{2\sigma^2}} + 2\sigma^2 Q\left(\frac{1-\varepsilon}{\sigma}\right) \leq \sigma \sqrt{\frac{2}{\pi}} \frac{1+\sigma^2}{1-\varepsilon} e^{-\frac{(1-\varepsilon)^2}{2\sigma^2}} \leq 4\sigma e^{-\frac{1}{8\sigma^2}},$$

and $\mathbb{E}[\|\theta \sigma \mathbf{y}\|_2] \leq 2e^{-\frac{1}{16\sigma^2}} d^{1/2}$. Meanwhile, $\mathbb{E} \sqrt{\sum_{i=1}^d |\theta \sigma \mathbf{y}(i)|^2} = \sqrt{d} \mathbb{E} \sqrt{\frac{\sum_{i=1}^d |\theta \sigma \mathbf{y}(i)|^2}{d}}$. It is not difficult to show²⁷ that $\mathbb{E} \sqrt{\frac{\sum_{i=1}^d |\theta \sigma \mathbf{y}(i)|^2}{d}} \rightarrow C'_\sigma$ for some constant $C'_\sigma > 0$, and so $\mathbb{E} \|\theta \sigma \mathbf{y}\|_2 \geq C'_\sigma d^{1/2}$. Since $f(\cdot) = \|\theta \sigma \cdot\|_2$ is 1-Lipschitz for $\sigma \leq 1$, $P[\|\theta \sigma \mathbf{y}\|_2 \geq 2\mathbb{E} \|\theta \sigma \mathbf{y}\|_2] \leq \exp(-8(\mathbb{E} \|\theta \sigma \mathbf{y}\|_2)^2/\pi^2)$ [3]. Plugging in the upper and lower bounds on $\mathbb{E} \|\theta \sigma \mathbf{y}\|_2$ yields the result. \blacksquare

C. Details of the Proof of Theorem 1

Proof: Consider the weak proportional growth setting $\text{WPG}_{\delta, \rho, C_0, \eta_0}$ with $\bar{\rho} < \delta$. We first consider a fixed, arbitrary sequence of signal supports $I \in \binom{[n]}{k_1}$. By Lemma 2, (I, J, σ) is ℓ^1 -recoverable if $\exists c \in (0, 1)$ such that

$$\|\mathbf{q}_0\|_2 + \frac{1}{1-\xi} \|\theta \mathbf{q}_0\|_2 \leq (1-\varepsilon)\sqrt{c\rho} = (1-\varepsilon)\sqrt{c(\bar{\rho} + \delta)} m^{1/2} + o(m^{1/2}), \quad (59)$$

where $\xi = \inf_{\|\mathbf{s}\|_0 \leq c\rho} \|\pi_{\mathcal{R}(G)} \mathbf{s}\|_2 / \|\mathbf{s}\|_2$. Choose c small enough that $\beta \doteq (\bar{\rho} + \delta)c$ satisfies $\beta < \min\left(\frac{\bar{\rho}}{1024}, \frac{\bar{\rho}}{64(1+2C_\mu \bar{\rho}^{-1/2})^2}\right)$ and $\bar{\rho}H(\beta/\bar{\rho}) + \delta H(\beta/\delta) < \frac{\bar{\rho}}{128\pi^2}$ (since in Lemma 3, $\|\mathbf{s}\|_0$ is a fraction of m , not p). Further suppose that $\nu < \min(\frac{1}{9}, \frac{1}{8(\sqrt{\delta}+1)}, (512/\delta)^{-1/4})$. Then by Lemma 3, $\xi < 1 - C_\xi \nu^8$, with probability $1 - e^{-Cm(1+o(1))}$.

Meanwhile, by Lemma 4, with probability at least $1 - e^{-Cm^{1-\eta_0/2}(1+o(1))}$, $\|\mathbf{q}_0\|_2 \leq \alpha_1 \nu m^{1/2} + o(m^{1/2})$ and $\|\theta \mathbf{q}_0\| \leq \alpha_2 \nu^{-8} e^{-\frac{1}{64\nu^2}}$. On the intersection of these three good events, the left hand side of (59) becomes

$$\|\mathbf{q}_0\|_2 + \frac{1}{1-\xi} \|\theta \mathbf{q}_0\| \leq \alpha_1 \nu m^{1/2} + \alpha_2 \nu^{-8} \exp\left(-\frac{1}{64\nu^2}\right) m^{1/2} + o(m^{1/2}). \quad (60)$$

²⁶And noting that $Q(z) \leq \frac{1}{z\sqrt{2\pi}} e^{-z^2/2}$.

²⁷Apply the strong law of large numbers to $d^{-1} \sum |\theta \sigma \mathbf{y}(i)|^2$ and Slutsky's theorem (Theorem 6 of [32]) to argue that $\mathbb{E} \sqrt{d^{-1} \sum |\theta \sigma \mathbf{y}(i)|^2} \rightarrow \sqrt{\mathbb{E} |\theta \sigma \mathbf{y}(i)|^2}$.

For ν sufficiently small, this is $\leq (1 - \varepsilon)\sqrt{c(\bar{\rho} + \delta)}m^{1/2} + o(m^{1/2})$, and hence, for m sufficiently large, on an event of probability $\geq 1 - \exp(-Cm^{1-\eta_0/2}(1 + o(1)))$, (I, J, σ) is ℓ^1 -recoverable. There are $\binom{m}{k_1} \leq \exp(m^{1-\eta_0} \log m)$ subsets I , and so the probability that (I, J, σ) is not ℓ^1 -recoverable for some I is bounded by

$$\exp\left(-Cm^{1-\eta_0/2}(1 + o(1))\right) \times \exp\left(m^{1-\eta_0} \log m\right) = \exp\left(-Cm^{1-\eta_0/2}(1 + o(1))\right) = o(1),$$

establishing the theorem. ■