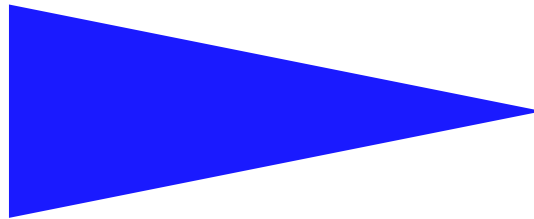


PUBLICATION
INTERNE
N° 1899



RESTRICTED ISOMETRY CONSTANTS WHERE ℓ^p SPARSE
RECOVERY CAN FAIL FOR $0 < p \leq 1$

MICHAEL E. DAVIES , RÉMI GRIBONVAL

Restricted Isometry Constants where ℓ^p sparse recovery can fail for $0 < p \leq 1$

Michael E. Davies^{*}, Rémi Gribonval^{**}

Systèmes cognitifs
Projet Metiss

Publication interne n° 1899 — Juillet 2008 — 22 pages

Abstract: We investigate conditions under which the solution of an underdetermined linear system with minimal ℓ^p norm, $0 < p \leq 1$, is guaranteed to be also the sparsest one. Our results highlight the pessimistic nature of sparse recovery analysis when recovery is predicted based on the restricted isometry constants (RIC) of the associated matrix.

We construct matrices with RIC δ_{2m} arbitrarily close to $1/\sqrt{2} \approx 0.717$ where sparse recovery with $p = 1$ fails for at least one m -sparse vector. This indicates that there is limited room for improving over the best known positive results of Foucart and Lai, which guarantee that ℓ^1 -minimisation recovers all m -sparse vectors for any matrix with $\delta_{2m} < 2(3 - \sqrt{2})/7 \approx 0.4531$. Another consequence of our construction is that recovery conditions expressed uniformly for all matrices in terms of RIC must require that all $2m$ -column submatrices are extremely well conditioned (condition numbers less than 2.5). In contrast, we also construct matrices with δ_{2m} arbitrarily close to one and $\delta_{2m+1} = 1$ where ℓ^1 -minimisation succeeds for any m -sparse vector. This illustrates the limits of RIC as a tool to predict the behaviour of ℓ^1 minimisation.

These constructions are a by-product of tight conditions for ℓ^p recovery ($0 \leq p \leq 1$) with matrices of unit spectral norm, which are expressed in terms of the minimal singular values of $2m$ -column submatrices. The results show that, compared to ℓ^1 -minimisation, ℓ^p -minimisation recovery failure is only slightly delayed in terms of the RIC values. Furthermore in this case the minimisation is nonconvex and it is important to consider the specific minimisation algorithm being used. We show that when ℓ^p optimisation is attempted using an iterative reweighted ℓ^1 scheme, failure can still occur for δ_{2m} arbitrarily close to $1/\sqrt{2}$.

Key-words: underdetermined linear system, sparse representation, overcomplete dictionary, compressed sensing, inverse problem, restricted isometry property, convex optimisation, nonconvex optimisation, iterative reweighted optimisation.

(Résumé : *tsvp*)

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

^{*} IDCOM & Joint Research Institute for Signal and Image Processing Edinburgh University, King's Buildings, Mayfield Road, Edinburgh EH9 3JL, Tel. +44 (0)131 650 5795, email: mike.davies@ed.ac.uk

^{**} remi.gribonval@irisa.fr

Sur les constantes d'isométrie restreintes et l'identification de représentation parcimonieuse par minimisation ℓ^p , $0 < p \leq 1$

Résumé : Nous nous intéressons aux conditions sous lesquelles la solution d'un système linéaire sous-déterminé de norme ℓ^p minimale, avec $0 < p \leq 1$, est aussi la plus parcimonieuse. Nos résultats mettent en lumière le caractère pessimiste de l'analyse de l'identification parcimonieuse lorsque l'identification est prédite en termes de constantes d'isométrie restreintes (CIR) de la matrice associée.

Nous construisons une matrice dont la CIR δ_{2m} est arbitrairement proche de $1/\sqrt{2} \approx 0.717$ pour laquelle il existe un vecteur à m composantes non nulles que la minimisation ℓ^1 ne permet pas d'identifier. Comparé au meilleur résultat positif connu de Foucart et Lai, qui garantit que la minimisation ℓ^1 identifie tous les vecteurs à m composantes non nulles pour toute matrice de CIR $\delta_{2m} < 2(3 - \sqrt{2})/7 \approx 0.4531$, notre construction indique que la marge possible d'amélioration du résultat positif est faible. Une autre conséquence de notre construction est que toute condition suffisante d'identification qui s'exprime en termes de CIR δ_{2m} d'une matrice doit imposer que toutes les sous-matrices à $2m$ colonnes soient extrêmement bien conditionnées (avec un conditionnement n'excédant pas 2.5). Nous illustrons plus avant les limites des CIR en construisant des matrices où δ_{2m} est arbitrairement proche de un et $\delta_{2m+1} = 1$ pour lesquelles la minimisation ℓ^1 identifie cependant tous les vecteurs à m composantes non nulles.

Nous exprimons enfin des résultats caractérisant précisément, pour toute matrice de norme spectrale unité, les conditions d'identification de représentations parcimonieuses par minimisation ℓ^p , ($0 \leq p \leq 1$). Nous remplaçons pour cela les CIR par les valeurs singulières minimales de sous-matrices du dictionnaire. Les résultats montrent que la mise en échec de la minimisation ℓ^p , $p < 1$ est à peine retardée en termes de CIR par rapport à la minimisation ℓ^1 . De plus, pour $p < 1$ la minimisation n'est plus convexe et il est important de tenir compte de l'algorithme de minimisation spécifiquement utilisé. Nous montrons qu'il existe des matrices de constante d'isométrie arbitrairement proche de $1/\sqrt{2}$ pour lesquels toute une classe d'algorithmes de minimisation ℓ^1 repondérée itérée –qui couvre plusieurs algorithmes proposés dans la littérature pour la minimisation ℓ^p , $p < 1$ – est mise en échec pour au moins un vecteur à m composantes non nulles.

Mots clés : système linéaire sous-déterminé, représentation parcimonieuse, dictionnaire redondant, acquisition compressée, problème inverse, propriété d'isométrie restreinte, optimisation convexe, optimisation non-convexe, optimisation itérative re-pondérée

Contents

1	Introduction and state of the art	4
1.1	Notations	4
1.2	Known conditions for ℓ^p sparse recovery	4
1.3	Contributions	6
2	Isometry measures for unit spectral norm dictionaries	7
3	Dictionaries with small δ_{2m} where ℓ^p can fail	8
4	Numerical studies of the σ_{2m}^2 and δ_{2m} conditions	15
4.1	Large dimensional ℓ^p failing dictionaries	15
4.2	Low dimensional examples	16
5	Rewighted ℓ^1 implementations for ℓ^p-optimisation	17
6	Discussion	19

1 Introduction and state of the art

This paper investigates conditions under which the solution $\hat{\mathbf{y}}$ of minimal ℓ^p norm, $0 < p \leq 1$, of an underdetermined linear system $\mathbf{x} = \Phi \mathbf{y}$ is guaranteed to be also the sparsest one. This is a central problem in sparse overcomplete signal representations, where \mathbf{x} is a vector representing some signal or image, Φ is an overcomplete signal dictionary, and \mathbf{y} is a sparse representation of the signal. This problem is also at the core of compressed sensing, where Φ is called a sensing matrix, \mathbf{x} is a collection of M linear measurements of some ideally sparse data \mathbf{y} . Although in both settings it is practically relevant to consider sparse *approximation* rather than exact sparse *representation*, most of the results of this paper are of a negative nature and naturally extend from the representation setting chosen here (for the sake of simplicity) to the approximation setting.

The proposed approach is twofold:

- we construct matrices (which we will call *dictionaries* from now on) Φ with “good” restricted isometry properties where sparse recovery with ℓ^p minimisation will nevertheless fail for at least one sparse vector.
- we construct dictionaries Φ with “bad” restricted isometry properties where sparse recovery with ℓ^p minimisation will nevertheless succeed for all (sufficiently) sparse vectors.

The goal is to understand how much improvement is possible over the best known positive results which relate restricted isometry constants to sparse ℓ^p recovery.

1.1 Notations

Given a vector $\mathbf{x} \in \mathbb{R}^M$ and a matrix $\Phi \in \mathbb{R}^{M \times N}$ with $M < N$, we are interested in sparse solutions to

$$\mathbf{x} = \Phi \mathbf{y} \quad (1)$$

We will denote by $\|\mathbf{y}\|_p$ the ℓ^p sparsity measure defined as:

$$\|\mathbf{y}\|_p \triangleq \left(\sum_{j=1}^N |y_j|^p \right)^{1/p} \quad (2)$$

where $0 < p \leq 1$. When $p = 0$, $\|\mathbf{y}\|_0$ denotes the ℓ^0 pseudo-norm that counts the number of non-zero elements of \mathbf{y} . The coefficient vector \mathbf{y} is said to be m -sparse if $\|\mathbf{y}\|_0 \leq m$.

We will use $\mathcal{N}(\Phi)$ for the null space of Φ . We will also make use of the subscript notation \mathbf{y}_T to denote a vector that is equal to some \mathbf{y} on the index set T and zero everywhere else. Denoting $|T|$ the cardinality of T , the vector \mathbf{y}_T is $|T|$ -sparse and we will say that the support of the vector \mathbf{y} lies within T whenever $\mathbf{y}_T = \mathbf{y}$.

1.2 Known conditions for ℓ^p sparse recovery

It has been shown in [13] that if:

$$\|\mathbf{z}_T\|_p < \|\mathbf{z}_{T^c}\|_p \quad (3)$$

holds for all $\mathbf{z} \in \mathcal{N}(\Phi)$ then any vector \mathbf{y}^* whose support lies within T , can be recovered by solving the following optimisation problem (which is non-convex for $0 \leq p < 1$):

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmin}} \|\mathbf{y}\|_p \text{ s.t. } \Phi \mathbf{y} = \Phi \mathbf{y}^*. \quad (4)$$

Irisa

Furthermore this condition, which is often referred to as the "null space property", is tight in that if the inequality (3) does not hold for some $\mathbf{z} \in \mathcal{N}(\Phi)$ then there exists a vector \mathbf{y}^* supported on T that is not the unique minimiser of (4). As a consequence, if (3) holds for all $\mathbf{z} \in \mathcal{N}(\Phi)$ and all index sets T of size m , then any m -sparse vector \mathbf{y}^* is recovered as the unique minimiser of (4). This condition is again tight, and it has been shown in [14, 15] that when it is satisfied for some $0 < p \leq 1$ it is also satisfied for all $0 \leq q \leq p$.

Using (4), particularly when $p = 1$, has become a popular means of solving for sparse representations. This is partly due to empirical evidence [5] that it often performs well and partly due to theoretical results [2, 3, 7, 13, 16]. An important concept in this regard that has been particularly influential in the emerging field of compressed sensing is the restricted isometry constant (RIC), δ_k . For a matrix Φ this is defined as the smallest number such that:

$$(1 - \delta_k) \leq \frac{\|\Phi \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} \leq (1 + \delta_k) \quad (5)$$

for every vector \mathbf{y} and every index set T with $|T| \leq k$. One weakness of the RIC is that the upper bound and the lower bound play fundamentally different roles and it is not preserved under a re-scaling of the dictionary [10] while recovery properties clearly are. One can, however, usually overcome the latter problem by considering an appropriately re-scaled dictionary such that both upper and lower bound is tight.

The RIC's importance can be linked with the following results:

1. Every m -sparse representation is unique if and only if [8]

$$\delta_{2m} < 1 \quad (6)$$

for an appropriately re-scaled dictionary. Furthermore almost every dictionary $\Phi \in \mathbb{R}^{M \times N}$ with $M \geq 2m$ satisfies this condition (again with appropriate re-scaling). Foucart and Lai [10] have also shown that for a given dictionary with $\delta_{2m+2} < 1$ there exists a sufficiently small p for which solving (4) is guaranteed to recover any m -sparse vector.

2. If

$$\delta_{2m} < 2(3 - \sqrt{2})/7 \approx 0.4531 \quad (7)$$

then every m -sparse representation can be exactly recovered using linear programming to solve (4) with $p = 1$, [10]. Furthermore most dictionaries $\Phi \in \mathbb{R}^{M \times N}$ (sampled from an appropriate probability model) will have an RIC $\delta_{2m} < \delta$ as long as: $M \geq C\delta^{-1}m \log(N/m)$, where C is some constant [1].

The RIC also bounds the condition number, κ , of submatrices, Φ_T , of a dictionary,

$$\kappa(\Phi_T) \leq \sqrt{\frac{1 + \delta_k}{1 - \delta_k}}, |T| \leq k \quad (8)$$

(indeed, Foucart and Lai [10] formulated their results in terms of the maximal submatrix condition number to avoid the re-scaling issues). This in turn bounds the Lipschitz constant of the inverse mapping resulting from solving the optimisation problem (4). In this regard the RIC also plays an important role in the noisy recovery problems [3, 10]: $\mathbf{x} = \Phi \mathbf{y} + \epsilon$ or $\mathbf{x} = \Phi(\mathbf{y} + \epsilon)$ where ϵ is an unknown but bounded noise term.

Note that when (7) holds all the $2m$ -submatrices have condition number $\kappa(\Phi_T) \leq 1.7$ when $|T| \leq 2m$, so they are extremely well behaved. In contrast, $\delta_{2m} < 1$ imposes no constraint on the condition number of the submatrices.

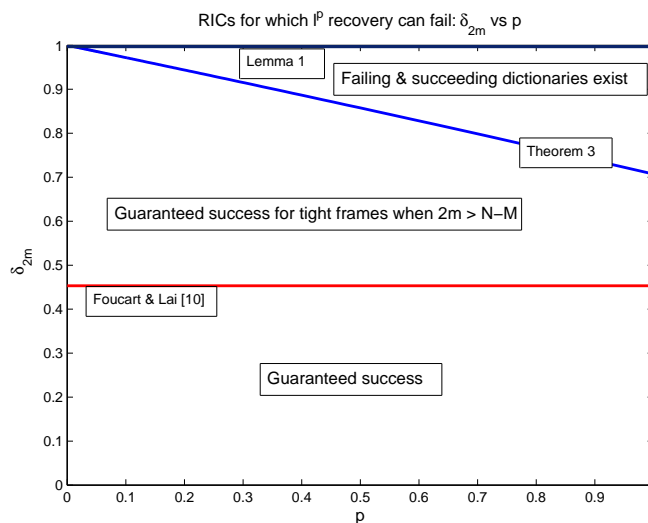


Figure 1: A summary of results now known ([10], Lemma 1 and Theorem 3) relating the restricted isometry constant to ℓ^p recovery.

1.3 Contributions

The bound (7) is an improvement over previous known bounds for ℓ^1 recovery [3]. However, in the proof of these bounds [3, 10], there are a number of estimates that are not tight. It is therefore an open question as to how much better we could expect to do, i.e. how large can we set $\delta \leq 1$ while still guaranteeing ℓ^1 recovery of any m -sparse vector for any dictionary with $\delta_{2m} < \delta$? This question is partially addressed by the following result:

Theorem 1. *For any $\epsilon > 0$ there exists an integer m and dictionary, Φ , with a restricted isometry constant $\delta_{2m} \leq 1/\sqrt{2} + \epsilon$ for which ℓ^1 recovery fails on some m -sparse vector.*

The proof is by an explicit construction which we will develop in the next section and is a by-product of some more general result concerning certain isometry conditions for which ℓ^p recovery fails, $0 < p \leq 1$. Indeed our complete results for RIC recovery conditions along with the result of [10] are summarised graphically in Figure 1.

The plot is divided up into three regions. Dictionaries in the bottom region [10] are guaranteed to succeed using any ℓ^p optimisation. In the top region there exist dictionaries (specifically minimally redundant unit norm tight frames) that are guaranteed to fail to recover at least one m -sparse vector y (Theorem 3). On the other hand, we can also find dictionaries (again minimally redundant unit norm tight frames) that are ℓ^p -succeeding for any $0 < p \leq 1$ with a RIC, δ_{2m} , arbitrarily close to one (Lemma 1).

Although there is a gap between the positive result of Foucart and Lai [10] for $p = 1$ and the negative result presented here, it is not a large one. For example, even if the positive result could be tightened to $\delta_{2m} < 1/\sqrt{2} - \epsilon$ – which would be the case if our negative results happened to be sharp (and the result *is* sharp for $2m > N - M$ with unit norm tight frames, see Corollary 1 below) – this would still require that the condition numbers of any $2m$ -column submatrix of Φ would have to be $\kappa(\Phi_T) \leq 2.5$, for $|T| \leq 2m$, which from any perspective is still extremely well conditioned.

The plot suggests that there might be some benefit in using $p \ll 1$ to improve sparse recovery. However in this case the optimisation problem is no longer convex and so we need to consider

algorithm specific recovery results. In this paper we examine the iterative reweighted ℓ^1 technique proposed in [3, 10] and present the following complement to Theorem 1.

Theorem 2. *For any $\epsilon > 0$ there exists an integer m and dictionary, Φ , with a restricted isometry constant $\delta_{2m} \leq 1/\sqrt{2} + \epsilon$ for which recovery using iteratively reweighted ℓ^1 fails on some m -sparse vector.*

This result does not necessarily imply that the uniform performance of iterative reweighted ℓ^1 is no better than ℓ^1 minimisation (although we suspect that the empirically observed benefits of such algorithms are more likely to be due to the presence of a range of coefficient scales). Instead the result highlights the danger of characterising sparse recovery uniformly in terms of the RIP.

The rest of the paper is structured as follows. In section 2 we introduce a variation on the classical RIC. We then develop our RIC results based upon an explicit minimally redundant unit spectral norm dictionary construction. In section 4 we explore our results numerically for both high dimensional dictionaries and a simple 1-sparse low dimensional example. Finally we examine the class of ℓ^p optimisation algorithms based upon iterative reweighted ℓ^1 . We conclude the paper with a discussion of implications of these results.

2 Isometry measures for unit spectral norm dictionaries

We will find it convenient to work with a slightly stronger condition than the usual restricted isometry property (RIP), one associated with **unit spectral norm dictionaries**, i.e. dictionaries such that

$$\|\Phi\| := \sup_{\mathbf{y} \neq 0} \frac{\|\Phi \mathbf{y}\|_2}{\|\mathbf{y}\|_2} = 1. \quad (9)$$

Definition 1 (asymmetric RIP). *Given a unit spectral norm dictionary $\Phi \in \mathbb{R}^{M \times N}$ let σ_k^2 be defined as:*

$$\sigma_k^2(\Phi) := \min_{\substack{\mathbf{y}_T \\ |T| \leq k}} \frac{\|\Phi \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} \quad (10)$$

We will usually drop the dependence on Φ when it is unambiguous. Clearly, as the maximum of any submatrix squared singular value is bounded by 1:

$$\max_{\substack{\mathbf{y}_T \\ |T| \leq k}} \frac{\|\Phi \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} \leq 1, \quad (11)$$

a unit spectral norm dictionary Φ with a given σ_k^2 implies the existence of a re-scaled dictionary, Ψ_k :

$$\Psi_k := \left(\frac{2}{1 + \sigma_k^2} \right)^{1/2} \Phi \quad (12)$$

with a RIC, δ_k :

$$\delta_k(\Psi_k) \leq \frac{(1 - \sigma_k^2(\Phi))}{(1 + \sigma_k^2(\Phi))} \quad (13)$$

The converse is not true since equality in (13) requires equality in (11). Under certain circumstances, however, equality can be assured.

Proposition 1 (Condition for equality in (13) with unit spectral norm tight frames). *Suppose that $\Phi \in \mathbb{R}^{M \times N}$ with $M \leq N$ is a unit spectral norm tight frame. Then for any $k > N - M$ we have:*

$$\delta_k(\Psi_k) = \frac{(1 - \sigma_k^2(\Phi))}{(1 + \sigma_k^2(\Phi))} \quad (14)$$

where Ψ_k is defined in (12). Moreover, Ψ_k is the optimal re-scaling of Φ with respect to the RIC δ_k in the sense that $\delta_k(\alpha\Phi) \geq \delta_k(\Psi_k)$ for any $\alpha > 0$. In particular for minimally redundant unit spectral norm tight frames (i.e., when $M = N - 1$) this is true for any $k \geq 2$.

Proof. Remember that by definition Φ is a frame [6] if there exists constants $0 < A \leq B < \infty$ such that for all \mathbf{x} ,

$$A\|\mathbf{x}\|_2^2 \leq \|\Phi^T \mathbf{x}\|_2^2 \leq B\|\mathbf{x}\|_2^2 \quad (15)$$

and a tight frame if the above holds with $A = B$. A unit spectral norm tight frame is therefore one for which $A = B = 1$, which is equivalently characterised by the condition $\Phi\Phi^T = \text{Id}$. For every vector $\mathbf{y} \in \mathbb{R}^N$, defining $\mathbf{x} := \Phi\mathbf{y}$ and $\mathbf{z} := \mathbf{y} - \Phi^T\Phi\mathbf{y}$ yields an orthogonal decomposition $\mathbf{y} = \Phi^T\mathbf{x} + \mathbf{z}$ hence

$$\|\Phi\mathbf{y}\|_2^2 = \|\mathbf{x}\|_2^2 = \|\Phi^T\mathbf{x}\|_2^2 = \|\mathbf{y}\|_2^2 - \|\mathbf{z}\|_2^2$$

and the upper bound in (11) is therefore achieved as long as we can find a \mathbf{y}_T that is in the range of Φ^T . For any T of size k , the dimension of the subspace spanned by all vectors of the form \mathbf{y}_T is k while the codimension of the range of Φ^T is $N - M$. Hence if $k > N - M$ there exists at least one nonzero vector in the intersection of these subspaces. The optimality of the re-scaled dictionary Ψ_k follows from the tightness of both upper and lower bounds in (5) for Ψ_k . \square

3 Dictionaries with small δ_{2m} where ℓ^p can fail

Our aim is to construct dictionaries Φ where sparse recovery will fail for at least one m -sparse vector $\mathbf{y} \in \mathbb{R}^N$. We consider the ℓ^p problem for any $0 < p \leq 1$ although we only provide closed form results for ℓ^1 . We are therefore looking for dictionaries that explicitly fail the ℓ^p recovery condition (3) while possessing small RIC δ_{2m} .

To find ' ℓ^p -failing dictionaries' (i.e., dictionaries for which ℓ^p minimisation fails to recover at least one m -sparse vector¹) with small RIC δ_{2m} , we will be looking for ℓ^p -failing dictionaries with largest possible σ_{2m}^2 . We will indeed prove somewhat more than Theorem 1, including tight results for ℓ^p -failure with unit spectral norm dictionaries in terms of asymmetric RIC σ_{2m}^2 , and tight results for ℓ^p -failure with unit spectral norm tight frames in terms of RIC δ_{2m} .

Theorem 3. *Consider $0 < p \leq 1$ and let $0 < \eta_p < 1$ be the unique positive solution to*

$$\eta_p^{2/p} + 1 = \frac{2}{p}(1 - \eta_p) \quad (16)$$

- *If $\Phi \in \mathbb{R}^{M \times N}$ is a unit spectral norm dictionary and $2m \leq M < N$ and*

$$\sigma_{2m}^2(\Phi) > 1 - \frac{2}{2-p}\eta_p \quad (17)$$

then all m -sparse vectors can be uniquely recovered by solving (4).

¹We will often omit the dependence on m when referring to ' ℓ^p -failing dictionaries'.

- For every $\epsilon > 0$, there exist integers $m \geq 1, N \geq 2m + 1$ and a minimally redundant unit spectral norm tight frame $\Phi \in \mathbb{R}^{(N-1) \times N}$ with:

$$\sigma_{2m}^2(\Phi) \geq 1 - \frac{2}{2-p}\eta_p - \epsilon \quad (18)$$

for which there exists an m -sparse vector which cannot be uniquely recovered by solving (4).

Whenever η_p is irrational the inequality in (17) can be replaced with \geq . Whenever η_p is rational, ϵ can be set to zero in (18).

Specialising to $p = 1$ we have $\eta_1^2 + 2\eta_1 - 2 = 0$, hence $\eta_1 = \sqrt{2} - 1$ and the right hand side in (17) is $3 - 2\sqrt{2}$. In terms of the standard RIP δ_{2m} for the re-scaled dictionary (12) with $k = 2m$ this means, using (13), that for any $\epsilon > 0$ there exists a dictionary Ψ with $\delta_{2m} < 1/\sqrt{2} + \epsilon$ where ℓ^1 recovery can fail, and Theorem 1 is proved.

Combining Theorem 3 with Proposition 1 above we get the following corollary:

Corollary 1. Assume that $\Phi \in \mathbb{R}^{M \times N}$ is a suitably re-scaled tight frame. If

$$N - M < 2m \leq M < N \quad (19)$$

and

$$\delta_{2m}(\Phi) < \frac{\eta_p}{2-p-\eta_p} \quad (20)$$

then all m -sparse vectors can be uniquely recovered by solving (4). Whenever η_p is irrational, the inequality in (20) can be replaced with \leq .

Strictly speaking the condition $2m \leq M$ is redundant with (20) since $2m > M$ implies $\delta_{2m} = 1$.

By the second part of Theorem 3, Corollary 1 is sharp in the sense that the right hand side in (20) cannot be weakened. This does not mean however that (20) is a necessary condition on the RIC for ℓ^p success, and there exist dictionaries with δ_{2m} arbitrarily close to one which recover every m -sparse vector, as expressed by the following lemma.

Lemma 1. For any $\epsilon > 0$, there exist integers m and N and a minimally redundant tight frame $\Phi_1 \in \mathbb{R}^{(N-1) \times N}$ along with re-scaled versions of Φ_1 , Φ_2 and Φ_3 , such that every m -sparse vector is recovered by solving (4) with any of Φ_1 , Φ_2 , Φ_3 and any $0 \leq p \leq 1$, yet

$$\sigma_{2m}^2(\Phi_1) \leq \epsilon \quad (21)$$

$$\sigma_{2m+1}^2(\Phi_1) = 0 \quad (22)$$

$$\delta_{2m}(\Phi_2) > 1 - \epsilon \quad (23)$$

$$\delta_{2m+1}(\Phi_3) = 1. \quad (24)$$

Theorem 3 will be proved by explicitly constructing the ℓ^p -failing unit spectral norm dictionaries with largest σ_{2m}^2 for a given pair (m, N) with $2m < N$, and a similar construction will be used to prove Lemma 1. We postpone the proofs and begin with a series of lemmatas.

Proposition 2 (Minimally redundant row orthonormal dictionaries are optimal among unit spectral norm dictionaries). Let $\Phi \in \mathbb{R}^{M \times N}$ be an arbitrary unit spectral norm dictionary which is ℓ^p -failing for some m -sparse vector with $M < N$. Then there exists a minimally redundant row orthonormal (unit spectral norm) dictionary $\Phi^* \in \mathbb{R}^{(N-1) \times N}$ which is ℓ^p -failing for some m -sparse vector such that for every k

$$\sigma_k^2(\Phi) \leq \sigma_k^2(\Phi^*). \quad (25)$$

Proof. We consider the singular value decomposition: $\Phi = V\Sigma U^T$ where $V \in \mathbb{R}^{M \times M}$ and $U^T \in \mathbb{R}^{M \times N}$ are row orthonormal, and $\Sigma \in \mathbb{R}^{M \times M}$ is diagonal. Since Φ has unit spectral norm $\|\Sigma\| = 1$ and we have for any \mathbf{y} and any T

$$\frac{\|\Phi \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} \leq \frac{\|U^T \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2}.$$

Since Φ is ℓ^p -failing for some m -sparse vector, by the characterisation (3), there exists some offending $\mathbf{z} \in \mathcal{N}(\Phi)$ and an index set I_m of size m such that

$$\|\mathbf{z}_{I_m}\|_p^p \geq \|\mathbf{z}_{I_m^c}\|_p^p \quad (26)$$

Now let $W \in \mathbb{R}^{N \times (N-M-1)}$ be an orthonormal basis over the orthogonal complement to $\{\mathbf{z}, U\}$, such that $\{\mathbf{z}, U, W\}$ forms an orthonormal basis over \mathbb{R}^N . We can then write any $\mathbf{y}_T \in \mathbb{R}^N$ as:

$$\mathbf{y}_T = \mathbf{z}a + U\mathbf{b} + W\mathbf{c}$$

for some $a \in \mathbb{R}$, $\mathbf{b} \in \mathbb{R}^M$ and $\mathbf{c} \in \mathbb{R}^{N-M-1}$. Define the minimally redundant row orthonormal dictionary $\Phi^* := [U, W]^T \in \mathbb{R}^{(N-1) \times N}$. First, for any \mathbf{y}_T we have:

$$\frac{\|\Phi \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} \leq \frac{\|U^T \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2} = \frac{\|\mathbf{b}\|_2^2}{a^2 + \|\mathbf{b}\|_2^2 + \|\mathbf{c}\|_2^2} \leq \frac{\|\mathbf{b}\|_2^2 + \|\mathbf{c}\|_2^2}{a^2 + \|\mathbf{b}\|_2^2 + \|\mathbf{c}\|_2^2} = \frac{\|\Phi^* \mathbf{y}_T\|_2^2}{\|\mathbf{y}_T\|_2^2}$$

therefore $\sigma_k^2(\Phi) \leq \sigma_k^2(\Phi^*)$. To conclude the proof we observe that by construction $\Psi \mathbf{z} = 0$, hence $\mathbf{z} \in \mathcal{N}(\Phi^*)$, which combined with (26) and the characterisation (3) shows that Ψ is ℓ^p -failing for at least one m -sparse vector. \square

The proposition above shows that ℓ^p -failing unit spectral norm dictionaries with largest σ_{2m}^2 can be searched within the restricted set of ℓ^p -failing minimally redundant row orthonormal dictionaries. We next evaluate the minimal singular values of the submatrices made of k columns of such Ψ .

Proposition 3 (Minimal singular values of submatrices are characterised by the null space). *Let $\Phi \in \mathbb{R}^{(N-1) \times N}$ be a minimally redundant row orthonormal dictionary, and let $\mathbf{z} \in \mathbb{R}^N$ with $\|\mathbf{z}\|_2 = 1$ be a vector which spans $\mathcal{N}(\Phi)$. Denoting I_k the set indexing the k largest components of \mathbf{z} we have for every k*

$$\sigma_k^2(\Phi) = 1 - \|\mathbf{z}_{I_k}\|_2^2. \quad (27)$$

Proof. Since $\Phi \mathbf{z} = 0$ and Φ is row orthonormal, $[\mathbf{z}, \Phi^T]$ forms an orthonormal basis in \mathbb{R}^N , and we can again write any vector \mathbf{y} as:

$$\mathbf{y} = \mathbf{z}a + \Phi^T \mathbf{b}$$

where $a \in \mathbb{R}$ and $\mathbf{b} \in \mathbb{R}^{N-1}$, and therefore $\|\Phi \mathbf{y}\|_2^2 = \|\mathbf{b}\|_2^2$. If \mathbf{y} has unit norm then

$$1 = \|\mathbf{y}\|_2^2 = a^2 + \|\mathbf{b}\|_2^2 = |\langle \mathbf{z}, \mathbf{y} \rangle|^2 + \|\Phi \mathbf{y}\|_2^2$$

To find the minimal singular value associated with the submatrix Φ_T we need to solve the problem

$$\begin{aligned} \sigma_k^2(\Phi) &= \min_{T, |T| \leq k} \min_{\substack{\mathbf{y}_T \\ \|\mathbf{y}_T\|_2 = 1}} \|\Phi \mathbf{y}_T\|_2^2 \\ &= 1 - \max_{T, |T| \leq k} \max_{\substack{\mathbf{y}_T \\ \|\mathbf{y}_T\|_2 = 1}} |\langle \mathbf{z}, \mathbf{y}_T \rangle|^2 \end{aligned}$$

i.e., we need to find the unit vector \mathbf{y}_T^* that is maximally correlated with \mathbf{z} . For a given T this is satisfied with $\mathbf{y}_T^* = \mathbf{z}_T / \|\mathbf{z}_T\|$, in which case $|\langle \mathbf{z}, \mathbf{y}_T \rangle|^2 = \|\mathbf{z}_T\|_2^2$. The best T is the one which captures the k largest components of \mathbf{z} , that is to say $T^* = I_k$. \square

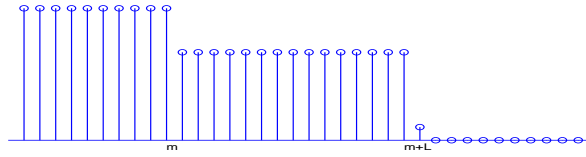


Figure 2: A stylised depiction of an optimal null vector for (28-31).

The proposition above shows that for minimally redundant row orthonormal dictionaries, $\sigma_k^2(\Phi)$ is completely determined by the unit vector \mathbf{z} which spans the null space $\mathcal{N}(\Phi)$. Our original problem was to select an ℓ^p -failing minimally redundant row orthonormal dictionary Φ with maximal $\sigma_k^2(\Phi)$ for $k = 2m$. This is now turned into an optimisation problem where we wish to select a unit norm vector \mathbf{z} that allows ℓ^p reconstruction failure for m -sparse vectors, while maximising σ_k^2 , i.e. minimising $\|\mathbf{z}_{I_k}\|_2^2$.

Without loss of generality, up to column permutation of Φ and sign changes, we may assume that $z_i \geq z_{i+1} \geq 0$, and the ℓ^p -failing assumption is that $\|\mathbf{z}_{I_m}\|_p^p \geq \|\mathbf{z}_{I_m^c}\|_p^p$. With a little manipulation the optimisation problem for finding a failing \mathbf{z} with maximal associated σ_k^2 can be written in the form of (28-31) below. The next lemma identifies the particularly simple form of the optimal null vectors which is also depicted in Figure 2.

Lemma 2 (Shape of the optimal vector \mathbf{z} of the null space). *Consider $k \geq 2m$ and let $\mathbf{z}^* \in \mathbb{R}^N$ be a solution to the following optimisation problem:*

$$\text{minimise: } J(\mathbf{z}) := \frac{\|\mathbf{z}_{\Lambda_0}\|_2^2 + \|\mathbf{z}_{\Lambda_1}\|_2^2}{\|\mathbf{z}_{\Lambda_2}\|_2^2} \quad (28)$$

$$\text{subject to: } \frac{\|\mathbf{z}_{\Lambda_1}\|_p^p + \|\mathbf{z}_{\Lambda_2}\|_p^p}{\|\mathbf{z}_{\Lambda_0}\|_p^p} \leq 1 \quad (29)$$

$$\|\mathbf{z}\|_2^2 = 1 \quad (30)$$

$$\text{and } z_i \geq z_{i+1} \geq 0 \quad (31)$$

where $\Lambda_0 = \{1, \dots, m\}$, $\Lambda_1 = \{m+1, \dots, k\}$ and $\Lambda_2 = \{k+1, \dots, N\}$. Then \mathbf{z}^* is piecewise flat, and has the form:

$$\mathbf{z}^* = [\underbrace{\alpha, \dots, \alpha}_m, \underbrace{\beta, \dots, \beta}_L, \gamma, 0, \dots, 0]^T \quad (32)$$

for some constants $\alpha > \beta > \gamma \geq 0$ and some L such that $k+1 \leq m+L \leq N$. Furthermore (29) holds with equality for \mathbf{z}^* .

Proof. We first note that, due to the continuity of $J(\mathbf{z})$ and the compactness of the constraint set, an optimum \mathbf{z}^* is guaranteed to exist. Then we prove by contradiction that \mathbf{z}^* must have the claimed form.

- $\mathbf{z}_{\Lambda_0}^*$ is flat. We know that

$$\|\mathbf{z}_{\Lambda_0}^*\|_2 \geq m^{1/2-1/p} \|\mathbf{z}_{\Lambda_0}^*\|_p$$

with equality only if $z_i^* = m^{-1/p} \cdot \|\mathbf{z}_{\Lambda_0}^*\|_p$ for all $i \in \Lambda_0$, i.e. if $\mathbf{z}_{\Lambda_0}^*$ is "flat". Therefore, if $\mathbf{z}_{\Lambda_0}^*$ is not flat, then we can find a \mathbf{z}' such that $\mathbf{z}'_{\Lambda_1 \cup \Lambda_2} = \mathbf{z}_{\Lambda_1 \cup \Lambda_2}^*$ and $z_i' = m^{-1/p} \cdot \|\mathbf{z}_{\Lambda_0}^*\|_p > \min_{j \in \Lambda_0} |z_j^*| \geq \|\mathbf{z}_{\Lambda_1 \cup \Lambda_2}^*\|_\infty$ for all $i \in \Lambda_0$. Now let $\mathbf{z}'' = \mathbf{z}' / \|\mathbf{z}'\|_2$. \mathbf{z}'' is feasible and $J(\mathbf{z}'') = J(\mathbf{z}') < J(\mathbf{z}^*)$ which contradicts the fact that \mathbf{z}^* is an optimum. Hence $\mathbf{z}_{\Lambda_0}^*$ must be flat.

- $\mathbf{z}_{\Lambda_1}^*$ is flat with all entries equal to $z_{k+1}^* = \|\mathbf{z}_{\Lambda_2}^*\|_\infty$. By contradiction, assume that $z_i^* \neq z_{k+1}^*$ for some $i \in \Lambda_1$. Then, we can construct a \mathbf{z}' with $\mathbf{z}'_{\Lambda_0 \cup \Lambda_2} = \mathbf{z}_{\Lambda_0 \cup \Lambda_2}$ and $z'_i = z_{k+1}^*$ for all $i \in \Lambda_1$. Again re-scale: $\mathbf{z}'' = \mathbf{z}' / \|\mathbf{z}'\|_2$. Thus \mathbf{z}'' is feasible and $J(\mathbf{z}'') = J(\mathbf{z}') < J(\mathbf{z}^*)$. Hence $\mathbf{z}_{\Lambda_1}^*$ must be flat with value z_{k+1}^* .
- **Shape of $\mathbf{z}_{\Lambda_2}^*$.** Now consider the index set Λ_2 . Suppose that there are two indices $k+1 < j < l < N$ such that $z_{k+1}^* = z_{j-1}^* > z_j^* \geq z_l^* > z_{l+1}^* = z_N^*$. We can then construct a \mathbf{z}' with non-increasing entries such that $z'_i = z_i$ for all $i \neq \{j, l\}$ and

$$\begin{aligned} z'_j &> z_j^* \\ z'_l &< z_l^* \\ |z'_j|^p + |z'_l|^p &= |z_j^*|^p + |z_l^*|^p \end{aligned}$$

Lemma 4 (in the Appendix) implies that $\|\mathbf{z}'_{\Lambda_2}\|_2 > \|\mathbf{z}_{\Lambda_2}^*\|_2$, hence $J(\mathbf{z}') < J(\mathbf{z}^*)$. Again we can re-scale to make the vector feasible. Hence we can conclude that $\mathbf{z}_{\Lambda_2}^*$ can only have one element not equal to z_{k+1}^* or z_N^* . A similar analysis shows that $z_N^* = 0$, and this concludes the proof that \mathbf{z}^* must have the form in (32) with $\alpha \geq \beta > \gamma \geq 0$ and $k+1 \leq m+L \leq N$. Moreover by (29) we have

$$m \cdot \alpha^p = \|\mathbf{z}_{\Lambda_0}^*\|_p^p \geq \|\mathbf{z}_{\Lambda_1}^*\|_p^p + \|\mathbf{z}_{\Lambda_2}^*\|_p^p \geq L \cdot \beta^p \geq (k+1-m) \cdot \beta^p > m \cdot \beta^p$$

hence $\alpha > \beta$.

- **Constraint (29) hold with equality for \mathbf{z}^* .** Suppose that the left hand side of (29) is strictly less than one for \mathbf{z}^* . Since $\alpha > \beta$, we could then find $a < 1$ such that \mathbf{z}' defined with $\mathbf{z}'_{\Lambda_1 \cup \Lambda_2} = \mathbf{z}_{\Lambda_1 \cup \Lambda_2}^*$ and $\mathbf{z}'_{\Lambda_0} = a\mathbf{z}_{\Lambda_0}^*$, properly re-scaled, simultaneously reduces the objective function (28) while still satisfying (29) and (31). Therefore (29) must hold with equality for any optimal \mathbf{z}^* . \square

Lemma 2 implies that we only have to consider a relatively simple form for \mathbf{z} , which is parameterised by $\alpha > \beta \geq \gamma \geq 0$ and m, L , where $k-m+1 \leq L \leq N-m$. Note that any zero elements in \mathbf{z} can be removed by simply reducing the dimension N of the dictionary. In order to calculate the largest σ_k^2 we need to evaluate optimal values for α, β, γ, m and L . In fact we will see that we can ignore γ , which comes from the fact that the optimal constructions will correspond to m and N very large. The following lemma is expressed for $k = 2m$ but straightforward modifications would make it possible to handle arbitrary $k \geq 2m$.

Lemma 3 (Calculating the largest σ_{2m}^2). *Consider $k = 2m < N$, $0 < p \leq 1$ and let η_p be the unique positive solution to (16). Let $\mathbf{z} \in \mathbb{R}^N$ be of the form (32) with $\alpha > \beta > \gamma \geq 0$ and $m+1 \leq L \leq N-m$, and assume that \mathbf{z} satisfies (29) with equality and (30). Then*

$$\|\mathbf{z}_{I_{2m}}\|_2^2 \geq \frac{2}{2-p} \eta_p. \quad (33)$$

If η_p is rational, equality is achieved for some \mathbf{z}^ . Otherwise, the inequality can be replaced with $>$, but one can get arbitrarily close to the lower bound with appropriate choices of $k = 2m < N$ and \mathbf{z} satisfying all the above conditions.*

Proof. Define

$$L' := L + (\gamma/\beta)^p \quad (34)$$

$$\eta := m/L'. \quad (35)$$

Since $\gamma < \beta$, we have $L \leq L' < L + 1$, and since $m + 1 \leq L$ we have $0 < \eta < 1$. The ℓ^p -failure equality constraint (29) reads $m\alpha^p = L\beta^p + \gamma^p = L'\beta^p$ hence

$$\beta = \eta^{1/p} \cdot \alpha < \alpha \quad (36)$$

Similarly by (30) we have $m\alpha^2 + L\beta^2 + \gamma^2 = 1$, and we let the reader check that this implies

$$m\alpha^2 + L'\beta^2 = 1 + (\gamma/\beta)^p\beta^2 - \gamma^2 \geq 1 \quad (37)$$

with equality when L' is integer (i.e. when $\gamma = 0$). Combining the two constraints we have:

$$m\alpha^2 \geq (1 + \eta^{2/p-1})^{-1} \quad (38)$$

and it follows that

$$\|\mathbf{z}_{I_{2m}}\|_2^2 = m\alpha^2 + m\beta^2 = m\alpha^2 (1 + \eta^{2/p}) \geq \frac{(1 + \eta^{2/p})}{(1 + \eta^{2/p-1})} \quad (39)$$

Differentiating the right hand side and equating to zero, we obtain that the value η_p that minimises the bound on $\|\mathbf{z}_{I_{2m}}\|_2^2$ for $0 < \eta < 1$ satisfies (16). Substituting this back into (39) gives:

$$\|\mathbf{z}_{I_{2m}}\|_2^2 \geq \frac{2}{2-p}\eta_p \quad (40)$$

Now that we have established the bound we discuss its tightness. First, one can check that for $0 < p \leq 1$, Equation (16) always has a unique solution in the region $\eta_p > 0$, though the solution does not appear to have a general closed form. Then, notice that by continuity, the right hand side in (39) can get arbitrarily close to the right hand side in (40) by choosing η sufficiently close to η_p . Moreover, by the density of the rational numbers in \mathbb{R} , we can always find integers m and L such that m/L gets arbitrarily close to η_p . For such integers, setting $\gamma = 0$ (so that $L' = L$ and $\eta = m/L$), choosing α to reach equality in (38), and setting β according to (36) yields a vector \mathbf{z}^* for which $\|\mathbf{z}_{I_{2m}}^*\|_2^2$ is arbitrarily close to the lower bound. If η_p is rational then equality is actually achieved. If η_p is irrational, then equality cannot be achieved. \square

We are now able to state the proof of Theorem 3.

Proof of Theorem 3. Consider a unit spectral norm dictionary Φ which satisfies (17). Assume that Φ is ℓ^p -failing for some m -sparse vector. Then, by Proposition 2, there exists a minimally redundant row orthonormal (unit spectral norm) dictionary $\Phi^* \in \mathbb{R}^{(N-1) \times N}$ which is ℓ^p -failing for some m -sparse vector such that

$$\sigma_{2m}^2(\Phi) \leq \sigma_{2m}^2(\Phi^*)$$

By Proposition 3,

$$\sigma_{2m}^2(\Phi^*) = 1 - \|\mathbf{z}_{I_{2m}}\|_2^2$$

where \mathbf{z} is a unit norm vector which spans the null space $\mathcal{N}(\Phi^*)$. Since Φ^* is ℓ^p -failing, \mathbf{z} (after proper reindexing and taking the absolute value) satisfies the constraints (29), (30) and (31), therefore by Lemma 2 and Lemma 3,

$$\|\mathbf{z}_{I_{2m}}\|_2^2 \geq \frac{2}{2-p}\eta_p. \quad (41)$$

We conclude that

$$\sigma_{2m}^2(\Phi) \leq 1 - \frac{2}{2-p}\eta_p.$$

By contraposition, if $\sigma_{2m}^2(\Phi) > 1 - \frac{2}{2-p}\eta_p$ then Φ cannot be ℓ^p -failing for any m -sparse vector. If η_p is irrational, the inequality in (41) can be replaced with $>$ hence it is sufficient to assume that $\sigma_{2m}^2(\Phi) \geq 1 - \frac{2}{2-p}\eta_p$.

Conversely, by the above Propositions and Lemmas, for every $\epsilon > 0$ there exists some \mathbf{z}^* satisfying the constraints (29), (30) and (31) for which

$$\|\mathbf{z}_{I_{2m}}^*\|_2^2 \leq \frac{2}{2-p}\eta_p + \epsilon, \quad (42)$$

yielding a (minimally redundant, row orthonormal) unit spectral norm dictionary Φ_p^* with

$$\sigma_{2m}^2(\Phi_p^*) \geq 1 - \frac{2}{2-p}\eta_p - \epsilon$$

which is ℓ^p -failing for some m -sparse vector. If η_p is rational, this is true for $\epsilon = 0$. \square

Let us proceed with the proof of Corollary 1.

Proof of Corollary 1. Since Φ is a tight frame, $\Phi = A \cdot \tilde{\Phi}$ for some unit spectral norm tight frame $\tilde{\Phi}$ and some real constant $0 < A < \infty$. For $N - M < 2m \leq M$, since Φ is a re-scaled version of $\tilde{\Phi}$, by Proposition 1 we have

$$\frac{1 - \sigma_{2m}^2(\tilde{\Phi})}{1 + \sigma_{2m}^2(\tilde{\Phi})} = \delta_{2m}(\tilde{\Psi}_{2m}) \leq \delta_{2m}(\Phi) < \frac{\eta_p}{2 - p - \eta_p}$$

with $\tilde{\Psi}_{2m}$ the optimally re-scaled version of $\tilde{\Phi}$ given by (12), hence

$$\sigma_{2m}^2(\tilde{\Phi}) > 1 - \frac{2}{2-p}\eta_p$$

and we can apply Theorem 3 to conclude. \square

We conclude this section with the proof of Lemma 1.

Proof of Lemma 1. Consider $\mathbf{z} = (\mathbf{z}_0, \mathbf{z}_1) \in \mathbb{R}^N$ where $N = 2m + 1$, $\mathbf{z}_0 \in \mathbb{R}^m$ is "flat" with entries $1/\sqrt{2m}$ and $\mathbf{z}_1 \in \mathbb{R}^{m+1}$ is "flat" with entries $1/\sqrt{2m+2}$. Check that \mathbf{z} has non-increasing entries, is ℓ^2 normalized and satisfies the ℓ^p -recovery condition for m -sparse vectors for $p = 0$ as well as for every $0 < p \leq 1$:

$$\|\mathbf{z}_0\|_p = m^{1/p-1/2} \cdot \|\mathbf{z}_0\|_2 = m^{1/p-1/2} \cdot \|\mathbf{z}_1\|_2 < (m+1)^{1/p-1/2} \cdot \|\mathbf{z}_1\|_2 = \|\mathbf{z}_1\|_p \quad (43)$$

Let $\Phi_1 \in \mathbb{R}^{(N-1) \times N}$ be a row orthonormal dictionary with null space spanned by \mathbf{z} : by the above properties, for every $0 \leq p \leq 1$, every m -sparse vector is recovered by the minimisation (4). By Proposition 3, for every k we have $\sigma_k^2(\Phi_1) = 1 - \|\mathbf{z}_{I_k}\|_2^2$, and in particular

$$\sigma_{2m}^2(\Phi_1) = \frac{1}{2m+2}; \sigma_{2m+1}^2(\Phi_1) = 0.$$

Moreover, since $2 \leq k = 2m \leq M$ and $2 \leq k' = 2m+1 \leq M$, by Proposition 1 we have

$$\delta_{2m}(\Phi_2) = \frac{2m+1}{2m+3}; \delta_{2m+1}(\Phi_3) = 1;$$

with $\Phi_2 = \Psi_{2m}$ and $\Phi_3 = \Psi_{2m+1}$ the appropriately re-scaled dictionaries. \square

4 Numerical studies of the σ_{2m}^2 and δ_{2m} conditions

We now take a brief look at numerical solutions for values of σ_{2m}^2 and δ_{2m} for which ℓ^p recovery can fail.

4.1 Large dimensional ℓ^p failing dictionaries

The analysis carried out so far, which relies on constructions for large dimensions m and N , shows that

$$\begin{aligned} \sup_{m, \Phi} \sigma_{2m}^2(\Phi) &= 1 - \frac{2\eta_p}{2-p} \\ \inf_{m, \Phi} \delta_{2m}(\Phi) &= \frac{\eta_p}{2-p-\eta_p} \end{aligned}$$

where the supremum is over integers m and unit spectral norm ℓ^p -failing dictionaries Φ , the infimum is over integers m and ℓ^p -failing tight frames $\Phi \in \mathbb{R}^{M \times N}$ with $2m \leq M < N < M + 2m$. This provides two curves $\sigma^2(p)$ and $\delta(p)$ for which there exists ℓ^p -failing dictionaries with σ_{2m}^2 above (respectively δ_{2m} below) or arbitrarily close to $\sigma^2(p)$ (resp. $\delta(p)$). To compute these curves we need to solve for η_p in Equation (16). For $p \in \{1, 2/3, 1/2\}$, this is a polynomial equation of degree $d = 2/p \in \{2, 3, 4\}$ which roots have algebraic expressions. In practice we rely on numerical solvers to compute η_p , $\sigma^2(p)$ and $\delta(p)$, which are displayed as a solid line on Figure 3 and Figure 4.

The work of [14, 15] showed that there is a whole family of sparsity measures including ℓ^p that span between ℓ^0 and ℓ^1 , and that solving (4) for $p < 1$ could offer gradually superior performance to ℓ^1 recovery when p decreases. The results in [10] provided quantitative ℓ^p -recovery conditions based on RIC. Here we see from Figure 4 that the offending RIC grows very gently as p shrinks. This implies that, at least in terms of worst case RIP analysis over all dictionaries, using a p slightly smaller than 1 does not provide a large benefit, and that one would need to rely on a $p \ll 1$ to expect a significant difference. However since solving (4) for $p < 1$ is non-convex such benefit will dependent on the specific choice of optimisation algorithm. For example we will see in the section 5 that iterative reweighted ℓ^1 techniques do not appear to provide uniform performance benefits beyond ℓ^1 minimisation.

These results may also seem at odds with a long history of empirical studies showing the benefits of ℓ^p optimisation for sparse recovery dating back to [12], however we note first that empirical results generally indicate an average performance bound rather than a uniform one and second the success of

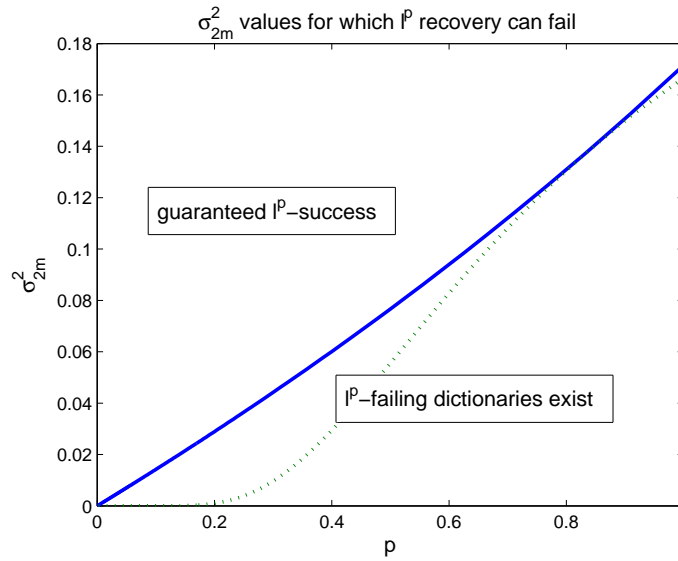


Figure 3: A plot of the σ_{2m}^2 values for which ℓ^p recovery with *unit spectral norm dictionaries* can fail (solid). This result is sharp in that for any (p, σ^2) above the line, a dictionary with $\sigma_{2m}^2 = \sigma^2$ is guaranteed to recover m -sparse representations by solving (4), while for any (p, σ^2) below the line we can find a dictionary with $\sigma_{2m}^2 = \sigma^2$ for which ℓ^p -recovery will fail on at least one m -sparse vector. The dashed line corresponds to the values for the best failing 2×3 dictionaries calculated in section 4.2.

ℓ^p optimisation seems to be predominantly associated with sparsity problems with a range of coefficient sizes, such as Gaussian distributed sparse coefficients, where the ℓ^p algorithm is able to pick off the larger coefficients first. Note that successful recovery in ℓ^1 optimisation is only a function of the sign of the coefficients [11] and thus is unable to exploit differences in coefficient size.

4.2 Low dimensional examples

Although our arguments above require $N \rightarrow \infty$ in order to approach the bound, in fact, it is very easy to construct a specific low dimensional example that is very close to it. Consider a $\Phi \in \mathbb{R}^{2 \times 3}$ for which ℓ^p minimisation just fails in the 1-sparse case. Select:

$$\mathbf{z} = \frac{1}{\sqrt{1 + 2^{1-2/p}}} \cdot \begin{bmatrix} 1 \\ -2^{-1/p} \\ -2^{-1/p} \end{bmatrix} \quad (44)$$

and generate any Φ such that Φ^T is the orthogonal complement to \mathbf{z} . For example we can have:

$$\Phi = \begin{bmatrix} \frac{1}{\sqrt{1+2^{2/p-1}}} & \frac{2^{1/p-1}}{\sqrt{1+2^{2/p-1}}} & \frac{2^{1/p-1}}{\sqrt{1+2^{2/p-1}}} \\ 0 & \frac{1}{\sqrt{2}} & \frac{-1}{\sqrt{2}} \end{bmatrix} \quad (45)$$

For the ℓ^1 case this gives:

$$\mathbf{z} = \frac{1}{\sqrt{6}} \begin{bmatrix} 2 \\ -1 \\ -1 \end{bmatrix} \quad (46)$$

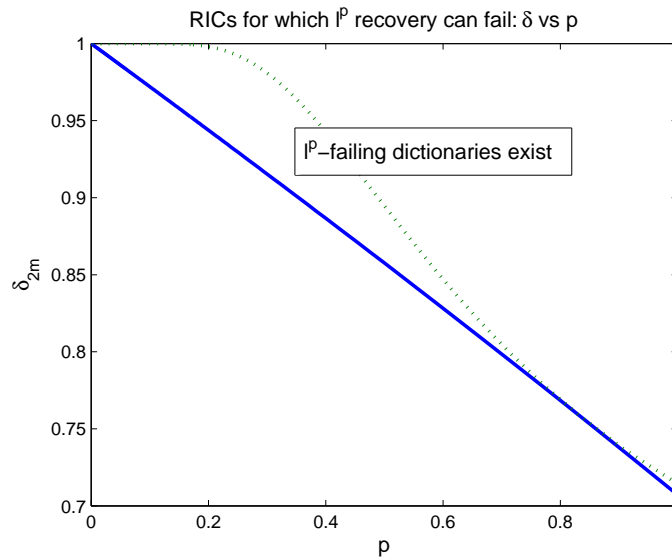


Figure 4: A plot of the RIC values, δ for which ℓ^p recovery can fail (solid) for *any* dictionary. Above the line, for any (p, δ) , we can find a dictionary with $\delta_{2m} = \delta$ for which ℓ^p -recovery will fail on at least one m -sparse vector. However, the result may not be sharp (except for the special case of tight frames with $2m > N - M$ – Corollary 1) since δ_{2m} and σ_{2m}^2 are only necessarily related through the inequality (13). Thus there may also exist failing dictionaries below the line. The dashed line corresponds to the RIC values for the re-scaled best failing 2×3 dictionaries in section 4.2.

and

$$\Phi = \frac{1}{\sqrt{6}} \begin{bmatrix} \sqrt{2} & \sqrt{2} & \sqrt{2} \\ 0 & \sqrt{3} & -\sqrt{3} \end{bmatrix} \quad (47)$$

The RIC, δ_2 , that can fail for this low dimensional example is also plotted as a dashed line in Figure 4. It was simply computed by considering the three 2×2 submatrices of Φ and computing their maximal and minimal singular values. Note that for ℓ^1 this is within 0.01 of the general condition for failure (due, no doubt, to the excellent engineering approximation of $\sqrt{2} \approx 3/2$: the offending \mathbf{z} in (46) has the shape (32) for $m = 1, L = 2$, i.e. with $\eta = 1/2$, while the optimum is for $\eta_1 = \sqrt{2} - 1$). The value of p for which $\eta_p = 1/2$ is optimal can be found by numerically solving $(1/2)^{2/p} + 1 = 1/p$. This gives $p \approx 0.839$ for which the 2×3 construction is actually optimum. Note that the two curves in Figure 4 touch at this value of p .

5 Reweighted ℓ^1 implementations for ℓ^p -optimisation

It is important to distinguish between optimisation functions and recovery algorithms. All the results in the previous sections have been derived for the recovery properties associated with the global minimum solutions for (4) without any regard for how these might be obtained. In practise, solving (4) for $p < 1$ is non-trivial. When $p < 1$ the cost function ceases to be convex and there are many local minima. One approach that has recently been proposed [4, 10] is to attempt to solve (4) by solving a sequence of reweighted ℓ^1 optimisation problems of the form:

$$\hat{\mathbf{y}}^{(n)} = \underset{\mathbf{y}}{\operatorname{argmin}} \|\mathbf{W}_n \mathbf{y}\|_1 \text{ s.t. } \Phi \mathbf{y} = \Phi \mathbf{y}^*. \quad (48)$$

where the initial weight matrix is set to the identity, $\mathbf{W}_1 = \mathbf{Id}$, and then subsequently \mathbf{W}_n is selected as a diagonal positive definite weight matrix that is a (possibly iteration dependent) function of the previous solution vector, $\mathbf{W}_n = f_n(\mathbf{y}^{(n-1)})$. At any step, the solution to the convex optimisation problem (48) can be characterised by the necessary and sufficient property

$$\forall \mathbf{z} \in \mathcal{N}(\Phi), |\langle \mathbf{W}_n \mathbf{z}, \text{sign}(\hat{\mathbf{y}}^{(n)}) \rangle| \leq \|(\mathbf{W}_n \mathbf{z})_{\Gamma_n}\|_1 \quad (49)$$

where Γ_n denotes the set indexing the *zero* entries in $\hat{\mathbf{y}}^{(n)}$.

In [4], as an approximation to the ℓ^0 minimisation problem, the following reweighting function was proposed:

$$W_n(k, k) = \left(\epsilon_n + |y_i^{(n-1)}| \right)^{p-1} \quad (50)$$

for $p = 0$ and some small $\epsilon_n > 0$. Here $W_n(k, k)$ denotes the k th diagonal element of $\mathbf{W}_n(k, k)$. In [10] the authors consider the same weighting function but including the full range of $0 \leq p < 1$. Note that the inclusion of the ϵ_n term is crucial as it keeps $W_n(k, k)$ bounded and ensures that a zero valued component y_i is able to become non-zero again at some subsequent iteration. Candès *et al.* [4] discuss using either a fixed ϵ_n or selecting it, while Foucart and Lai [10] argue that, at least in terms of the associated cost functions, letting $\epsilon_n \rightarrow 0$ converges to a solution for (4). It is also noted in [4] that there are various other reweighting strategies that could be deployed, some of which may not even be associated with a specific cost function.

A natural question to ask is: *what is the guaranteed performance of such algorithms?* In order to consider the widest possible set of reweighting schemes we define the following that we consider to encompass all ‘reasonable’ reweighting schemes.

Definition 2 (Admissible reweighting schemes). *A reweighting scheme is considered to be admissible if, $\mathbf{W}_1 = \mathbf{Id}$ and if, for each n , there exists a $w_{\max}^n < \infty$ such that for all k , $0 \leq W_n(k, k) \leq w_{\max}^n$ and $W_n(k, k) = w_{\max}^n \Leftrightarrow \hat{y}_k^{(n)} = 0$.*

The next two proposals shed some light on what performance guarantees we might expect from such schemes.

Proposition 4 (Iteratively reweighted ℓ^1 is not worse than ℓ^1). *Let Φ be an arbitrary dictionary and T an arbitrary support set. If ℓ^1 recovery is successful for all vectors with support set T , then recovery using iteratively reweighted ℓ^1 with any admissible reweighting scheme is also successful for all vectors with support T .*

Proof. Assume that T is a support set for which ℓ^1 is guaranteed to succeed: i.e., $\|\mathbf{z}_T\|_1 \leq \|\mathbf{z}_{T^c}\|_1, \forall \mathbf{z} \in \mathcal{N}(\Phi)$. Since $\mathbf{W}_1 = \mathbf{Id}$, for any \mathbf{y}^* supported in T , $\hat{\mathbf{y}}^{(1)}$ is the ℓ^1 minimiser therefore $\hat{\mathbf{y}}^{(1)} = \mathbf{y}^*$. As a result, $T^c \subset \Gamma_1$, and for $k \in T^c$, $W_2(k, k) = w_{\max}^2$, therefore

$$\forall \mathbf{z} \in \mathcal{N}(\Phi), |\langle \mathbf{W}_2 \mathbf{z}, \text{sign}(\hat{\mathbf{y}}^{(1)}) \rangle| \leq w_{\max}^2 \|\mathbf{z}_T\|_1 \leq w_{\max}^2 \|\mathbf{z}_{T^c}\|_1 \leq \|(\mathbf{W}_2 \mathbf{z})_{\Gamma_1}\|_1$$

It follows that $\hat{\mathbf{y}}^{(2)} = \hat{\mathbf{y}}^{(1)}$ and iteratively one gets $\hat{\mathbf{y}}^{(n)} = \mathbf{y}^*$ for all n . \square

Proposition 4 indicates that the reweighting strategy cannot damage an already successful solution. However we also have the following negative result.

Proposition 5 (Iteratively reweighted ℓ^1 is not uniformly better than ℓ^1). *Let $\Phi \in \mathbb{R}^{(N-1) \times N}$ be a minimally redundant dictionary of maximal rank $N - 1$. Let T be a support set for which ℓ^1 recovery fails. Then iteratively reweighted ℓ^1 with any admissible reweighting scheme will also fail for some vector \mathbf{y} with support T .*

Proof. Let $\Phi \in \mathbb{R}^{N-1 \times N}$ be a minimally redundant dictionary with maximal rank and let $\mathbf{z} \in \mathcal{N}(\Phi)$ be an arbitrary generator of its null space. Consider any set T for which ℓ^1 recovery can fail, i.e., $\|\mathbf{z}_T\|_1 \geq \|\mathbf{z}_{T^c}\|_1$. Let $\mathbf{y}^* = \mathbf{z}_T$. Because of the dimensionality of the null space, any representation satisfying $\Phi \mathbf{y} = \Phi \mathbf{y}^*$ takes the form $\mathbf{y} = \mathbf{z}_T - \alpha \mathbf{z} = (1 - \alpha)\mathbf{z}_T - \alpha \mathbf{z}_{T^c}$. For any weight

$$\|\mathbf{W}_n \mathbf{y}\|_1 = |1 - \alpha| \cdot \|\mathbf{W}_n \mathbf{z}_T\|_1 + |\alpha| \cdot \|\mathbf{W}_n \mathbf{z}_{T^c}\|_1, \alpha \in \mathbb{R}$$

hence there are only two possible unique solutions to (48), corresponding to $\alpha = 0$ and $\alpha = 1$. Since ℓ^1 fails to recover \mathbf{y}^* , we have $\hat{\mathbf{y}}^{(1)} = -\mathbf{z}_{T^c}$, therefore $T \subset \Gamma_1$ and $\mathbf{W}_2(k, k) = w_{\max}^2, k \in T$.² It follows that

$$|\langle \mathbf{W}_2 \mathbf{z}, \text{sign}(\hat{\mathbf{y}}^{(1)}) \rangle| \leq w_{\max}^2 \|\mathbf{z}_{T^c}\|_1 \leq w_{\max}^2 \|\mathbf{z}_T\|_1 \leq \|(\mathbf{W}_2 \mathbf{z})_{\Gamma_1}\|_1$$

and we obtain that $\hat{\mathbf{y}}^{(n)} = -\mathbf{z}_{T^c}$ for all n . \square

Combining this with the results from section 3 immediately gives Theorem 2.

6 Discussion

In this paper we have quantified values of the RIC, δ_{2m} for which there exist dictionaries where minimization of (4) for some $0 < p \leq 1$ will fail to recover at least one m -sparse vector. This result is in some sense complementary to existing positive results [3, 10] and leaves limited room for improvement. Indeed for the special case of appropriately re-scaled tight frames our negative result becomes sharp when $2m > N - M$.

On the other hand we have also shown that there exist minimally redundant tight frames with RIC, δ_{2m} arbitrarily close to one for which ℓ^p recovery is successful for any p .³ This should not be that surprising, RIP recovery conditions (be they for ℓ^1 or ℓ^p) come from a worst case analysis with respect to several parameters: worst case over all coefficients for a given sign pattern; worst case over all sign patterns for a given support; worst case over all supports of a given size; and worst case over all dictionaries with a given RIC. Our results emphasize the pessimism of such a worst case analysis.

In the context of compressed sensing [7, 3], there is also the desire to characterize the degree of undersampling (M/N) that is possible while still achieving exact recovery. Here RIP can be used to show that certain random matrices with high probability are guaranteed exact recovery with an undersampling of the order $(m/N) \log(N/m)$. However this result is indirect, firstly due to the worst case analysis discussed above and then secondly through the application of the concentration of measure [1]. A more direct approach, characterizing the phase transition between exact recovery and undersampling for classes of random matrices, seems to provide a much clearer indication of the relationship between undersampling and recovery [9]. Of course, deriving expressions for such phase transitions when $p \neq 1$ is likely to be a very challenging problem. Interestingly, the ‘strong’ phase transition of Donoho and Tanner [9] indicates that as $M/N \rightarrow 1$ most minimally redundant tight frames will fail when $m/M \approx 0.18$. In contrast, our result for the ℓ^1 -failing minimally redundant tight frame with the smallest RIC is associated with $m/M \rightarrow 1/(\sqrt{(2)} + 2) \approx 0.29$ and so is clearly not indicative of the boundary behaviour.

²If the solution to (48) is not unique then all values of α between 0 and 1 result in valid solutions and the algorithm has no means for determining the correct one. We therefore make the pessimistic assumption that the algorithm will select the incorrect representation associated with $\alpha = 1$.

³When the dictionary is not tight it is trivial to find such dictionaries by post-multiplying any ℓ^p -successful dictionary with a matrix $A \in \mathbb{R}^{M \times M}$ that introduces the required ill-conditioning (i.e. $\Phi \rightarrow A\Phi$) to make $\delta_{2m} > 1 - \epsilon$. As the null space is unaffected by this action ℓ^p -recovery is still maintained.

Foucart and Lai [10] have also presented guaranteed recovery results for general ℓ^p minimisation with $0 < p \leq 1$. These results are couched in terms of δ_{2m+2} rather than δ_{2m} and are also explicitly dependent upon m . In contrast, the general result in Theorem 3 is independent of m though this could be refined to include m -dependence. Indeed for small m and p the positive result in [10] actually exceeds the negative bound computed from Theorem 3. However, we also note that for fixed p the m -dependent results rapidly converge to the m -independent result of $\delta_{2m+2} < 2(3 - \sqrt{2})/7$, which is slightly weaker than their ℓ^1 recovery result since $\delta_{2m+2} \geq \delta_{2m}$. Theorem 3 seems to suggest that, at least in terms of worst case RIP analysis, there is limited value in reducing p a little below one.

Reducing $p < 1$ also introduces other issues. As the cost function is no longer convex the performance of the ℓ^p optimisation will be a function of the minimisation algorithm used. Our analysis of the iterative reweighted ℓ^1 algorithm (Theorem 2) shows that in terms of worst case RIP analysis there appears to be no gain in using this over unweighted ℓ^1 .

Empirical evidence with iterative reweighting suggests that there can be substantial improvement over unweighted ℓ^1 . However, while this might again be put down to the pessimistic nature of the worst case RIP analysis, we also suspect that the benefits of such algorithms do stem from typically having a range of coefficient scales and that the performance of iterative reweighted ℓ^1 algorithms is probably highly coefficient dependent. Such non-uniform performance cannot be captured by a worst case performance analysis.

Although we have not explicitly considered it here, the RIP also plays a role in quantifying the robustness of ℓ^p recovery to observation noise [3, 10], i.e. when $\mathbf{x} = \Phi\mathbf{y} + \epsilon$. However, as noted here and in [10] exact recovery is independent of dictionary scaling, $\Phi \rightarrow c\Phi$, while robustness to noise is directly related to the scale of the dictionary. It is possible to define the error relative to the isometry constants as in [10], however it could be argued that a fairer measure of robustness would be in terms of absolute error when the dictionary is also constrained to have some physically reasonable property. For example, one might require that the dictionary or ‘sensing matrix’ cannot amplify observations, in other words $\|\Phi\| \leq 1$. Interestingly, in this case, the notion of *asymmetric RIP* that we introduced in section 2 becomes the relevant measure. When viewed in this regard the existing robustness results for random matrices [3, 10] become significantly more pessimistic. This is because such matrices (e.g. random unit spectral norm orthoprojectors) typically shrink sparse vectors by a factor of $\sqrt{M/N}$ and to obtain an appropriate RIC requires re-scaling. However, this in turn implies that typically $\|\Phi\| \approx \sqrt{N/M}$. Hence the robustness of *unit spectral norm* random matrices to observation error scales inversely proportional to the square root of the degree of undersampling.

We finally note that there are a couple of straight forward extensions that we have not pursued in order to keep the paper reasonably concise. First, it would be possible to extend the results in Proposition 1 and Corollary 1 to include the factor A/B associated with non-tight frame bounds. Second, our main results are derived in terms of σ_k^2 and δ_k for $k = 2m$. However, there are a number of positive results based on RICs associated with larger index sets (as in [10]), $k > 2m$. Results similar to Lemma 3 and consequently Theorem 3 in terms of such sets should also be straight forward.

Acknowledgements

This research was partly supported by EPSRC grant D000246/1. The authors would like to thank Thomas Blumensath and Jared Tanner for many interesting discussions on ℓ^1 recovery and discussion of [3]. MED acknowledges support of his position from the Scottish Funding Council and their support of the Joint Research Institute with the Heriot-Watt University as a component part of the Edinburgh Research Partnership.

Appendix

Lemma 4. Let $0 < p < 2$ and $u_1 > v_1 \geq v_2 > u_2 \geq 0$ such that $u_1^p + u_2^p = v_1^p + v_2^p$. Then $u_1^2 + u_2^2 > v_1^2 + v_2^2$.

Proof. Let $J = u_1^2 + u_2^2$ and $u_1^p + u_2^p = c$ for some constant $c > 0$. It is sufficient to show that $\partial J / \partial u_1 > 0$ whenever $u_1 > u_2$.

$$\frac{\partial J}{\partial u_1} = 2u_1 + 2u_2 \frac{\partial u_2}{\partial u_1} = 2u_1 - u_2 \left(\frac{u_2}{u_1} \right)^{1-p} = 2u_1 \left(1 - \left(\frac{u_2}{u_1} \right)^{2-p} \right)$$

which is strictly positive if $u_1 > u_2 \geq 0$ and $p < 2$. □

References

- [1] R. Baraniuk, M. Davenport, R. DeVore, and M. Wakin, "A simple proof of the restricted isometry property for random matrices." *To appear in Constructive Approximation*, 2008.
- [2] E. Candès, J. Romberg, and T. Tao, "Robust uncertainty principles: Exact signal reconstruction from highly incomplete frequency information." *IEEE Trans. Info. Theory*, vol. 52, pp. 489–509, Feb 2006.
- [3] E. Candès, "The Restricted Isometry Property and its implications for Compressed Sensing." *Compte Rendus de l'Academie des Sciences, Paris, Serie I*, 346, 589–592, 2008.
- [4] E. J. Candès, M. B. Wakin and S. P. Boyd, "Enhancing sparsity by reweighted ℓ_1 minimization", to appear in *J. Fourier Anal. Appl.*, 2008.
- [5] S. Chen, D. L. Donoho, and M. A. Saunders, "Atomic decomposition by basis pursuit," *SIAM J. Sci Comp.*, vol. 20, no. 1, pp. 33-61, 1999.
- [6] O. Christensen, "An introduction to frames and Riesz bases", Birkhauser, Boston, MA, 2003.
- [7] D. Donoho, "Compressed Sensing." *IEEE Trans. Info. Theory*, vol. 52, no. 4, pp 1289-1306, April, 2006.
- [8] D. L. Donoho and M. Elad, "Optimally sparse representation from overcomplete dictionaries via ℓ^1 norm minimization," *Proc. Natl. Acad. Sci. USA*, vol. 100, no. 5, pp. 2197-2002, Mar. 2002.
- [9] D. Donoho and J. Tanner, "Counting faces of randomly-projected polytopes when the projection radically lowers dimension", to appear in *Journal of the AMS*, 2008.
- [10] S. Foucart and M.-J. Lai, "Sparsest solutions of underdetermined linear systems via ℓ_q -minimization for $0 < q \leq 1$." *Submitted to Applied and Computational Harmonic Analysis*, 2008.
- [11] J.-J. Fuchs, "On Sparse Representations in Arbitrary Redundant Bases", *IEEE Trans. Inform. Theory*, Vol. 50, No. 6, pp. 1341–1344, June 2004.

- [12] I.F. Gorodnitsky and B.D. Rao, “Sparse signal reconstruction from limited data using FOCUSS: A re-weighted norm minimization algorithm.” *IEEE Trans. on Signal Processing*, vol. 45, pp. 600–616, March 1997.
- [13] R. Gribonval and M. Nielsen, “Sparse decompositions in unions of bases.” *IEEE Trans. Info. Theory*, vol. 49, no. 12, pp 3320-3325, Dec 2003.
- [14] R. Gribonval and M. Nielsen, “On the strong uniqueness of highly sparse expansions from redundant dictionaries.” *In Proc. Int Conf. Independent Component Analysis (ICA’04)*, Sep 2004.
- [15] R. Gribonval and M. Nielsen, “Highly sparse representations from dictionaries are unique and independent of the sparseness measure.” *Applied and Computational Harmonic Analysis*, vol. 22, no. 3, pp 335–355, May 2007. [Technical report October 2003].
- [16] J. Tropp, “Just relax: Convex programming methods for identifying sparse signals.” *IEEE Trans. Info. Theory*, vol. 51, no. 3, pp. 1030-1051, Mar 2006.