

Stochastic Geometry and Random Matrix Theory in Compressed Sensing

Jared Tanner
University of Edinburgh

Institut Henri Poincaré
June 20th to the 22nd, 2011

Sensing: the information acquisition step

- ▶ Digital data: $x \in \mathbb{R}^n$. Measure vectors by point sensing

$$\langle x, e_i \rangle = x_i \quad \text{for } i = 1, 2, \dots, n$$

Same for images and other digital data.

Cannot improve if x_i independent of each other.

- ▶ Analog data: $f(x) \in B_\sigma$. Bandlimited model popular in EE

$$f(x) \in B_\sigma \quad \Leftrightarrow \quad f(x) := \frac{1}{\sqrt{2\pi}} \int_{-\sigma}^{\sigma} \hat{f}(w) e^{2\pi i w x} dw$$

where $\hat{f}(w) = 0$ for $|w| > \sigma$. Shannon Sampling Theorem:

$$f(x) = \sum_{k \in \mathbb{Z}} f(kT) \psi(x - kT) \quad \text{if } T \leq 1/2\sigma$$

Can exactly recovery $f(x)$ from its point samples $\{f(kT)\}_{k \in \mathbb{Z}}$.

“Nyquist sampling rate” is $T = 1/2\sigma$, largest for $f \in B_\sigma$.

After we acquire the data, then we compress

- ▶ On a computer all data is stored digital
For vectors $x \in \mathbb{R}^n$ already in digital form.
For $f(x) \in B_\sigma$ store $\{f(kT)\}_{\mathbb{Z}}$ and use Shannon Theorem
- ▶ Most vectors have entries that are very dependent and
 $f(x) \in B_\sigma$ are smooth (analytic) with nearby entries related.
- ▶ Do we really need to store all of these entries if highly dependent?

After we acquire the data, then we compress

- ▶ On a computer all data is stored digital
For vectors $x \in \mathbb{R}^n$ already in digital form.
For $f(x) \in B_\sigma$ store $\{f(kT)\}_{\mathbb{Z}}$ and use Shannon Theorem
- ▶ Most vectors have entries that are very dependent and
 $f(x) \in B_\sigma$ are smooth (analytic) with nearby entries related.
- ▶ Do we really need to store all of these entries if highly dependent? Certainly not!
Welcome to the wonderful world of approximation theory.
- ▶ Smoothness (even piecewise smooth) implies compressibility
The options are endless:
Fourier series, orthogonal polynomials, wavelets, curvelets, shearlets, and any other “let” you can imagine.
(More on these to come...)

Compression

- ▶ Discrete data:

Let the columns of Φ span \mathbb{R}^n so that $x = \Phi z$ for some z .

Let $H_k(z)$ be the k -term *hard threshold*, setting all but the largest k entries (in magnitude) to zero.

x compressible in representation Φ if $\|x - \Phi H_k(z)\| \ll \|x\|$.

- ▶ Polynomial decay for large problems:

$$\|x - \Phi H_k(z)\| \leq \text{Const.} k^{-p}$$

- ▶ Analog functions: truncated series in representation $\{\psi_\ell(x)\}_\ell$

$$S_N f(x) := \sum_{\ell=0}^N \hat{f}_\ell \psi_\ell(x)$$

Polynomial or exponential decay $\|f(x) - S_N f(x)\| \leq C \cdot \tau^{-N}$

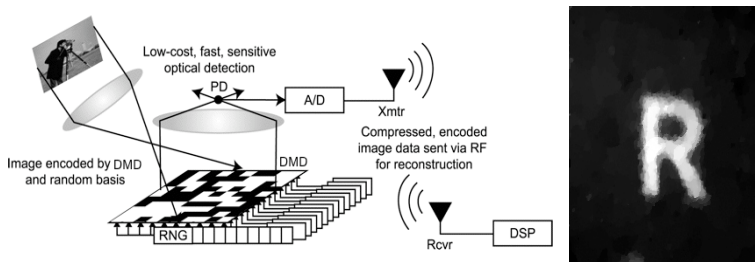
- ▶ Compression is ubiquitous, essentially always performed.
- ▶ Full sensing and then compression is very wasteful

The advantages of Compressed Sensing

- ▶ If k coefficients are sufficient to accurately approximate the data, why measure it all in the first place?
- ▶ Move compression into acquisition: Compressed Sensing (CS)
- ▶ There is a cost associated with CS, use when sensing is costly
- ▶ A few applications:
 - MRI Scanner – length of time in device, through-put
 - UAV imaging – time of flight over target
 - Nuclear Medicine (CT/SPECT/PET) – radiation dosage
 - Genomic sequencing – through-put
 - Satellite – limited communications and battery
 - (lets see a few pictures...)

Single Photo-diode digital camera

- ▶ Proof of concept for compressed sensing: Baraniuk and Kelly



- ▶ 2% measurements compared to number of pixels in recon.
- ▶ Savings, measurement time, simple device, power of device, ...
- ▶ Multi-spectral variants have been constructed.

Magnetic Force Resonance Microscopy (A. Hero, M. Ting)

- Non-linear sparsity exploiting reconstruction algorithms:

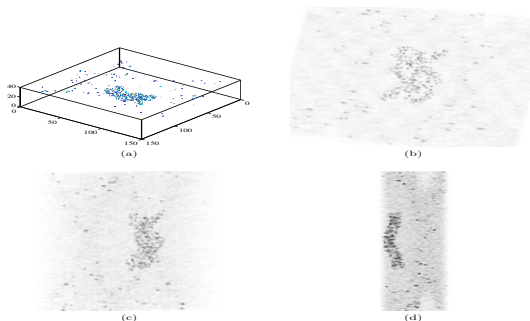
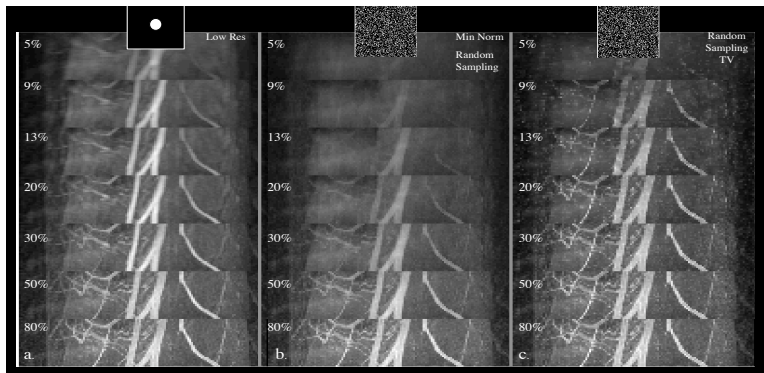


Figure 5.28: Three dimensional visualization of the MAP2 reconstruction of 103D's hydrogen atoms with $g^* = (\sqrt{2})^{-1}$ at an SNR of 6.02 dB. Different viewing angles are shown. The helical structure of 103D is apparent.

103D DNA Molecule - 272 Hydrogen Atoms

MRI - Angiography

- Stanford MRI Lab: T. Cucker, M. Lustig, and D. G. Nishimura



Astronomy applications: Herschel

► CIRM (France): J.L. Starck

HERSCHEL



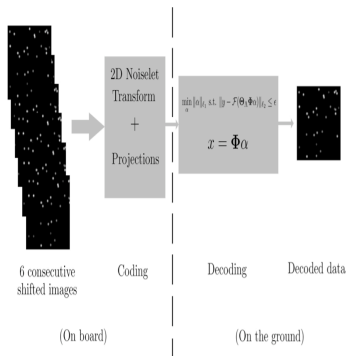
This space telescope has been designed to observe in the far-infrared and sub-millimeter wavelength range. Its launch is scheduled for the beginning of 2009. The shortest wavelength band, 57-210 microns, is covered by PACS (Photodetector Array Camera and Spectrometer).

Herschel data transfer problem:
-no time to do sophisticated data compression on board.
-a compression ratio of 6 must be achieved.

=> solution: averaging of six successive images on board

CS may offer another alternative.

The proposed Herschel compression scheme

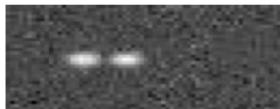


Astronomy applications: Hershel

► CIRM (France): J.L. Starck

Resolution: CS versus Mean

Simulated image



Mean of six images

Simulated noisy image with flat and dark



Compressed sensing reconstructed images

Resolution limit versus SNR

SNR	-17.3	-9.35	-3.3	0.21	2.7	4.7	6.2	7.6	8.7
Intensity	900	2250	4500	6750	9000	11250	13500	15750	18000
MO6	3	3	3	3	3	3	3	3	3
CS	2.33	2.33	2	2	2	2	2	2	2

THE CS-BASED COMPRESSION ENTAILS A RESOLUTION GAIN EQUAL TO A 30% OF THE SPATIAL RESOLUTION PROVIDED BY MO6.

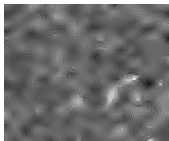
Astronomy applications: Herschel

- CIRM (France): J.L. Starck

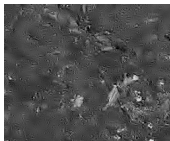
JPEG2000 Versus Compressed Sensing

Compression Rate: 25

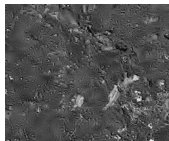
One observation



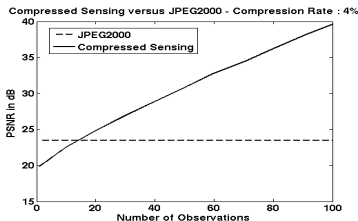
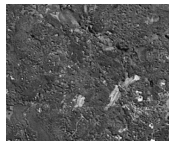
10 observations



20 observations



100 observations



Matrix completion : Inpainting

- Duke: L. Carin lab



80% of RGB Voxels Missing at Random



Recovered Image via Beta Process and

Matrix completion : Inpainting

► Duke: L. Carin lab



80% of RGB Voxels Missing at Random



Recovered Image via Beta Process and

		1954 Pitchers			2008 Pitchers		
		D. Mossi	B. Turley	R. Narleski	C. Marmol	G. Balfour	B. Morrow
1954 Batters	D. Mueller	0.2829 ± 0.0166	0.2340 ± 0.0171	0.2543 ± 0.0176	0.1472 ± 0.0192	0.1755 ± 0.0217	0.1627 ± 0.0211
	S. Burgess	0.2619 ± 0.0147	0.2277 ± 0.0131	0.2391 ± 0.0152	0.1338 ± 0.0163	0.1550 ± 0.0195	0.1502 ± 0.0154
	B. Skowron	0.2657 ± 0.0135	0.2141 ± 0.0133	0.2339 ± 0.0144	0.1254 ± 0.0153	0.1530 ± 0.0200	0.1407 ± 0.0193
	C. Zambrano	0.2652 ± 0.0217	0.2215 ± 0.0225	0.2384 ± 0.0219	0.1313 ± 0.0235	0.1566 ± 0.0276	0.1475 ± 0.0267
2008 Batters	P. Sandoval	0.3570 ± 0.0355	0.3316 ± 0.0400	0.3402 ± 0.0386	0.2612 ± 0.0529	0.2792 ± 0.0500	0.2740 ± 0.0522
	R. Furcal	0.2971 ± 0.0140	0.2599 ± 0.0141	0.2748 ± 0.0144	0.1721 ± 0.0158	0.1948 ± 0.0187	0.1876 ± 0.0189

Matrix completion : Inpainting

► Duke: L. Carin lab

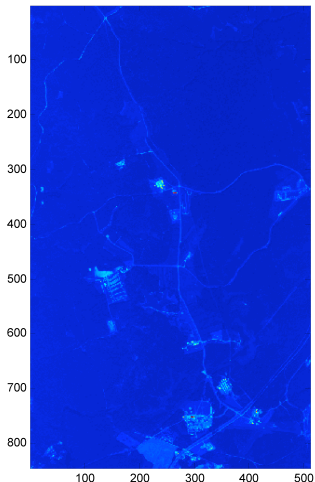
- Name: HyMapAPHill (NGA)
- Image size: 845 by 512
- Total Channels: 106

Original Scene:

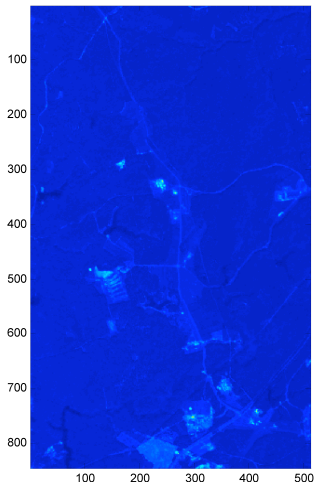


Duke: L. Carin lab

- 2% of Hyperspectral datacube at random, band 1



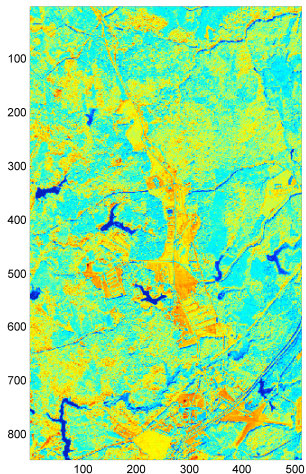
Original image



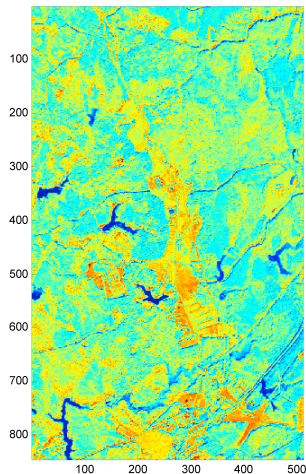
Restored image

Duke: L. Carin lab

- 2% of Hyperspectral datacube at random, band 50



Original image



Restored image

Matrix Completion, segmentation, video

► Stanford: E. Candes

Candes's model, noise free, scene2

original frame



Low rank component



Sparse component



original frame



Low rank component



Sparse component



original frame



Low rank component



Sparse component



E. J. Candès, X. Li, Y. Ma, and J. Wright, "Robust Principal Component Analysis?"

Back to compression: Fourier Series

Definition (C^s -periodic)

The $C^s[-\pi, \pi]$ seminorm is defined as

$$\|f\|_{C^s[-\pi, \pi]} := \int_{-\pi}^{\pi} |f^{(s)}(x)| dx$$

where $f^{(s)}(x)$ denotes the s^{th} derivative of $f(x)$. A function is said to be in $C^s[-\pi, \pi]$ if $\|f\|_{C^s[-\pi, \pi]} < \infty$. We refer to a function as being C^s -periodic over $[-\pi, \pi]$ if it is in $C^{(s)}[-\pi, \pi]$ and $f^{(j)}(\pi) = f^{(j)}(-\pi)$ for $j = 0, \dots, s-1$.

Definition (Fourier series)

Let $f(x)$ be in $L^2[-\pi, \pi]$ and be C^s -periodic over $[-\pi, \pi]$. Then, it can be represented in the Fourier orthonormal basis as

$$f(x) = (2\pi)^{-1/2} \sum_{k \in \mathbb{Z}} \hat{f}_k e^{ikx} \quad \text{with} \quad \hat{f}_k := (2\pi)^{-1/2} \int_{-\pi}^{\pi} f(x) e^{-ikx} dx.$$

Fourier Series compression rate

Theorem (Truncated Fourier Approximation)

Let $f(x)$ be in $L^2[-\pi, \pi]$ and be C^s -periodic for $s \geq 2$. Then

$$\|f - S_N f\|_{L^\infty[-\pi, \pi]} \leq \left(\frac{2}{\pi}\right)^{1/2} (s-1)^{-1} \|f\|_{C^s[-\pi, \pi]} \cdot N^{-s+1}.$$

Proof.

Integrate by parts and triangle inequality

$$\hat{f}_k = (2\pi)^{-1/2} (-ik)^{-s} \int_{-\pi}^{\pi} f^{(s)}(x) e^{-ikx} dx,$$

$$\max_{x \in [-\pi, \pi]} |f(x) - S_N f(x)| = \max_{x \in [-\pi, \pi]} \left| \sum_{|k| > N} \hat{f}_k e^{ikx} \right| \leq \sum_{|k| > N} |\hat{f}_k|,$$

$$\sum_{k=N+1}^{\infty} k^{-s} \leq \int_N^{\infty} k^{-s} dk \text{ for } s \geq 2$$

□

Generalized Fourier Series

- ▶ Fourier series superb for smooth periodic functions.
- ▶ Smooth non-periodic functions via orthogonal polynomials
- ▶ *Global* bases have difficulty for non-smooth functions.
Gibbs' Phenomenon can be overcome through edge detection and postprocessing, but does not work well for noisy data.
- ▶ Localized expansions allow better qualitative understanding
- ▶ Haar system is composed of the scaling function

$$\phi(x) = \begin{cases} 1 & x \in [0, 1] \\ 0 & x \notin [0, 1] \end{cases}$$

and translation and dilations of the mother wavelet

$$\psi(x) = \begin{cases} 1 & x \in [0, 1/2] \\ -1 & x \in (1/2, 1) \\ 0 & x \notin [0, 1] \end{cases}$$

Let $\psi_{n,k}(x) := 2^{n/2}\psi(2^n x - k)$.

Haar Wavelets

Definition (Haar Wavelet)

The Haar system

$$\phi(x) \cup \psi_{n,k}(x)_{n \in \mathbb{N}, 0 \leq k < 2^n}$$

is an orthonormal basis for $L^2[0, 1]$. Define the Haar coefficients as

$$f_0 := \int_0^1 f(x) \phi(x) dx \quad \text{and} \quad f_{n,k} := \int_0^1 f(x) \psi_{n,k}(x) dx$$

and the truncated Haar approximation of $f(x)$ as

$$W_M f(x) := f_0 + \sum_{n=0}^M \sum_{k=0}^{2^n-1} f_{n,k} \psi_{n,k}(x).$$

The truncated Haar expansion converges to the original function in $L^2[0, 1]$,

$$\lim_{M \rightarrow \infty} \|f - W_M f\|_{L^2[0,1]} \rightarrow 0.$$

Wavelet convergence rates: Vanishing moments

- ▶ Convergence rate of truncated Wavelet approximations dictated by the decay rate of coefficients $f_{n,k}$ for n large.
- ▶ Consider Haar as example. $\text{supp}(\psi_{n,k}(x)) = 2^{-n}[k, k+1)$
Taylor series $f(x)$ about $2^{-n}(k+1/2)$

$$\begin{aligned} f_{n,k} &= \int_{2^{-n}k}^{2^{-n}(k+1)} \psi_{n,k}(x) [f(x_0) + (x - x_0)f'(x_0) + \cdots] dx \\ &= 2^{-3n/2} 4f'(x_0) + \mathcal{O}(2^{-5n/2}). \end{aligned}$$

- ▶ If $f(x)$ piecewise smooth with $\mathcal{O}(\ell)$ discontinuities then ℓ of $f_{n,k} \sim 2^{-n/2}$ and $2^n - \mathcal{O}(\ell)$ are of size $2^{-3n/2}$.
- ▶ Overall decay rate $2^{-n/2}$ dictated by discontinuities.
- ▶ Appears exponential, but needs $N = 2^n$ coefficients.
- ▶ Decay in N is a slow $N^{-1/2}$ rate if a linear approximation, but at exponential rate $2^{-n/2}$ if only large entries kept.

Other wavelets and higher dimensional “lets”

- ▶ Wavelets beyond Haar, are they better?

Other wavelets and higher dimensional “lets”

- ▶ Wavelets beyond Haar, are they better? Yes and No.
Convergence rate of all 1D wavelets can be viewed similarly.
Wavelets that cross discontinuities have “large” $\mathcal{O}(2^{-n/2})$ coefficients, and other coefficients size dictated by number of vanishing moments, $\mathcal{O}(2^{-(2p+1)n/2})$ for order p wavelet.
- ▶ Higher order wavelets have faster convergence, with “wider” wavelets and more crossing the discontinuities
- ▶ Time-frequency tiling:
Wavelets use translation and dilation
Gabor atoms use translation and modulation
Multi-dimensional variants use other operators such as rotation and shear
- ▶ A few examples to see how they work, discrete case

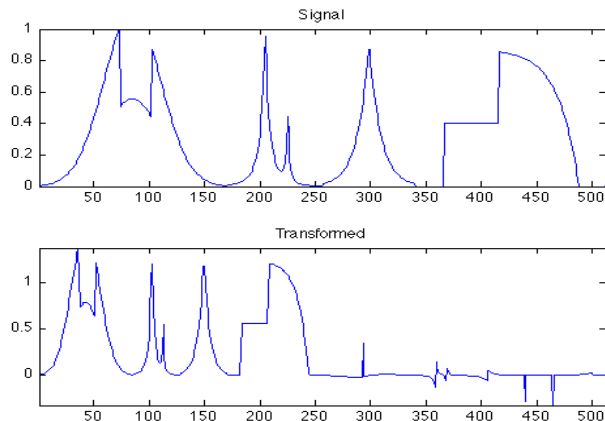
The filterbank viewpoint for discrete data

- ▶ Convolve discrete vector $f \in \mathbb{R}^{2n}$ with two vectors
 h a “low pass filter” that approximates f and
 g a “high pass filter” that captures $f - h$
- ▶ Downsample $a = (f \star h) \downarrow 2$ and $d = (f \star g) \downarrow 2$,
where $(u \downarrow 2)[k] = u[2k]$ to keep $2n$ entries (same as f)
- ▶ If h and g are designed properly then f can be recovered.
Upsample a and d by adding a zero after each entry, $(u \uparrow 2)$
Convolve upsampled vectors with the reverse order of h and g

$$f = (a \uparrow 2) \star \tilde{h} + (d \uparrow 2) \star \tilde{g}$$

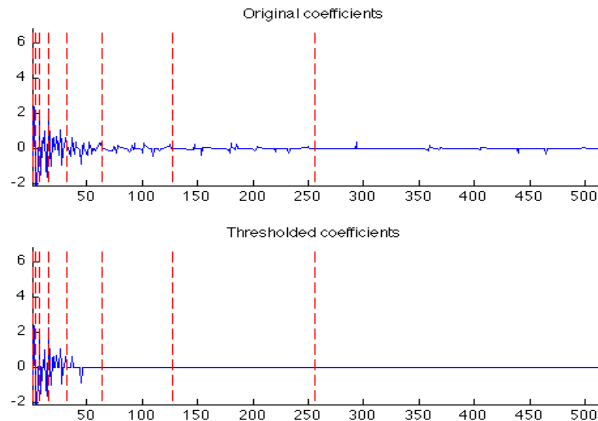
- ▶ Example: $h = [1 \ 1]$ and $g = [-1 \ 1]$
 $a[1] = f[1] + f[2]$ and $d[1] = -f[1] + f[2]$
 $(a \uparrow 2) \star \tilde{h}$ has two entries both equal to $a[1]$
 $(d \uparrow 2) \star \tilde{g}$ has first entry $-d[1]$ and second entry $d[1]$

Haar example, one step [Peyre]

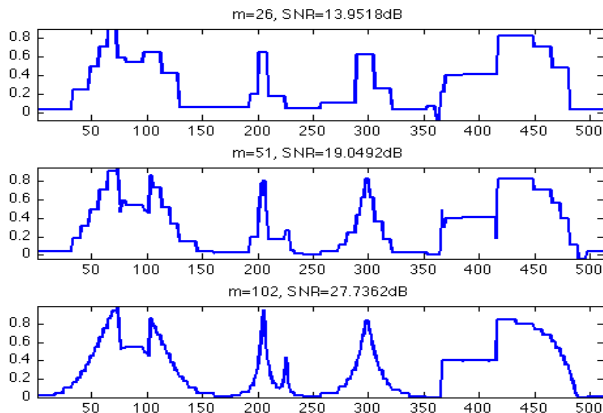


- ▶ First half is a and second half is d
- ▶ Repeat process on a portion

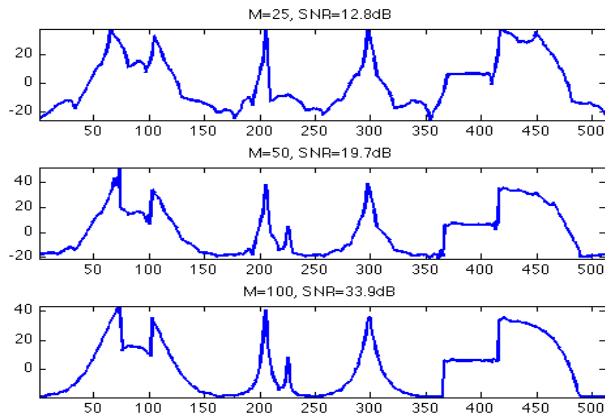
Haar example, full transform [Peyre]



Haar m term approximation [Peyre]



Daubechies4 m term approximation [Peyre]



Searching for simplicity (sparsity)

- Sparse solutions to underdetermined systems of equations

$$\min_x \|x\|_0 \quad \text{subject to} \quad \|y - Ax\| \leq \tau$$

- Basis pursuit: find best set of columns of A for y
- design algorithms that find sparse solutions and hope...

Searching for simplicity (sparsity)

- Sparse solutions to underdetermined systems of equations

$$\min_x \|x\|_0 \quad \text{subject to} \quad \|y - Ax\| \leq \tau$$

- Basis pursuit: find best set of columns of A for y
- design algorithms that find sparse solutions and hope...
- Simple solutions to under determined systems of equations

$$x \quad \text{such that} \quad y = Ax \quad \text{and} \quad \alpha_i \leq x_i \leq \beta_i$$

- if enough of x_i are equal to α_i or β_i is it unique?

Searching for simplicity (sparsity)

- ▶ Sparse solutions to underdetermined systems of equations

$$\min_x \|x\|_0 \quad \text{subject to} \quad \|y - Ax\| \leq \tau$$

- Basis pursuit: find best set of columns of A for y
- design algorithms that find sparse solutions and hope...
- ▶ Simple solutions to under determined systems of equations

$$x \quad \text{such that} \quad y = Ax \quad \text{and} \quad \alpha_i \leq x_i \leq \beta_i$$

- if enough of x_i are equal to α_i or β_i is it unique?
- ▶ Low rank matrix approximation (matrix completion)

$$\min_M \text{rank}(M) \quad \text{subject to} \quad \|y - \mathcal{A}(M)\| \leq \tau$$

- unknown representation in which M is simple

Coherence

Let A be the sensing matrix and a_i its i^{th} column

$$\mu_2(A) := \max_{i \neq j} |a_i^* a_j|$$

Coherence

Let A be the sensing matrix and a_i its i^{th} column

$$\mu_2(A) := \max_{i \neq j} |a_i^* a_j|$$

► Pros:

- Easy to calculate!
- Easy to use to prove pretty good results
- A general tool for any algorithm (wide usage)

Coherence

Let A be the sensing matrix and a_i its i^{th} column

$$\mu_2(A) := \max_{i \neq j} |a_i^* a_j|$$

- ▶ Pros:
 - Easy to calculate!
 - Easy to use to prove pretty good results
 - A general tool for any algorithm (wide usage)
- ▶ Cons:
 - A general tool for any algorithm (bad results)
 - Worst case results are limited to “sqrt” proportionality

Coherence

Let A be the sensing matrix and a_i its i^{th} column

$$\mu_2(A) := \max_{i \neq j} |a_i^* a_j|$$

- ▶ Pros:
 - Easy to calculate!
 - Easy to use to prove pretty good results
 - A general tool for any algorithm (wide usage)
- ▶ Cons:
 - A general tool for any algorithm (bad results)
 - Worst case results are limited to “sqrt” proportionality

Use coherence analyze: Thresholding, Matching Pursuit,
Orthogonal Matching Pursuit, and ℓ^1 -regularization

* for the moment assume solution is unique

One step thresholding

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set Λ the index set of the $k \leq m$ largest in $|A_{m,n}^* y|$
Output the n -vector x whose entries are

$$x_\Lambda = (A_\Lambda^* A_\Lambda)^{-1} A_\Lambda^* y \quad \text{and} \quad x(i) = 0 \text{ for } i \notin \Lambda.$$

One step thresholding

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set Λ the index set of the $k \leq m$ largest in $|A_{m,n}^* y|$
Output the n -vector x whose entries are

$$x_\Lambda = (A_\Lambda^* A_\Lambda)^{-1} A_\Lambda^* y \quad \text{and} \quad x(i) = 0 \text{ for } i \notin \Lambda.$$

Theorem

Let $y = A_{m,n} x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, and

$$\|x_0\|_0 < \frac{1}{2} (\nu_\infty(x_0) \cdot \mu_2(A_{m,n})^{-1} + 1),$$

then the Thresholding decoder with $k = \|x_0\|_0$ will return x_0 , with $\nu_p(x) := \min_{j \in \text{supp}(x)} |x(j)| / \|x\|_p$.

One step thresholding (proof)

Proof.

With $y = A_{m,n}x_0$, denote $w = A_{m,n}^*y = A_{m,n}^*A_{m,n}x_0$.

The i^{th} entry in w is equal to $w_i = \sum_{j \in \text{supp}(x_0)} x_0(j) a_i^* a_j$.

One step thresholding (proof)

Proof.

With $y = A_{m,n}x_0$, denote $w = A_{m,n}^*y = A_{m,n}^*A_{m,n}x_0$.

The i^{th} entry in w is equal to $w_i = \sum_{j \in \text{supp}(x_0)} x_0(j) a_i^* a_j$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{j \in \text{supp}(x_0)} |x_0(j)| \cdot |a_i^* a_j| \leq k \mu_2(A_{m,n}) \|x_0\|_\infty.$$

One step thresholding (proof)

Proof.

With $y = A_{m,n}x_0$, denote $w = A_{m,n}^*y = A_{m,n}^*A_{m,n}x_0$.

The i^{th} entry in w is equal to $w_i = \sum_{j \in \text{supp}(x_0)} x_0(j) a_i^* a_j$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{j \in \text{supp}(x_0)} |x_0(j)| \cdot |a_i^* a_j| \leq k \mu_2(A_{m,n}) \|x_0\|_\infty.$$

For $i \in \text{supp}(x_0)$ the magnitude of w_i is bounded below as:

$$\begin{aligned} |w_i| &\geq |x_0(i)| - \left| \sum_{j \in \text{supp}(x_0), j \neq i} x_0(j) a_i^* a_j \right| \\ &\geq |x_0(i)| - (k-1) \mu_2(A_{m,n}) \|x_0\|_\infty. \end{aligned}$$

One step thresholding (proof)

Proof.

With $y = A_{m,n}x_0$, denote $w = A_{m,n}^*y = A_{m,n}^*A_{m,n}x_0$.

The i^{th} entry in w is equal to $w_i = \sum_{j \in \text{supp}(x_0)} x_0(j) a_i^* a_j$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{j \in \text{supp}(x_0)} |x_0(j)| \cdot |a_i^* a_j| \leq k \mu_2(A_{m,n}) \|x_0\|_\infty.$$

For $i \in \text{supp}(x_0)$ the magnitude of w_i is bounded below as:

$$\begin{aligned} |w_i| &\geq |x_0(i)| - \left| \sum_{j \in \text{supp}(x_0), j \neq i} x_0(j) a_i^* a_j \right| \\ &\geq |x_0(i)| - (k-1) \mu_2(A_{m,n}) \|x_0\|_\infty. \end{aligned}$$

Recovery if $\max_{i \notin \text{supp}(x_0)} |w_i| < \min_{i \in \text{supp}(x_0)} |w_i|$. □

Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$, and let $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$ and define
 $x^{j+1} = x^j + (a_i^* r^j) e_i$ where e_i is the i^{th} coordinate vector.

Output x^j when a termination criteria is obtained.

Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$, and let $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$ and define $x^{j+1} = x^j + (a_i^* r^j) e_i$ where e_i is the i^{th} coordinate vector.

Output x^j when a termination criteria is obtained.

Theorem

Let $y = A_{m,n} x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, and

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then Matching Pursuit will have $\operatorname{supp}(x^j) \subseteq \operatorname{supp}(x_0)$ for all j .

Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$, and let $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$ and define $x^{j+1} = x^j + (a_i^* r^j) e_i$ where e_i is the i^{th} coordinate vector.

Output x^j when a termination criteria is obtained.

Theorem

Let $y = A_{m,n} x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, and

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then Matching Pursuit will have $\operatorname{supp}(x^j) \subseteq \operatorname{supp}(x_0)$ for all j .

* Preferable over one step thresholding: no dependence on $\nu_p(x_0)$.

Matching Pursuit (proof)

Proof.

Assume $\text{supp}(x^j) \subset \text{supp}(x_0)$ for some j , which is true for $j = 0$.

Let $r^j = y - A_{m,n}x^j$, and $w_i = \sum_{\ell \in \text{supp}(x_0)} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell$.

Matching Pursuit (proof)

Proof.

Assume $\text{supp}(x^j) \subset \text{supp}(x_0)$ for some j , which is true for $j = 0$.

Let $r^j = y - A_{m,n}x^j$, and $w_i = \sum_{\ell \in \text{supp}(x_0)} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{\ell \in \text{supp}(x_0)} |(x_0 - x^j)(\ell)| \cdot |a_i^* a_\ell| \leq k\mu_2(A_{m,n}) \|x_0 - x^j\|_\infty.$$

Matching Pursuit (proof)

Proof.

Assume $\text{supp}(x^j) \subset \text{supp}(x_0)$ for some j , which is true for $j = 0$.

Let $r^j = y - A_{m,n}x^j$, and $w_i = \sum_{\ell \in \text{supp}(x_0)} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{\ell \in \text{supp}(x_0)} |(x_0 - x^j)(\ell)| \cdot |a_i^* a_\ell| \leq k\mu_2(A_{m,n}) \|x_0 - x^j\|_\infty.$$

For $i \in \text{supp}(x_0)$ the magnitude of w_i is bounded below as:

$$\begin{aligned} |w_i| &\geq |(x_0 - x^j)(i)| - \left| \sum_{\ell \in \text{supp}(x_0), \ell \neq i} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell \right| \\ &\geq |(x_0 - x^j)(i)| - (k-1)\mu_2(A_{m,n}) \|x_0 - x^j\|_\infty. \end{aligned}$$

Matching Pursuit (proof)

Proof.

Assume $\text{supp}(x^j) \subset \text{supp}(x_0)$ for some j , which is true for $j = 0$.

Let $r^j = y - A_{m,n}x^j$, and $w_i = \sum_{\ell \in \text{supp}(x_0)} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell$.

For $i \notin \text{supp}(x_0)$ the magnitude of w_i is bounded above as:

$$|w_i| \leq \sum_{\ell \in \text{supp}(x_0)} |(x_0 - x^j)(\ell)| \cdot |a_i^* a_\ell| \leq k\mu_2(A_{m,n}) \|x_0 - x^j\|_\infty.$$

For $i \in \text{supp}(x_0)$ the magnitude of w_i is bounded below as:

$$\begin{aligned} |w_i| &\geq |(x_0 - x^j)(i)| - \left| \sum_{\ell \in \text{supp}(x_0), \ell \neq i} (x_0 - x^j)(\ell) \cdot a_i^* a_\ell \right| \\ &\geq |(x_0 - x^j)(i)| - (k-1)\mu_2(A_{m,n}) \|x_0 - x^j\|_\infty. \end{aligned}$$

Recovery if $\max_{i \in \text{supp}(x_0)} |w_i| > \max_{i \notin \text{supp}(x_0)} |w_i|$.

□

Orthogonal Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$ and Λ^0 to be the empty set, and set $j = 0$.

Let $r^j := y - Ax^j$, $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$, and $\Lambda^{j+1} = i \cup \Lambda^j$.

Set $x_{\Lambda^{j+1}}^{j+1} = (A_{\Lambda^{j+1}}^* A_{\Lambda^{j+1}})^{-1} A_{\Lambda^{j+1}}^* y$

and $x^{j+1}(\ell) = 0$ for $\ell \notin \Lambda^{j+1}$, and set $j = j + 1$.

Output x^j when a termination criteria is obtained.

Orthogonal Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$ and Λ^0 to be the empty set, and set $j = 0$.

Let $r^j := y - Ax^j$, $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$, and $\Lambda^{j+1} = i \cup \Lambda^j$.

Set $x_{\Lambda^{j+1}}^{j+1} = (A_{\Lambda^{j+1}}^* A_{\Lambda^{j+1}})^{-1} A_{\Lambda^{j+1}}^* y$

and $x^{j+1}(\ell) = 0$ for $\ell \notin \Lambda^{j+1}$, and set $j = j + 1$.

Output x^j when a termination criteria is obtained.

Theorem

Let $y = A_{m,n} x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, and

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then after $\|x_0\|_{\ell^0}$ steps, Orthogonal Matching Pursuit recovers x_0 .

Orthogonal Matching Pursuit [Tr05]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Let $r^j := y - Ax^j$.

Set $x^0 = 0$ and Λ^0 to be the empty set, and set $j = 0$.

Let $r^j := y - Ax^j$, $i := \operatorname{argmax}_{\ell} |a_{\ell}^* r^j|$, and $\Lambda^{j+1} = i \cup \Lambda^j$.

Set $x_{\Lambda^{j+1}}^{j+1} = (A_{\Lambda^{j+1}}^* A_{\Lambda^{j+1}})^{-1} A_{\Lambda^{j+1}}^* y$

and $x^{j+1}(\ell) = 0$ for $\ell \notin \Lambda^{j+1}$, and set $j = j + 1$.

Output x^j when a termination criteria is obtained.

Theorem

Let $y = A_{m,n}x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, and

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then after $\|x_0\|_{\ell^0}$ steps, Orthogonal Matching Pursuit recovers x_0 .

* Proof, same as Matching Pursuit. Finite number of steps.

ℓ^1 -regularization [Tr05]

Input: y and $A_{m,n}$.

“Algorithm”: Return $\operatorname{argmin} \|x\|_1$ subject to $y = Ax$.

Theorem

Let $y = A_{m,n}x_0$, with

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then the solution of ℓ^1 -regularization is x_0 .

ℓ^1 -regularization [Tr05]

Input: y and $A_{m,n}$.

“Algorithm”: Return $\operatorname{argmin} \|x\|_1$ subject to $y = Ax$.

Theorem

Let $y = A_{m,n}x_0$, with

$$\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1),$$

then the solution of ℓ^1 -regularization is x_0 .

* Preferable over OMP: faster if use good ℓ^1 solver.

ℓ^1 -regularization (proof, page 1)

Proof.

Let $\Lambda_0 := \text{supp}(x_0)$ and $\Lambda_1 := \text{supp}(x_1)$ with $y = A_{m,n}x_0 = A_{m,n}x_1$, and $\exists i$ with $i \in \Lambda_1$ with $i \notin \Lambda_0$. Note that because $y = A_{\Lambda_0}x_0 = A_{\Lambda_1}x_1$,

$$\begin{aligned}\|x_0\|_1 &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* A_{\Lambda_0} x_0\|_1 \\ &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* y\|_1 \\ &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* A_{\Lambda_1} x_1\|_1.\end{aligned}$$

Establish bounds on $(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i$.

ℓ^1 -regularization (proof, page 1)

Proof.

Let $\Lambda_0 := \text{supp}(x_0)$ and $\Lambda_1 := \text{supp}(x_1)$ with $y = A_{m,n}x_0 = A_{m,n}x_1$, and $\exists i$ with $i \in \Lambda_1$ with $i \notin \Lambda_0$. Note that because $y = A_{\Lambda_0}x_0 = A_{\Lambda_1}x_1$,

$$\begin{aligned}\|x_0\|_1 &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* A_{\Lambda_0} x_0\|_1 \\ &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* y\|_1 \\ &= \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* A_{\Lambda_1} x_1\|_1.\end{aligned}$$

Establish bounds on $(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i$.

To establish proof need bounds for $i \in \Lambda$ and $i \notin \Lambda$.

$$\begin{aligned}\text{For } i \in \Lambda_0: \| (A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i \|_1 \\ = \| (A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* A_{\Lambda_0} e_i \|_1 = \| e_i \|_1 = 1\end{aligned}$$

□

ℓ^1 -regularization (proof, page 2)

Proof.

For any $i \notin \Lambda_0$ we establish the bound in two parts; first,

$$\|A_{\Lambda_0}^* a_i\|_1 \leq \sum_{\ell \in \Lambda_0} |a_\ell^* a_i| \leq k\mu_2(A_{m,n}).$$

Noting $A_{\Lambda_0}^* A_{\Lambda_0} = I_{k,k} + B$ where $B_{i,i} = 0$ and $|B_{i,j}| \leq \mu_2(A_{m,n})$, then

$$\|(I_{k,k} + B)^{-1}\|_1 = \left\| \sum_{\ell=0}^{\infty} (-B)^\ell \right\|_1 \leq \sum_{\ell=0}^{\infty} \|B\|_1^\ell = \frac{1}{1 - \|B\|_1} \leq \frac{1}{1 - (k-1)\mu_2(A_{m,n})}.$$

Therefore, for $i \notin \Lambda_0$:

$$\|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 \leq \frac{k\mu_2(A_{m,n})}{(1 - (k-1)\mu_2(A_{m,n}))} < 1$$

□

ℓ^1 -regularization (proof, page 3)

Proof.

Proof concludes through triangle inequality and use that:

- For $i \in \Lambda_0$: $\|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 = 1$
- For $i \notin \Lambda_0$: $\|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 < 1$
- And $\exists i$ with $i \in \Lambda_1$ and $i \notin \Lambda_0$.

ℓ^1 -regularization (proof, page 3)

Proof.

Proof concludes through triangle inequality and use that:

- For $i \in \Lambda_0$: $\|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 = 1$
- For $i \notin \Lambda_0$: $\|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 < 1$
- And $\exists i$ with $i \in \Lambda_1$ and $i \notin \Lambda_0$.

Then,

$$\begin{aligned}\|x_0\|_1 &= \left\| \sum_{i \in \Lambda_1} (A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i x_1(i) \right\|_1 \\ &\leq \sum_{i \in \Lambda_1} |x_1(i)| \cdot \|(A_{\Lambda_0}^* A_{\Lambda_0})^{-1} A_{\Lambda_0}^* a_i\|_1 \\ &< \sum_{i \in \Lambda_1} |x_1(i)| = \|x_1\|_1.\end{aligned}$$



But, is the solution even unique?

The sparsity of the sparsest vector in the nullspace of A ,

$$\text{spark}(A) := \min_z \|z\|_{\ell^0} \quad \text{subject to} \quad Az = 0.$$

Theorem (Spark and Coherence)

$$\text{spark}(A_{m,n}) \geq \min(m+1, \mu_2(A_{m,n})^{-1} + 1)$$

If $\|x_0\| < (\mu_2(A_{m,n})^{-1} + 1)/2$ unique satisfying $y = A_{m,n}x_0$.

Proof.

Gershgorin disc theorem for $A_\Lambda^* A_\Lambda$ with $|\Lambda| = k$:

1 on diagonal, off diagonals bounded by $\mu_2(A_{m,n})$.

If $k < \mu_2(A_{m,n})^{-1} + 1$, smallest singular value of $A_\Lambda^* A_\Lambda$ is > 0 \square

How to interpret these results, is better possible?

- ▶ When is $\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1)$?

Grassman Frames: $\mu_2(A_{m,n}) \geq \left(\frac{n-m}{m(n-1)} \right)^{1/2} \sim m^{-1/2}$

“Sqrt bottleneck” $\|x_0\|_{\ell^0} \lesssim \sqrt{m}$

How to interpret these results, is better possible?

- ▶ When is $\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1)$?

Grassman Frames: $\mu_2(A_{m,n}) \geq \left(\frac{n-m}{m(n-1)} \right)^{1/2} \sim m^{-1/2}$

“Sqrt bottleneck” $\|x_0\|_{\ell^0} \lesssim \sqrt{m}$

- ▶ Is better possible? (not without more)

Fourier & Dirac: $A_{m,n} = [F \ I]$ for m the square of an integer:

Let $\Lambda = [\sqrt{m}, 2\sqrt{m}, \dots, m]$, then

$$\sum_{j \in \Lambda} e_j = \sum_{j \in \Lambda} f_j \implies \text{spark}(A_{m,n}) = 2\sqrt{m}.$$

How to interpret these results, is better possible?

- ▶ When is $\|x_0\|_{\ell^0} < \frac{1}{2} (\mu_2(A_{m,n})^{-1} + 1)$?

Grassman Frames: $\mu_2(A_{m,n}) \geq \left(\frac{n-m}{m(n-1)} \right)^{1/2} \sim m^{-1/2}$

“Sqrt bottleneck” $\|x_0\|_{\ell^0} \lesssim \sqrt{m}$

- ▶ Is better possible? (not without more)

Fourier & Dirac: $A_{m,n} = [F \ I]$ for m the square of an integer:

Let $\Lambda = [\sqrt{m}, 2\sqrt{m}, \dots, m]$, then

$$\sum_{j \in \Lambda} e_j = \sum_{j \in \Lambda} f_j \implies \text{spark}(A_{m,n}) = 2\sqrt{m}.$$

- ▶ Slightly more accurate sometimes with cumulative coherence:
 $\max_{i \in \Lambda} \max_{\Lambda'} \sum_{j \in \Lambda'} a_i^* a_j$
- ▶ To avoid pathological cases introduce randomness

One step thresholding: average sign pattern [ScVa07]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set Λ the index set of the $k \leq m$ largest in $|A_{m,n}^* y|$

Output the n -vector x whose entries are

$$x_\Lambda = (A_\Lambda^* A_\Lambda)^{-1} A_\Lambda y \quad \text{and} \quad x(i) = 0 \text{ for } i \notin \Lambda.$$

One step thresholding: average sign pattern [ScVa07]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set Λ the index set of the $k \leq m$ largest in $|A_{m,n}^* y|$
Output the n -vector x whose entries are

$$x_\Lambda = (A_\Lambda^* A_\Lambda)^{-1} A_\Lambda y \quad \text{and} \quad x(i) = 0 \text{ for } i \notin \Lambda.$$

Theorem

Let $y = A_{m,n} x_0$, with the columns of $A_{m,n}$ having unit ℓ^2 norm, the sign of the nonzeros in x_0 selected randomly from ± 1 independent of $A_{m,n}$, and

$$\|x_0\|_{\ell^0} < (128 \log(2n/\epsilon))^{-1} \nu_\infty^2(x_0) \mu_2^{-2}(A_{m,n}),$$

then, with probability greater than $1 - \epsilon$, the Thresholding decoder with $k = \|x_0\|_{\ell^0}$ will return x_0 .

One step thresholding: average sign pattern (proof, pg. 1)

Theorem (Rademacher concentration)

Fix a vector α . Let ϵ be a Rademacher series, vector with entries drawn uniformly from ± 1 , of the same length as α , then

$$\text{Prob} \left(\left| \sum_i \epsilon_i \alpha_i \right| > t \right) \leq 2 \exp \left(\frac{-t^2}{32 \|\alpha\|_2^2} \right)$$

Let $\Lambda := \text{supp}(x_0)$. Thresholding fail to recover x_0 if

$$\max_{i \notin \Lambda} |a_i^* y| > \min_{i \in \Lambda} |a_i^* y|.$$

$$\text{Prob} \left(\max_{i \notin \Lambda} |a_i^* y| > p \quad \text{and} \quad \min_{i \in \Lambda} |a_i^* y| < p \right) \leq$$

$$\text{Prob} \left(\max_{i \notin \Lambda} |a_i^* y| > p \right) + \text{Prob} \left(\min_{i \in \Lambda} |a_i^* y| < p \right) =: P_1 + P_2$$

One step thresholding: average sign pattern (proof, pg. 2)

$$\begin{aligned}P_1 &= \text{Prob} \left(\max_{i \notin \Lambda} |a_i^* y| > p \right) \\&\leq \sum_{i \notin \Lambda} \text{Prob} (|a_i^* y| > p) \\&= \sum_{i \notin \Lambda} \text{Prob} \left(\left| \sum_{j \in \Lambda} x_0(j) (a_i^* a_j) \right| > p \right) \\&\leq 2 \sum_{i \notin \Lambda} \exp \left(\frac{-p^2}{32 \sum_{j \in \Lambda} |x_0(j)|^2 |a_i^* a_j|^2} \right) \\&\leq 2(n - k) \exp \left(\frac{-p^2}{32k \|x_0\|_\infty^2 \mu_2^2(A_{m,n})} \right).\end{aligned}$$

One step thresholding: average sign pattern (proof, pg. 3)

$$\begin{aligned}P_2 &= \text{Prob} \left(\min_{i \in \Lambda} |a_i^* y| < p \right) \\&\leq \text{Prob} \left(\min_{i \in \Lambda} |x_0(i)| - \max_{i \in \Lambda} \left| \sum_{j \in \Lambda, j \neq i} x_0(j) (a_i^* a_j) \right| < p \right) \\&\leq \sum_{i \in \Lambda} \text{Prob} \left(\left| \sum_{j \in \Lambda, j \neq i} x_0(j) (a_i^* a_j) \right| > \min_{i \in \Lambda} |x_0(i)| - p \right) \\&\leq 2 \sum_{i \in \Lambda} \exp \left(\frac{-(\min_{i \in \Lambda} |x_0(i)| - p)^2}{32 \sum_{j \in \Lambda, j \neq i} |x_0(j)|^2 |a_i^* a_j|^2} \right) \\&\leq 2k \exp \left(\frac{-(\min_{i \in \Lambda} |x_0(i)| - p)^2}{32k \|x_0\|_\infty^2 \mu_2^2(A_{m,n})} \right).\end{aligned}$$

One step thresholding: average sign pattern (proof, pg. 4)

Balance P_1 and P_2 by setting $p := \min_{i \in \Lambda} |x_0(i)|/2$:

$$P_1 + P_2 \leq 2n \exp \left(\frac{-(\min_{i \in \Lambda} |x_0(i)|)^2}{128k \|x_0\|_\infty^2 \mu_2^2(A_{m,n})} \right) \leq 2n \exp \left(\frac{-\nu_\infty(x_0)^2}{128k \mu_2^2(A_{m,n})} \right).$$

Setting this bound on the probability of failure equal to ϵ and solving for k yields the conclusion of the proof. \square

- ▶ Similar work for matching pursuit by Schnass, ℓ^1 by Tropp, and in Statistical RICs
- ▶ Stronger uniform statements we need more than coherence.

Restricted Isometry Constants

The set of k -sparse vectors

$$\chi^n(k) := \{x \in \mathbb{R}^n : \|x\|_{\ell^0} \leq k\}.$$

Upper and lower RICs of A , U_k and L_k respectively, are defined as

$$U_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 + c)\|x\|_2^2 \geq \|Ax\|_2^2 \quad \forall x \in \chi^n(k).$$

$$L_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 - c)\|x\|_2^2 \leq \|Ax\|_2^2, \quad \forall x \in \chi^n(k);$$

Restricted Isometry Constants

The set of k -sparse vectors

$$\chi^n(k) := \{x \in \mathbb{R}^n : \|x\|_{\ell^0} \leq k\}.$$

Upper and lower RICs of A , U_k and L_k respectively, are defined as

$$U_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 + c)\|x\|_2^2 \geq \|Ax\|_2^2 \quad \forall x \in \chi^n(k).$$

$$L_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 - c)\|x\|_2^2 \leq \|Ax\|_2^2, \quad \forall x \in \chi^n(k);$$

► Pros:

- Easy to use to prove optimal order results
- A general tool for any algorithm (wide usage)

► Cons:

- Don't know how to calculate it
- A general tool for any algorithm (bad results)

Restricted Isometry Constants

The set of k -sparse vectors

$$\chi^n(k) := \{x \in \mathbb{R}^n : \|x\|_{\ell^0} \leq k\}.$$

Upper and lower RICs of A , U_k and L_k respectively, are defined as

$$U_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 + c)\|x\|_2^2 \geq \|Ax\|_2^2 \quad \forall x \in \chi^n(k).$$

$$L_k := \min_{c \geq 0} c \quad \text{subject to} \quad (1 - c)\|x\|_2^2 \leq \|Ax\|_2^2, \quad \forall x \in \chi^n(k);$$

► Pros:

- Easy to use to prove optimal order results
- A general tool for any algorithm (wide usage)

► Cons:

- Don't know how to calculate it
 - A general tool for any algorithm (bad results)
- No known matrix with bounded RICs for $k \sim m \sim n$
- Coherence for $k \sim m^2$ or random matrices used

The first RIC bounds (Gaussian): [CaTa05]

Let $\sigma^{\max}(B)$ and $\sigma^{\min}(B)$ be the largest and smallest singular values of B respectively. Then,

$$\begin{aligned}\text{Prob}(\sigma^{\max}(A_k) > 1 + \sqrt{k/m} + o(1) + t) &\leq \exp(-mt^2/2) \\ \text{Prob}(\sigma^{\min}(A_k) < 1 - \sqrt{k/m} + o(1) - t) &\leq \exp(-mt^2/2),\end{aligned}$$

where $o(1)$ denotes a quantity that tends to zero as $m \rightarrow \infty$.

Definition

Set $\delta = m/n$ and $\rho = k/m$ with $(\delta, \rho) \in (0, 1)^2$ and define:

$$\mathcal{U}^{CT}(\delta, \rho) := \left[1 + \sqrt{\rho} + (2\delta^{-1}H(\delta\rho))^{1/2} \right]^2 - 1$$

$$\mathcal{L}^{CT}(\delta, \rho) := 1 - \max \left\{ 0, \left[1 - \sqrt{\rho} - (2\delta^{-1}H(\delta\rho))^{1/2} \right]^2 \right\},$$

where Shannon Entropy $H(p) := -p \log p - (1 - p) \log(1 - p)$

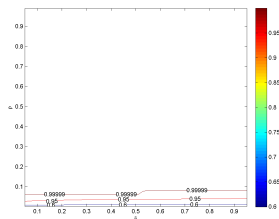
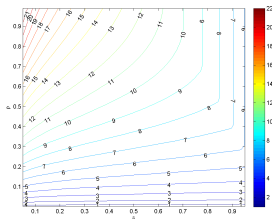
RIC bounds (Gaussian): [CaTa05]

Theorem

$A_{m,n}$ entries are drawn i.i.d. from the Gaussian normal $\mathcal{N}(0, 1/m)$. Let $\delta_m = m/n$ and $\rho_m = k/m$. For any fixed $\epsilon > 0$, in the limit as $\delta_m \rightarrow \delta \in (0, 1)$ and $\rho_m \rightarrow \rho \in (0, 1)$ as $m \rightarrow \infty$,

$$\mathbf{P}(L_k < \mathcal{L}^{CT}(\delta, \rho) - \epsilon) \rightarrow 0 \quad \text{and} \quad \mathbf{P}(U_k < \mathcal{U}^{CT}(\delta, \rho) + \epsilon) \rightarrow 0$$

exponentially in m .



$\mathcal{U}^{CT}(\delta, \rho)$ (left panel) and $\mathcal{L}^{CT}(\delta, \rho)$ (right panel).

RIC bounds (Gaussian): [CaTa05] (proof)

Proof.

$$\begin{aligned} & \text{Prob} \left(\max_{K \subset \Omega, |K|=k} \sigma^{\max}(A_K) > (1 + \sqrt{k/m}) + o(1) + t \right) \\ & \leq \sum_{K \subset \Omega, |K|=k} \text{Prob} \left(\sigma^{\max}(A_K) > (1 + \sqrt{k/m}) + o(1) + t \right) \\ & \leq \binom{n}{k} \exp(-mt^2/2) \leq \text{poly}(n) \cdot \exp \left(m \left[\delta^{-1} H(\rho\delta) - t^2/2 \right] \right), \end{aligned}$$

Use smallest t such that prob goes to zero.

Solve for the zero of the exponent: $t = [2\delta^{-1}H(\rho\delta)]^{1/2}$.

This corresponds to an upper bound on

$$\begin{aligned} & \text{Prob} \left(\max \sigma^{\max}(A_K) > (1 + \sqrt{k/m}) + [2\delta^{-1}H(\rho\delta)]^{1/2} + \epsilon + o(1) \right) \\ & \leq \text{poly}(n) \cdot e^{-m(\epsilon+o(1))} \end{aligned}$$

which converges to zero exponentially with m .

□

Iterative Hard Thresholding [Fu10]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set $x^0 = 0$ and $j = 0$.

While $\|y - A_{m,n}x^j\|_2 < Tol$ repeat the following steps:

set $v^j := x^j + A_{m,n}^*(y - A_{m,n}x^j)$, and $x^{j+1} = H_k(v^j)$.

Output x^j .

Iterative Hard Thresholding [Fu10]

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set $x^0 = 0$ and $j = 0$.

While $\|y - A_{m,n}x^j\|_2 < Tol$ repeat the following steps:

set $v^j := x^j + A_{m,n}^*(y - A_{m,n}x^j)$, and $x^{j+1} = H_k(v^j)$.

Output x^j .

Theorem

Let $y = A_{m,n}x_0 + e$ for x_0 k -sparse and $A_{m,n}$ in General Position.

Set $\mu^{iht} := 2 \max(L_{3k}, U_{3k})$ and $\xi^{iht} := 2(1 + U_{2k})^{1/2}$.

With k used for the hard thresholding function, IHT satisfy the inequality

$$\|x^j - x_0\|_2 \leq (\mu^{iht})^j \|x_0\| + \frac{\xi^{iht}}{1 - \mu^{iht}} \|e\|_2.$$

For $\mu^{iht} < 1$ convergence of x^j to approximation of x_0 .

Iterative Hard Thresholding (proof, pg. 1)

Proof.

$H_k(\cdot)$ returns the k -sparse closest in the ℓ^2 norm, for instance

$$\|v^j - H_k(v^j)\|_2 = \|v^j - x^{j+1}\|_2 \leq \|v^j - x_0\|_2. \quad (1)$$

Note that

$$\begin{aligned} \|v^j - x^{j+1}\|_2^2 &= \|(v^j - x_0) + (x_0 - x^{j+1})\|_2^2 = \\ &\|v^j - x_0\|_2^2 + \|x_0 - x^{j+1}\|_2^2 + 2\operatorname{Re}((v^j - x_0)^*(x_0 - x^{j+1})) \end{aligned}$$

where $\operatorname{Re}(c)$ denotes the real part of c .

Bounding the above expression using (1) and canceling the $\|v^j - x_0\|_2^2$ term yields

$$\|x^{j+1} - x_0\|_x^2 \leq 2\operatorname{Re}((v^j - x_0)^*(x^{j+1} - x_0)).$$

Iterative Hard Thresholding (proof, pg. 2)

Consider the $3k$ sparse set

$$\Lambda = \text{supp}(x_0) \cup \text{supp}(x^j) \cup \text{supp}(x^{j+1}):$$

$$\begin{aligned}\|x^{j+1} - x_0\|_2^2 &\leq 2\text{Re}((x^j - x_0)^*(x^{j+1} - x_0)) \\ &= 2\text{Re}\left(\left((I - A_{m,n}^* A_{m,n})(x^j - x_0)\right)^*(x^{j+1} - x_0)\right) \\ &\quad + 2\text{Re}(e^* A_{m,n}(x^{j+1} - x_0)) \\ &= 2\text{Re}\left(\left((I - A_{\Lambda}^* A_{\Lambda})(x^j - x_0)_{\Lambda}\right)^*(x^{j+1} - x_0)_{\Lambda}\right) \\ &\quad + 2\text{Re}(e^* A_{m,n}(x^{j+1} - x_0)) \\ &\leq 2\|I - A_{\Lambda}^* A_{\Lambda}\|_2 \cdot \|x^j - x_0\|_2 \cdot \|x^{j+1} - x_0\|_2 \\ &\quad + 2\|e\|_2 \cdot \|A_{m,n}(x^{j+1} - x_0)\|_2\end{aligned}$$

Iterative Hard Thresholding (proof, pg. 3)

RIC bounds $\|I - A_{\Lambda}^* A_{\Lambda}\|_2 \leq \max(U_{3k}, L_{3k})$ and $\|A_{m,n}(x^{j+1} - x_0)\|_2 \leq (1 + U_{2k})^{1/2} \|x^{j+1} - x_0\|_2$ then dividing by $\|x^{j+1} - x_0\|_2$ yields

$$\|x^{j+1} - x_0\|_2 \leq 2 \max(L_{3k}, U_{3k}) \cdot \|x^j - x_0\|_2 + 2(1 + U_{2k})^{1/2} \|e\|_2$$

Let $\mu^{iht} := 2 \max(L_{3k}, U_{3k})$ and $\xi^{iht} := 2(1 + U_{2k})^{1/2}$.

Error at step j in terms of initial error $\|x^0 - x_0\|_2 = \|x_0\|_2$

$$\|x^j - x_0\|_2 \leq (\mu^{iht})^J \cdot \|x_0\|_2 + \xi^{iht} \|e\|_2 \sum_{\ell=0}^{j-1} (\mu^{iht})^{\ell}$$

Replacing final sum with bound $1/(1 - \mu^{iht})$ completes the proof.



ℓ^1 -regularization [Ca08]

Input: y , $A_{m,n}$, and tolerance ϵ .

“Algorithm”: Return $x^* = \operatorname{argmin} \|x\|_1$ subject to $\|y - Ax\|_2 \leq \epsilon$.

Theorem

Let $y = A_{m,n}x_0 + e$ for x_0 k -sparse, $\|e\|_2 \leq \epsilon$, and $A_{m,n}$ in General Position.

Set $\mu^{\ell^1} := 2^{-1/2}(U_{2k} + L_{2k})/(1 - L_{2k})$ and

$\xi^{\ell^1} := 2^{3/2}(1 + U_{2k})^{1/2}/(1 - L_{2k})$

With $x^* = \operatorname{argmin} \|x\|_{\ell}^1$ subject to $\|y - A_{m,n}x\|_2 \leq \epsilon$
and $\mu^{iht} < 1$ then

$$\|x_0 - x^*\|_2 < \frac{\xi^{\ell^1}}{1 - \mu^{\ell^1}} \cdot \|e\|_2$$

ℓ^1 -regularization (proof, pg. 1)

Proof.

Let $h := x^* - x_0$. The goal is to show $\|h\|_2 \leq \text{Const.} \|e\|_2$

ℓ^1 -regularization (proof, pg. 1)

Proof.

Let $h := x^* - x_0$. The goal is to show $\|h\|_2 \leq \text{Const.} \|e\|_2$

Let $\Lambda_0 := \text{supp}(x_0)$. Partition the rest of $1, 2, \dots, n$ into k sets

Let Λ_1 be the support of the largest k entries of $|h_{\Lambda_0^c}|$,

Λ_2 the support set of the next largest k entries in $|h_{\Lambda_0^c}|$, etc...

ℓ^1 -regularization (proof, pg. 1)

Proof.

Let $h := x^* - x_0$. The goal is to show $\|h\|_2 \leq \text{Const.} \|e\|_2$

Let $\Lambda_0 := \text{supp}(x_0)$. Partition the rest of $1, 2, \dots, n$ into k sets

Let Λ_1 be the support of the largest k entries of $|h_{\Lambda_0^c}|$,

Λ_2 the support set of the next largest k entries in $|h_{\Lambda_0^c}|$, etc...

Show that $\|h\|_2 \leq \text{Const.} \|e\|_2$ small by considering

$h_{(\Lambda_0 \cup \Lambda_1)^c}$ and $h_{(\Lambda_0 \cup \Lambda_1)^c}$

$\Lambda_{01} := (\Lambda_0 \cup \Lambda_1)$ contains support of x_0 and where h is largest

Λ_{01}^c contains the rest of the n -vector

ℓ^1 -regularization (proof, pg. 2: Show $\|h_{\Lambda_0^c}\|_2$ “small”)

For vectors satisfying $\|y - A_{m,n}x\|_2 \leq \epsilon$, x^* has smallest ℓ^1 norm

$$\|x\|_1 \geq \|x^*\|_1 = \|x + h\|_1 \geq \|x_{\Lambda_0}\|_1 - \|h_{\Lambda_0}\|_1 + \|h_{\Lambda_0^c}\|_1$$

which implies that $\|h_{\Lambda_0^c}\|_1 \leq \|h_{\Lambda_0}\|_1$.

ℓ^1 -regularization (proof, pg. 2: Show $\|h_{\Lambda_{01}^c}\|_2$ “small”)

For vectors satisfying $\|y - A_{m,n}x\|_2 \leq \epsilon$, x^* has smallest ℓ^1 norm

$$\|x\|_1 \geq \|x^*\|_1 = \|x + h\|_1 \geq \|x_{\Lambda_0}\|_1 - \|h_{\Lambda_0}\|_1 + \|h_{\Lambda_0^c}\|_1$$

which implies that $\|h_{\Lambda_0^c}\|_1 \leq \|h_{\Lambda_0}\|_1$.

By construction, largest entry in h_{Λ_j} smaller than average in $h_{\Lambda_{j-1}}$

$$\|h_{\Lambda_j}\|_2 \leq k^{1/2} \|h_{\Lambda_j}\|_\infty \leq k^{1/2} (k^{-1} \|h_{\Lambda_{j-1}}\|_1) = k^{-1/2} \|h_{\Lambda_{j-1}}\|_1$$

Use above bound and triangle inequality to obtain

$$\|h_{\Lambda_{01}^c}\|_2 = \left\| \sum_{j \geq 2} h_{\Lambda_j} \right\|_2 \leq \sum_{j \geq 2} \|h_{\Lambda_j}\|_2 \leq \sum_{j \geq 1} k^{-1/2} \|h_{\Lambda_j}\|_1 = k^{-1/2} \|h_{\Lambda_0^c}\|_1$$

With the above, $\|h_{\Lambda_0^c}\|_1 \leq \|h_{\Lambda_0}\|_1$, and Cauchy Schwartz

$$\|h_{\Lambda_{01}^c}\|_2 \leq k^{-1/2} \|h_{\Lambda_0}\|_1 \leq k^{-1/2} \left(k^{1/2} \|h_{\Lambda_0}\|_2 \right) = \|h_{\Lambda_0}\|_2 \leq \|h_{\Lambda_{01}}\|_2$$

ℓ^1 -regularization (proof, pg. 3: a few notes)

For any $j \neq k$

$$|(Ah_{\Lambda_j})^*(Ah_{\Lambda_k})| \leq \frac{U_{2k} + L_{2k}}{2} \|h_{\Lambda_j}\|_2 \cdot \|h_{\Lambda_k}\|_2$$

Proof.

Let u and v be unit norm k -sparse with disjoint support I and J then $\|Au \pm Av\|_2 = \|A_{I \cup J}(u \pm v)\|_2$ and using RIC bounds for $2k$

$$(1 - L_{2k})\|u + v\|_2^2 \leq \|A_{I \cup J}(u + v)\|_2^2 \leq (1 + U_{2k})\|u + v\|_2^2$$

with u and v disjoint unit norm we have $\|u + v\|_2^2 = 2$.

Substituting the above upper and lower bounds into the following

$$|(Au)^*Av| = \frac{1}{4} \left| \|Au + Av\|_2^2 - \|Au - Av\|_2^2 \right|$$

□

ℓ^1 -regularization (proof, pg. 4: Show $\|h_{\Lambda_{01}}\|_2$ “small”)

First note that:

$$\|Ah\|_2 = \|A(x^* - x_0)\|_2 \leq \|Ax^* - y\|_2 + \|y - Ax_0\|_2 \leq 2\|e\|$$

Bound $\|h_{\Lambda_{01}}\|_2$ through upper and lower bounds on $\|Ah_{\Lambda_{01}}\|_2^2$

Begin with the upper bound:

$$\begin{aligned}\|Ah_{\Lambda_{01}}\|_2^2 &= (Ah_{\Lambda_{01}})^* \left(Ah - \sum_{j \geq 2} Ah_{\Lambda_j} \right) \\ &\leq \|Ah_{\Lambda_{01}}\|_2 \cdot \|Ah\|_2 \\ &\quad + \sum_{j \geq 2} [(Ah_{\Lambda_0})^* Ah_{\Lambda_j} + (Ah_{\Lambda_1})^* Ah_{\Lambda_j}] \\ &\leq (1 + U_{2k})^{1/2} \|h_{\Lambda_{01}}\|_2 \cdot 2\|e\| \\ &\quad + \frac{U_{2k} + L_{2k}}{2} (\|h_{\Lambda_0}\|_2 + \|h_{\Lambda_1}\|_2) \sum_{j \geq 2} \|h_{\Lambda_j}\|_2\end{aligned}$$

ℓ^1 -regularization (proof, pg. 5: Show $\|h_{\Lambda_{01}}\|_2$ “small”)

Continue upper bound using $\|h_{\Lambda_0}\|_2 + \|h_{\Lambda_1}\|_2 \leq \sqrt{2}\|h_{\Lambda_{01}}\|_2$

$$\begin{aligned}\|Ah_{\Lambda_{01}}\|_2^2 &\leq 2(1 + U_{2k})^{1/2}\|h_{\Lambda_{01}}\|_2 \cdot \|e\| \\ &\quad + \frac{\sqrt{2}}{2}(U_{2k} + L_{2k})\|h_{\Lambda_{01}}\|_2 k^{-1/2}\|h_{\Lambda_0^c}\|_1\end{aligned}\quad (2)$$

Lower bound $\|Ah_{\Lambda_{01}}\|_2^2$ using simple RIP bound

$$(1 - L_{2k})\|h_{\Lambda_{01}}\|_2^2 \leq \|Ah_{\Lambda_{01}}\|_2^2$$

Stating lower and upper bound of $\|Ah_{\Lambda_{01}}\|_2^2$ and divide by $\|h_{\Lambda_{01}}\|_2$

$$(1 - L_{2k})\|h_{\Lambda_{01}}\|_2 \leq 2(1 + U_{2k})\|e\|_2 + \frac{\sqrt{2}}{2}(U_{2k} + L_{2k})k^{-1/2}\|h_{\Lambda_0^c}\|_1$$

ℓ^1 -regularization (proof, pg. 6: Show $\|h_{\Lambda_{01}}\|_2$ “small”)

Recall

$$\|h_{\Lambda_0^c}\|_1 \leq \|h_{\Lambda_0}\|_1 \leq k^{1/2} \|h_{\Lambda_0}\|_2 \leq k^{1/2} \|h_{\Lambda_{01}}\|_2$$

and substitute into bound of $\|h_{\Lambda_{01}}\|_2$ from prior slide gives

$$(1 - L_{2k}) \|h_{\Lambda_{01}}\|_2 \leq 2(1 + U_{2k}) \|e\|_2 + \frac{\sqrt{2}}{2} (U_{2k} + L_{2k}) \|h_{\Lambda_{01}}\|_2$$

If $1 - L_{2k} < (U_{2k} + L_{2k})2^{-1/2}$, solving for $\|h_{\Lambda_{01}}\|_2$ gives bound

$$\|h_{\Lambda_{01}}\|_2 \leq \left(1 - \mu^{\ell^1}\right)^{-1} \frac{2(1 + U_{2k})^{1/2}}{1 - L_{2k}} \cdot \|e\|_2$$

where $\mu^{\ell^1} := 2^{-1/2}(U_{2k} + L_{2k})/(1 - L_{2k}) < 1$

ℓ^1 -regularization (proof, pg. 7: putting it all together)

The goal was to bound $\|x^* - x_0\|_2^2 = \|h\|_2^2 = \|h_{\Lambda_{01}}\|_2^2 + \|h_{\Lambda_{01}^c}\|_2^2$

Using $\|h_{\Lambda_{01}^c}\|_2^2 \leq \|h_{\Lambda_{01}}\|_2^2$ and bound on $\|h_{\Lambda_{01}}\|_2$ obtain

$$\|x^* - x_0\|_2 \leq \sqrt{2} \left(1 - \mu^{\ell^1}\right)^{-1} \frac{2(1 + U_{2k})^{1/2}}{1 - L_{2k}} \cdot \|e\|_2$$

Let $\xi^{\ell^1} := 2^{3/2}(1 + U_{2k})^{1/2}/(1 - L_{2k})$ and have standard form

$$\|x^* - x_0\|_2 \leq \frac{\xi^{\ell^1}}{1 - \mu^{\ell^1}} \|e\|_2$$

recovery guarantee provided $\mu^{\ell^1} < 1$.



How to interpret this result, should we be happy?

- ▶ Optimal order if L_{3k}, U_{3k} bounded for $k \sim m$ and $m \sim n$
- ▶ There are random matrices what w.h.p. have L_k, U_k bounded!

How to interpret this result, should we be happy?

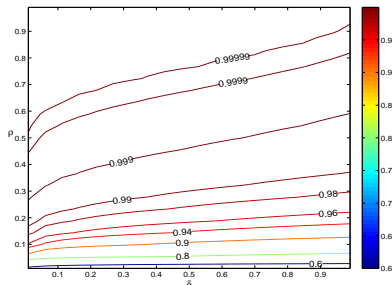
- ▶ Optimal order if L_{3k}, U_{3k} bounded for $k \sim m$ and $m \sim n$
- ▶ There are random matrices what w.h.p. have L_k, U_k bounded!
- ▶ When is $\mu^{iht} := 2 \max(L_{3k}, U_{3k}) < 1$
- ▶ Many algorithms with bounds of this form, which to use?

How to interpret this result, should we be happy?

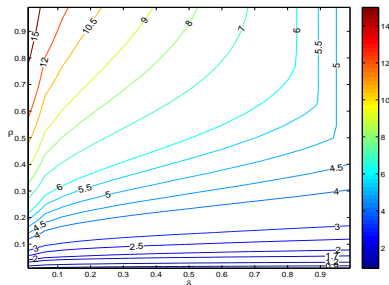
- ▶ Optimal order if L_{3k}, U_{3k} bounded for $k \sim m$ and $m \sim n$
- ▶ There are random matrices what w.h.p. have L_k, U_k bounded!
- ▶ When is $\mu^{iht} := 2 \max(L_{3k}, U_{3k}) < 1$
- ▶ Many algorithms with bounds of this form, which to use?
- ▶ To answer these questions need to have bounds on the RICs.
- ▶ Previous CaTa05 bounds insufficient for reasonable k

RIC bounds for Gaussian $\mathcal{N}(0, m^{-1})$ [BaTa10, BICaTa09]

$$(1 - L(\delta, \rho))\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + U(\delta, \rho))\|x\|_2^2$$



$L(\delta, \rho)$

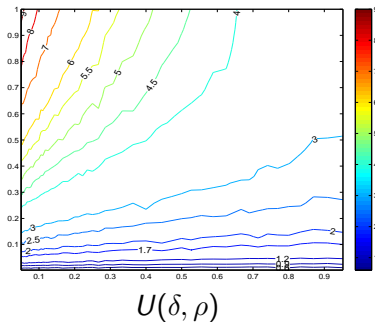
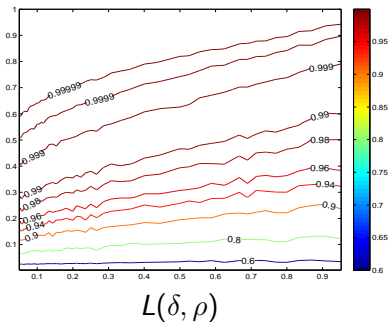


$U(\delta, \rho)$

- ▶ Using Wishart Distributions and groupings
- ▶ Less than 1.57 times empirically observed values

RIC bounds for Gaussian $\mathcal{N}(0, m^{-1})$ [BaTa10, BICaTa09]

$$(1 - L(\delta, \rho))\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + U(\delta, \rho))\|x\|_2^2$$



- ▶ Empirical draw with $n = 400$, consistent with $n = 200, 800$
- ▶ Local searches for local extremal singular values: algorithms of Richtarik (U) and Dossal et al (L).

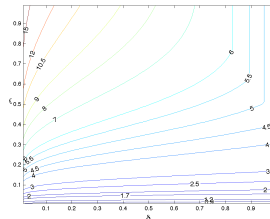
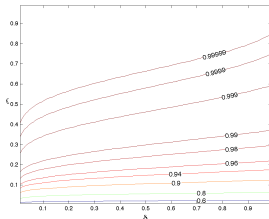
RIC bounds for Gaussian $\mathcal{N}(0, m^{-1})$: [BICaTa09]

Theorem

$A_{m,n}$ entries are drawn i.i.d. from the Gaussian normal $\mathcal{N}(0, 1/m)$. Let $\delta_m = m/n$ and $\rho_m = k/m$. For any fixed $\epsilon > 0$, in the limit as $\delta_m \rightarrow \delta \in (0, 1)$ and $\rho_m \rightarrow \rho \in (0, 1)$ as $m \rightarrow \infty$,

$$\mathbf{P}(L_k < \mathcal{L}^{\text{BCT}}(\delta, \rho) - \epsilon) \rightarrow 1 \quad \text{and} \quad \mathbf{P}(U_k < \mathcal{U}^{\text{BCT}}(\delta, \rho) + \epsilon) \rightarrow 1$$

exponentially in m .



$\mathcal{L}^{\text{BCT}}(\delta, \rho)$ (left panel) and $\mathcal{U}^{\text{BCT}}(\delta, \rho)$ (right panel).

Definition of BCT bounds

Let $H(p) := p \log(1/p) + (1 - p) \log(1/(1 - p))$ denote the usual Shannon Entropy with base e logarithms, and let

$$\psi_{\min}(\lambda, \rho) := H(\rho) + \frac{1}{2} [(1 - \rho) \log \lambda + 1 - \rho + \rho \log \rho - \lambda],$$

$$\psi_{\max}(\lambda, \rho) := \frac{1}{2} [(1 + \rho) \log \lambda + 1 + \rho - \rho \log \rho - \lambda].$$

Define $\lambda^{\min}(\delta, \rho)$ and $\lambda^{\max}(\delta, \rho)$ as the solution to (3) and (4), respectively:

$$\delta \psi_{\min}(\lambda^{\min}(\delta, \rho), \rho) + H(\rho \delta) = 0 \quad \text{for} \quad \lambda^{\min}(\delta, \rho) \leq 1 - \rho \quad (3)$$

$$\delta \psi_{\max}(\lambda^{\max}(\delta, \rho), \rho) + H(\rho \delta) = 0 \quad \text{for} \quad \lambda^{\max}(\delta, \rho) \geq 1 + \rho. \quad (4)$$

Define $\mathcal{L}^{\text{BCT}}(\delta, \rho)$ and $\mathcal{U}^{\text{BCT}}(\delta, \rho)$ as

$$\mathcal{L}^{\text{BCT}}(\delta, \rho) := 1 - \lambda^{\min}(\delta, \rho) \quad \text{and} \quad \mathcal{U}^{\text{BCT}}(\delta, \rho) := \min_{\nu \in [\rho, 1]} \lambda^{\max}(\delta, \nu) - 1.$$

RIC bounds for $\mathcal{N}(0, m^{-1})$ (proof, pg. 1: largest)

Begin with behaviour of largest singular value[Edelman88]

Let A_Λ be a matrix of size $m \times k$ whose entries are drawn i.i.d from $\mathcal{N}(0, m^{-1})$. Let $f_{\max}(k, m; \lambda)$ denote the probability density function for the largest eigenvalue of the Wishart matrix $A_\Lambda^T A_\Lambda$ of size $k \times k$. Then $f_{\max}(k, m; \lambda)$ satisfies:

$$f_{\max}(k, m; \lambda) \leq \left[(2\pi)^{1/2} (m\lambda)^{-3/2} \left(\frac{m\lambda}{2} \right)^{(m+k)/2} \frac{1}{\Gamma(\frac{k}{2})\Gamma(\frac{m}{2})} \right] \cdot e^{-m\lambda/2}$$

Large deviation (large k and m) behavior of f_{\max} , apply $m^{-1} \log(\cdot)$

$$\frac{1}{2} \left[(1 + \rho_m) \log \lambda - \left(\rho_m - \frac{1}{m} \right) \log \rho_m + \frac{2}{m} \log \frac{m}{2} + 1 + \rho_m - \lambda \right].$$

Large m limit gives exponential behaviour $\psi_{\max}(\lambda, \rho)$

RIC bounds for $\mathcal{N}(0, m^{-1})$ (proof, pg. 2: smallest)

Smallest singular value [Edelman88] similarly

Let $f_{\min}(k, m; \lambda)$ denote the probability density function for the smallest eigenvalue of the Wishart matrix $A_{\Lambda}^T A_{\Lambda}$ of size $k \times k$.

Then $f_{\min}(k, m; \lambda)$ bounded above by:

$$\leq \left(\frac{\pi}{2m\lambda}\right)^{1/2} \cdot e^{-m\lambda/2} \left(\frac{m\lambda}{2}\right)^{(m-k)/2} \cdot \left[\frac{\Gamma(\frac{m+1}{2})}{\Gamma(\frac{k}{2})\Gamma(\frac{m-k+1}{2})\Gamma(\frac{m-k+2}{2})} \right]$$

Large deviation (large k and m) behavior of f_{\min} , apply $m^{-1} \log(\cdot)$

$$\psi_{\min}(\lambda, \rho_m) := H(\rho_m) + \frac{1}{2} [(1 - \rho_m) \log \lambda + \rho_m \log \rho_m + 1 - \rho_m - \lambda].$$

Large m limit gives $f_{\min}(k, m; \lambda) \leq \exp[m \cdot \psi_{\min}(\lambda, \rho)]$

RIC bounds for $\mathcal{N}(0, m^{-1})$ (proof, pg. 3: union bound)

Have bounds on PDFs of largest and smallest eigenvalues:

$$f_{\max}(k, m; \lambda) \leq \exp[m \cdot \psi_{\max}(\lambda, \rho)] \quad f_{\min}(k, m; \lambda) \leq \exp[m \cdot \psi_{\min}(\lambda, \rho)]$$

with

$$\psi_{\min}(\lambda, \rho) := H(\rho) + \frac{1}{2} [(1 - \rho) \log \lambda + 1 - \rho + \rho \log \rho - \lambda],$$

$$\psi_{\max}(\lambda, \rho) := \frac{1}{2} [(1 + \rho) \log \lambda + 1 + \rho - \rho \log \rho - \lambda].$$

Note: $\lim_{\lambda \downarrow 0} \psi_{\min}(\lambda, \rho) \rightarrow -\infty$ and $\lim_{\lambda \uparrow \infty} \psi_{\max}(\lambda, \rho) \rightarrow -\infty$

RIC bounds for $\mathcal{N}(0, m^{-1})$ (proof, pg. 3: union bound)

Have bounds on PDFs of largest and smallest eigenvalues:

$$f_{\max}(k, m; \lambda) \leq \exp[m \cdot \psi_{\max}(\lambda, \rho)] \quad f_{\min}(k, m; \lambda) \leq \exp[m \cdot \psi_{\min}(\lambda, \rho)]$$

with

$$\psi_{\min}(\lambda, \rho) := H(\rho) + \frac{1}{2} [(1 - \rho) \log \lambda + 1 - \rho + \rho \log \rho - \lambda],$$

$$\psi_{\max}(\lambda, \rho) := \frac{1}{2} [(1 + \rho) \log \lambda + 1 + \rho - \rho \log \rho - \lambda].$$

Note: $\lim_{\lambda \downarrow 0} \psi_{\min}(\lambda, \rho) \rightarrow -\infty$ and $\lim_{\lambda \uparrow \infty} \psi_{\max}(\lambda, \rho) \rightarrow -\infty$

Apply union bound over $\binom{n}{k} \sim \exp(n \cdot H(\delta\rho))$ sets

Solve zero level curve of exponent to get $\lambda^{\min}(\delta, \rho)$ and $\lambda^{\max}(\delta, \rho)$

□

Improve bounds further through grouping

Set $r = \binom{n}{k} \binom{p}{k}^{-1}$ and draw $u := rn$ sets M_i each of cardinality p , drawn uniformly at random from the $\binom{n}{p}$ possible p -sets.

Let G be the union of all u groups,

$$\text{Prob} \left(|G| < \binom{n}{k} \right) < C(k/n) n^{-1/2} e^{-n(1-\ln 2)}$$

where $C(z) \leq \frac{5}{4} (2\pi z(1-z))^{(-1/2)}$.

Improve bounds further through grouping

Set $r = \binom{n}{k} \binom{p}{k}^{-1}$ and draw $u := rn$ sets M_i each of cardinality p , drawn uniformly at random from the $\binom{n}{p}$ possible p -sets. Let G be the union of all u groups,

$$\text{Prob} \left(|G| < \binom{n}{k} \right) < C(k/n) n^{-1/2} e^{-n(1-\ln 2)}$$

where $C(z) \leq \frac{5}{4} (2\pi z(1-z))^{(-1/2)}$. **Proof.**

Select one set $K \subset 1, 2, \dots, N$ of cardinality $|K| = k$, draw of the sets M_i . The probability that K is not contained in M_i is $1/r$.

Probability K is not in any of the u sets M_i is $(1 - r^{-1})^u \leq e^{-u/r}$.

Applying a union bound over all $\binom{n}{k}$ sets K bounds

$$\text{Prob} \left(|G| < \binom{n}{k} \right) < \binom{n}{k} e^{-u/r}.$$

Stirling's Inequality: $\binom{n}{zn} \leq \frac{5}{4} (2\pi z(1-z)n)^{(-1/2)} e^{nH(z)}$

Note that $H(z) \leq \ln 2$ for $z \in [0, 1]$, and substituting u . □

Algorithms for Sparse Approximation

Input: A , y , and possibly tuning parameters

- ▶ ℓ^q -regularization (for $q \in (0, 1]$):

$$\min_x \|x\|_{\ell^q} \quad \text{subject to} \quad \|Ax - y\|_2 \leq \tau$$

- ▶ Simple Iterated Thresholding:

$$x^{t+1} = H_k(x^t + \kappa A^T(y - Ax^t))$$

- ▶ Two-Stage Thresholding (Subspace Pursuit, CoSaMP):

$$v^{t+1} = x^{t+1} = H_{\alpha k}(x^t + \kappa A^T(y - Ax^t))$$

$$I_t = \text{supp}(v^t) \cup \text{supp}(x^t) \quad \text{Join supp. sets}$$

$$w_{I_t} = (A_{I_t}^T A_{I_t})^{-1} A_{I_t}^T y \quad \text{Least squares fit}$$

$$x^{t+1} = H_{\beta k}(w^t) \quad \text{Second threshold}$$

When does RIP guarantee they work?

Phase transition (lower bounds) implied by RIP

Theorem: Let $y = Ax + e$ for any k -sparse x and with A $\mathcal{N}(0, m^{-1})$ iid. Define $\rho_S^{alg}(\delta)$ as the solution to $\mu^{alg}(\delta, \rho) = 1$.

For any $\epsilon > 0$, as $(k, m, n) \rightarrow \infty$ with $m/n \rightarrow \delta \in (0, 1)$ and $k/m \rightarrow \rho < (1 - \epsilon)\rho_S^{alg}(\delta)$, there is an exponentially high probability on the draw of A that after l iterations, the algorithm output \hat{x} approximates x within the bound

$$\|x - \hat{x}\|_2 \leq \left[\mu^{alg}(\delta, \rho) \right]^l \|x\|_2 + \frac{\xi^{alg}(\delta, \rho)}{1 - \mu^{alg}(\delta, \rho)} \|e\|_2.$$

Moreover, if $e = 0$, algorithm recovers x exactly in no more than

$$\ell_{max}^{alg}(x) := \left\lceil \frac{\log \nu_{\infty}(x)}{\log \mu^{alg}(\delta, \rho)} + 1 \right\rceil$$

iterates.

Lemma to show, naive replacement of bounds is ok

For some $\tau < 1$, define the set $\mathcal{Z} := (0, \tau)^p \times (0, \infty)^q$ and let $F : \mathcal{Z} \rightarrow \mathbb{R}$ be continuously differentiable on \mathcal{Z} . Let A be $m \times n$ with aRIP constants $L(\cdot, m, n)$, $U(\cdot, m, n)$ and let $L(\delta, \cdot)$, $U(\delta, \cdot)$ be their bounds. Define $\mathbf{1}$ to be the vector of all ones, and

$$z(k, m, n) := [L(k, m, n), \dots, L(pk, m, n), U(k, m, n), \dots, U(qk, m, n)]$$

$$z(\delta, \rho) := [L(\delta, \rho), \dots, L(\delta, p\rho), U(\delta, \rho), \dots, U(\delta, q\rho)].$$

Suppose, for all $t \in \mathcal{Z}$, $(\nabla F[t])_i \geq 0$ for all $i = 1, \dots, p + q$ and for any $v \in \mathcal{Z}$ we have $\nabla F[t] \cdot v > 0$. Then for any $c\epsilon > 0$, as $(k, m, n) \rightarrow \infty$ with $m/n \rightarrow \delta$, $k/n \rightarrow \rho$, there is an exponentially high probability on the draw of the matrix A that

$$\text{Prob}(F[z(k, m, n)] < F[z(\delta, \rho) + 1c\epsilon]) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

- ▶ Let F be μ or $\frac{\mu}{1-\xi}$
- ▶ Can replace (k, m, n) by (δ, ρ) bound and $\mathcal{O}(\epsilon)$

Lemma to show, naive replacement of bounds is ok

For some $\tau < 1$, define the set $\mathcal{Z} := (0, \tau)^p \times (0, \infty)^q$ and let $F : \mathcal{Z} \rightarrow \mathbb{R}$ be continuously differentiable on \mathcal{Z} . Let A be $m \times n$ with a RIP constants $L(\cdot, m, n)$, $U(\cdot, m, n)$ and let $L(\delta, \cdot)$, $U(\delta, \cdot)$ be their bounds. Define $\mathbf{1}$ to be the vector of all ones, and

$$z(k, m, n) := [L(k, m, n), \dots, L(pk, m, n), U(k, m, n), \dots, U(qk, m, n)]$$

$$z(\delta, \rho) := [L(\delta, \rho), \dots, L(\delta, p\rho), U(\delta, \rho), \dots, U(\delta, q\rho)].$$

Suppose, for all $t \in \mathcal{Z}$, $(\nabla F[t])_i \geq 0$ for all $i = 1, \dots, p + q$ and there exists $j \in \{1, \dots, p\}$ such that $(\nabla F[t])_j > 0$. Then there exists $c \in (0, 1)$ depending only on F, δ , and ρ such that for any $\epsilon \in (0, 1)$

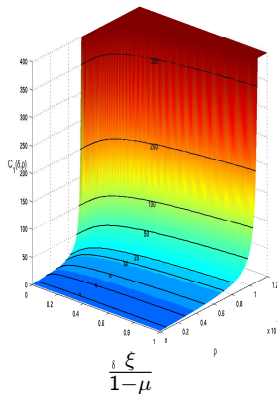
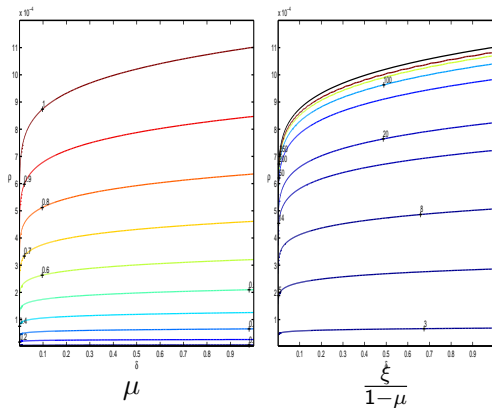
$$F[z(\delta, \rho) + \mathbf{1}c\epsilon] < F[z(\delta, (1 + \epsilon)\rho)],$$

and so there is an exponentially high probability on the draw of A that

$$\text{Prob}(F[z(k, n, N)] < F[z(\delta, (1 + \epsilon)\rho)]) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

► Can absorb the $\mathcal{O}(\epsilon)$ inside ρ component.

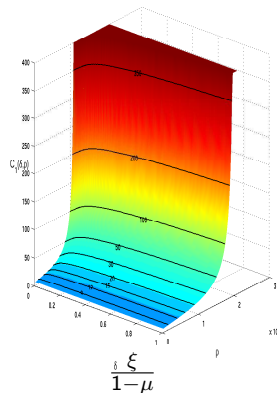
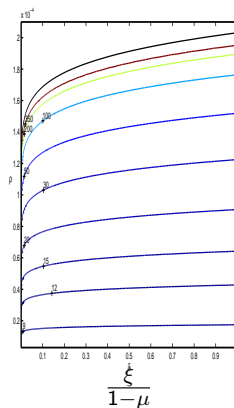
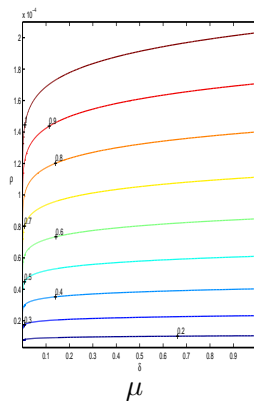
Iterated Hard Thresholding



- Success can only be guaranteed below $\mu(\delta, \rho) < 1$.

Bounding stability and complexity gives yet lower thresholds.

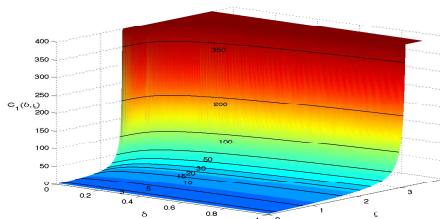
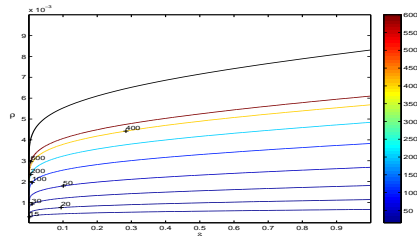
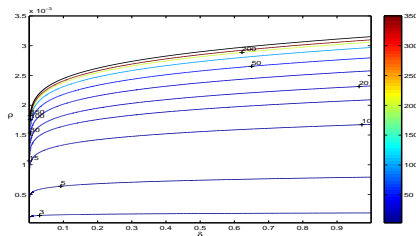
CoSaMP



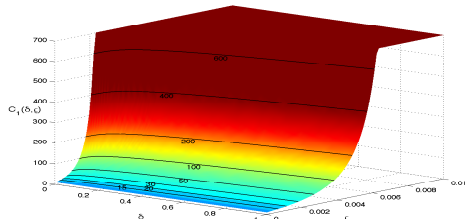
- Success can only be guaranteed below $\mu(\delta, \rho) < 1$.

Bounding stability and complexity gives yet lower thresholds.

ℓ^q -regularization, $\mu/(1 - \xi)$

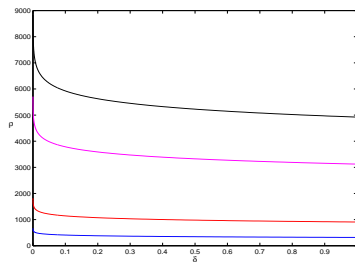
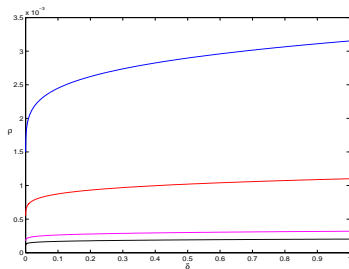


$q = 1$



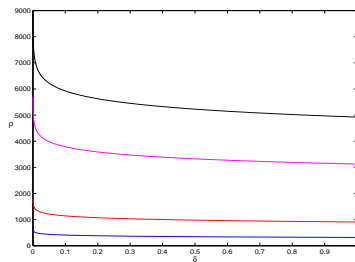
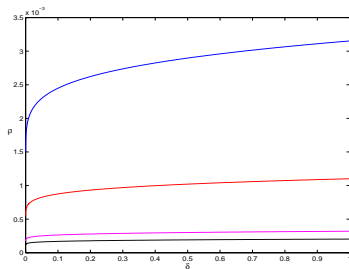
$q = 1/2$

Best known bounds implied by RIP [BlCaTaTh09]



- ▶ Lower bounds on the Strong exact recovery phase transition for Gaussian random matrices for the algorithms ℓ^1 -regularization, IHT, SP, and CoSaMP (black).
 - Unfortunately recovery thresholds are impractically low.
 $m > 317k$, $m > 907k$, $m > 3124k$, $m > 4925k$
- ▶ Coherence and RICs of structured encoders go to zero.
- ▶ Targeted techniques give more precise results, $m > 5.9k$.

Best known bounds implied by RIP, asymptotic [BaTa11]



- ▶ Lower bounds on the Strong exact recovery phase transition for Gaussian random matrices for the algorithms ℓ^1 -regularization, IHT, SP, and CoSaMP (black).
 - Asymptotic recovery condition for $m > \gamma k \log(n/m)$
 $\gamma = 36$, $\gamma = 93$, $\gamma = 272$, $\gamma = 365$
- ▶ RIP analysis of OMP yields $m > 6k^2 \log(n/k)$, seems sharp.
- ▶ Targeted techniques give more precise results, $\gamma = 2e$ for ℓ^1 .

Theory inadequate for many algorithms, experiment!

- ▶ Experimental testing of universality for $\rho_W(\delta, C)$ and $\rho_W(\delta, T)$ via embarrassingly parallel on 1400 node cluster.
- ▶ HPC specific GPUs are a major advance in computing power, c2050 1 TeraFlop/s, top UK computer in 2002 was 2TF/s
- ▶ Many core, 448 on c2050, requires careful use of parallelism

Theory inadequate for many algorithms, experiment!

- ▶ Experimental testing of universality for $\rho_W(\delta, C)$ and $\rho_W(\delta, T)$ via embarrassingly parallel on 1400 node cluster.
- ▶ HPC specific GPUs are a major advance in computing power, c2050 1 TeraFlop/s, top UK computer in 2002 was 2TF/s
- ▶ Many core, 448 on c2050, requires careful use of parallelism
- ▶ NIHT experimental setup:
 - ▶ Single precision, matrix-vector multiplication via DCT
 - ▶ Fast support set detection via linear binning
Not all bins counted in initial steps, effective k smaller initially
Avoid counting bins for small values to avoid long queues
Avoid rebinning when support set couldn't have changed.

Computing environment

CPU:

- ▶ Intel Xeon 5650 (released March 2010)
- ▶ 6 core, 2.66 GHz
- ▶ 12 GB of DDR2 PC3-1066, 6.4 GT/s
- ▶ Matlab 2010a, 64 bit (inherent multi-core threading)

GPU:

- ▶ NVIDIA Tesla c2050 (release April 2010)
- ▶ 448 Cores, peak performance 1.03 Tflop/s
- ▶ 3GB GDDR5 (on device memory)
- ▶ Error-correction

Check that NIHT still performs similarly

Sparsity k	GPU iters	CPU1 iters	CPU2 iters
512	54	54	54
1024	56	56	56
2048	60	60	60
4096	66	66	66
8192	75	75	75
16384	106	106	106
32768	382	374	377
65536	1000	1000	1000

Table: Iterations: $N = 1,048,576$ and $n = 262,144$, i.e. $\delta = .25$.

- ▶ GPU: GPU NIHT
- ▶ CPU1: same as used on GPU, but on CPU in matlab
- ▶ CPU2: standard CPU matlab implementation

Timings of NIHT

Sparsity k	GPU time (s)	CPU1 time (s)	CPU2 time (s)
512	0.6790	28.4154	25.9532
1024	0.6963	29.9603	25.1112
2048	0.7481	33.0999	28.3294
4096	0.8272	36.5142	33.7068
8192	0.9120	41.3536	37.2952
16384	1.2025	48.0748	46.3941
32768	3.4911	198.5140	183.1571
65536	9.7955	548.3572	475.3499

Table: Timings: $N = 1,048,576$ and $n = 262,144$, i.e. $\delta = .25$.

- ▶ GPU: GPU NIHT
- ▶ CPU1: same as used on GPU, but on CPU in matlab
- ▶ CPU2: standard CPU matlab implementation

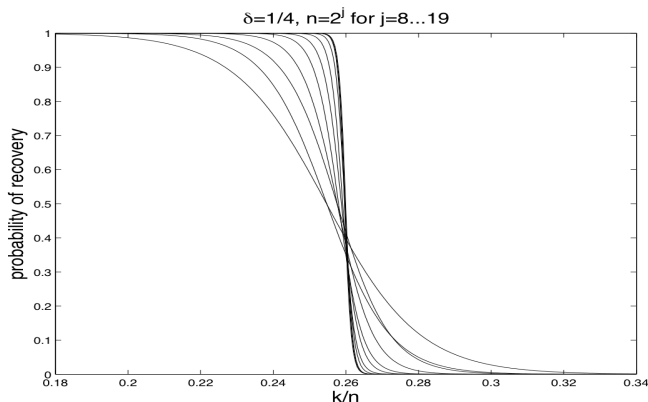
Acceleration of NIHT

Sparsity k	CPU1/GPU acceleration	CPU2/GPU acceleration
512	41.8460	38.2200
1024	43.0305	36.0660
2048	44.2444	37.8677
4096	44.1419	40.7480
8192	45.3417	40.8919
16384	39.9789	38.5813
32768	56.8625	52.4637
65536	55.9806	48.5274

Table: Acceleration: $N = 1,048,576$ and $n = 262,144$, i.e. $\delta = .25$.

- ▶ GPU workstation (4 card) equivalent to ≈ 1000 node cluster
- ▶ Computing resources allows large scale testing of algorithms
- ▶ Empirical investigation of phase transition and other properties

Empirical analysis of NIHT, $\delta = 0.25$



- ▶ Logit fit, $\frac{\exp(\beta_0 + \beta_1 k)}{1 + \exp(\beta_0 + \beta_1 k)}$, of data collected of about 10^5 tests
- ▶ $\rho_W^{niht}(1/4) \approx 0.25967$
- ▶ Transition width proportional to $m^{-1/2}$
- ▶ Can also extract iterations, time, convergence rate...

The polytope model and face survival

There are three high dimensional regular polytopes.

Each can be used to model compressed sensing questions

- ▶ Crosspolytope $C^n := \|x\|_1 \leq 1$
models ℓ^1 -regularization
- ▶ Simplex $T^{n-1} := \sum_{i=1}^n x_i \leq 1$ with $x_i \geq 0$ for all i
models ℓ^1 -regularization with sign prior
- ▶ Hypercube $H^n := \|x\|_\infty \leq 1$
models bound constraints, different notion of simplicity

The polytope model and face survival

There are three high dimensional regular polytopes.

Each can be used to model compressed sensing questions

- ▶ Crosspolytope $C^n := \|x\|_1 \leq 1$
models ℓ^1 -regularization
- ▶ Simplex $T^{n-1} := \sum_{i=1}^n x_i \leq 1$ with $x_i \geq 0$ for all i
models ℓ^1 -regularization with sign prior
- ▶ Hypercube $H^n := \|x\|_\infty \leq 1$
models bound constraints, different notion of simplicity

Lemma

F a k-face of the polytope or polyhedral cone Q and x_0 a vector in $\text{relint}(F)$. For $m \times n$ matrix A the following are equivalent:

$$\begin{aligned} (\text{Survive}(A, F, Q)): & \quad AF \text{ is a } k\text{-face of } AQ, \\ (\text{Transverse}(A, x_0, Q)): & \quad \mathcal{N}(A) \cap \text{Feas}_{x_0}(Q) = \{0\}. \end{aligned}$$

Explaining the models

- Crosspolytope $C^n := \{x \mid \|x\|_1 \leq 1\}$

If $x_0 \in \mathbb{R}^n$ is on a k -face of C^n and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(C^n) = \{0\}$ then x_0 has the minimum ℓ^1 norm and $y = Ax_0$.

Explaining the models

- ▶ Crosspolytope $C^n := \{x \mid \|x\|_1 \leq 1\}$
If $x_0 \in \mathbb{R}^n$ is on a k -face of C^n and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(C^n) = \{0\}$ then x_0 has the minimum ℓ^1 norm and $y = Ax_0$.
- ▶ Simplex $T^{n-1} := \{x \mid \sum_{i=1}^n x_i \leq 1, x_i \geq 0 \text{ for all } i\}$
If $x_0 \in \mathbb{R}_+^n$ is on a k -face of T^{n-1} and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(T^{n-1}) = \{0\}$ then x_0 has the minimum ℓ^1 norm with nonnegative prior and $y = Ax_0$.

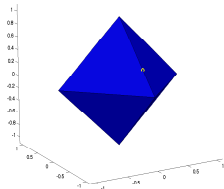
Explaining the models

- ▶ Crosspolytope $C^n := \|x\|_1$
If $x_0 \in \mathbb{R}^n$ is on a k -face of C^n and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(C^n) = \{0\}$ then x_0 has the minimum ℓ^1 norm and $y = Ax_0$.
- ▶ Simplex $T^{n-1} := \sum_{i=1}^n x_i \leq 1$ with $x_i \geq 0$ for all i
If $x_0 \in \mathbb{R}_+^n$ is on a k -face of T^{n-1} and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(T^{n-1}) = \{0\}$ then x_0 has the minimum ℓ^1 norm with nonnegative prior and $y = Ax_0$.
- ▶ Hypercube $H^n := \|x\|_\infty \leq 1$
If $x_0 \in \mathbb{R}^n$ is on a k -face of H^n and $\mathcal{N}(A) \cap \text{Feas}_{x_0}(H^n) = \{0\}$ then x_0 has the minimum ℓ^∞ norm, and is the unique vector satisfying H^n bounds and $y = Ax_0$.

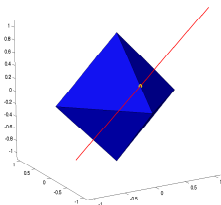
Graphical representation for ℓ^1 -regularization and Crosspolytope

Geometry of ℓ^1 -regularization, \mathbb{R}^n

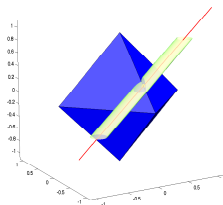
- Sparsity: $x_0 \in \mathbb{R}^n$ with $k < m$ nonzeros on $k - 1$ face of C^n .
- Null space of A intersects C^n at only x_0 , or pierces C^n



$$\ell^1 \text{ ball} \in \mathbb{R}^n$$



$$x_0 + \mathcal{N}(A)$$

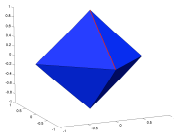


$$\|A(x - h)\| \leq \eta$$

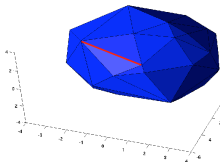
- If $\{x_0 + \mathcal{N}(A)\} \cap C^n = x_0$, ℓ^1 minimization recovers x_0
- Faces pierced by $x_0 + \mathcal{N}(A)$ do not recover k sparse x_0

Geometry of ℓ^1 -regularization, \mathbb{R}^m

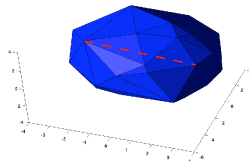
- ▶ Sparsity: $x_0 \in \mathbb{R}^n$ with $k < m$ nonzeros on $k - 1$ face of C^n .
- ▶ Matrix A projects face of ℓ^1 ball either onto or into $\text{conv}(\pm A)$.



ℓ^1 ball $\in \mathbb{R}^n$



edge onto AC^n

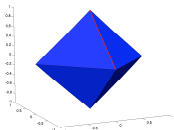


edge into AC^n

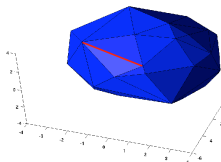
- ▶ Survived faces are sparsity patterns in x where $\ell^1 \rightarrow \ell^0$
- ▶ Faces which fall inside AC^n are not solutions to ℓ^1

Geometry of ℓ^1 -regularization, \mathbb{R}^m

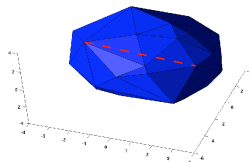
- ▶ Sparsity: $x_0 \in \mathbb{R}^n$ with $k < m$ nonzeros on $k - 1$ face of C^n .
- ▶ Matrix A projects face of ℓ^1 ball either onto or into $\text{conv}(\pm A)$.



ℓ^1 ball $\in \mathbb{R}^n$



edge onto AC^n

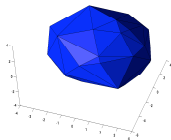


edge into AC^n

- ▶ Survived faces are sparsity patterns in x where $\ell^1 \rightarrow \ell^0$
- ▶ Faces which fall inside AC^n are not solutions to ℓ^1
- ▶ Neighborliness of random polytopes [Affentranger & Schneider]
- ▶ Exact recoverability of k sparse signals by “counting faces”

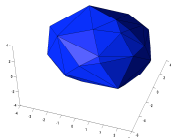
Stochastic geometry and uniform recovery

- ▶ Convex hull of n points in m dimensions
- ▶ $a_i \in \mathbb{R}^m$: $i = 1, 2, \dots, n$
- ▶ $P = \text{conv}(A)$



Stochastic geometry and uniform recovery

- ▶ Convex hull of n points in m dimensions
- ▶ $a_i \in \mathbb{R}^m$: $i = 1, 2, \dots, n$
- ▶ $P = \text{conv}(A)$
- ▶ Definition of A being k -neighborly:
 - ▶ Every a_i is a vertex of $\text{conv}(A)$
 - ▶ Every pair (a_i, a_j) span an edge of $\text{conv}(A)$
 - ▶ Every k -tuple of A span a $k - 1$ face of $\text{conv}(A)$
- ▶ Cyclic Polytopes are maximally $\lfloor m/2 \rfloor$ -neighborly, Vandermonde
- ▶ Gale (1956) suggested most polytopes are neighborly

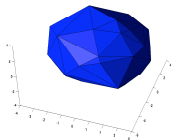


Classical Result - m fixed

- ▶ Convex hull of n points in m dimensions
- ▶ $a_i \in \mathbb{R}^m$: $i = 1, 2, \dots, n$
- ▶ $P = \text{conv}(A)$
- ▶ Classically: a_i i.i.d. Gaussian $N(0, \Sigma)$, m fixed

$$\# \text{ vert}(P) \sim c_m \log^{(m-1)/2} n, \quad n \rightarrow \infty.$$

- ▶ Not even 0-neighborly
- ▶ Renyi-Sulanke (1963), Efron (1965),
Raynaud (1971), Hueter (1998)
- ▶ Is this the typical structure for a random polytope?



Proportional growth

- ▶ Modern high-dimensional setting:
 $a_i \in \mathbb{R}^n$ iid Gaussian $N(0, \Sigma)$
 $\delta = m/n \in (0, 1)$, m and n large

Proportional growth

- ▶ Modern high-dimensional setting:
 $a_i \in \mathbb{R}^n$ iid Gaussian $N(0, \Sigma)$
 $\delta = m/n \in (0, 1)$, m and n large
- ▶ **Surprise** - neighborliness proportional to m is typical

$$\text{Prob}\{\text{conv}(A) \text{ is } k\text{-neighborly}\} \rightarrow 1, \text{ as } m, n \rightarrow \infty$$

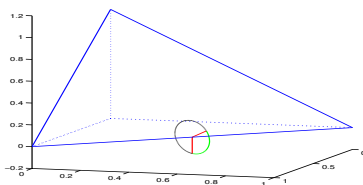
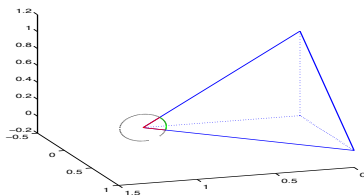
for $k < \rho_S(m/n; T) \cdot m$, [DoTa05].

- ▶ What is $\rho_S(m/n; T)$?
- ▶ Similarly for C^n (central neighborliness) and H^n (zonotope)
- ▶ For C^n known that $\rho_S(m/n; C) \leq 1/3$, unknown construction
- ▶ Nice model, but how do we calculate $\rho_S(m/n; Q)$?

Expected number of faces, random ortho-projector

$$f_k(Q) - \mathcal{E}f_k(AQ) = 2 \sum_{s \geq 0} \sum_{F \in \mathcal{F}_k(Q)} \sum_{G \in \mathcal{F}_{m+1+2s}(Q)} \beta(F, G) \gamma(G, Q)$$

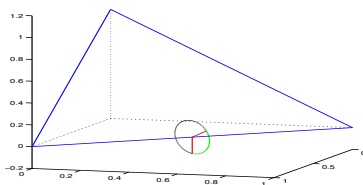
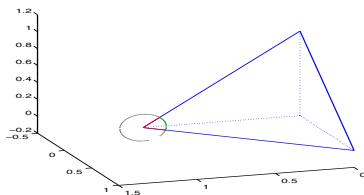
where β and γ are internal and external angles respectively
[Affentranger, Schneider]



Expected number of faces, random ortho-projector

$$f_k(Q) - \mathcal{E}f_k(AQ) = 2 \sum_{s \geq 0} \sum_{F \in \mathcal{F}_k(Q)} \sum_{G \in \mathcal{F}_{m+1+2s}(Q)} \beta(F, G) \gamma(G, Q)$$

where β and γ are internal and external angles respectively
[Affentranger, Schneider]



► Hypercube is easily to calculate angles, others less so

$$\gamma(T^\ell, T^{m-1}) = \sqrt{\frac{\ell+1}{\pi}} \int_0^\infty e^{-(\ell+1)x^2} \left(\frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy \right)^{n-\ell-1} dx.$$

Hypercube angles and face counting [DoTa08]

- ▶ Faces of H^n are all hypercubes
- ▶ $\beta(H^k, H^\ell) = 2^{-(\ell-k)}$ for all $H^k \in H^\ell$
- ▶ $\gamma(H^\ell, H^n) = 2^{-(n-\ell)}$ for all $H^\ell \in H^n$
- ▶ For a given H^k , the number of $H^k \in H^\ell \in H^n$ is $\binom{n-k}{\ell-k}$.

$$f_k(H^n) - \mathcal{E}f_k(AH^n) = 2 \sum_{s \geq 0} \sum_{F \in \mathcal{F}_k(H^n)} 2^{-(n-k)} \binom{n-k}{m+1+2s-k}$$

Hypercube angles and face counting [DoTa08]

- ▶ Faces of H^n are all hypercubes
- ▶ $\beta(H^k, H^\ell) = 2^{-(\ell-k)}$ for all $H^k \in H^\ell$
- ▶ $\gamma(H^\ell, H^n) = 2^{-(n-\ell)}$ for all $H^\ell \in H^n$
- ▶ For a given H^k , the number of $H^k \in H^\ell \in H^n$ is $\binom{n-k}{\ell-k}$.

$$f_k(H^n) - \mathcal{E}f_k(AH^n) = 2 \sum_{s \geq 0} \sum_{F \in \mathcal{F}_k(H^n)} 2^{-(n-k)} \binom{n-k}{m+1+2s-k}$$

- ▶ There are $2^{n-k} \binom{n}{k}$ different k -faces of H^n

$$f_k(H^n) - \mathcal{E}f_k(AH^n) = 2 \binom{n}{k} \sum_{s \geq 0} \binom{n-k}{m+1+2s-k}$$

- ▶ Compare $s = 0$ with $f_k(H^n) = 2^{n-k} \binom{n}{k}$, most faces survive

Hypercube weak and strong phase transitions [DoTa08]

- ▶ Weak phase transitions separate when *most* k -faces survive

$$\frac{f_k(H^n) - \mathcal{E}f_k(AH^n)}{f_k(H^n)} = 2^{-(n-k-1)} \sum_{s \geq 0} \binom{n-k}{m+1+2s-k}$$

- ▶ Main effect from $s = 0$ (bound by n times $s = 0$ factor)
- ▶ When is $2^{-(n-k)} \binom{n-k}{m-k}$ exponentially small?

Hypercube weak and strong phase transitions [DoTa08]

- ▶ Weak phase transitions separate when *most* k -faces survive

$$\frac{f_k(H^n) - \mathcal{E}f_k(AH^n)}{f_k(H^n)} = 2^{-(n-k-1)} \sum_{s \geq 0} \binom{n-k}{m+1+2s-k}$$

- ▶ Main effect from $s = 0$ (bound by n times $s = 0$ factor)
- ▶ When is $2^{-(n-k)} \binom{n-k}{m-k}$ exponentially small?
Combinatorial term largest at $m - k = \frac{n-k}{2}$, then $= 2^{(n-k)}$
- ▶ Weak phase transitions $\rho_W(\delta; H) := \max(0, 2 - \delta^{-1})$
- ▶ No strong phase transition (proof to come)
- ▶ Hypercube is sufficiently simple we can say much more, later

Large $k \sim m \sim n$ behavior of C^n and T^{n-1} angles

Exemplify through one angle, the external angle between simplices

$$\gamma(T^\ell, T^{m-1}) = \sqrt{\frac{\ell+1}{\pi}} \int_0^\infty e^{-(\ell+1)x^2} \left(\frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy \right)^{n-\ell-1} dx.$$

Define internal (dy) integral as $\Phi(x) := \frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy$, then

$$\left(\frac{2}{\sqrt{\pi}} \int_0^x e^{-y^2} dy \right)^{n-\ell-1} = \exp[(n-\ell-1) \ln(\Phi(x))]$$

Full integral then given by

$$\gamma(T^\ell, T^{m-1}) = \sqrt{\frac{\ell+1}{\pi}} \int_0^\infty e^{-(\ell+1)x^2 + (n-\ell-1) \ln(\Phi(x))} dx$$

Integrand maximized at $-2\ell x + (n-\ell)\Phi_x(x)/\Phi(x) = 0$

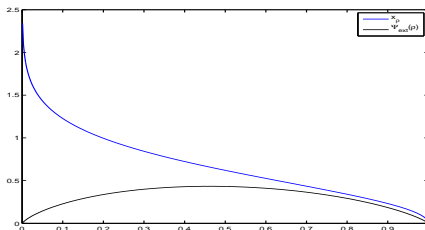
Let $\nu := \ell/n$ and x_ν satisfies $2x_\nu = (\nu^{-1} - 1)\Phi_x(x_\nu)/\Phi(x_\nu)$

Large deviation exponent of external simplex angle

- ▶ Bound using dominant exponential behavior

$$\gamma(T^\ell, T^{m-1}) = \sqrt{\frac{\ell+1}{\pi}} e^{-n[\nu x_\nu^2 + (1-\nu)\ln(\Phi(x_\nu))]} \int_0^\infty e^{-x^2 - \ln(\Phi(x))} dx$$

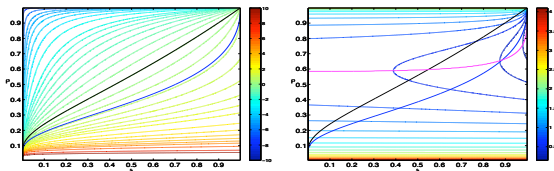
- ▶ $\sqrt{(\ell+1)/\pi}$ and remaining integral have small effect,
- ▶ Dominant effect in sum given at $\nu = \rho$
- ▶ Large deviation exponent $\Psi_{\text{ext}}(\rho) := \nu x_\nu^2 + (1-\nu)\ln(\Phi(x_\nu))$



- ▶ Other angles and combinatorial terms similarly

Probability exponents for C^n and T^{n-1}

$$f_k(Q) - \mathcal{E}f_k(AQ) = 2 \sum_{s \geq 0} \sum_{F \in \mathcal{F}_k(Q)} \sum_{G \in \mathcal{F}_{m+1+2s}(Q)} \beta(F, G) \gamma(G, Q)$$



- Strong Phase transitions (uniform bounds)

$$f_k(Q) - \mathcal{E}f_k(AQ) \leq \text{poly}(m, n) \cdot \exp(-n\Psi_{\text{net}}(\delta, \rho; Q))$$

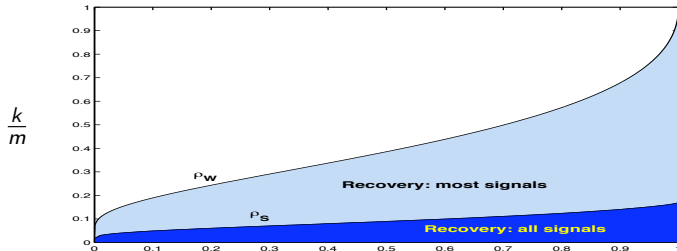
- Weak Phase transitions (average performance)

$$\frac{f_k(Q) - \mathcal{E}f_k(AQ)}{f_k(Q)} \leq \text{poly}(m, n) \cdot \exp(-n(\Psi_{\text{net}} - \Psi_{\text{face}})(\delta, \rho; Q))$$

- Widths of phase transitions: Strong m^{-1} and Weak $m^{-1/2}$

Phase Transition: ℓ^1 ball, C^n [Do05]

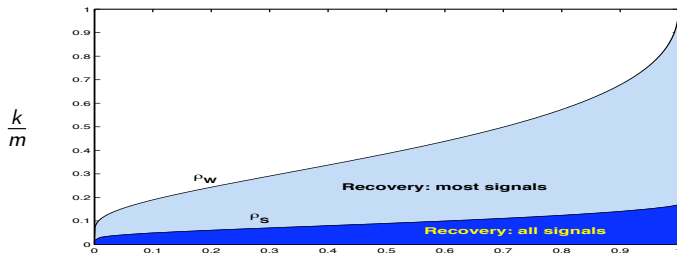
- With overwhelming probability on measurements $A_{m,n}$:
for any $\epsilon > 0$, as $(k, m, n) \rightarrow \infty$
- All k -sparse signals if $k/m \leq \rho_S(m/n, C^n)(1 - \epsilon)$
 - Most k -sparse signals if $k/m \leq \rho_W(m/n, C^n)(1 - \epsilon)$
 - Failure typical if $k/m \geq \rho_W(m/n, C^n)(1 + \epsilon)$



$$\delta = m/n$$

Phase Transition: ℓ^1 ball, C^n [Do05]

- ▶ With overwhelming probability on measurements $A_{m,n}$:
for any $\epsilon > 0$, as $(k, m, n) \rightarrow \infty$
 - All k -sparse signals if $k/m \leq \rho_S(m/n, C^n)(1 - \epsilon)$
 - Most k -sparse signals if $k/m \leq \rho_W(m/n, C^n)(1 - \epsilon)$
 - Failure typical if $k/m \geq \rho_W(m/n, C^n)(1 + \epsilon)$

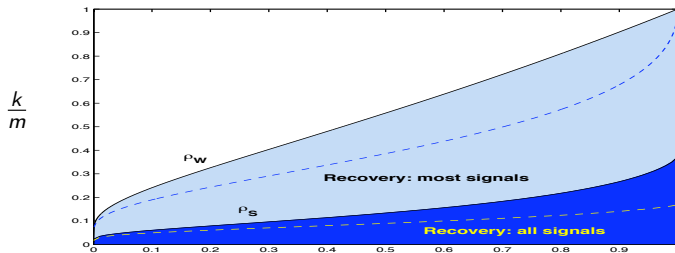


$$\delta = m/n$$

- ▶ Finite n sampling theorems proven, empirical agreement
- ▶ For $m \ll n$ requires $m > 2(e)k \cdot \log(n/m)$

Phase Transition: Simplex, T^{n-1} , $x \geq 0$ [DoTa05]

- ▶ With overwhelming probability on measurements $A_{m,n}$:
for any $\epsilon > 0$, $x \geq 0$, as $(k, m, n) \rightarrow \infty$
 - All k -sparse signals if $k/m \leq \rho_S(m/n, T^{n-1})(1 - \epsilon)$
 - Most k -sparse signals if $k/m \leq \rho_W(m/n, T^{n-1})(1 - \epsilon)$
 - Failure typical if $k/m \geq \rho_W(m/n, T^{n-1})(1 + \epsilon)$

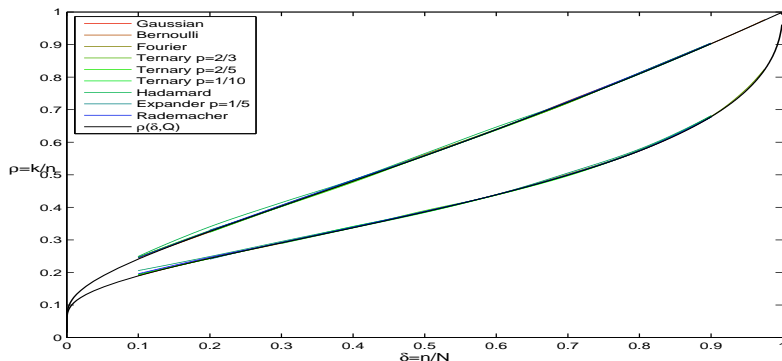


$$\delta = m/n$$

- ▶ Finite n sampling theorems proven, empirical agreement
- ▶ For $m \ll n$ requires $m > 2(e)k \cdot \log(n/m)$

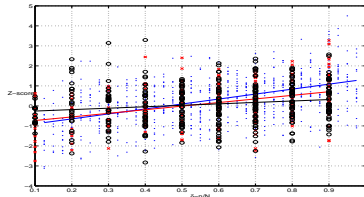
Weak Phase Transitions: Observed Universality [DoTa09]

- ▶ Black: Weak phase transition: $x \geq 0$ (top), x signed (bot.)
- ▶ Empirical evidence of 50% success rate, $n = 1600$,

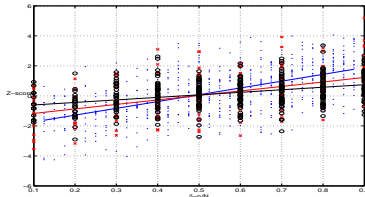


- ▶ Rigorous statistical testing of non-Gaussian vs. Gaussian
- ▶ Over 7 cpu years of data collected

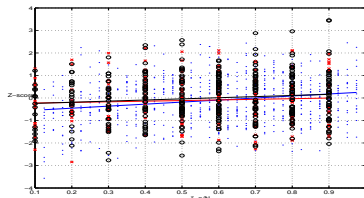
Bulk Z-scores: signed



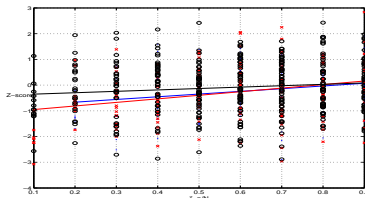
(a) Bernoulli



(b) Fourier



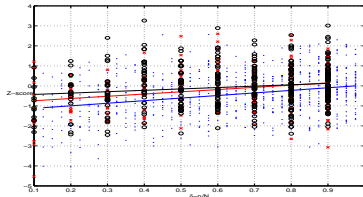
(c) Ternary (1/3)



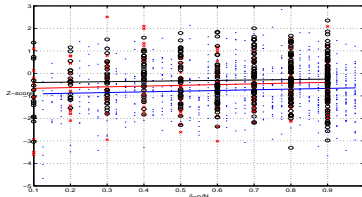
(d) Rademacher

- ▶ $n = 200$, $n = 400$ and $n = 1600$
- ▶ Linear trend with $\delta = m/n$, decays at rate $m^{-1/2}$

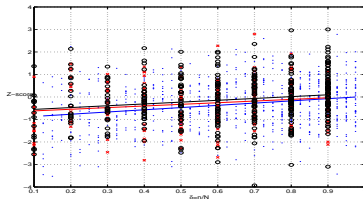
Bulk Z-scores: nonnegative



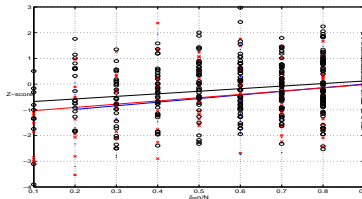
(a) Bernoulli



(b) Fourier



(c) Ternary (1/3)



(d) Rademacher

- ▶ $n = 200$, $n = 400$ and $n = 1600$
- ▶ Linear trend with $\delta = m/n$, decays at rate $m^{-1/2}$

Hypercube: universality result [DoTa09]

Theorem

Let A be an $m \times n$ matrix in general position. Then

$$f_k(AH^n) = (1 - P_{n-m, n-k})f_k(H^n)$$

where

$$P_{q,Q} = 2^{-Q+1} \sum_{\ell=0}^{q-1} \binom{Q-1}{\ell}.$$

- ▶ Universal: for every general position matrix (worse if not g.p.)
- ▶ Finite dimensional and exact

Hypercube: universality result [DoTa09]

Theorem

Let A be an $m \times n$ matrix in general position. Then

$$f_k(AH^n) = (1 - P_{n-m, n-k})f_k(H^n)$$

where

$$P_{q,Q} = 2^{-Q+1} \sum_{\ell=0}^{q-1} \binom{Q-1}{\ell}.$$

- ▶ Universal: for every general position matrix (worse if not g.p.)
- ▶ Finite dimensional and exact

Theorem (Cover & Winder)

A set of Q hyperplanes in general position in \mathbb{R}^q , all passing through a common point, divides the space into $2^Q P_{q,Q}$ regions.

Hypercube: universality result (proof)

- ▶ Consider a k -set Λ . For each k -face whose entries are not at bounds on Λ , translate (without rotation) $Feas_{F(\Lambda)}(H^n)$ so that its “spine” is at the origin.
- ▶ The union of these $Feas_{F(\Lambda)}(H^n)$ is a covering of \mathbb{R}^n with $n - k$ hyperplanes used to partition it

Hypercube: universality result (proof)

- ▶ Consider a k -set Λ . For each k -face whose entries are not at bounds on Λ , translate (without rotation) $Feas_{F(\Lambda)}(H^n)$ so that its “spine” is at the origin.
- ▶ The union of these $Feas_{F(\Lambda)}(H^n)$ is a covering of \mathbb{R}^n with $n - k$ hyperplanes used to partition it
- ▶ There are $\binom{n}{k}$ of these Λ coverings
- ▶ The $\mathcal{N}(A)$ is $n - m$ dimensional passing through the origin and is bisected by the $n - k$ planes $\binom{n}{k}$ times

Hypercube: universality result (proof)

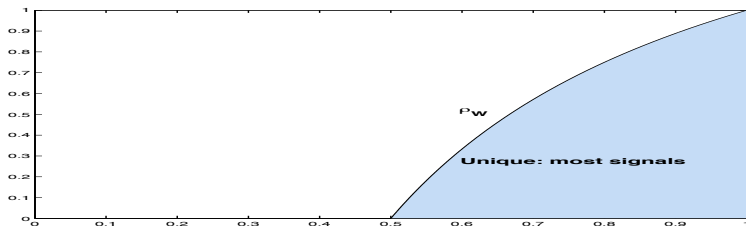
- ▶ Consider a k -set Λ . For each k -face whose entries are not at bounds on Λ , translate (without rotation) $Feas_{F(\Lambda)}(H^n)$ so that its “spine” is at the origin.
- ▶ The union of these $Feas_{F(\Lambda)}(H^n)$ is a covering of \mathbb{R}^n with $n - k$ hyperplanes used to partition it
- ▶ There are $\binom{n}{k}$ of these Λ coverings
- ▶ The $\mathcal{N}(A)$ is $n - m$ dimensional passing through the origin and is bisected by the $n - k$ planes $\binom{n}{k}$ times
- ▶ Each region of $\mathcal{N}(A)$ corresponds to a k -face where $Feas_{F(\Lambda)}(H^n) \cap \mathcal{N}(A) \neq \emptyset$, a lost k -face

$$f_k(H^n) - f_k(AH^n) = \binom{n}{k} 2^{n-k} P_{n-m, n-k} = f_k(H^n) P_{n-m, n-k}$$



Phase Transition [DoTa09]: Hypercube, H^n

- ▶ Let $-1 \leq x \leq 1$ have k entries $\neq -1, 1$ and form $y = Ax$.
- ▶ Are there other $z \in H^n[-1, 1]$ such that $Az = y$, $z \neq x$?
- ▶ As $m, n \rightarrow \infty$, Typically No provided $k/m < \rho_W(\delta; H)$



- ▶ Unlike R , T and C : no strong phase transition, $f_k(H^n)$ large
- ▶ Universal: A need only be in general position
- ▶ Simplicity beyond sparsity: Hypercube k -faces correspond to vectors with only k entries away from bounds (not -1 or 1).

Orthant: centro-symmetric result

Theorem

Let A be an $m \times n$ matrix in general position with a centro-symmetric nullspace and exchangeable columns. Then

$$\mathcal{E}f_k(A\mathbb{R}_+^n) = (1 - P_{n-m, n-k})f_k(\mathbb{R}_+^n)$$

- Similar to hypercube, but in expectation

Orthant: centro-symmetric result

Theorem

Let A be an $m \times n$ matrix in general position with a centro-symmetric nullspace and exchangeable columns. Then

$$\mathcal{E}f_k(A\mathbb{R}_+^n) = (1 - P_{n-m,n-k})f_k(\mathbb{R}_+^n)$$

- Similar to hypercube, but in expectation

Theorem (Wendel)

Let Q points in \mathbb{R}^q be drawn i.i.d. from a centro-symmetric distribution such that the points are in general position, then the probability that all the points fall in some half space is $P_{q,Q}$.

Orthant: centro-symmetric result (proof)

- ▶ Let $x \geq 0$ have k entries $x_i > 0$ and form $y = Ax$; $x_{\Lambda^c} > 0$ for $|\Lambda^c| = k$, $x_{\Lambda} = 0$.

Orthant: centro-symmetric result (proof)

- ▶ Let $x \geq 0$ have k entries $x_i > 0$ and form $y = Ax$; $x_{\Lambda^c} > 0$ for $|\Lambda^c| = k$, $x_{\Lambda} = 0$.
- ▶ Not unique if $\exists z \in \mathcal{N}(A)$ with $z_{\Lambda} \geq 0$

Orthant: centro-symmetric result (proof)

- ▶ Let $x \geq 0$ have k entries $x_i > 0$ and form $y = Ax$; $x_{\Lambda^c} > 0$ for $|\Lambda^c| = k$, $x_{\Lambda} = 0$.
- ▶ Not unique if $\exists z \in \mathcal{N}(A)$ with $z_{\Lambda} \geq 0$
- ▶ Let $B \in \mathbb{R}^{n, n-k}$ be a basis for $\mathcal{N}(A)$, then $z = Bc$ for some c .
- ▶ Not unique if $(B^T c)_{\Lambda} \geq 0$ where $|\Lambda| = n - k$.

Orthant: centro-symmetric result (proof)

- ▶ Let $x \geq 0$ have k entries $x_i > 0$ and form $y = Ax$; $x_{\Lambda^c} > 0$ for $|\Lambda^c| = k$, $x_{\Lambda} = 0$.
- ▶ Not unique if $\exists z \in \mathcal{N}(A)$ with $z_{\Lambda} \geq 0$
- ▶ Let $B \in \mathbb{R}^{n, n-k}$ be a basis for $\mathcal{N}(A)$, then $z = Bc$ for some c .
- ▶ Not unique if $(B^T c)_{\Lambda} \geq 0$ where $|\Lambda| = n - k$.
- ▶ Geometrically, not unique if $n - k$ row of B^T fall in some half-space of \mathbb{R}^{n-m} .

Orthant: centro-symmetric result (proof)

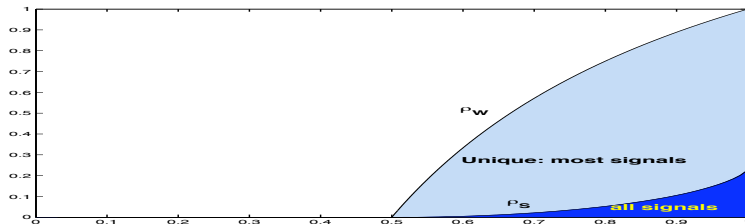
- ▶ Let $x \geq 0$ have k entries $x_i > 0$ and form $y = Ax$; $x_{\Lambda^c} > 0$ for $|\Lambda^c| = k$, $x_{\Lambda} = 0$.
- ▶ Not unique if $\exists z \in \mathcal{N}(A)$ with $z_{\Lambda} \geq 0$
- ▶ Let $B \in \mathbb{R}^{n,n-k}$ be a basis for $\mathcal{N}(A)$, then $z = Bc$ for some c .
- ▶ Not unique if $(B^T c)_{\Lambda} \geq 0$ where $|\Lambda| = n - k$.
- ▶ Geometrically, not unique if $n - k$ row of B^T fall in some half-space of \mathbb{R}^{n-m} .
- ▶ For rows of B drawn iid from centro-symmetric, row exchangeable, in general position: [Wendel's Theorem](#)
- ▶ Probability of failure is

$$2^{-n+k+1} \sum_{\ell=0}^{n-m-1} \binom{n-k-1}{\ell}$$

- ▶ Probability of failure $\rightarrow 0$ if $n - m - 1 < (n - k - 1)/2$.

Projected Orthant [DoTa09]

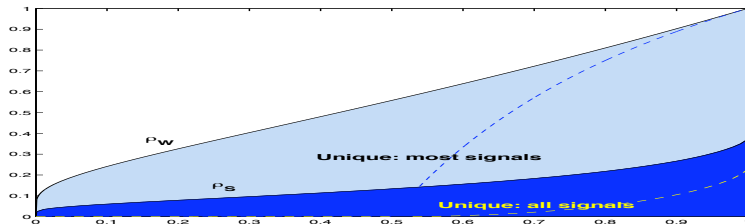
- ▶ Let $x \geq 0$ be k -sparse and form $y = Ax$.
- ▶ Are there other $z \in \mathbb{R}^n$ such that $Az = y$, $z \geq 0$, $z \neq x$?
- ▶ As $m, n \rightarrow \infty$, Typically No provided $k/m < \rho_W(\delta; \mathbb{R}_+)$



- ▶ Universal: A an ortho-complement of $B \in \mathbb{R}^{n-m \times n}$ with entries selected i.i.d. from a symmetric distribution
- ▶ For $k/m < \rho_W(\delta, H^n) := [2 - 1/\delta]_+$ and $x \geq 0$, any “feasible” method will work.

Projected Orthant, matrix design [DoTa09]

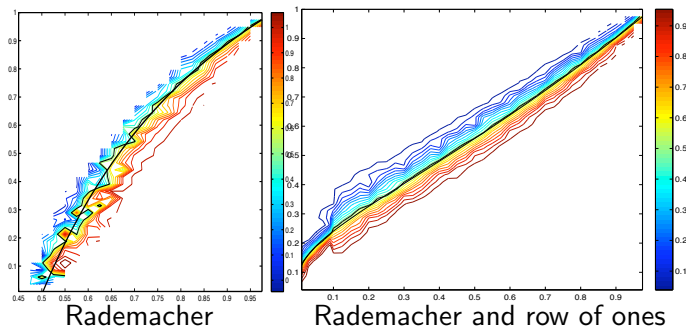
- ▶ Let $x \geq 0$ be k -sparse and form $y = Ax$.
- ▶ Are there other $z \in \mathbb{R}^n$ such that $Az = y$, $z \geq 0$, $z \neq x$?
- ▶ As $m, n \rightarrow \infty$, Typically No provided $k/m < \rho_W(\delta; \mathbb{R}_+)$



- ▶ Gaussian and measuring the mean (row of ones):
 $\rho_W(m/n; \mathbb{R}_+) \rightarrow \rho_W(m/n; T)$
- ▶ Simple modification of A makes profound difference
Unique even for $m/n \rightarrow 0$ with $m > 2(e)k \log(n/m)$

Orthant matrix design, it's really true

- ▶ Let $x \geq 0$ be k -sparse and form $y = Ax$.
- ▶ Not ℓ^1 , but: $\max_y \|x - z\|$ subject to $Az = Ax$ and $z \geq 0$
- ▶ Good empirical agreement for $n = 200$.



Simplicity as low rank

- ▶ Sparse approximation considers sparsity or bound simplicity
- ▶ Matrix completion considers low rank simplicity
- ▶ Main innovation isn't low rank simplicity, but unknown space

Simplicity as low rank

- ▶ Sparse approximation considers sparsity or bound simplicity
- ▶ Matrix completion considers low rank simplicity
- ▶ Main innovation isn't low rank simplicity, but unknown space
- ▶ Matrices that have low rank representation in a known basis

Definition. A matrix M has a k -sparse representation in the matrix dictionary $\Psi := \{\Psi_j\}_{j=1}^n$ if

$$M = \sum_j x_0(j) \Psi_j \quad \text{with} \quad \|\mathbf{x}_0\|_{\ell^0} = k.$$

- ▶ How should we sense M ?

Simplicity as low rank

- ▶ Sparse approximation considers sparsity or bound simplicity
- ▶ Matrix completion considers low rank simplicity
- ▶ Main innovation isn't low rank simplicity, but unknown space
- ▶ Matrices that have low rank representation in a known basis

Definition. A matrix M has a k -sparse representation in the matrix dictionary $\Psi := \{\Psi_j\}_{j=1}^n$ if

$$M = \sum_j x_0(j) \Psi_j \quad \text{with} \quad \|\mathbf{x}_0\|_{\ell^0} = k.$$

- ▶ How should we sense M ?
- ▶ Let M model a channel and h a known “pilot vector”
- ▶ Sense channel M by sending h , recover M from Mh and h

Matrices with known sparse representation [PfRaTa08]

- Sense channel M by sending h , recover M from Mh and h

$$\begin{aligned}Mh &= \left(\sum_{j=1}^N x_0(j) \Psi_j \right) h = \sum_{j=1}^n x_0(j) (\Psi_j h) \\ &= (\Psi_1 h \mid \Psi_2 h \mid \dots \mid \Psi_n h) x =: (\Psi h) x_0\end{aligned}$$

where $(\Psi h) = (\Psi_1 h \mid \Psi_2 h \mid \dots \mid \Psi_n h)$.

- Let $y = Mh$ and $A = \Psi$ and we are back to usual CS

Matrices with known sparse representation [PfRaTa08]

- ▶ Sense channel M by sending h , recover M from Mh and h

$$\begin{aligned}Mh &= \left(\sum_{j=1}^N x_0(j) \Psi_j \right) h = \sum_{j=1}^n x_0(j) (\Psi_j h) \\ &= (\Psi_1 h \mid \Psi_2 h \mid \dots \mid \Psi_n h) x =: (\Psi h) x_0\end{aligned}$$

where $(\Psi h) = (\Psi_1 h \mid \Psi_2 h \mid \dots \mid \Psi_n h)$.

- ▶ Let $y = Mh$ and $A = \Psi h$ and we are back to usual CS
- ▶ Exemplar applications: wireless communication and sonar
Seeking channel for detection or repair channel corruption
Model channel as a few dominant translations (delays) and modulations (reflections/dopler) and let $\Psi_n h$ be Gabor

Matrix completion oracle recovery

- ▶ Sensing of matrices $M \in \mathbb{R}^{n_1 \times n_2}$ with $\text{rank}(M) = r$.
- ▶ Then $M = U\Sigma V^T$ for $U \in \mathbb{R}^{n_1 \times r}$ with orthonormal columns, $\Sigma \in \mathbb{R}^{r \times r}$ diagonal, and $V \in \mathbb{R}^{n_2 \times r}$ with orthonormal columns

Matrix completion oracle recovery

- ▶ Sensing of matrices $M \in \mathbb{R}^{n_1 \times n_2}$ with $\text{rank}(M) = r$.
- ▶ Then $M = U\Sigma V^T$ for $U \in \mathbb{R}^{n_1 \times r}$ with orthonormal columns, $\Sigma \in \mathbb{R}^{r \times r}$ diagonal, and $V \in \mathbb{R}^{n_2 \times r}$ with orthonormal columns
- ▶ What is the dimensionality of a rank r matrix?
There are $n_1 r + n_2 r + r$ values in U , Σ , and V
Orthogonality of columns in U and V impose $r^2 + r$ constraints
- ▶ Dimensionality of rank r matrices is $r(n_1 + n_2 - r)$, not $n_1 n_2$
- ▶ If $r(n_1 + n_2 - r) \ll n_1 n_2$ then maybe can exploit low dimensionality for a form of compressed sensing
- ▶ Need at least $m \geq \min(r(n_1 + n_2 - r), n_1 n_2)$ measurements
- ▶ How can we sense and recover matrices with optimal order?

Sensing in matrix completion

- ▶ Sensing in compressed sensing via inner products (vectors):
 - good idea - vectors that do not have sparse representation in the same basis as the vector being sensed
 - bad idea - point sensing a k -sparse vector

Sensing in matrix completion

- ▶ Sensing in compressed sensing via inner products (vectors):
 - good idea - vectors that do not have sparse representation in the same basis as the vector being sensed
 - bad idea - point sensing a k -sparse vector
- ▶ Matrix completion is no different, inner products (matrices):
 - good idea - matrices that do not have low rank representation in the same U and V column and row space
 - bad idea - point sensing a matrix that is sparse, low rank in point entries
- ▶ Designate \mathcal{A} the linear sensing operator from $\mathbb{R}^{n_1 \times n_2} \rightarrow \mathbb{R}^m$
- ▶ Measurements $y = \mathcal{A}(M)$ where $y_p = \sum_{i,j} A(p)_{i,j} M_{i,j}$
- ▶ Standard choices for $A(p)$: point sensing via one nonzero or dense sensing via i.i.d. centro-symmetric distribution

Algorithms for matrix completion

- ▶ Given \mathcal{A} and $y = \mathcal{A}(M_0)$, how to recover M_0
- ▶ Convex relaxation notion
 - Compressed sensing replaced $\min \|x\|_0$ s.t. $y = Ax$ with smallest convex relaxation $\min \|x\|_1$ s.t. $y = Ax$.
 - Matrix completion uses the obvious same replacement of $\min \text{rank}(M)$ s.t. $y = \mathcal{A}(M)$ with smallest convex relaxation $\min \|M\|_*$ s.t. $y = \mathcal{A}(M)$.

Algorithms for matrix completion

- ▶ Given \mathcal{A} and $y = \mathcal{A}(M_0)$, how to recover M_0
- ▶ Convex relaxation notion
 - Compressed sensing replaced $\min \|x\|_0$ s.t. $y = Ax$ with smallest convex relaxation $\min \|x\|_1$ s.t. $y = Ax$.
 - Matrix completion uses the obvious same replacement of $\min \text{rank}(M)$ s.t. $y = \mathcal{A}(M)$ with smallest convex relaxation $\min \|M\|_*$ s.t. $y = \mathcal{A}(M)$.
- ▶ Iterative hard thresholding
 - CS used steepest descent on $\|y - Ax\|_2$, restrict $\|x\|_0 = k$
 - Matrix completion uses steepest descent on $\|y - \mathcal{A}(M)\|_F$ then restrict to $\text{rank}(M) = r$
- ▶ Any of the algorithmic ideas from CS can be extended to Matrix completion using the obvious related property

Analysis of matrix completion algorithms: coherence

- Coherence μ if all three satisfied [Candés and Recht]
(let $n := \max(n_1, n_2)$)

$$\max_i \sum_{j=1}^r U_{i,j}^2 \leq \mu \frac{r}{n}$$

$$\max_i \sum_{j=1}^r V_{i,j}^2 \leq \mu \frac{r}{n}$$

$$\max_{i,k} \left| \sum_{j=1}^r U_{i,j} V_{k,j} \right| \leq \mu \frac{r}{n}$$

Theorem

If given $p \geq c \cdot \mu r n^{6/5} \log n$ entries of M then with high probability on M , nuclear (Schatten) norm recovers M .

Analysis of matrix completion algorithms: RICs

- ▶ Matrix completion version of RICs, for all M with $\text{rank}(M) = r$

$$(1 - R_r(\mathcal{A}))\|M\|_F \leq \|\mathcal{A}(M)\|_2 \leq (1 + R_r(\mathcal{A}))\|M\|_F$$

Analysis of matrix completion algorithms: RICs

- ▶ Matrix completion version of RICs, for all M with $\text{rank}(M) = r$

$$(1 - R_r(\mathcal{A}))\|M\|_F \leq \|\mathcal{A}(M)\|_2 \leq (1 + R_r(\mathcal{A}))\|M\|_F$$

Theorem

Let $\text{rank}(M_0) \leq r$, $y = \mathcal{A}(M_0)$, and $R_{2r}(\mathcal{A}) < 1$, then M_0 is the matrix of minimum rank satisfying $y = \mathcal{A}(M)$, and is the minimizer of the minimum rank decoder.

Analysis of matrix completion algorithms: RICs

- ▶ Matrix completion version of RICs, for all M with $\text{rank}(M) = r$

$$(1 - R_r(\mathcal{A}))\|M\|_F \leq \|\mathcal{A}(M)\|_2 \leq (1 + R_r(\mathcal{A}))\|M\|_F$$

Theorem

Let $\text{rank}(M_0) \leq r$, $y = \mathcal{A}(M_0)$, and $R_{2r}(\mathcal{A}) < 1$, then M_0 is the matrix of minimum rank satisfying $y = \mathcal{A}(M)$, and is the minimizer of the minimum rank decoder.

Theorem (Recht, Fazel, Parrilo)

Let $\text{rank}(M_0) = r$, $y = \mathcal{A}(M_0)$, and $R_{5r}(\mathcal{A}) < 1/10$, then

$$M_0 = \underset{M}{\operatorname{argmin}} \|M\|_* \quad \text{subject to} \quad y = \mathcal{A}(M).$$

Nuclear norm recovery guarantee via RIC (proof, pg. 1)

- ▶ Follow proof for ℓ^1 -regularization, but with matrices
- ▶ Need to decompose null-space matrix

Lemma

Let A and B be matrices with the same dimensions. There exist matrices B_1 and B_2 with $B = B_1 + B_2$, $AB_2^ = 0$ and $A^*B_2 = 0$, and $\langle B_1, B_2 \rangle = 0$, and $\text{rank}(B_1) \leq 2\text{rank}(A)$.*

Proof. Let A have a full SVD $A = U\Sigma V^*$. Let $\hat{B} = U^*BV$ and partition it into blocks

$$\hat{B} = \begin{bmatrix} \hat{B}_{11} & \hat{B}_{12} \\ \hat{B}_{21} & \hat{B}_{22} \end{bmatrix}$$

with \hat{B}_{11} square of size $\text{rank}(A)$, then

$$B_1 := U \begin{bmatrix} \hat{B}_{11} & \hat{B}_{12} \\ \hat{B}_{21} & 0 \end{bmatrix} V^* \quad \text{and} \quad B_2 := U \begin{bmatrix} 0 & 0 \\ 0 & \hat{B}_{22} \end{bmatrix} V^*$$

Nuclear norm recovery guarantee via RIC (proof, pg. 2)

Let $X^* = \operatorname{argmin}_X \|X\|_*$ s.t. $y = \mathcal{A}(X)$

Let $R = X^* - X_0$. By X^* being the argmin:

$$\|X_0\|_* \geq \|X_0 + R\|_* \geq \|X_0 + R_c\|_* - \|R_0\|_* = \|X_0\|_* + \|R_c\|_* - \|R_0\|_*$$

which yields $\|R_0\|_* \geq \|R_c\|_*$. (analogous to NSP)

Partition R_c into matrices of rank R_1, R_2, \dots , with R_1 having the largest $3r$ singular values of R_c , R_2 the next largest $3r$ singular values...

Compare largest singular value in set $i + 1$ with average in set i

$$\max(\sigma(R_{i+1})) \leq \frac{1}{3r} \sum \sigma(R_i) \implies \|R_{i+1}\|_F^2 \leq \frac{1}{3r} \|R_i\|_*^2.$$

Nuclear norm recovery guarantee via RIC (proof, pg. 3)

Use norm relations and $\text{rank}(R_0) \leq 2r$ to derive bound

$$\sum_{j \geq 2} \|R_j\|_F \leq \frac{1}{\sqrt{3r}} \sum_{j \geq 1} \|R_j\|_* = \frac{1}{\sqrt{3r}} \|R_c\|_* \leq \frac{1}{\sqrt{3r}} \|R_0\|_* \leq \frac{\sqrt{2r}}{\sqrt{3r}} \|R_0\|_F$$

Use RICs with bound from below

$$\begin{aligned} \|\mathcal{A}(R)\|_2 &\geq \|\mathcal{A}(R_0 + R_1)\|_2 - \sum_{j \geq 2} \|\mathcal{A}(R_j)\| \\ &\geq (1 - R_{5r}) \|R_0 + R_1\|_F - (1 + R_{3r}) \sum_{j \geq 2} \|R_j\|_F \\ &\geq ((1 - R_{5r}) - \sqrt{\frac{2}{3}}(1 + R_{3r})) \|R_0\|_F \end{aligned} \quad (5)$$

If the factor multiplying $\|R_0\|_F$ is positive, and by construction $\mathcal{A}(R) = 0$ we must have $R = 0$. □

Matrices having bounded RICs

Concentration of measure bounds analogous to before

If the entries in \mathcal{A} is a map from $\mathbb{R}^{n_1 \times n_1} \rightarrow \mathbb{R}^p$ with entries drawn from a distribution that is mean zero and has a finite fourth moment then for all $0 < \epsilon < 1$

$$\text{Prob}(|\|\mathcal{A}(M)\|_2^2 - \|M\|_F^2| \geq \epsilon \|M\|_F^2) \leq 2 \exp(-p(\epsilon^2/2 - \epsilon^3/3)/2).$$

Theorem. If \mathcal{A} is a near isometry, then for every $1 \leq r \leq m$, there exists constants c such that with exponentially high probability the RICs remain bounded whenever $p \geq cr(m+n) \log(mn)$.

- ▶ The story of matrix completion parallels that of compressed sensing, but with fewer quantitative statements and more open problems.
- ▶ There is also a “polytope” style analysis for the convex relaxation, nuclear norm, algorithm for matrix completion

Structured sparsity

- ▶ Standard CS model: $\Sigma_k(n) := \{x \mid \|x\|_0 \leq k\}$
The union of $\binom{n}{k}$ subspaces.
- ▶ A reasonable model due to the prevalence of compressibility

Structured sparsity

- ▶ Standard CS model: $\Sigma_k(n) := \{x \mid \|x\|_0 \leq k\}$
The union of $\binom{n}{k}$ subspaces.
- ▶ A reasonable model due to the prevalence of compressibility
- ▶ Wavelet transforms convert piecewise smooth signal to coefficients that decay at rate, j^{th} coefficient $\sim j^{-p}$ or τ^{-j}
- ▶ Decay of wavelet coefficients indicate k largest coefficients gives faithful approximation

Structured sparsity

- ▶ Standard CS model: $\Sigma_k(n) := \{x \mid \|x\|_0 \leq k\}$
The union of $\binom{n}{k}$ subspaces.
- ▶ A reasonable model due to the prevalence of compressibility
- ▶ Wavelet transforms convert piecewise smooth signal to coefficients that decay at rate, j^{th} coefficient $\sim j^{-p}$ or τ^{-j}
- ▶ Decay of wavelet coefficients indicate k largest coefficients gives faithful approximation
- ▶ Randomly permute where the k largest coefficients, compute inverse wavelet transform, looks like noise not piecewise smooth.
- ▶ k term wavelet approximation to find has structure.
Not all $\binom{n}{k}$ of the k -sparse vectors likely, don't look for them.

Structured sparsity model and RIC [BaCeDuHe08]

- Need not specify structure of the model yet, just number p
- Definition. [Model sparsity]** For any p distinct support sets Λ_j with $|\Lambda_j| = k \ \forall j$, let

$$\mathcal{M}_k := \{x \mid \text{supp}(x) \in \Lambda_j \text{ for some } j\}.$$

We refer to \mathcal{M}_k as a model based sparsity space.

Structured sparsity model and RIC [BaCeDuHe08]

- Need not specify structure of the model yet, just number p
- Definition. [Model sparsity]** For any p distinct support sets Λ_j with $|\Lambda_j| = k \ \forall j$, let

$$\mathcal{M}_k := \{x \mid \text{supp}(x) \in \Lambda_j \text{ for some } j\}.$$

We refer to \mathcal{M}_k as a model based sparsity space.

- Need a method of analysis, no gain for coherence, use RICs
- Definition. [Model RICs]** Given matrix $A \in \mathbb{R}^{m \times n}$, let $R_{\mathcal{M}_k}$ be the smallest constant that satisfies

$$(1 - R_{\mathcal{M}_k})\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + R_{\mathcal{M}_k})\|x\|_2^2 \quad \forall x \in \mathcal{M}_k$$

Structured sparsity model and RIC [BaCeDuHe08]

- ▶ Need not specify structure of the model yet, just number p
- Definition. [Model sparsity]** For any p distinct support sets Λ_j with $|\Lambda_j| = k \quad \forall j$, let

$$\mathcal{M}_k := \{x \mid \text{supp}(x) \in \Lambda_j \text{ for some } j\}.$$

We refer to \mathcal{M}_k as a model based sparsity space.

- ▶ Need a method of analysis, no gain for coherence, use RICs
- Definition. [Model RICs]** Given matrix $A \in \mathbb{R}^{m \times n}$, let $R_{\mathcal{M}_k}$ be the smallest constant that satisfies

$$(1 - R_{\mathcal{M}_k})\|x\|_2^2 \leq \|Ax\|_2^2 \leq (1 + R_{\mathcal{M}_k})\|x\|_2^2 \quad \forall x \in \mathcal{M}_k$$

- ▶ Letting $p = \binom{n}{k}$ recovers the usual $\Sigma_k(n)$ and RICs
- ▶ Could improve results if $R_{\mathcal{M}_k}$ replaced with asymmetric
- ▶ The essential improvement: select p such that $p \sim e^{\alpha \cdot k}$ without any n dependence. (Normal case $p \sim e^{n \cdot H(k/n)}$.)

A subtle point

- ▶ If $x, y \in \Sigma_k(n)$ then $x + y \in \Sigma_{2k}(n)$, use RICs of order $2k$
 - ▶ If $x, y \in \mathcal{M}_k$ then $x + y$ may be $2k$ sparse, but model changes
- Definition. [Union of model sparsity]** Define \mathcal{M}_k^r as the union of r possibly different model sparsity sets:

$$\mathcal{M}_k^r := \left\{ x \mid \sum_{\ell=1}^r x^{(\ell)} \quad \text{where} \quad x^{(\ell)} \in \mathcal{M}_k \right\}.$$

- ▶ \mathcal{M}_k^r can be thought of as \mathcal{M}_{rk} with p modified to $\sim p^r$

A subtle point

- ▶ If $x, y \in \Sigma_k(n)$ then $x + y \in \Sigma_{2k}(n)$, use RICs of order $2k$
 - ▶ If $x, y \in \mathcal{M}_k$ then $x + y$ may be $2k$ sparse, but model changes
- Definition.** [Union of model sparsity] Define \mathcal{M}_k^r as the union of r possibly different model sparsity sets:

$$\mathcal{M}_k^r := \{x \mid \sum_{\ell=1}^r x^{(\ell)} \quad \text{where} \quad x^{(\ell)} \in \mathcal{M}_k\}.$$

- ▶ \mathcal{M}_k^r can be thought of as \mathcal{M}_{rk} with p modified to $\sim p^r$
- ▶ Same algorithms work, with restriction to \mathcal{M}_k at each step
- ▶ Analysis for IHT as example

Iterative Hard Model Thresholding

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set $x^0 = 0$ and $j = 0$.

While $\|y - A_{m,n}x^j\|_2 < Tol$ repeat the following steps:

set $v^j := x^j + A_{m,n}^*(y - A_{m,n}x^j)$, and

$x^{j+1} = H_k(v^j)$ where $H_k(\cdot)$ thresholds to best \mathcal{M}_k .

Output x^j .

Iterative Hard Model Thresholding

Input: y , $A_{m,n}$ and k (number of nonzeros in output vector).

Algorithm: Set $x^0 = 0$ and $j = 0$.

While $\|y - A_{m,n}x^j\|_2 < Tol$ repeat the following steps:

set $v^j := x^j + A_{m,n}^*(y - A_{m,n}x^j)$, and

$x^{j+1} = H_k(v^j)$ where $H_k(\cdot)$ thresholds to best \mathcal{M}_k .

Output x^j .

Theorem

Let $y = A_{m,n}x_0 + e$ for $x_0 \in \mathcal{M}_k$ and $A_{m,n}$ in General Position.

Set $\mu^{iht} := 2R_{\mathcal{M}_k^3}$ and $\xi^{iht} := 2(1 + R_{\mathcal{M}_k^2})^{1/2}$.

With k used for the hard thresholding function, IHT satisfy the inequality

$$\|x^j - x_0\|_2 \leq (\mu^{iht})^j \|x_0\| + \frac{\xi^{iht}}{1 - \mu^{iht}} \|e\|_2.$$

For $\mu^{iht} < 1$ convergence of x^j to approximation of x_0 .

Iterative Hard Model Thresholding (proof, pg. 1)

Proof.

$H_k(\cdot)$ returns the vector in \mathcal{M}_k closest in the ℓ^2 norm, for instance

$$\|v^j - H_k(v^j)\|_2 = \|v^j - x^{j+1}\|_2 \leq \|v^j - x_0\|_2. \quad (6)$$

Note that

$$\begin{aligned} \|v^j - x^{j+1}\|_2^2 &= \|(v^j - x_0) + (x_0 - x^{j+1})\|_2^2 = \\ &= \|v^j - x_0\|_2^2 + \|x_0 - x^{j+1}\|_2^2 + 2\operatorname{Re}((v^j - x_0)^*(x_0 - x^{j+1})) \end{aligned}$$

where $\operatorname{Re}(c)$ denotes the real part of c .

Bounding the above expression using (6) and canceling the $\|v^j - x_0\|_2^2$ term yields

$$\|x^{j+1} - x_0\|_x^2 \leq 2\operatorname{Re}((v^j - x_0)^*(x^{j+1} - x_0)).$$

Iterative Hard Model Thresholding (proof, pg. 2)

Consider \mathcal{M}_k^3 model set from joining models for x_0 , x_j and x_{j+1} .

$$\begin{aligned}\|x^{j+1} - x_0\|_2^2 &\leq 2\operatorname{Re}((x^j - x_0)^*(x^{j+1} - x_0)) \\&= 2\operatorname{Re}\left(((I - A_{m,n}^* A_{m,n})(x^j - x_0))^*(x^{j+1} - x_0)\right) \\&\quad + 2\operatorname{Re}(e^* A_{m,n}(x^{j+1} - x_0)) \\&= 2\operatorname{Re}\left(\left((I - A_{\mathcal{M}_k^3}^* A_{\mathcal{M}_k^3})(x^j - x_0)\right)^*(x^{j+1} - x_0)\right) \\&\quad + 2\operatorname{Re}(e^* A_{m,n}(x^{j+1} - x_0)) \\&\leq 2\|I - A_{\mathcal{M}_k^3}^* A_{\mathcal{M}_k^3}\|_2 \cdot \|x^j - x_0\|_2 \cdot \|x^{j+1} - x_0\|_2 \\&\quad + 2\|e\|_2 \cdot \|A_{m,n}(x^{j+1} - x_0)\|_2\end{aligned}$$

Iterative Hard Model Thresholding (proof, pg. 3)

Model RIC bounds $\|I - A_{\mathcal{M}_k^3}^* A_{\mathcal{M}_k^3}\|_2 \leq R_{\mathcal{M}_k^3}$ and
 $\|A_{m,n}(x^{j+1} - x_0)\|_2 \leq (1 + R_{\mathcal{M}_k^2})^{1/2} \|x^{j+1} - x_0\|_2$ then
dividing by $\|x^{j+1} - x_0\|_2$ yields

$$\|x^{j+1} - x_0\|_2 \leq 2R_{\mathcal{M}_k^3} \cdot \|x^j - x_0\|_2 + 2(1 + R_{\mathcal{M}_k^2})^{1/2} \|e\|_2$$

Let $\mu^{iht} := 2R_{\mathcal{M}_k^3}$ and $\xi^{iht} := 2(1 + R_{\mathcal{M}_k^2})^{1/2}$.

Error at step j in terms of initial error $\|x^0 - x_0\|_2 = \|x_0\|_2$

$$\|x^j - x_0\|_2 \leq (\mu^{iht})^j \cdot \|x_0\|_2 + \frac{\xi^{iht}}{1 - \mu^{iht}} \|e\|_2$$

□

- It looks like nothing has changed, but when is $R_{\mathcal{M}_k^3} < 1/2$?

Wavelet tree model

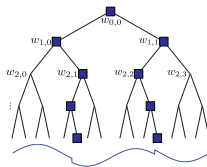
- ▶ Piecewise smooth functions (and their vector samples) is likely the most encompassing versatile model for signals and images.
- ▶ Wavelets have (rapid) polynomial decay in smooth regions, only lack decay for wavelets that interact with discontinuity.

Wavelet tree model

- ▶ Piecewise smooth functions (and their vector samples) is likely the most encompassing versatile model for signals and images.
- ▶ Wavelets have (rapid) polynomial decay in smooth regions, only lack decay for wavelets that interact with discontinuity.
- ▶ Convention of narrower wavelets as scale (label i) increases, coef. (i, j) large suggests coefficient $(i - 1, \lfloor j/2 \rfloor)$ also large.
- ▶ Connected subtree model: if (i, j) coefficient is kept, then so is $(i - 1, \lfloor j/2 \rfloor)$ up to top scale

Wavelet tree model

- ▶ Piecewise smooth functions (and their vector samples) is likely the most encompassing versatile model for signals and images.
- ▶ Wavelets have (rapid) polynomial decay in smooth regions, only lack decay for wavelets that interact with discontinuity.
- ▶ Convention of narrower wavelets as scale (label i) increases, coef. (i, j) large suggests coefficient $(i - 1, \lfloor j/2 \rfloor)$ also large.
- ▶ Connected subtree model: if (i, j) coefficient is kept, then so is $(i - 1, \lfloor j/2 \rfloor)$ up to top scale



Wavelet tree model

- ▶ Piecewise smooth functions (and their vector samples) is likely the most encompassing versatile model for signals and images.
- ▶ Wavelets have (rapid) polynomial decay in smooth regions, only lack decay for wavelets that interact with discontinuity.
- ▶ Convention of narrower wavelets as scale (label i) increases, coef. (i, j) large suggests coefficient $(i - 1, \lfloor j/2 \rfloor)$ also large.
- ▶ Connected subtree model: if (i, j) coefficient is kept, then so is $(i - 1, \lfloor j/2 \rfloor)$ up to top scale
- ▶ If there are k nonzeros kept in a subtree, there are $p = \text{const.} (2e)^k$ different subtrees to consider
- ▶ This helps in controlling the size of the Model RICs for $m \ll n$

Wavelet model RIC bounds

- Use basic concentration of measure bound

$$\text{Prob}(\sigma^{\max}(A_k) > 1 + \sqrt{k/m} + o(1) + t) \leq \exp(-mt^2/2)$$

$$\text{Prob}(\sigma^{\min}(A_k) < 1 - \sqrt{k/m} + o(1) - t) \leq \exp(-mt^2/2),$$

and union bound over $p = \text{const.}$ $(2e)^k$ sets

$$\text{Prob}\left(\max_{K \in \mathcal{M}_k} \sigma^{\max}(A_K) > 1 + \sqrt{\rho} + t\right) \leq c \cdot \exp(m[\rho \log(2e) - t^2/2])$$

- To have probability going to zero solve zero level curve,
 $t^* := \sqrt{2\rho \log(2e)}$
- Note, only depends on ρ , not δ
- $R_{\mathcal{M}_k}(\rho) := [1 + \sqrt{\rho} + \sqrt{2\rho \log(2e)}]^2 - 1$
- For any α , there is a ρ such that $R_{\mathcal{M}_k} < \alpha$ is satisfied
- $R_{\mathcal{M}_k^3} < 1/2$ corresponds to $m \geq 43k$, independent of n .

The impact of including a model [BaCeDuHe08]

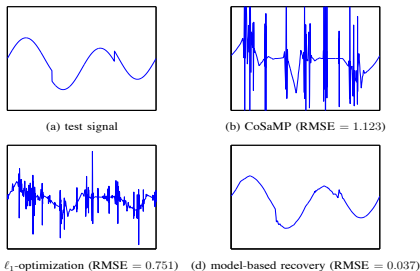


Fig. 1. Example performance of model-based signal recovery. (a) Piecewise-smooth Heavisine test signal of length $N = 1024$. This signal is compressible under a connected wavelet tree model. Signal recovered from $M = 80$ random Gaussian measurements using (b) the iterative recovery algorithm CoSaMP, (c) standard ℓ_1 linear programming, and (d) the wavelet tree-based CoSaMP algorithm from Section V. In all figures, root mean-squared error (RMSE) values are normalized with respect to the ℓ_2 norm of the signal.

- Comes with a cost. Parallel IHT has about 40% time cost for $H_k(\cdot)$ when using fast matrix vector products. Use dynamic programming to find model greatly increases the computational burden.
- If using model based, use more sophisticated (costly) decoder

Summary

- ▶ Most signals/data that we are interested in in practise has some underlying simplicity such as: compressibility, known bounds, inherent lower dimensionality
- ▶ Can move knowledge of this simplicity into the acquisition step
- ▶ Simple linear measurement processes have optimal rate, with reasonable constants, no need for learning
- ▶ Most of the contributions are on design and analysis for algorithms to recover vectors/matrices from their compressed measurements
- ▶ Methods of analysis: coherence, RICs, convex geometry
- ▶ Much is known, and there is much to be done
 - accurate understanding of average case performance
 - effect of imposing more prior information
 - extensions to other models of simplicity such as low rank

Thank you for your time and attention