
Structured Signal Processing



Joel A. Tropp

Computing + Mathematical Sciences

California Institute of Technology

jtropp@cms.caltech.edu

Amuse Bouche

The Compressed Sensing Problem

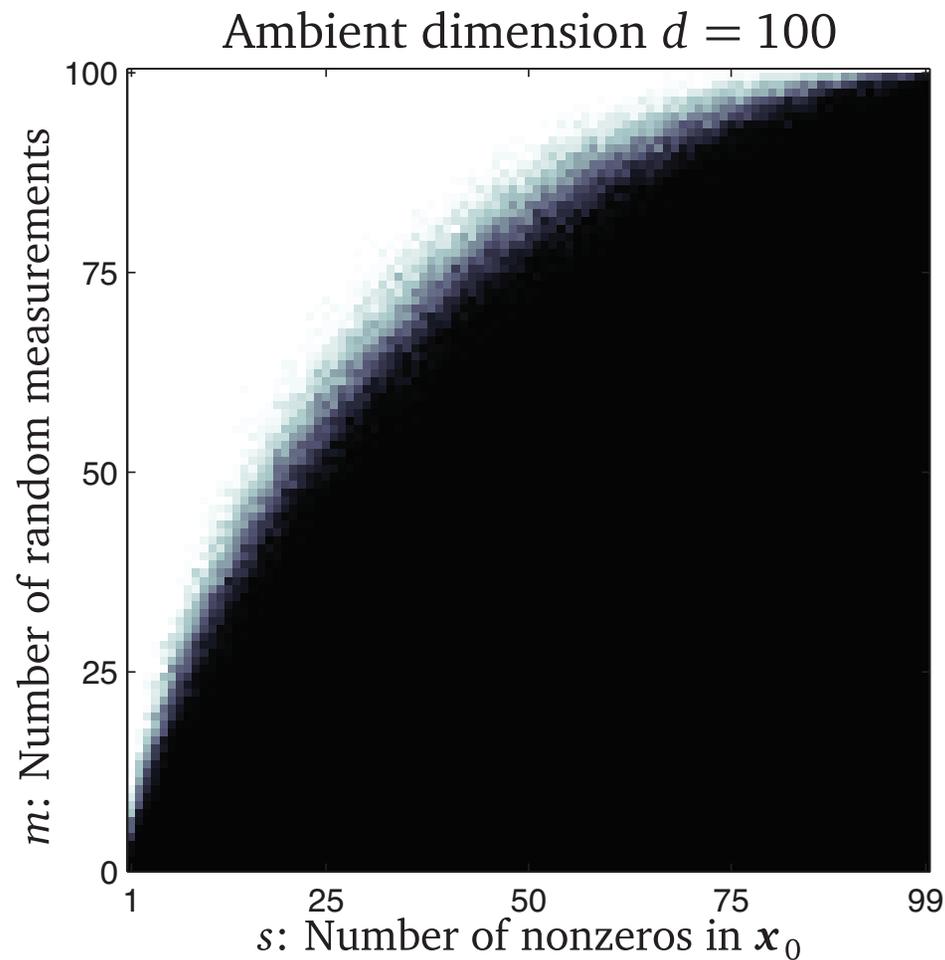
- Suppose $\mathbf{x}_\natural \in \mathbb{R}^d$ has s nonzero entries
- Let $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ be a standard normal matrix
- Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_\natural \in \mathbb{R}^m$
- Find estimate $\hat{\mathbf{x}}$ by solving convex program

$$\text{minimize } \|\mathbf{x}\|_{\ell_1} \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

- **Hope:** $\hat{\mathbf{x}} = \mathbf{x}_\natural$

Sources: Donoho 2006; Candès & Tao 2006.

Empirical Performance of Compressed Sensing

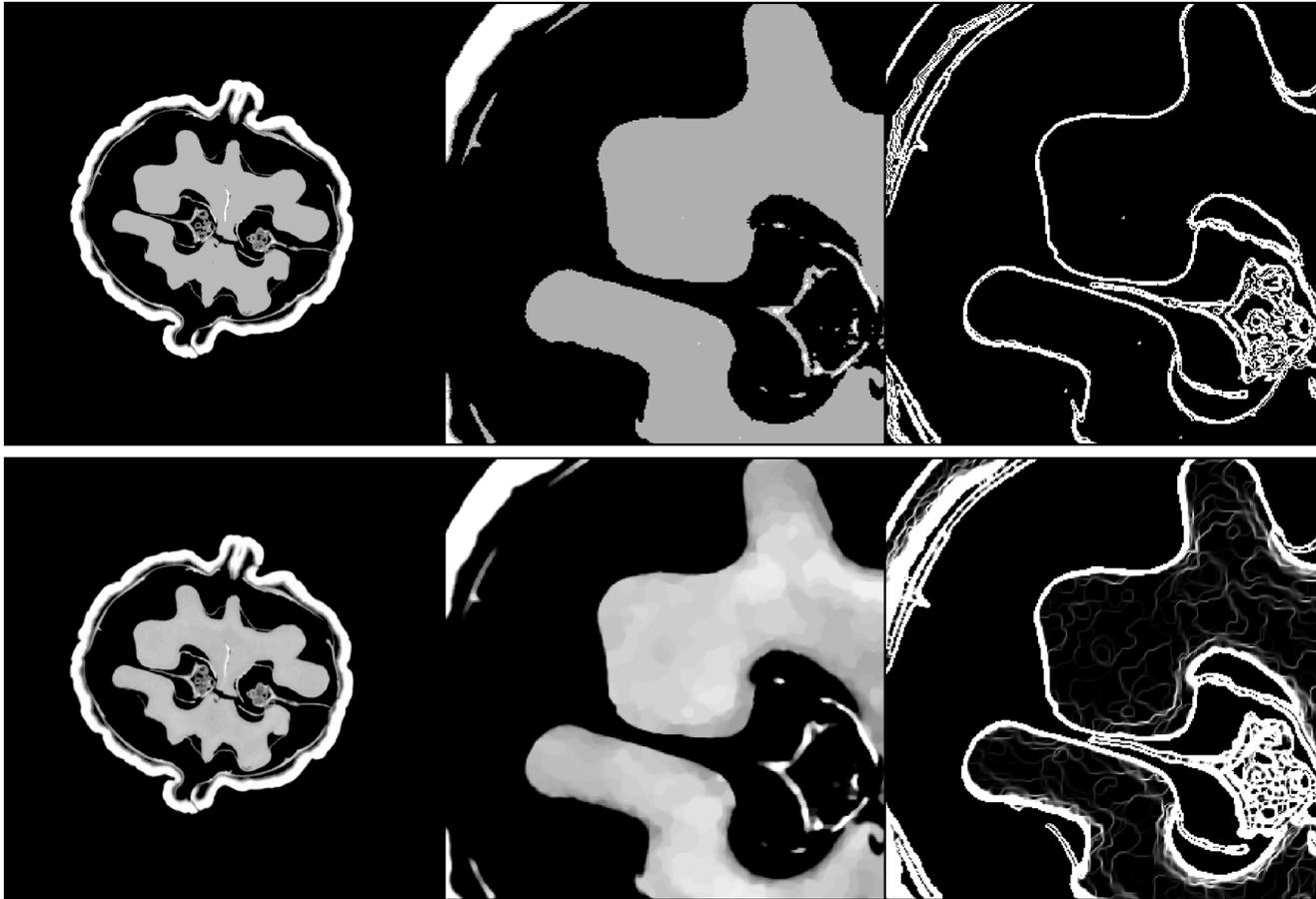


What's Going on Here?!?

- 🐼 What is the probability of success as a function of (s, m, d) ?
- 🐼 Does a phase transition exist?
- 🐼 Can we locate the phase transition?
- 🐼 How wide is the transition region?
- 🐼 Is there a geometric explanation for this phenomenon?
- 🐼 Can we export this reasoning to understand other problems?
- 🐼 Who cares?

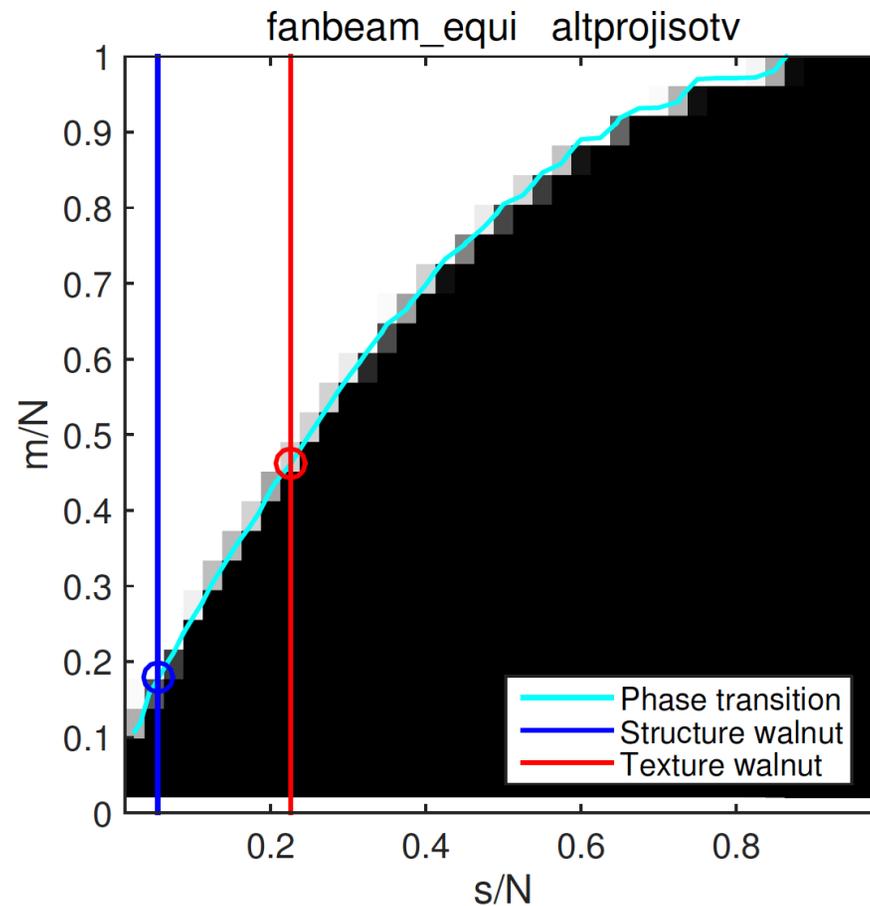
Sources: Donoho 2006; Donoho & Tanner 2009; Stojnic 2009, 2013; McCoy & Tropp 2013, 2014; Thrampoulidis et al. 2013–2016; Amelunxen et al. 2014; Foygel & Mackey 2014; Goldstein et al. 2016.

Case Study: Walnut Phantoms



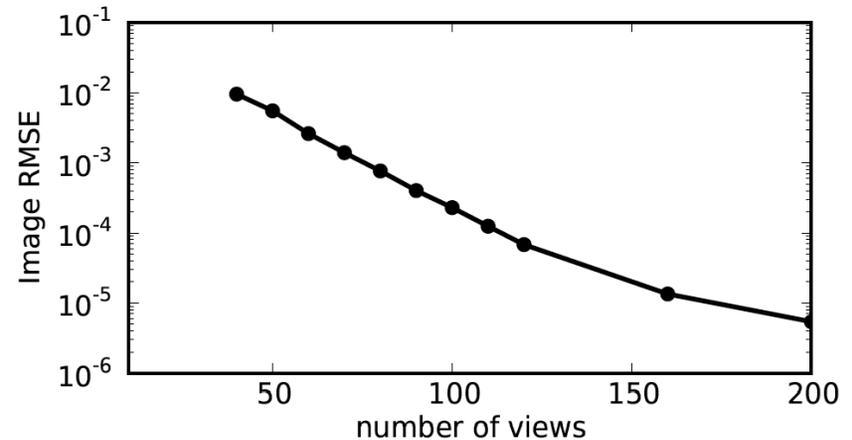
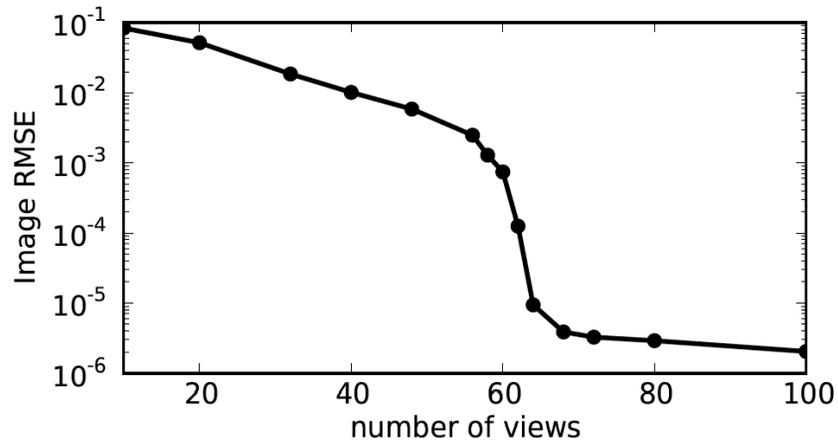
Source: Jørgensen & Sidky 2014.

Case Study: Walnut Phantoms



Source: Jørgensen & Sidky 2014.

Case Study: Walnut Phantoms

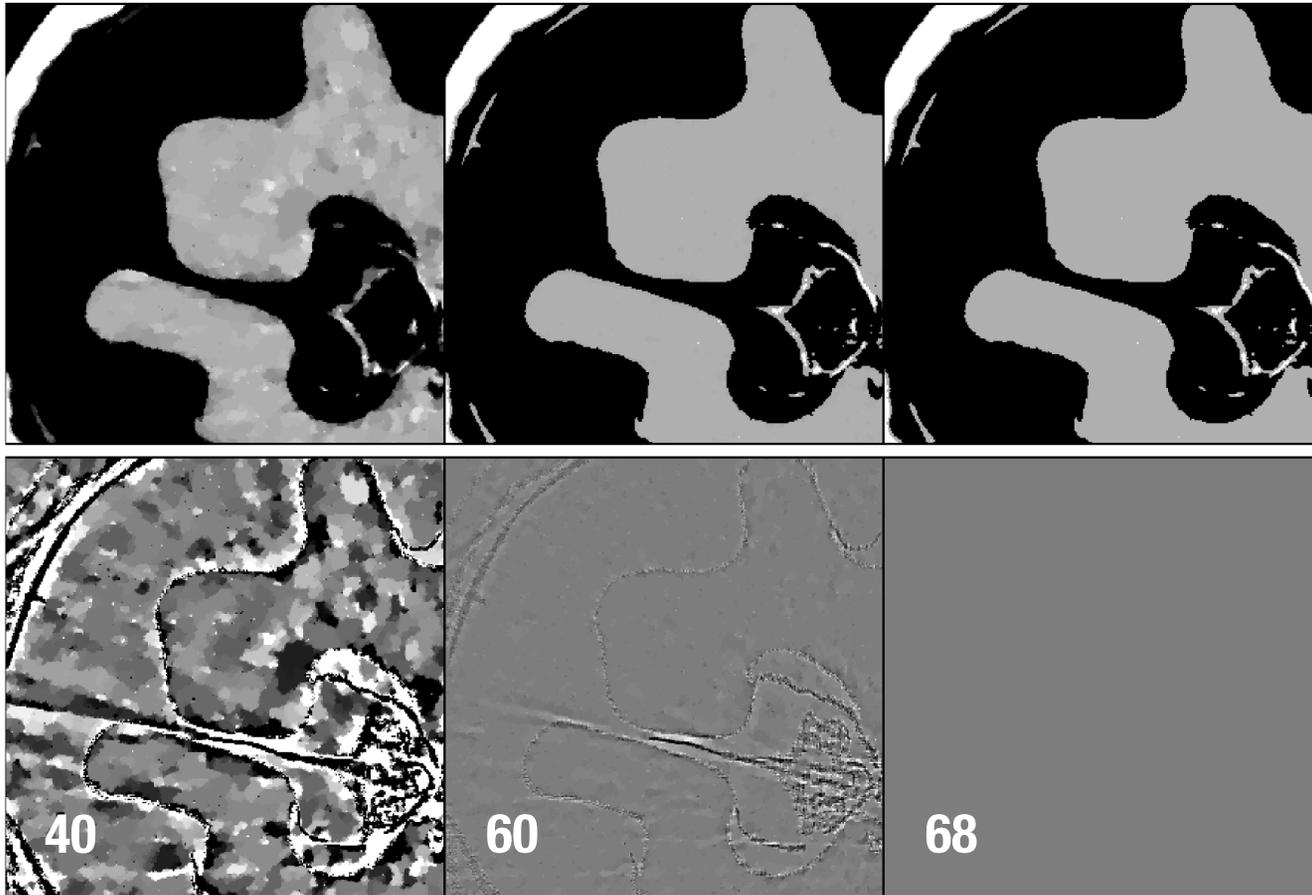


Walnut image	Gradient sparsity	Recovered at	DT prediction	ALMT prediction
Structure	45,074	68	69.3	71.7
Texture	186,306	?	188.7	185.8

Table 1: Walnut test images with gradient-domain sparsity levels, number of projections at which recovery is observed, and DT and ALMT phase-diagram predictions of critical sampling levels. A reference point of full sampling is $N_v \geq 403$ projections, where the system matrix has more rows than columns.

Source: Jørgensen & Sidky 2014.

Case Study: Walnut Phantoms



Source: Jørgensen & Sidky 2014.

Atomic Decomposition

Atoms and Dictionaries

Definition 1. Consider a compact collection \mathcal{A} of vectors:

$$\mathcal{A} = \{\mathbf{a}_\xi : \xi \in \Xi\} \subset \mathbb{R}^d$$

The collection \mathcal{A} is called a **dictionary**, and the elements \mathbf{a}_ξ are called **atoms**.

- Atoms are “elementary structures” that compose signals of interest
- Closely related to definition from nonlinear approximation (1990s)
- Terminology motivated by atomic decomposition in harmonic analysis (1970s)
- Generalizes the concept of a frame in signal processing (1980s)

Sources: Duffin & Schaeffer 1952; Coifman 1974; Daubechies et al. 1986; Mallat & Zhang 1993; Davis et al. 1994; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Temlyakov 2002; Donoho 2005; Fuchs 2005; Chandrasekaran et al. 2012.

Example: Astronomical Image

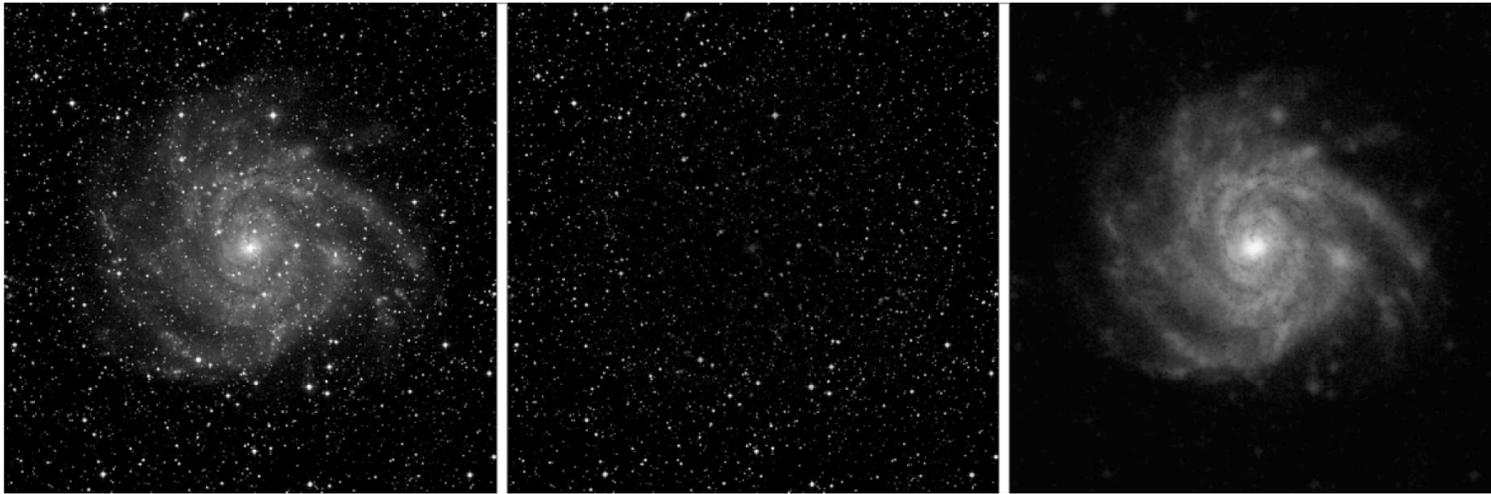


Image credit: NASA

Observation z_0

Sparse component x_0

DCT-sparse component y_0

Source: Starck et al. 2003; McCoy & Tropp 2013; McCoy et al. 2013; Amelunxen et al. 2014.

Sparsity with Respect to a Dictionary

Definition 2. For a natural number s , define the set of vectors of the form

$$\mathcal{A}_s = \left\{ \mathbf{x} = \sum_{\omega \in \Omega} c_\omega \mathbf{a}_\omega : c_\omega \geq 0 \text{ and } |\Omega| \leq s \right\}$$

The members of \mathcal{A}_s are said to be s -sparse with respect to the dictionary

- \mathcal{A}_s are increasing approximation classes
- Sparsity s parameterizes complexity of signals
- A “semiparametric” model with wide application

Sources: Stechkin 1955; Miller 1989; Jones 1992; Barron 1993; Mallat & Zhang 1993; Davis et al. 1994; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Donoho & Huo 2001; Temlyakov 2002.

Sparse Approximation Problems

👉 Given a vector \mathbf{x}_h and a dictionary \mathcal{A} ...

👉 Minimize the error subject to a bound on the sparsity:

$$\text{minimize } \|\mathbf{x} - \mathbf{x}_h\|^2 \quad \text{subject to } \mathbf{x} \in \mathcal{A}_s$$

👉 Minimize the sparsity subject to a bound on the error:

$$\text{minimize } s \quad \text{subject to } \mathbf{x} \in \mathcal{A}_s \quad \text{and} \quad \|\mathbf{x} - \mathbf{x}_h\|^2 \leq \varepsilon$$

👉 **Fact (Natarajan 1995):** In general, sparse approximation is NP-hard

$\|\cdot\|$ always denotes the Euclidean norm

Sources: Schmidt 1908; Stechkin 1955; Friedman & Stuetzle 1981; Miller 1989; Jones 1992; Barron 1993; Mallat & Zhang 1993; Davis et al. 1994; Natarajan 1995; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Donoho & Huo 2001; Temlyakov 2002;

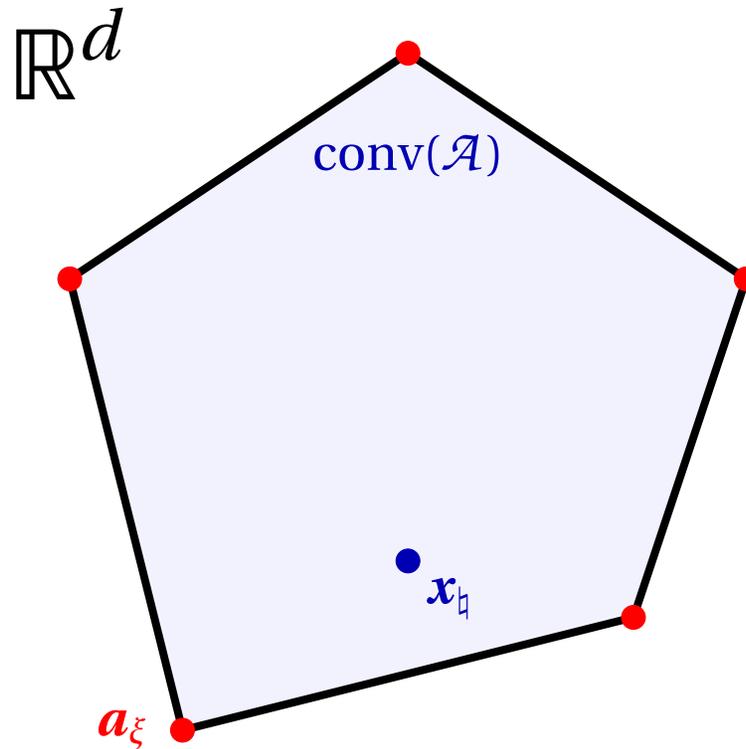
Advice for Research & Life

If at first you don't succeed...

Lower your standards!



Relax: Approximation in Convex Hulls



- 🐼 **Carathéodory:** Can write x_h as a convex combination of $d + 1$ atoms
- 🐼 How well can we approximate x_h with s atoms?

Approximate Carathéodory

Theorem 3 (Maurey 1970s). *Assume*

• \mathcal{A} is a dictionary

• $\mathbf{x}_h \in \text{conv}(\mathcal{A})$

Then there exists an s -sparse vector $\mathbf{x}_s \in \mathcal{A}_s$ with

$$\|\mathbf{x}_s - \mathbf{x}_h\| \leq \text{diam}(\mathcal{A}) \cdot \frac{1}{\sqrt{s}}$$

where $\text{diam}(\mathcal{A}) = \max_{\omega, \xi \in \Xi} \|\mathbf{a}_\omega - \mathbf{a}_\xi\|$

Sources: Maurey 1970s; Pisier 1980; Carl 1985; DeVore & Temlyakov 1996.

Proof of Maurey's Theorem: The Empirical Method

☞ Carathéodory implies

$$\mathbf{x}_h = \sum_{\omega \in \Omega} p_\omega \mathbf{a}_\omega \quad \text{where} \quad p_\omega \geq 0, \quad \sum_{\omega \in \Omega} p_\omega = 1, \quad |\Omega| \leq d + 1$$

☞ Define random vector \mathbf{z} that takes values $\mathbf{z} = \mathbf{a}_\omega$ with probability p_ω

☞ Observe: $\mathbb{E} \mathbf{z} = \mathbf{x}_h$

☞ Set $\bar{\mathbf{z}} = \frac{1}{s} \sum_{i=1}^s \mathbf{z}_i \in \mathcal{A}_s$ where $\mathbf{z}_1, \dots, \mathbf{z}_s$ are iid copies of \mathbf{z}

☞ Using orthogonality,

$$\mathbb{E} \|\bar{\mathbf{z}} - \mathbf{x}_h\|^2 = \frac{1}{s^2} \mathbb{E} \left\| \sum_{i=1}^s (\mathbf{z}_i - \mathbb{E} \mathbf{z}_i) \right\|^2 = \frac{1}{s^2} \sum_{i=1}^s \mathbb{E} \|\mathbf{z}_i - \mathbb{E} \mathbf{z}_i\|^2 = \frac{1}{s} \mathbb{E} \|\mathbf{z} - \mathbf{x}_h\|^2 \leq \frac{\text{diam}^2(\mathcal{A})}{s}$$

☞ The probabilistic method yields the desired $\mathbf{x}_s \in \mathcal{A}_s$

Atomic Gauges

Definition 4. Let \mathcal{A} be a dictionary. The **atomic gauge** is defined as

$$\begin{aligned}\|\mathbf{x}\|_{\mathcal{A}} &= \min \{t \geq 0 : \mathbf{x} \in t \cdot \text{conv}(\mathcal{A})\} \\ &= \min \left\{ \sum_{\xi \in \Xi} c_{\xi} : \mathbf{x} = \sum_{\xi \in \Xi} c_{\xi} \mathbf{a}_{\xi} \text{ and } c_{\xi} \geq 0 \right\}\end{aligned}$$

The gauge is nonnegative, positively homogeneous, and convex.

- ☞ Gauges are norm-like functions
- ☞ They are linear on rays from the origin
- ☞ They can take the values zero or $+\infty$ on an entire ray

Sources: Stechkin 1955; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Temlyakov 2002; Donoho 2005; Fuchs 2005; Chandrasekaran et al. 2012.

Atomic Gauges and Sparse Approximation

Corollary 5. For any vector \mathbf{x}_h , there is an s -sparse vector $\mathbf{x}_s \in \mathcal{A}_s$ that achieves

$$\|\mathbf{x}_s - \mathbf{x}_h\| \leq \text{diam}(\mathcal{A}) \cdot \frac{\|\mathbf{x}_h\|_{\mathcal{A}}}{\sqrt{s}}$$

- Atomic gauge always controls the quality of sparse approximations
- Evidence that atomic gauge is a reasonable proxy for complexity with respect to dictionary
- Bound is optimal for worst-case \mathbf{x}_h
- Very poor bound for exactly sparse \mathbf{x}_h

Sources: Stechkin 1955; Jones 1992; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Temlyakov 2002; Jaggi 2013.

Examples of Dictionaries and Atomic Gauges

Signal Type	Dictionary (\mathcal{A})	Atomic Gauge ($\ \cdot\ _{\mathcal{A}}$)
Sparse vector	$\{\pm \mathbf{e}_i\}$	$\ \cdot\ _{\ell_1}$
Frequency-sparse	$\{\pm \mathbf{f}_i\}$	$\ \mathcal{F}(\cdot)\ _{\ell_1}$
Spikes + sines	$\{\pm \mathbf{e}_i\} \cup \{\pm \mathbf{f}_i\}$	$\ \cdot\ _{\ell_1} \square \ \mathcal{F}(\cdot)\ _{\ell_1}$
Sparse gradient	$\{\pm(\mathbf{e}_{i+1} - \mathbf{e}_i)\}$	$\ \cdot\ _{\text{TV}}$
Sparse + nonnegative	$\{\mathbf{e}_i\}$	$\ \cdot\ _{\ell_1}$ if nn, else $+\infty$
Saturated	$\{\pm 1\}^d$	$\ \cdot\ _{\ell_\infty}$
Row-sparse matrix	$\{\mathbf{e}_i \mathbf{u}^* : \ \mathbf{u}\ = 1\}$	$\sum_i \ (\cdot)_i\ $
Low-rank matrix	$\{\mathbf{u} \mathbf{v}^* : \ \mathbf{u}\ = \ \mathbf{v}\ = 1\}$	$\ \cdot\ _{S_1}$
Low rank + psd	$\{\mathbf{u} \mathbf{u}^* : \ \mathbf{u}\ = 1\}$	$\ \cdot\ _{S_1}$ if psd, else $+\infty$
Orthogonal	$\{\mathbf{U} : \mathbf{U} \mathbf{U}^* = \mathbf{U}^* \mathbf{U} = \mathbf{I}\}$	$\ \cdot\ _{S_\infty}$

Sources: Rudin et al. 1992; Mallat & Zhang 1993; DeVore & Temlyakov 1996; Chen et al. 1997, 2001; Donoho & Huo 2001; Temlyakov 2002; Fazel 2002; Tropp 2006; Chandrasekaran et al. 2012; Jaggi 2011, 2013;

Atomic Regularization for Sparse Approximation

👉 Given a vector \mathbf{x}_b and a dictionary \mathcal{A} ...

👉 Minimize the error subject to a bound on the atomic gauge:

$$\text{minimize } \|\mathbf{x} - \mathbf{x}_b\|^2 \quad \text{subject to } \|\mathbf{x}\|_{\mathcal{A}} \leq \alpha$$

👉 Minimize the atomic gauge subject to a bound on the error:

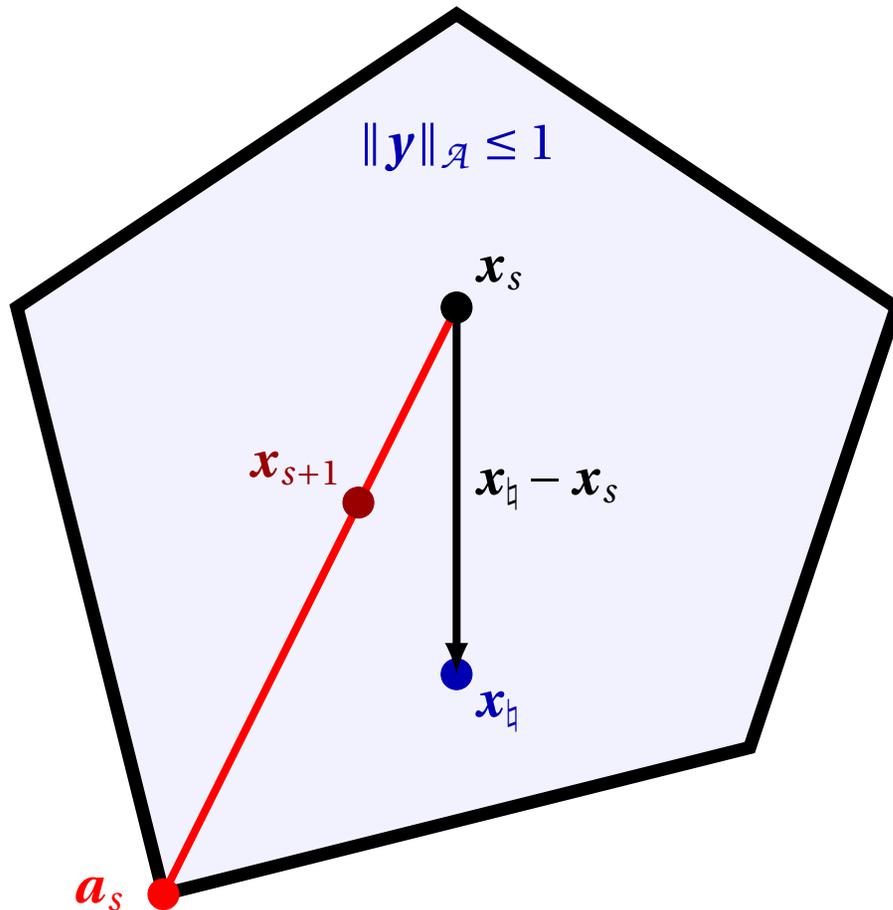
$$\text{minimize } \|\mathbf{x}\|_{\mathcal{A}} \quad \text{subject to } \|\mathbf{x} - \mathbf{x}_b\|^2 \leq \varepsilon$$

👉 Tradeoff the value of the error and the atomic gauge:

$$\text{minimize } \|\mathbf{x} - \mathbf{x}_b\|^2 + \lambda \|\mathbf{x}\|_{\mathcal{A}}$$

Sources: Chen et al. 1997, 2001; Donoho & Huo 2001; Gribonval & Nielson 2002; Donoho & Elad 2004; Tropp 2006; Fuchs 2006; Donoho et al. 2006; Chandrasekaran et al. 2012;

Conditional Gradient Method (CGM)



Initialize $\mathbf{x}_0 = \mathbf{0}$; $s = 0$

$$\mathbf{a}_s = \arg \max_{\|\mathbf{y}\|_{\mathcal{A}} \leq 1} \langle \mathbf{y}, \mathbf{x}_b - \mathbf{x}_s \rangle$$

$$\mathbf{x}_{s+1} = (1 - \theta_s)\mathbf{x}_s + \theta_s \mathbf{a}_s$$

$$\theta_s = (s + 1)^{-1}$$

$$\min_{\|\mathbf{x}\|_{\mathcal{A}} \leq 1} \|\mathbf{x} - \mathbf{x}_b\|^2$$

Sparse Approximation via CGM

Theorem 6 (CGM Convergence). *Assume*

• \mathcal{A} is a dictionary

• \mathbf{x}_h is a vector with $\|\mathbf{x}_h\|_{\mathcal{A}} \leq \alpha$, where α is known

Then, after s iterations, CGM produces an s -sparse vector \mathbf{x}_s with

$$\|\mathbf{x}_s - \mathbf{x}_h\| \leq \text{diam}(\mathcal{A}) \cdot \frac{\alpha}{\sqrt{s+1}}$$

• If α is misspecified, CGM converges to best approximation of \mathbf{x}_h with such α

• CGM also known as Frank–Wolfe, Relaxed Greedy Algorithm, or Relaxed Matching Pursuit

Sources: Frank & Wolfe 1956; Levitin & Poljak 1967; Jones 1992; DeVore & Temlyakov 1996; Temlyakov 2002; Hazan 2008; Clarkson 2010; Jaggi 2011, 2013.

Proof of CGM Convergence

Without loss, assume $\alpha = 1$

Since $\mathbf{x}_{s+1} = (1 - \theta_s)\mathbf{x}_s + \theta_s \mathbf{a}_s$,

$$\|\mathbf{x}_{s+1} - \mathbf{x}_h\|^2 = \|\mathbf{x}_s - \mathbf{x}_h\|^2 - 2\theta_s \langle \mathbf{x}_h - \mathbf{x}_s, \mathbf{a}_s - \mathbf{x}_s \rangle + \theta_s^2 \|\mathbf{a}_s - \mathbf{x}_s\|^2$$

By construction of \mathbf{a}_s ,

$$\langle \mathbf{x}_h - \mathbf{x}_s, \mathbf{a}_s - \mathbf{x}_s \rangle = \max_{\|\mathbf{y}\|_{\mathcal{A}} \leq 1} \langle \mathbf{x}_h - \mathbf{x}_s, \mathbf{y} - \mathbf{x}_s \rangle \geq \langle \mathbf{x}_h - \mathbf{x}_s, \mathbf{x}_h - \mathbf{x}_s \rangle = \|\mathbf{x}_s - \mathbf{x}_h\|^2$$

Therefore,

$$\|\mathbf{x}_{s+1} - \mathbf{x}_h\|^2 \leq (1 - 2\theta_s) \|\mathbf{x}_s - \mathbf{x}_h\|^2 + \theta_s^2 \text{diam}^2(\mathcal{A})$$

Since $\theta_s(1 - \theta_s) < \theta_{s+1}$, induction yields

$$\|\mathbf{x}_s - \mathbf{x}_h\|^2 \leq \theta_s \text{diam}^2(\mathcal{A})$$

Exact Recovery Theorems

Theorem 7 (Donoho & Elad 2004; Fuchs 2004; Tropp 2004). *Assume*

• \mathcal{A} is a standardized dictionary ($\|\mathbf{a}_\xi\| = 1$)

• the coherence $\mu = \max_{\xi \neq \omega} |\langle \mathbf{a}_\xi, \mathbf{a}_\omega \rangle|$

• $\mathbf{x}_\natural \in \mathcal{A}_s$ where $s \leq \frac{1}{2} \sqrt{\mu^{-1} + 1}$

Then we can obtain an *s-sparse representation* of \mathbf{x}_\natural by solving the linear program

$$\text{minimize } \sum_{\xi \in \Xi} c_\xi \quad \text{subject to } \mathbf{x}_\natural = \sum_{\xi \in \Xi} c_\xi \mathbf{a}_\xi \quad \text{and} \quad c_\xi \geq 0$$

• Can sometimes obtain optimal error bounds for sparse approximation!

Sources: Donoho & Huo 2001; Gribonval & Nielsen 2002; Gilbert et al. 2003; Fuchs 2004–2006; Gribonval & Vandergheynst 2004; Donoho & Elad 2004; Tropp 2004, 2006.

Applications of Sparse Models...

- 🐼 Denoising and estimation
- 🐼 Signal recovery, regression, and compressed sensing
- 🐼 Simultaneous sparse approximation and group sparsity
- 🐼 Demixing and morphological component analysis
- 🐼 Matrix completion and phase retrieval
- 🐼 Superresolution and line spectral estimation
- 🐼 Blind deconvolution and self-calibration
- 🐼 ...

Sources: Donoho & Johnstone 1992; Mallat & Zhang 1993; Chen et al. 1997, 2001; Fazel 2002; Starck et al. 2003; Tropp et al. 2006; Recht et al. 2009, 2010; Bodmann et al. 2009; Jaggi 2011, 2013; Bhaskar et al. 2012; Fernandez-Granda 2013; Romberg et al. 2013; McCoy & Tropp 2013; Amelunxen et al. 2014; Ling & Strohmer 2015;

Statistical Dimension

Regularized Denoising

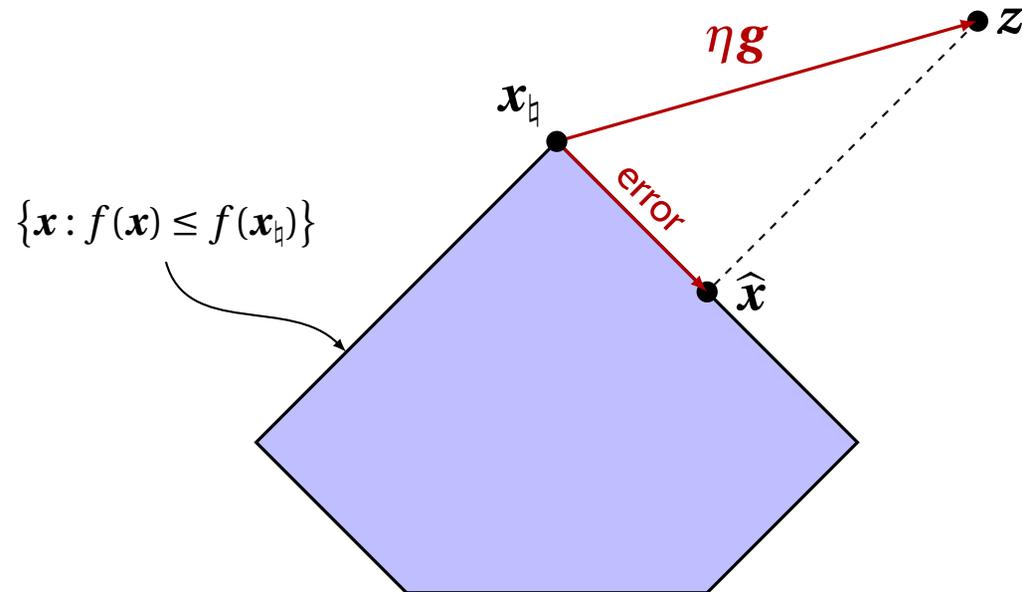
- ☞ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex “structural” penalty (e.g., an atomic gauge)
- ☞ Let \mathbf{x}_\dagger be “structured” but unknown
- ☞ Observe $\mathbf{z} = \mathbf{x}_\dagger + \eta \mathbf{g}$ where $\mathbf{g} \sim \text{NORMAL}(0, \mathbf{I})$
- ☞ Remove noise by solving the convex program

$$\text{minimize } \|\mathbf{z} - \mathbf{x}\|^2 \quad \text{subject to } f(\mathbf{x}) \leq f(\mathbf{x}_\dagger)$$

- ☞ **Hope:** The minimizer $\hat{\mathbf{x}}$ approximates \mathbf{x}_\dagger
- ☞ **Remark:** Other formulations more practical, but this is easier to analyze

Sources: Donoho et al. 2009, 2013; Bhaskar et al. 2012; Chandrasekaran & Jordan 2013; Oymak & Hassibi 2013; Amelunxen et al. 2014.

Geometry of Atomic Denoising I



Cones and Projections

Definition 8. A **convex cone** is a convex set K that satisfies $K \subseteq \tau K$ for $\tau \geq 0$.

Definition 9. Let K be a convex cone. The **polar** is the closed convex cone $K^\circ = \{\mathbf{y} : \langle \mathbf{y}, \mathbf{x} \rangle \leq 0 \text{ for all } \mathbf{x} \in K\}$.

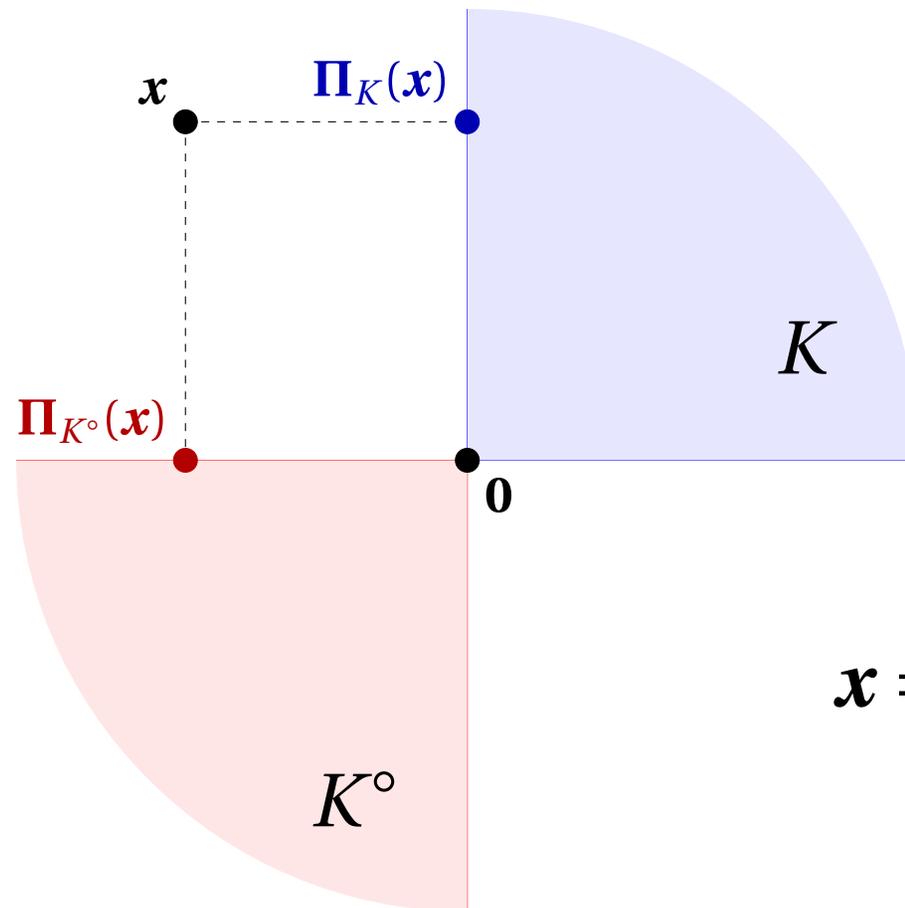
Definition 10. Let K be a closed convex cone. The **Euclidean projection** onto K is

$$\mathbf{\Pi}_K(\mathbf{x}) = \arg \min_{\mathbf{y} \in K} \|\mathbf{y} - \mathbf{x}\|^2$$

For a general convex cone C , define $\mathbf{\Pi}_C = \mathbf{\Pi}_{\text{closure}(C)}$

Sources: Rockafellar 1970; Rockafellar & Wets 1997; Hiriart-Urruty & Lemaréchal 2002.

Moreau's Theorem



$$\mathbf{x} = \Pi_K(\mathbf{x}) + \Pi_{K^\circ}(\mathbf{x})$$

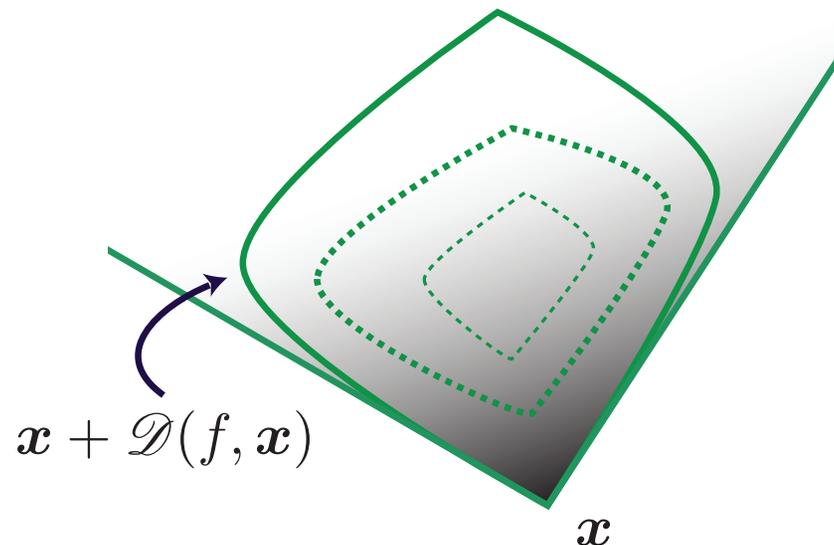
$$\Pi_K(\mathbf{x}) \perp \Pi_{K^\circ}(\mathbf{x})$$

Sources: Moreau 1965; Rockafellar 1970; Rockafellar & Wets 1997; Hiriart-Urruty & Lemaréchal 2002.

Descent Cones

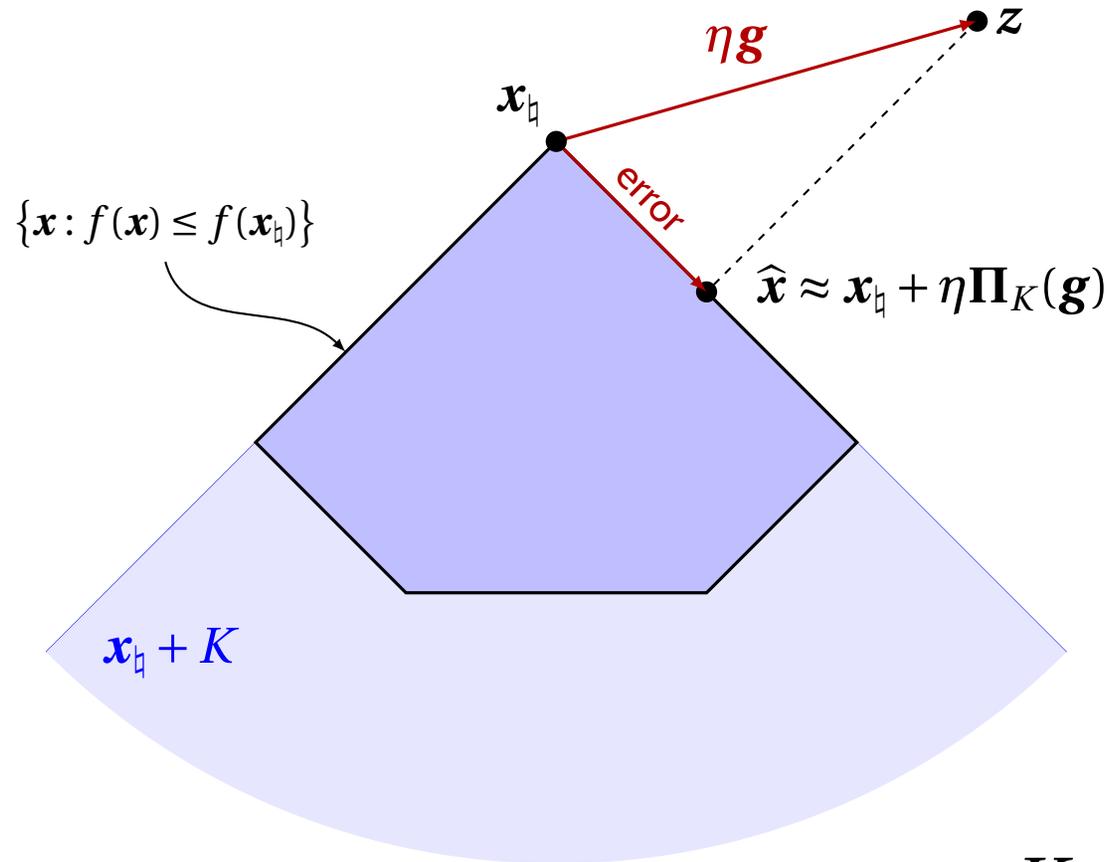
Definition 11. For convex f and a point \mathbf{x} , the **descent cone** is the convex cone

$$\mathcal{D}(f, \mathbf{x}) = \{\mathbf{u} : f(\mathbf{x} + \varepsilon \mathbf{u}) \leq f(\mathbf{x}) \text{ for some } \varepsilon > 0\}$$



Sources: Rockafellar 1970; Rockafellar & Wets 1997; Hiriart-Urruty & Lemaréchal 2002.

Geometry of Regularized Denoising II



$$K = \mathcal{D}(f, x_q)$$

Analysis of Regularized Denoising

Theorem 12 (Oymak & Hassibi 2013). *Assume*

- f is a convex function and $\mathbf{x}_q \in \text{dom}(f)$
- Observe $\mathbf{z} = \mathbf{x}_q + \eta \mathbf{g}$ where $\mathbf{g} \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } \|\mathbf{x} - \mathbf{z}\|^2 \quad \text{subject to } f(\mathbf{x}) \leq f(\mathbf{x}_q)$$

Then

$$\sup_{\eta > 0} \frac{\mathbb{E} \|\hat{\mathbf{x}} - \mathbf{x}_q\|^2}{\eta^2} = \lim_{\eta \downarrow 0} \frac{\mathbb{E} \|\hat{\mathbf{x}} - \mathbf{x}_q\|^2}{\eta^2} = \mathbb{E} \|\mathbf{\Pi}_K(\mathbf{g})\|^2$$

where $K = \mathcal{D}(f, \mathbf{x}_q)$

Sources: Donoho et al. 2009, 2013; Oymak & Hassibi 2013; Chandrasekaran & Jordan 2013; Amelunxen et al. 2014.

Regularized Denoising: Upper Bound

☞ Change variables: $\mathbf{u} = \mathbf{x} - \mathbf{x}_q$

$$\text{minimize } \|\mathbf{u} - \eta \mathbf{g}\|^2 \quad \text{subject to } f(\mathbf{x}_q + \mathbf{u}) \leq f(\mathbf{x}_q)$$

☞ First-order optimality for $\hat{\mathbf{u}} = \hat{\mathbf{x}} - \mathbf{x}_q$:

$$\langle \mathbf{v} - \hat{\mathbf{u}}, \hat{\mathbf{u}} - \eta \mathbf{g} \rangle \geq 0 \quad \text{for all feasible } \mathbf{v}$$

☞ Choose $\mathbf{v} = \mathbf{0}$ and use the fact $\hat{\mathbf{u}} \in \mathcal{D}(f, \mathbf{x}_q) = K$:

$$\|\hat{\mathbf{u}}\|^2 \leq \langle \eta \mathbf{g}, \hat{\mathbf{u}} \rangle = \eta \langle \Pi_K(\mathbf{g}) + \Pi_{K^c}(\mathbf{g}), \hat{\mathbf{u}} \rangle \leq \eta \langle \Pi_K(\mathbf{g}), \hat{\mathbf{u}} \rangle \leq \eta \|\hat{\mathbf{u}}\| \|\Pi_K(\mathbf{g})\|$$

☞ Rearrange and take expectation:

$$\eta^{-2} \mathbb{E} \|\hat{\mathbf{u}}\|^2 \leq \mathbb{E} \|\Pi_K(\mathbf{g})\|^2$$

☞ The lower bound is technical

Statistical Dimension

Definition 13. Let K be a convex cone in \mathbb{R}^d . The **statistical dimension** is

$$\delta(K) = \mathbb{E} \|\mathbf{\Pi}_K(\mathbf{g})\|^2 \quad \text{where} \quad \mathbf{g} \sim \text{NORMAL}(\mathbf{0}, \mathbf{I})$$

An extension of linear dimension to cones:

☛ $0 \leq \delta(K) \leq d$

☛ $C \subset K$ implies $\delta(C) \leq \delta(K)$

☛ If L is a subspace, $\delta(L) = \dim(L)$

Sources: Amelunxen et al. 2014; McCoy & Tropp 2014.

Properties of Statistical Dimension

- Let C, K be convex cones in \mathbb{R}^d
- Gaussian formulation:** $\delta(K) = \mathbb{E} \|\Pi_K(\mathbf{g})\|^2$
 - Nonnegativity:** $\delta(K) \geq 0$
 - Subspaces:** $\delta(L) = \dim(L)$ for a subspace L
 - Rotational invariance:** $\delta(K) = \delta(\mathbf{Q}K)$ for any orthogonal \mathbf{Q}
- Complementarity:** $\delta(K) + \delta(K^\circ) = d$
 - Upper bound:** $\delta(K) \leq d$
- Polar formulation:** $\delta(K) = \mathbb{E} \text{dist}^2(\mathbf{g}, K^\circ)$
- Mean-squared-width formulation:** $\delta(K) = \mathbb{E} \left(\sup_{\|\mathbf{u}\| \leq 1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2$
 - Monotonicity:** $C \subset K$ implies $\delta(C) \leq \delta(K)$

Sources: Chandrasekaran et al. 2012; Amelunxen et al. 2014; McCoy & Tropp 2014.

Proof of Statistical Dimension Properties

🐼 Moreau's Theorem + basic facts about standard normal vectors

🐼 Complementarity:

$$d = \mathbb{E} \|\mathbf{g}\|^2 = \mathbb{E} \|\mathbf{\Pi}_K(\mathbf{g})\|^2 + \mathbb{E} \|\mathbf{\Pi}_{K^\circ}(\mathbf{g})\|^2 = \delta(K) + \delta(K^\circ)$$

🐼 Polar formulation:

$$\delta(K) = \mathbb{E} \|\mathbf{\Pi}_K(\mathbf{g})\|^2 = \mathbb{E} \|\mathbf{g} - \mathbf{\Pi}_{K^\circ}(\mathbf{g})\|^2 = \mathbb{E} \text{dist}^2(\mathbf{g}, K^\circ)$$

🐼 Mean-squared-width formulation:

$$\sup_{\substack{\|\mathbf{u}\| \leq 1 \\ \mathbf{u} \in K}} \langle \mathbf{g}, \mathbf{u} \rangle = \sup_{\substack{\|\mathbf{u}\| \leq 1 \\ \mathbf{u} \in K}} \langle \mathbf{\Pi}_K(\mathbf{g}) + \mathbf{\Pi}_{K^\circ}(\mathbf{g}), \mathbf{u} \rangle \leq \sup_{\substack{\|\mathbf{u}\| \leq 1 \\ \mathbf{u} \in K}} \langle \mathbf{\Pi}_K(\mathbf{g}), \mathbf{u} \rangle = \|\mathbf{\Pi}_K(\mathbf{g})\|$$

$$\sup_{\substack{\|\mathbf{u}\| \leq 1 \\ \mathbf{u} \in K}} \langle \mathbf{g}, \mathbf{u} \rangle \geq \langle \mathbf{g}, \mathbf{\Pi}_K(\mathbf{g}) / \|\mathbf{\Pi}_K(\mathbf{g})\| \rangle = \|\mathbf{\Pi}_K(\mathbf{g})\|$$

Gaussian Width and Statistical Dimension

Proposition 14 (Amelunxen et al. 2014). Assume K is a convex cone in \mathbb{R}^d and $\mathbf{g} \in \mathbb{R}^d$ is standard normal. Then

$$\delta(K) - 1 \leq \left(\mathbb{E} \sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \leq \delta(K)$$

👉 **Upper bound** (Jensen):

$$\left(\mathbb{E} \sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \leq \mathbb{E} \left(\sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \leq \mathbb{E} \left(\sup_{\|\mathbf{u}\| \leq 1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 = \delta(K)$$

👉 **Lower bound** (Poincaré):

$$\begin{aligned} 1 + \left(\mathbb{E} \sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 &\geq \mathbb{E} \left(\sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \\ &\geq \mathbb{E} \left[\left(\sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \cdot \mathbb{1}_{\mathbf{g} \in \mathbb{R}^d \setminus K^\circ} \right] = \mathbb{E} \left[\left(\sup_{\|\mathbf{u}\| \leq 1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \right)^2 \right] = \delta(K) \end{aligned}$$

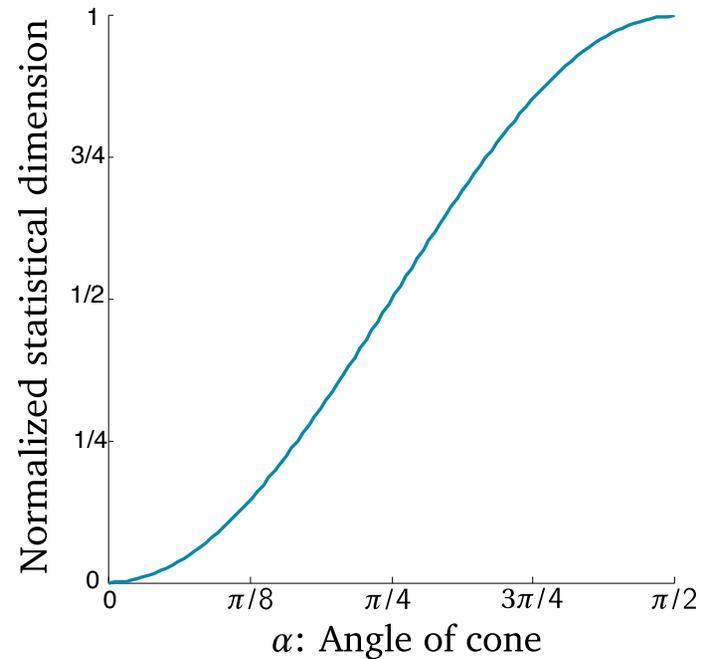
Source: Amelunxen et al. 2014.

Basic Examples of Statistical Dimension

Cone	Notation	Statistical Dimension
Subspace	L_j	j
Nonnegative orthant	\mathbb{R}_+^d	$\frac{1}{2}d$
Second-order cone	\mathbb{L}^{d+1}	$\frac{1}{2}(d+1)$
Real psd cone	\mathbb{S}_+^d	$\frac{1}{4}d(d-1)$
Complex psd cone	\mathbb{H}_+^d	$\frac{1}{2}d^2$

Sources: Chandrasekaran et al. 2012; Amelunxen et al. 2014; McCoy & Tropp 2014.

Circular Cones



• For $\alpha \in (0, \pi/2)$, define $\text{Circ}_d(\alpha) = \{\mathbf{x} \in \mathbb{R}^d : x_1 \geq \|\mathbf{x}\| \cos(\alpha)\}$

• $\delta(\text{Circ}_d(\alpha)) = d \sin^2(\alpha) + \cos(2\alpha) + o(1)$

Sources: Amelunxen et al. 2014; McCoy & Tropp 2014.

Polar Form of a Descent Cone

Definition 15. The subdifferential ∂f of a convex function f at a point \mathbf{x} is

$$\partial f(\mathbf{x}) = \{\mathbf{u} : f(\mathbf{y}) \geq f(\mathbf{x}) + \langle \mathbf{u}, \mathbf{y} - \mathbf{x} \rangle \text{ for all } \mathbf{y}\}$$

Fact 16. Assume that $\partial f(\mathbf{x})$ is nonempty, compact, and does not contain the origin.

Then

$$\mathcal{D}(f, \mathbf{x})^\circ = \bigcup_{\tau \geq 0} \tau \cdot \partial f(\mathbf{x})$$

Remark: Compactness is not essential.

Sources: Rockafellar 1970; Rockafellar & Wets 1998; Hiriart-Urruty & Lemaréchal 2002.

The Descent Cone Recipe

Proposition 17 (Amelunxen et al. 2014). Assume that f is a convex function whose subdifferential $\partial f(\mathbf{x})$ is nonempty and does not contain the origin. Then

$$\delta(\mathcal{D}(f, \mathbf{x})) \leq \inf_{\tau \geq 0} \mathbb{E} \text{dist}^2(\mathbf{g}, \tau \cdot \partial f(\mathbf{x})) = \inf_{\tau \geq 0} J(\tau)$$

where \mathbf{g} is standard normal

🐼 Assume $\partial f(\mathbf{x})$ is compact for simplicity

🐼 Calculate:

$$\begin{aligned} \delta(\mathcal{D}(f, \mathbf{x})) &= \mathbb{E} \text{dist}^2(\mathbf{g}, \mathcal{D}(f, \mathbf{x})^\circ) = \mathbb{E} \text{dist}^2\left(\mathbf{g}, \bigcup_{\tau \geq 0} \tau \cdot \partial f(\mathbf{x})\right) \\ &= \mathbb{E} \inf_{\tau \geq 0} \text{dist}^2(\mathbf{g}, \tau \cdot \partial f(\mathbf{x})) \leq \inf_{\tau \geq 0} \mathbb{E} \text{dist}^2(\mathbf{g}, \tau \cdot \partial f(\mathbf{x})) \end{aligned}$$

Sources: Stojnic 2009, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Foygel & Mackey 2014.

The Descent Cone Recipe: Error Estimate

Theorem 18 (Amelunxen et al. 2014). Assume that f is a norm. For nonzero \mathbf{x} ,

$$0 \leq \inf_{\tau \geq 0} J(\tau) - \delta(\mathcal{D}(f, \mathbf{x})) \leq \frac{2 \max\{\|\mathbf{u}\| : \mathbf{u} \in \partial f(\mathbf{x})\}}{f(\mathbf{x} / \|\mathbf{x}\|)}$$

- 🐼 **Idea:** Linearize $\tau \mapsto \text{dist}^2(\mathbf{g}, \tau \cdot \partial f(\mathbf{x}))$ at minimizer τ_* of J
- 🐼 Related result of Foygel & Mackey based on other ideas
- 🐼 Still room for improvement

Sources: Amelunxen et al. 2014; Foygel & Mackey 2014.

Example: ℓ_1 Statistical Dimension I

🐼 **Goal:** Compute $\delta(\mathcal{D}(\|\cdot\|_{\ell_1}, \mathbf{x}_s))$ for vector $\mathbf{x}_s \in \mathbb{R}^d$ with s nonzero entries

🐼 By symmetry, can assume $\mathbf{x}_s = (\mathbf{1}_s, \mathbf{0}_{d-s})$

🐼 Subdifferential: $\partial f(\mathbf{x}_s) = \{(\mathbf{1}_s, \mathbf{v}) \in \mathbb{R}^d : \|\mathbf{v}\|_\infty \leq 1\}$

🐼 Distance to scaled subdifferential:

$$\begin{aligned} J_s(\tau) &= \mathbb{E} \text{dist}^2(\mathbf{g}, \tau \cdot \partial \|\mathbf{x}_s\|_{\ell_1}) = \sum_{i=1}^s \mathbb{E}(g_i - \tau)^2 + \sum_{i=s+1}^d \mathbb{E} \min_{|v_i| \leq 1} (g_i - \tau v_i)^2 \\ &= s \cdot \mathbb{E}(g - \tau)^2 + (d - s) \cdot \mathbb{E}(|g| - \tau)_+^2 \leq s \cdot (1 + \tau^2) + 2(d - s) \cdot e^{-\tau^2/2} \end{aligned}$$

🐼 Choose $\tau^2 = 2 \log(d/s)$ to reach

$$\delta(\mathcal{D}(\|\cdot\|_{\ell_1}, \mathbf{x}_s)) \leq \inf_{\tau \geq 0} J_s(\tau) \leq 2s(1 + \log(d/s))$$

Sources: Donoho 2006; Donoho & Tanner 2009; Stojnic 2009, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Foygel & Mackey 2014; Goldstein et al. 2014.

Example: ℓ_1 Statistical Dimension II

• $\mathbf{x}_s \in \mathbb{R}^d$ has s nonzero entries

• Observe $\mathbf{z} = \mathbf{x}_s + \eta \mathbf{g}$

• $\hat{\mathbf{x}}$ solves

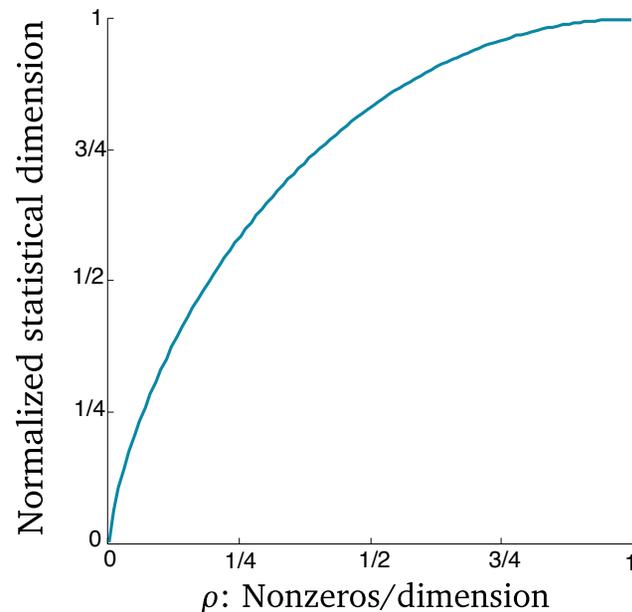
$$\text{minimize } \|\mathbf{x} - \mathbf{z}\|^2 \quad \text{subject to } \|\mathbf{x}\|_{\ell_1} \leq \|\mathbf{x}_s\|_{\ell_1}$$

• Then

$$\mathbb{E} \|\hat{\mathbf{x}} - \mathbf{x}_s\|^2 \leq 2s(1 + \log(d/s)) \cdot \eta^2$$

• Almost achieve same MSE as if we knew $\text{supp}(\mathbf{x}_s)$!

Example: ℓ_1 Statistical Dimension III

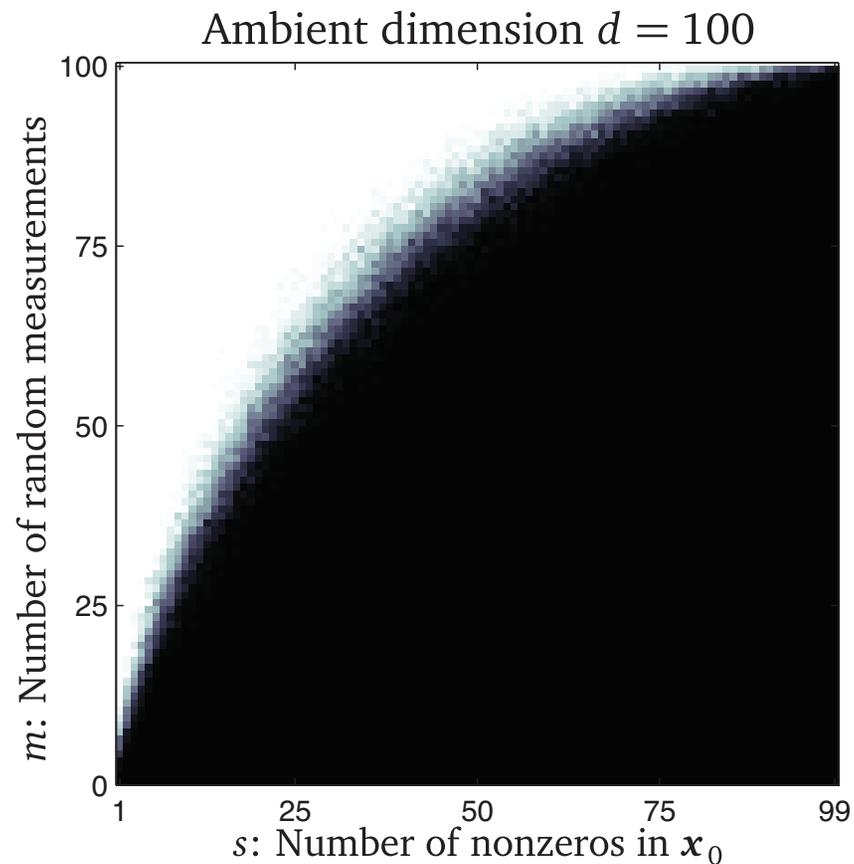


• With $\rho = s/d \in [0, 1]$,

$$\frac{\delta(\mathcal{D}(f, \mathbf{x}_s))}{d} \leq \inf_{\tau \geq 0} [\rho(1 + \tau^2) + (1 - \rho) \mathbb{E}(|g| - \tau)_+^2] \leq \frac{\delta(\mathcal{D}(f, \mathbf{x}_s))}{d} + \frac{1}{d\sqrt{\rho}}$$

Sources: Affentranger & Schneider 1992; Betke & Henk 1993; Böröczky & Henk 1999; Donoho 2006; Donoho & Tanner 2009; Stojnic 2009, 2013; McCoy & Tropp 2013; Amelunxen et al. 2014; Foygel & Mackey 2014; Goldstein et al. 2014.

Compressed Sensing: Hmmm...



Sources: Donoho et al. 2009–2013.

Example: S_1 Statistical Dimension I

• $\mathbf{X}_r \in \mathbb{R}^{d_1 \times d_2}$ has rank r

• Observe $\mathbf{Z} = \mathbf{X}_r + \eta \mathbf{G}$

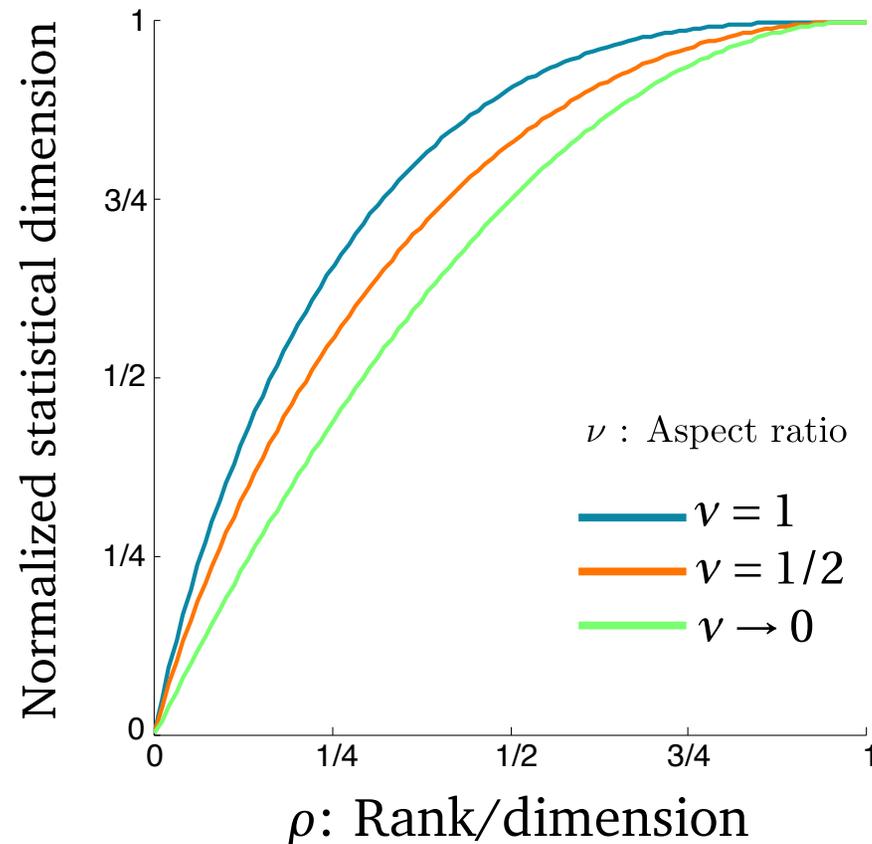
• $\hat{\mathbf{X}}$ solves

$$\text{minimize } \|\mathbf{X} - \mathbf{Z}\|_{\text{F}}^2 \quad \text{subject to } \|\mathbf{X}\|_{S_1} \leq \|\mathbf{X}_r\|_{S_1}$$

• Then

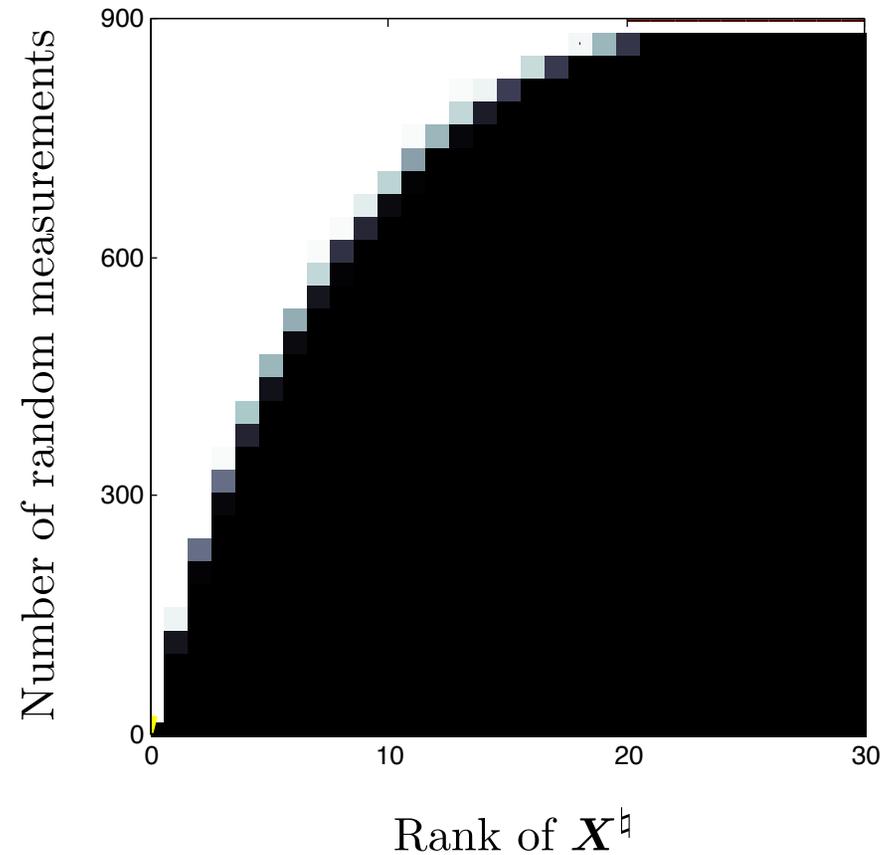
$$\mathbb{E} \|\hat{\mathbf{X}} - \mathbf{X}_r\|_{\text{F}}^2 \leq \delta(\mathcal{D}(\|\cdot\|_{S_1}, \mathbf{X}_r)) \cdot \eta^2 \leq 3r(d_1 + d_2 - r) \cdot \eta^2$$

Example: S_1 Statistical Dimension II



Sources: Oymak et al. 2010; Chandrasekaran et al. 2012; Amelunxen et al. 2014; McCoy & Tropp 2013, 2014; Goldstein et al. 2014.

Matrix Compressed Sensing: Hmmm...



Sources: Recht et al. 2010; Oymak et al. 2010; Donoho et al. 2013.

A Conjecture...

The phase transition of matrix recovery from Gaussian measurements matches the minimax MSE of matrix denoising

David L. Donoho^{a,1}, Matan Gavish^a, and Andrea Montanari^{a,b}

Departments of ^aStatistics and ^bElectrical Engineering, Stanford University, Stanford, CA 94305

Contributed by David L. Donoho, April 3, 2013 (sent for review February 5, 2013)

Let X_0 be an unknown M by N matrix. In matrix recovery, one takes $n < MN$ linear measurements y_1, \dots, y_n of X_0 , where $y_i = \text{Tr}(A_i^T X_0)$ and each A_i is an M by N matrix. A popular approach for matrix recovery is nuclear norm minimization (NNM): solving the convex optimization problem $\min \|X\|_*$ subject to $y_i = \text{Tr}(A_i^T X)$ for all $1 \leq i \leq n$, where $\|\cdot\|_*$ denotes the nuclear norm, namely, the sum of singular values. Empirical work reveals a phase transition curve, stated in terms of the undersampling fraction $\delta(n, M, N) = n/(MN)$, rank fraction $\rho = \text{rank}(X_0)/\min\{M, N\}$, and aspect ratio $\beta = M/N$. Specifically when the measurement matrices A_i have independent standard Gaussian random entries, a curve $\delta^*(\rho) = \delta^*(\rho; \beta)$ exists such that, if $\delta > \delta^*(\rho)$, NNM typically succeeds for large M, N , whereas if $\delta < \delta^*(\rho)$, it typically fails. An apparently quite different problem is matrix denoising in Gaussian noise, in which an unknown M by N matrix X_0 is to be estimated based on direct noisy measurements $Y = X_0 + Z$, where the matrix Z has independent and identically distributed Gaussian entries. A popular

measurement operator \mathcal{A} , then the solution $X_1 = X_1(y)$ to (P_{nuc}) is precisely X_0 . Such incoherence can be obtained by letting \mathcal{A} be random, for instance if $\mathcal{A}(X_0)_i = \text{Tr}(A_i^T X_0)$ with $A_i \in \mathbb{R}^{m \times n}$ having independent and identically distributed (i.i.d) Gaussian entries. In this case we speak of “matrix recovery from Gaussian measurements” (1).

A key phrase from the previous paragraph is “if X_0 is sufficiently low rank.” Clearly, there must be a quantitative tradeoff between the rank of X_0 and the number of measurements required, such that higher-rank matrices require more measurements. In the Gaussian measurements model, with N, M sufficiently large, empirical work by Recht et al. (1, 7, 8), Tanner and Wei (9), and Oymak and Hassibi (10) documents a phase transition phenomenon. For matrices of a given rank, there is a fairly precise number of required samples, in the sense that a transition from nonrecovery to complete recovery takes place sharply as the number of samples varies across this value. For example, in Table S1 we present data from

Gaussian Comparisons

Slepian's Lemma

Theorem 19 (Slepian 1962). *Assume* $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ be centered, jointly Gaussian vectors whose covariance structures satisfy

$$\begin{cases} \mathbb{E} X_i X_j \leq \mathbb{E} Y_i Y_j & \text{for all } i \neq j \\ \mathbb{E} X_i^2 = \mathbb{E} Y_i^2 & \text{for all } i \end{cases}$$

Then, for all choices of $\lambda_i \in \mathbb{R}$,

$$\mathbb{P} \left(\bigcup_{i=1}^N \{Y_i > \lambda_i\} \right) \leq \mathbb{P} \left(\bigcup_{i=1}^N \{X_i > \lambda_i\} \right)$$

In particular,

$$\mathbb{E} \max_i Y_i \leq \mathbb{E} \max_i X_i$$

Sources: Slepian 1962; Sudakov 1969, 1971; Marcus & Shepp 1970, 1972; Fernique 1974; Joag-Dev et al. 1983; Kahane 1986; Ledoux & Talagrand (Cor. 3.12).

Chevet's Theorem

Corollary 20 (Chevet 1977, Gordon 1985). **Assume**

- $U \subset \mathbb{R}^m$ and $V \subset \mathbb{R}^n$ are compact subsets of the unit sphere
- $\Gamma \in \mathbb{R}^{n \times m}$ and $\mathbf{g} \in \mathbb{R}^m$ and $\mathbf{h} \in \mathbb{R}^n$ are independent standard normal

Then

$$\mathbb{E} \max_{\mathbf{u} \in U, \mathbf{v} \in V} \langle \Gamma \mathbf{u}, \mathbf{v} \rangle \leq \mathbb{E} \max_{\mathbf{u} \in U, \mathbf{v} \in V} [\langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle]$$

Sources: Chevet 1977; Gordon 1985, 1988; Ledoux & Talagrand (Thm. 3.20); Davidson & Szarek 2001; Stojnic 2013; Amelunxen et al. 2014; Thrampoulidis et al. 2015; Amelunxen & Lotz 2015.

Example: Spectral Norm of a Gaussian Matrix

🐼 **Goal:** For standard normal $\Gamma \in \mathbb{R}^{n \times m}$, bound expectation of

$$\|\Gamma\| = \max_{\|u\|=\|v\|=1} \langle \Gamma u, v \rangle$$

🐼 Apply Chevet's Theorem:

$$\begin{aligned} \mathbb{E} \|\Gamma\| &= \mathbb{E} \max_{\|u\|=\|v\|=1} \langle \Gamma u, v \rangle \\ &\leq \mathbb{E} \max_{\|u\|=\|v\|=1} [\langle \mathbf{g}, u \rangle + \langle \mathbf{h}, v \rangle] \\ &= \mathbb{E} \|\mathbf{g}\| + \mathbb{E} \|\mathbf{h}\| \\ &\leq \sqrt{m} + \sqrt{n} \end{aligned}$$

🐼 **Result is sharp, including constants!**

Sources: Marchenko & Pastur 1967; Chevet 1977; Gordon 1985, 1988; Yin et al. 1988; Bai et al. 1988; Edelman 1988; Ledoux & Talagrand (Thm. 3.20); Davidson & Szarek 2001.

Proof of Chevet's Theorem

☞ Let $\gamma \in \mathbb{R}$ be an independent standard normal

☞ Define independent Gaussian processes on $\{(\mathbf{u}, \mathbf{v}) : \|\mathbf{u}\| = \|\mathbf{v}\| = 1\}$:

$$Y_{\mathbf{u}\mathbf{v}} = \langle \mathbf{\Gamma} \mathbf{u}, \mathbf{v} \rangle + \gamma \quad \text{and} \quad X_{\mathbf{u}\mathbf{v}} = \langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle$$

☞ Compute the covariances:

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}} Y_{\mathbf{u}'\mathbf{v}'} = \langle \mathbf{u}, \mathbf{u}' \rangle \langle \mathbf{v}, \mathbf{v}' \rangle + 1 \quad \text{and} \quad \mathbb{E} X_{\mathbf{u}\mathbf{v}} X_{\mathbf{u}'\mathbf{v}'} = \langle \mathbf{u}, \mathbf{u}' \rangle + \langle \mathbf{v}, \mathbf{v}' \rangle$$

☞ Comparison:

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}} Y_{\mathbf{u}'\mathbf{v}'} - \mathbb{E} X_{\mathbf{u}\mathbf{v}} X_{\mathbf{u}'\mathbf{v}'} = (1 - \langle \mathbf{u}, \mathbf{u}' \rangle)(1 - \langle \mathbf{v}, \mathbf{v}' \rangle) \geq 0$$

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}}^2 - \mathbb{E} X_{\mathbf{u}\mathbf{v}}^2 = 0$$

☞ Apply Slepian's Lemma (on finite subsets of U, V):

$$\mathbb{E} \max_{\|\mathbf{u}\|=\|\mathbf{v}\|=1} \langle \mathbf{\Gamma} \mathbf{u}, \mathbf{v} \rangle = \mathbb{E} \max_{\|\mathbf{u}\|=\|\mathbf{v}\|=1} Y_{\mathbf{u}\mathbf{v}} \leq \mathbb{E} \max_{\|\mathbf{u}\|=\|\mathbf{v}\|=1} X_{\mathbf{u}\mathbf{v}} = \mathbb{E} \max_{\|\mathbf{u}\|=\|\mathbf{v}\|=1} [\langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle]$$

Gordon's Theorem

Theorem 21 (Gordon 1985). **Assume** $\mathbf{X} \in \mathbb{R}^{M \times N}$ and $\mathbf{Y} \in \mathbb{R}^{M \times N}$ are centered, jointly Gaussian matrices whose covariance structures satisfy

$$\begin{cases} \mathbb{E} X_{ij} X_{kl} \leq \mathbb{E} Y_{ij} Y_{kl} & \text{for all } i \neq k \text{ and all } j, \ell \\ \mathbb{E} X_{ij} X_{il} \geq \mathbb{E} Y_{ij} Y_{il} & \text{for all } i \text{ and } j \neq \ell \\ \mathbb{E} X_{ij}^2 = \mathbb{E} Y_{ij}^2 & \text{for all } i, j \end{cases}$$

Then, for all choices of $\lambda_{ij} \in \mathbb{R}$,

$$\mathbb{P} \left(\bigcap_{i=1}^M \bigcup_{j=1}^N \{Y_{ij} > \lambda_{ij}\} \right) \geq \mathbb{P} \left(\bigcap_{i=1}^M \bigcup_{j=1}^N \{X_{ij} > \lambda_{ij}\} \right)$$

In particular,

$$\mathbb{E} \min_i \max_j Y_{ij} \geq \mathbb{E} \min_i \max_j X_{ij}$$

Sources: Joag-Dev et al. 1983; Gordon 1985, 1988; Kahane 1986; Ledoux & Talagrand (Cor. 3.13).

The Gaussian Minimax Theorem

Corollary 22 (Gordon 1985). **Assume**

- $U \subset \mathbb{R}^m$ and $V \subset \mathbb{R}^n$ are compact subsets of the unit sphere
- $\Gamma \in \mathbb{R}^{n \times m}$ and $\mathbf{g} \in \mathbb{R}^m$ and $\mathbf{h} \in \mathbb{R}^n$ are independent standard normal

Then

$$\mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} \langle \Gamma \mathbf{u}, \mathbf{v} \rangle \geq \mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} [\langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle]$$

Sources: Chevet 1977; Gordon 1985, 1988; Ledoux & Talagrand (Thm. 3.20); Davidson & Szarek 2001; Stojnic 2013; Amelunxen et al. 2014; Thrampoulidis et al. 2015; Amelunxen & Lotz 2015.

Example: Minimum Singular Value of a Gaussian Matrix

🐼 **Goal:** For standard normal $\Gamma \in \mathbb{R}^{n \times m}$, bound

$$\mathbb{E} \sigma_{\min}(\Gamma) = \min_{\|u\|=1} \max_{\|v\|=1} \langle \Gamma u, v \rangle$$

🐼 **Gaussian Minimax Theorem:**

$$\begin{aligned} \mathbb{E} \sigma_{\min}(\Gamma) &= \min_{\|u\|=1} \max_{\|v\|=1} \langle \Gamma u, v \rangle \\ &\geq \mathbb{E} \min_{\|u\|=1} \max_{\|v\|=1} [\langle \mathbf{g}, u \rangle + \langle \mathbf{h}, v \rangle] \\ &= \mathbb{E} \|\mathbf{h}\| - \mathbb{E} \|\mathbf{g}\| \\ &\geq \sqrt{n-1} - \sqrt{m} \end{aligned}$$

🐼 **Result is sharp, including constants!**

🐼 **Remark:** Can replace $\sqrt{n-1}$ with \sqrt{n} if you work hard enough

Sources: Marchenko & Pastur 1967; Chevet 1977; Silverstein 1985; Gordon 1985, 1988; Ledoux & Talagrand (Thm. 3.20); Szarek 1991; Bai & Yin 1993; Davidson & Szarek 2001.

Proof of Gaussian Minimax Theorem

Define independent Gaussian processes on $\{(\mathbf{u}, \mathbf{v}) : \mathbf{u} \in U, \mathbf{v} \in V\}$:

$$Y_{\mathbf{u}\mathbf{v}} = \langle \Gamma \mathbf{u}, \mathbf{v} \rangle + \gamma \quad \text{and} \quad X_{\mathbf{u}\mathbf{v}} = \langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle$$

Compute the covariances:

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}} Y_{\mathbf{u}'\mathbf{v}'} = \langle \mathbf{u}, \mathbf{u}' \rangle \langle \mathbf{v}, \mathbf{v}' \rangle + 1 \quad \text{and} \quad \mathbb{E} X_{\mathbf{u}\mathbf{v}} X_{\mathbf{u}'\mathbf{v}'} = \langle \mathbf{u}, \mathbf{u}' \rangle + \langle \mathbf{v}, \mathbf{v}' \rangle$$

Comparison:

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}} Y_{\mathbf{u}'\mathbf{v}'} - \mathbb{E} X_{\mathbf{u}\mathbf{v}} X_{\mathbf{u}'\mathbf{v}'} = (1 - \langle \mathbf{u}, \mathbf{u}' \rangle)(1 - \langle \mathbf{v}, \mathbf{v}' \rangle) \geq 0$$

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}} Y_{\mathbf{u}\mathbf{v}'} - \mathbb{E} X_{\mathbf{u}\mathbf{v}} X_{\mathbf{u}\mathbf{v}'} = 0$$

$$\mathbb{E} Y_{\mathbf{u}\mathbf{v}}^2 - \mathbb{E} X_{\mathbf{u}\mathbf{v}}^2 = 0$$

Apply Gordon's Theorem (on finite subsets of U, V):

$$\begin{aligned} \mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} \langle \Gamma \mathbf{u}, \mathbf{v} \rangle &= \mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} Y_{\mathbf{u}\mathbf{v}} \\ &\geq \mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} X_{\mathbf{u}\mathbf{v}} = \mathbb{E} \min_{\mathbf{u} \in U} \max_{\mathbf{v} \in V} [\langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle] \end{aligned}$$

Kahane's Approach to Gaussian Comparison

Theorem 23 (Kahane 1986). **Assume** $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ are centered, jointly Gaussian vectors. For sets $A, B \subset \{1, \dots, N\}^2$, assume that the covariance structures satisfy

$$\begin{cases} \mathbb{E} X_i X_j \leq \mathbb{E} Y_i Y_j & \text{for } (i, j) \in A \\ \mathbb{E} X_i X_j \geq \mathbb{E} Y_i Y_j & \text{for } (i, j) \in B \\ \mathbb{E} X_i X_j = \mathbb{E} Y_i Y_j & \text{for } (i, j) \notin A \cup B \end{cases}$$

Let $f : \mathbb{R}^N \rightarrow \mathbb{R}$ be a function whose second (distributional) derivative satisfies

$$\begin{cases} \partial_{ij} f \geq 0 & \text{for } (i, j) \in A \\ \partial_{ij} f \leq 0 & \text{for } (i, j) \in B \end{cases}$$

Then

$$\mathbb{E} f(\mathbf{x}) \leq \mathbb{E} f(\mathbf{y})$$

Sources: Joag-Dev et al. 1983; Kahane 1986; Ledoux & Talagrand (Thm. 3.11).

Gaussian Integration by Parts I

Lemma 24 (Univariate Gaussian IBP). *Let $\gamma \in \mathbb{R}$ be a standard normal random variable. For “any” function $f : \mathbb{R} \rightarrow \mathbb{R}$,*

$$\mathbb{E}[\gamma f(\gamma)] = \mathbb{E}[f'(\gamma)]$$

🐼 Calculate:

$$\mathbb{E}[\gamma f(\gamma)] = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} u f(u) e^{-u^2/2} du = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} f'(u) e^{-u^2/2} du = \mathbb{E}[f'(\gamma)]$$

🐼 Sufficient that f absolutely continuous and $f' \in L_1(d\gamma)$

Gaussian Integration by Parts II

Lemma 25 (Gaussian IBP). Let $\mathbf{x} \in \mathbb{R}^n$ be a centered, jointly Gaussian vector with covariance Σ . For “any” function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\mathbb{E}[X_i f(\mathbf{x})] = \sum_{j=1}^n (\Sigma)_{ij} \mathbb{E}[(\partial_j f)(\mathbf{x})]$$

☛ Since $\mathbf{x} \sim \Sigma^{1/2} \mathbf{z}$ for standard normal $\mathbf{z} \in \mathbb{R}^n$,

$$\mathbb{E}[X_i f(\mathbf{x})] = \sum_{k=1}^n (\Sigma^{1/2})_{ik} \mathbb{E}[Z_k f(\Sigma^{1/2} \mathbf{z})] = \sum_{k=1}^n (\Sigma^{1/2})_{ik} \mathbb{E}[Z_k g(\mathbf{z})]$$

☛ By univariate Gaussian IBP,

$$\mathbb{E}[Z_k g(\mathbf{z})] = \mathbb{E}[(\partial_k g)(\mathbf{z})] = \sum_{j=1}^n (\Sigma^{1/2})_{kj} \mathbb{E}[(\partial_j f)(\Sigma^{1/2} \mathbf{z})] = \sum_{j=1}^n (\Sigma^{1/2})_{kj} \mathbb{E}[(\partial_j f)(\mathbf{x})]$$

☛ Therefore,

$$\mathbb{E}[X_i f(\mathbf{x})] = \sum_{j,k=1}^n (\Sigma^{1/2})_{ik} (\Sigma^{1/2})_{kj} \mathbb{E}[(\partial_j f)(\mathbf{x})] = \sum_{j=1}^n \Sigma_{ij} \mathbb{E}[(\partial_j f)(\mathbf{x})]$$

Gaussian Interpolation

Lemma 26 (Gaussian Interpolation). Let $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in \mathbb{R}^n$ be independent, centered, jointly Gaussian vectors with covariances $\Sigma_{\mathbf{x}}$ and $\Sigma_{\mathbf{y}}$. Define

$$\mathbf{z}(t) = \sqrt{t} \mathbf{x} + \sqrt{1-t} \mathbf{y} \quad \text{for } t \in [0, 1]$$

For “any” function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\frac{d}{dt} \mathbb{E}[f(\mathbf{z}(t))] = \frac{1}{2} \sum_{i,j=1}^n ((\Sigma_{\mathbf{x}})_{ij} - (\Sigma_{\mathbf{y}})_{ij}) \mathbb{E}[(\partial_{ij} f)(\mathbf{z}(t))]$$

☞ Calculate:

$$\frac{d}{dt} \mathbb{E}[f(\mathbf{z}(t))] = \sum_{i=1}^n \mathbb{E}[(\partial_i f)(\mathbf{z}(t))(\partial_i \mathbf{z})(t)] = \frac{1}{2} \sum_{i=1}^n \mathbb{E} \left[(\partial_i f)(\mathbf{z}(t)) \left(\frac{1}{\sqrt{t}} X_i - \frac{1}{\sqrt{1-t}} Y_i \right) \right]$$

☞ Apply Gaussian IBP to each term; e.g.,

$$\frac{1}{\sqrt{t}} \mathbb{E}[(\partial_i f)(\mathbf{z}(t)) X_i] = \sum_{j=1}^n (\Sigma_{\mathbf{x}})_{ij} \mathbb{E}[(\partial_{ij} f)(\mathbf{z}(t))]$$

Proof of Kahane's Theorem

🐼 Observe that $f(\mathbf{z}(0)) = f(\mathbf{y})$ and $f(\mathbf{z}(1)) = f(\mathbf{x})$

🐼 By Gaussian interpolation,

$$\frac{d}{dt} \mathbb{E}[f(\mathbf{z}(t))] = \frac{1}{2} \sum_{i,j=1}^n ((\boldsymbol{\Sigma}_x)_{ij} - (\boldsymbol{\Sigma}_y)_{ij}) \mathbb{E}[(\partial_{ij} f)(\mathbf{z}(t))]$$

🐼 By hypothesis,

$$\text{for } (i, j) \in A, \quad (\boldsymbol{\Sigma}_x)_{ij} \leq (\boldsymbol{\Sigma}_y)_{ij} \quad \text{and} \quad \partial_{ij} f \geq 0$$

$$\text{for } (i, j) \in B, \quad (\boldsymbol{\Sigma}_x)_{ij} \geq (\boldsymbol{\Sigma}_y)_{ij} \quad \text{and} \quad \partial_{ij} f \leq 0$$

$$\text{otherwise,} \quad (\boldsymbol{\Sigma}_x)_{ij} = (\boldsymbol{\Sigma}_y)_{ij}$$

🐼 Thus,

$$\frac{d}{dt} \mathbb{E}[f(\mathbf{z}(t))] \leq 0$$

🐼 Conclude: $\mathbb{E} f(\mathbf{x}) \leq \mathbb{E} f(\mathbf{y})$

Proof of Slepian's Lemma

☞ Let $\mathbf{x} \in \mathbb{R}^N$ and $\mathbf{y} \in \mathbb{R}^N$ be centered, jointly Gaussian vectors with

$$\begin{cases} \mathbb{E} X_i X_j \leq \mathbb{E} Y_i Y_j & \text{for all } i \neq j \\ \mathbb{E} X_i^2 = \mathbb{E} Y_i^2 & \text{for all } i \end{cases}$$

☞ Set $A = \{(i, j) : i \neq j\}$ and $B = \emptyset$. Define the function

$$f(\mathbf{w}) = \prod_{i=1}^N \mathbb{1}\{w_i \leq \lambda_i\}$$

☞ For $(i, j) \in A$, compute second derivative:

$$(\partial_i f)(\mathbf{w}) = -\mathbb{1}\{w_i = \lambda_i\} \prod_{j \neq i} \mathbb{1}\{w_j \leq \lambda_j\}$$

$$(\partial_{ij} f)(\mathbf{w}) = \mathbb{1}\{w_i = \lambda_i, w_j = \lambda_j\} \prod_{k \notin \{i, j\}} \mathbb{1}\{w_k \leq \lambda_k\} \geq 0$$

☞ Apply Kahané's Theorem:

$$\mathbb{P}\left(\bigcap_{i=1}^N \{X_i \leq \lambda_i\}\right) = \mathbb{E} f(\mathbf{x}) \leq \mathbb{E} f(\mathbf{y}) = \mathbb{P}\left(\bigcap_{i=1}^N \{Y_i \leq \lambda_i\}\right)$$

Proof of Gordon's Theorem I

• Let $\mathbf{X} \in \mathbb{R}^{M \times N}$ and $\mathbf{Y} \in \mathbb{R}^{M \times N}$ be centered, jointly Gaussian matrices with

$$\begin{cases} \mathbb{E} X_{ij} X_{k\ell} \leq \mathbb{E} Y_{ij} Y_{k\ell} & \text{for all } i \neq k \text{ and all } j, \ell \\ \mathbb{E} X_{ij} X_{i\ell} \geq \mathbb{E} Y_{ij} Y_{i\ell} & \text{for all } i \text{ and } j \neq \ell \\ \mathbb{E} X_{ij}^2 = \mathbb{E} Y_{ij}^2 & \text{for all } i, j \end{cases}$$

• Set $A = \{((i, j), (k, \ell)) : i \neq k\}$ and $B = \{((i, j), (k, \ell)) : i = k, j \neq \ell\}$. Define

$$f(\mathbf{W}) = \prod_{i=1}^N \left[1 - \prod_{j=1}^M \mathbb{1}\{w_{ij} \leq \lambda_{ij}\} \right]$$

• Compute first derivative:

$$(\partial_{(i,j)} f)(\mathbf{W}) = \mathbb{1}\{w_{ij} = \lambda_{ij}\} \prod_{j' \neq j} \mathbb{1}\{w_{ij'} \leq \lambda_{ij'}\} \prod_{i' \neq i} \left[1 - \prod_{j'} \mathbb{1}\{w_{i'j'} \leq \lambda_{i'j'}\} \right]$$

Proof of Gordon's Theorem II

☛ For $((i, j), (k, \ell)) \in B$, compute second derivative:

$$\begin{aligned}
 (\partial_{(i,j),(k,\ell)} f)(\mathbf{W}) &= -\mathbb{1}\{w_{ij} = \lambda_{ij}, w_{i\ell} = \lambda_{i\ell}\} \\
 &\quad \times \prod_{j' \notin \{j, \ell\}} \mathbb{1}\{w_{ij'} \leq \lambda_{ij'}\} \prod_{i' \neq i} \left[1 - \prod_{j'} \mathbb{1}\{w_{i'j'} \leq \lambda_{i'j'}\} \right] \leq 0
 \end{aligned}$$

☛ For $((i, j), (k, \ell)) \in A$, compute second derivative:

$$\begin{aligned}
 (\partial_{(i,j),(k,\ell)} f)(\mathbf{W}) &= \mathbb{1}\{w_{ij} = \lambda_{ij}, w_{k\ell} = \lambda_{k\ell}\} \\
 &\quad \times \prod_{j' \neq j} \mathbb{1}\{w_{ij'} \leq \lambda_{ij'}\} \prod_{j' \neq \ell} \mathbb{1}\{w_{kj'} \leq \lambda_{kj'}\} \prod_{i' \notin \{i, k\}} \left[1 - \prod_{j'} \mathbb{1}\{w_{i'j'} \leq \lambda_{i'j'}\} \right] \geq 0
 \end{aligned}$$

☛ Apply Kahané's Theorem:

$$\mathbb{P} \left(\bigcap_{i=1}^N \left(\bigcap_{j=1}^M \{X_{ij} \leq \lambda_{ij}\} \right)^c \right) = \mathbb{E} f(\mathbf{x}) \leq \mathbb{E} f(\mathbf{y}) = \mathbb{P} \left(\bigcap_{i=1}^N \left(\bigcap_{j=1}^M \{Y_{ij} \leq \lambda_{ij}\} \right)^c \right)$$

Gaussian Concentration

Theorem 27. Assume

☞ $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a Lipschitz function: $|f(\mathbf{x}) - f(\mathbf{y})| \leq L \cdot \|\mathbf{x} - \mathbf{y}\|$ for all \mathbf{x}, \mathbf{y}

☞ $Z = f(\mathbf{g})$ where \mathbf{g} is standard normal

Then

$$\text{Var}[Z] = \mathbb{E}(Z - \mathbb{E}Z)^2 \leq L^2 \quad (\text{Poincaré})$$

$$\mathbb{P}\{|Z - \mathbb{E}Z| \geq t\} \leq e^{-t^2/(2L^2)} \quad \text{for all } t \geq 0 \quad (\text{concentration})$$

Sources: Ledoux & Talagrand (Sec. 1.1); Bogachev (Sec. 1.7); Boucheron et al. (Sec. 3.7, 5.4).

Pisier's Approach to Gaussian Concentration

Theorem 28 (Pisier 1986). *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is an L -Lipschitz function
- $\mathbf{x} \in \mathbb{R}^d$ and $\gamma \in \mathbb{R}$ are standard normal
- $\Phi : \mathbb{R} \rightarrow \mathbb{R}$ is convex
- $t \mapsto \Phi(t) + \Phi(-t)$ is increasing on \mathbb{R}_+

Then

$$\mathbb{E} \Phi(f(\mathbf{x}) - \mathbb{E} f(\mathbf{x})) \leq \mathbb{E} \Phi\left(\frac{1}{2}\pi L \gamma\right)$$

In particular,

$$\begin{aligned} \Phi(t) = t^2 : & \quad \text{Var}[f(\mathbf{x})] \leq \frac{1}{4}\pi^2 L^2 \\ \Phi(t) = e^{\theta t} : & \quad \log \mathbb{E} e^{\theta(f(\mathbf{x}) - \mathbb{E} f(\mathbf{x}))} \leq \frac{1}{8}\pi^2 L^2 \theta^2 \quad \text{for } \theta \in \mathbb{R} \end{aligned}$$

Sources: Pisier 1986; Ledoux & Talagrand (Eqn. (1.5)).

Proof of Pisier's Theorem

- Let \mathbf{x}, \mathbf{y} be independent standard normal variables
- Define $\mathbf{z}(\alpha) = \cos(\alpha)\mathbf{x} + \sin(\alpha)\mathbf{y}$ and $\mathbf{z}'(\alpha) = -\sin(\alpha)\mathbf{x} + \cos(\alpha)\mathbf{y}$
- For fixed α , the vectors $\mathbf{z}(\alpha)$ and $\mathbf{z}'(\alpha)$ are independent standard normal
- Calculate:

$$\begin{aligned}\mathbb{E}\Phi(f(\mathbf{x}) - \mathbb{E}f(\mathbf{y})) &\leq \mathbb{E}\Phi(f(\mathbf{x}) - f(\mathbf{y})) \\ &= \mathbb{E}\Phi\left(\int_0^{\pi/2} \langle \nabla f(\mathbf{z}(\alpha)), \mathbf{z}'(\alpha) \rangle d\alpha\right) \\ &\leq \frac{2}{\pi} \int_0^{\pi/2} \mathbb{E}\Phi\left(\frac{\pi}{2} \langle \nabla f(\mathbf{z}(\alpha)), \mathbf{z}'(\alpha) \rangle\right) d\alpha \\ &= \mathbb{E}\Phi\left(\frac{\pi}{2} \langle \nabla f(\mathbf{x}), \mathbf{y} \rangle\right) \\ &= \mathbb{E}\Phi\left(\frac{\pi}{2} \|\nabla f(\mathbf{x})\| \cdot \gamma\right) \\ &\leq \mathbb{E}\Phi\left(\frac{\pi}{2} L \cdot \gamma\right)\end{aligned}$$

Signal Recovery

Convex Signal Recovery

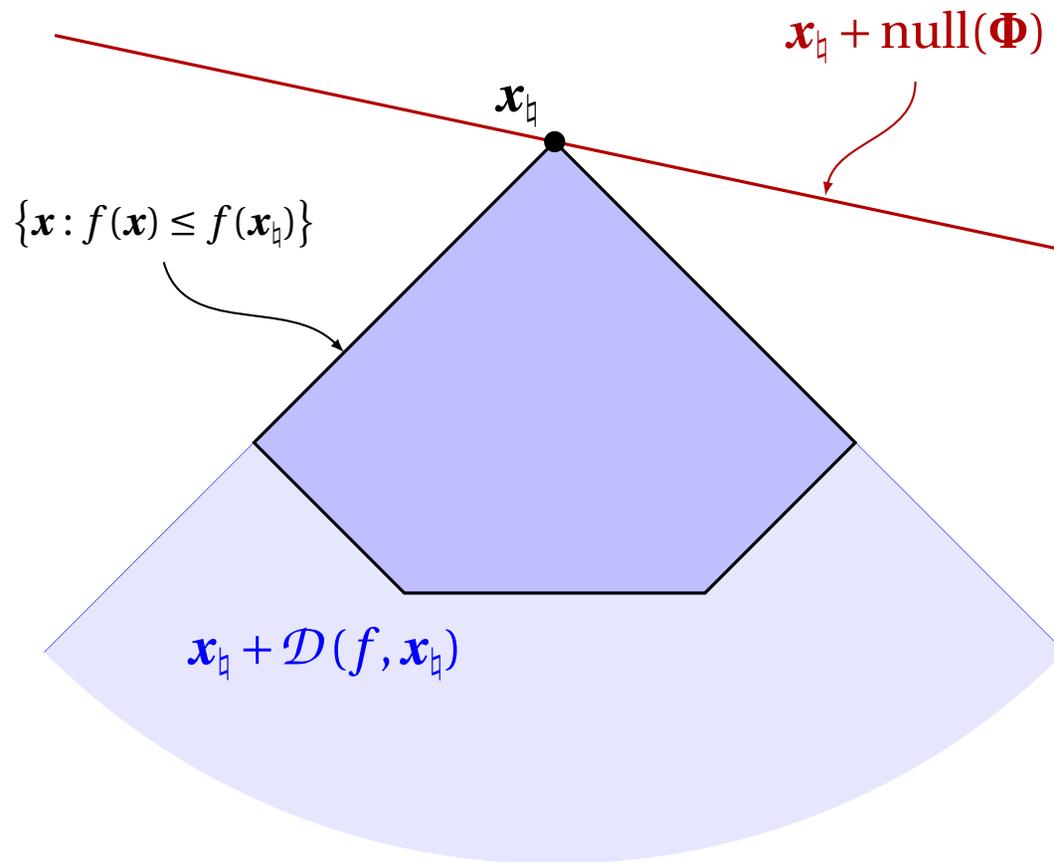
- ☞ Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be a convex structural penalty (e.g., an atomic gauge)
- ☞ Let $\mathbf{x}_\natural \in \mathbb{R}^d$ be “structured” but unknown
- ☞ Let $\Phi \in \mathbb{R}^{m \times d}$ be a known measurement matrix
- ☞ Observe $\mathbf{z} = \Phi \mathbf{x}_\natural \in \mathbb{R}^m$
- ☞ Find estimate $\hat{\mathbf{x}}$ by solving convex program

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \Phi \mathbf{x} = \mathbf{z}$$

☞ **Hope:** $\hat{\mathbf{x}} = \mathbf{x}_\natural$

Sources: Chen et al. 1997, 2001; Chandrasekaran et al. 2012, McCoy & Tropp 2013; Oymak et al. 2013; Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Geometry of Convex Signal Recovery



Analysis of Convex Signal Recovery: Success

Proposition 29 (Geometric Formulation). *Convex signal recovery **succeeds** ($\hat{\mathbf{x}} = \mathbf{x}_q$) if and only if*

$$\mathcal{D}(f, \mathbf{x}_q) \cap \text{null}(\Phi) = \{\mathbf{0}\}$$

Proposition 30 (Analytic Condition for Success). *Convex signal recovery **succeeds** if*

$$\sigma_{\min}(\Phi; K) = \inf_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x} \in K}} \|\Phi \mathbf{x}\| > 0 \quad (\text{“Minimum conic singular value”})$$

where $K = \mathcal{D}(f, \mathbf{x}_q)$

Sources: Rudelson & Vershynin 2006; Stojnic 2009, 2013; Oymak et al. 2010; Chandrasekaran et al. 2012; McCoy & Tropp 2013; Oymak et al. 2013; Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Gaussian Measurements

Proposition 31 (Gaussian Measurements). **Assume** $\Gamma \in \mathbb{R}^{m \times d}$ is a standard normal matrix. **Then** $\text{null}(\Gamma)$ is a uniformly distributed subspace of \mathbb{R}^d with codimension $m \wedge d$, almost surely.

Sources: Donoho 2006; Candès & Tao 2006; Rudelson & Vershynin 2006; Stojnic 2009, 2013; Donoho & Tanner 2009; Recht et al. 2010; Oymak et al. 2010; Chandrasekaran et al. 2012; McCoy & Tropp 2013; Oymak et al. 2013; Amelunxen et al. 2014; Goldstein et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Minimum Conic Singular Value of a Gaussian Matrix

Proposition 32. Assume

- ☛ K is a convex cone in \mathbb{R}^d
- ☛ $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal

Then

$$\mathbb{E} \sigma_{\min}(\mathbf{\Gamma}; K) \geq \sqrt{m-1} - \sqrt{\delta(K)}$$

- ☛ Write $\sigma_{\min}(\mathbf{\Gamma}; K) = \inf_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \max_{\|\mathbf{v}\|=1} \langle \mathbf{\Gamma} \mathbf{u}, \mathbf{v} \rangle$

- ☛ Gaussian Minimax Theorem:

$$\begin{aligned} \mathbb{E} \sigma_{\min}(\mathbf{\Gamma}; K) &\geq \mathbb{E} \inf_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \max_{\|\mathbf{v}\|=1} [\langle \mathbf{g}, \mathbf{u} \rangle + \langle \mathbf{h}, \mathbf{v} \rangle] \\ &\geq \sqrt{m-1} - \mathbb{E} \sup_{\|\mathbf{u}\|=1, \mathbf{u} \in K} \langle \mathbf{g}, \mathbf{u} \rangle \geq \sqrt{m-1} - \sqrt{\delta(K)} \end{aligned}$$

Sources: Gordon 1985, 1988; Rudelson & Vershynin 2006; Stojnic 2009, 2013; Oymak et al. 2010; Chandrasekaran et al. 2012; Oymak et al. 2013; Amelunxen et al. 2014, Thrampoulidis et al. 2014–2016; Tropp 2015.

Concentration of Minimum Conic Singular Value

Proposition 33. *Assume* K is a convex cone, and let Γ be standard normal. *Then*

$$\mathbb{P}\{\sigma_{\min}(\Gamma; K) < \mathbb{E}\sigma_{\min}(\Gamma; K) - t\} \leq e^{-t^2/2}$$

🐼 Bound the Lipschitz constant of $\sigma_{\min}(\cdot; K)$:

$$\inf_{\|u\|=1, u \in K} \|\Gamma u\| - \inf_{\|u\|=1, u \in K} \|\Gamma' u\| \leq \|\Gamma u'\| - \|\Gamma' u'\| \leq \|(\Gamma - \Gamma')u'\| \leq \|\Gamma - \Gamma'\|_F$$

where u' is a near-minimizer of the second term

🐼 Apply Gaussian concentration with $L = 1$

Sources: Rudelson & Vershynin 2006; Stojnic 2009, 2013; Oymak et al. 2010; Chandrasekaran et al. 2012; Oymak et al. 2013; Amelunxen et al. 2014, Thrampoulidis et al. 2013–2016; Tropp 2015.

Gaussian Measurements: Success

Theorem 34 (Chandrasekaran et al. 2012). *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_\dagger \in \text{dom}(f)$
- Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_\dagger$ where $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

Then

$$m \geq \delta(K) + C\sqrt{d} \quad \text{implies} \quad \hat{\mathbf{x}} = \mathbf{x}_\dagger \quad \text{with high probability}$$

where $K = \mathcal{D}(f, \mathbf{x}_\dagger)$

Sources: Rudelson & Vershynin 2006; Stojnic 2009, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Gaussian Measurements: Success Proof

☞ Let $K = \mathcal{D}(f, \mathbf{x}_q)$

☞ Combine last two results:

$$\begin{aligned} e^{-t^2/2} &\geq \mathbb{P}\{\sigma_{\min}(\mathbf{\Gamma}; K) < \mathbb{E}\sigma_{\min}(\mathbf{\Gamma}; K) - t\} \\ &\geq \mathbb{P}\left\{\sigma_{\min}(\mathbf{\Gamma}; K) < \sqrt{m-1} - \sqrt{\delta(K)} - t\right\} \end{aligned}$$

☞ Set $t = 3$ to achieve probability less than 2%

☞ Success ($\sigma_{\min}(\mathbf{\Gamma}; K) > 0$): $\sqrt{m-1} - \sqrt{\delta(K)} - 3 > 0$

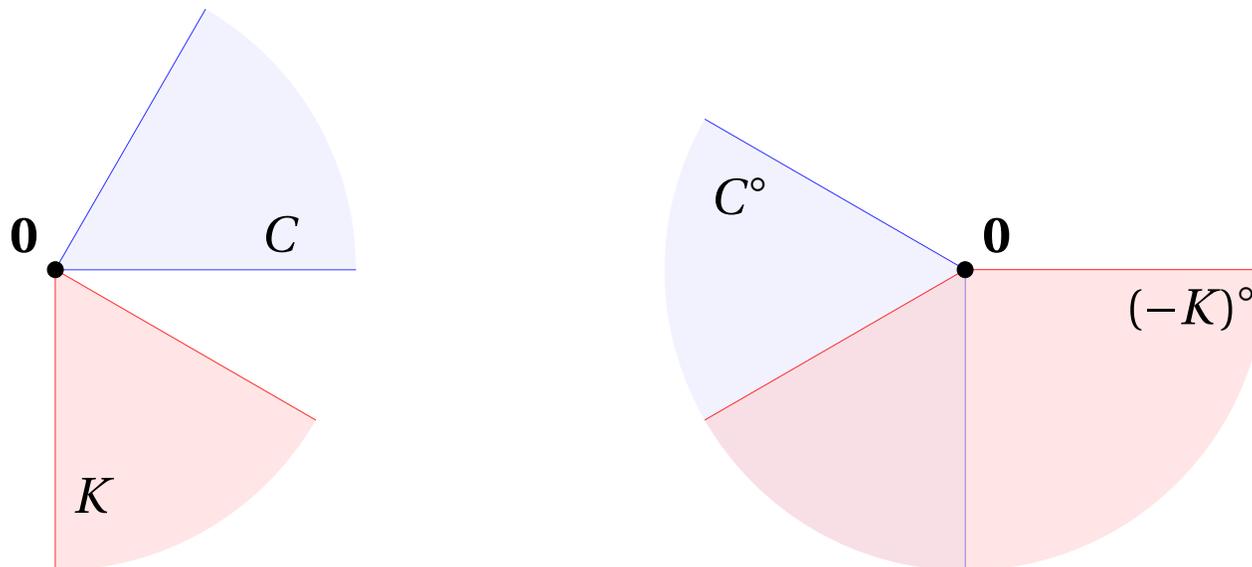
☞ Equivalently, $m > \delta(K) + 6\sqrt{\delta(K)} + 10$

☞ Use upper bound $\delta(K) \leq d$

Separation of Convex Cones

Theorem 35 (Klee 1955). **Assume** C and K are convex cones in \mathbb{R}^d , one of which is not a subspace. **Then**

$$C \cap K = \{\mathbf{0}\} \quad \text{implies} \quad C^\circ \cap (-K)^\circ \neq \{\mathbf{0}\}$$



Source: Klee 1955.

Analysis of Convex Signal Recovery: Failure

Proposition 36 (Geometric Formulation). *Convex signal recovery **fails** (i.e., \mathbf{x}_\dagger is not the unique solution) if and only if*

$$\mathcal{D}(f, \mathbf{x}_\dagger) \cap \text{null}(\Phi) \neq \{\mathbf{0}\}$$

If $\mathcal{D}(f, \mathbf{x}_\dagger)$ is not a subspace, a sufficient condition for failure is

$$(\mathcal{D}(f, \mathbf{x}_\dagger))^\circ \cap \text{null}(\Phi)^\circ = \{\mathbf{0}\}$$

Proposition 37 (Analytic Condition for Failure). *Let Ψ be a matrix with $\text{null}(\Psi) = \text{null}(\Phi)^\circ$. Convex signal recovery **fails** if*

$$\sigma_{\min}(\Psi; K^\circ) = \min_{\substack{\|\mathbf{x}\|=1 \\ \mathbf{x} \in K^\circ}} \|\Psi \mathbf{x}\| > 0$$

where $K = \mathcal{D}(f, \mathbf{x}_\dagger)$

Sources: Stojnic 2013; McCoy & Tropp 2013; Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Gaussian Measurements: Failure

Theorem 38 (Amelunxen et al. 2014). *Assume*

- ☞ $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_q \in \text{dom}(f)$
- ☞ Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_q$ where $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal
- ☞ The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

Then

$m \leq \delta(K) - C\sqrt{d}$ implies \mathbf{x}_q is not the unique solution with high prob.

where $K = \mathcal{D}(f, \mathbf{x}_q)$

Sources: Stojnic 2013; Oymak et al. 2013; Amelunxen et al. 2014; Foygel & Mackey 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Gaussian Measurements: Failure Proof

- ☛ Can assume $m < d$ or else failure with probability zero
- ☛ If $\mathcal{D}(f, \mathbf{x}_\dagger)$ is a subspace, just count dimensions
- ☛ Let $\mathbf{\Gamma}^\circ \in \mathbb{R}^{(d-m) \times d}$ be standard normal
- ☛ $\text{null}(\mathbf{\Gamma}^\circ)$ has same distribution as $\text{null}(\mathbf{\Gamma})^\circ$ (a unif. rdm subspace, codim m)
- ☛ Let $K = \mathcal{D}(f, \mathbf{x}_\dagger)$. As before,

$$\begin{aligned}
 e^{-t^2/2} &\geq \mathbb{P} \{ \sigma_{\min}(\mathbf{\Gamma}^\circ; K^\circ) < \mathbb{E} \sigma_{\min}(\mathbf{\Gamma}^\circ; K^\circ) - t \} \\
 &\geq \mathbb{P} \left\{ \sigma_{\min}(\mathbf{\Gamma}^\circ; K^\circ) < \sqrt{d-m-1} - \sqrt{\delta(K^\circ)} - t \right\}
 \end{aligned}$$

☛ Failure ($\sigma_{\min}(\mathbf{\Gamma}^\circ; K^\circ) > 0$): $\sqrt{d-m-1} - \sqrt{\delta(K^\circ)} - 3 > 0$

☛ Equivalently, $d - m > \delta(K^\circ) + 6\sqrt{\delta(K^\circ)} + 10$

☛ Use facts $\delta(K^\circ) = d - \delta(K) \leq d$

Gaussian Measurements: Summary

Theorem 39 (Chandrasekaran et al. 2012; Amelunxen et al. 2014). **Assume**

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_\dagger \in \text{dom}(f)$
- Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_\dagger$ where $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

Then

$m \leq \delta(K) - C\sqrt{d}$ implies **success** with high probability

$m \geq \delta(K) + C\sqrt{d}$ implies **failure** with high probability

where $K = \mathcal{D}(f, \mathbf{x}_\dagger)$

Sources: Rudelson & Vershynin 2006; Stojnic 2009, 2013; Oymak et al. 2010, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Foygel & Mackey 2014; Thrampoulidis et al. 2014–2016; Tropp 2015.

Gaussian Measurements: Improved

Theorem 40 (Amelunxen et al. 2014). *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_q \in \text{dom}(f)$
- Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_q$ where $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

Then

$m \leq \delta(K) - C\sqrt{\delta(K) \wedge \delta(K^\circ)}$ implies **success** with high prob.

$m \geq \delta(K) + C\sqrt{\delta(K) \wedge \delta(K^\circ)}$ implies **failure** with high prob.

where $K = \mathcal{D}(f, \mathbf{x}_q)$

Sources: Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Goldstein et al. 2017.

Example: ℓ_1 Minimization

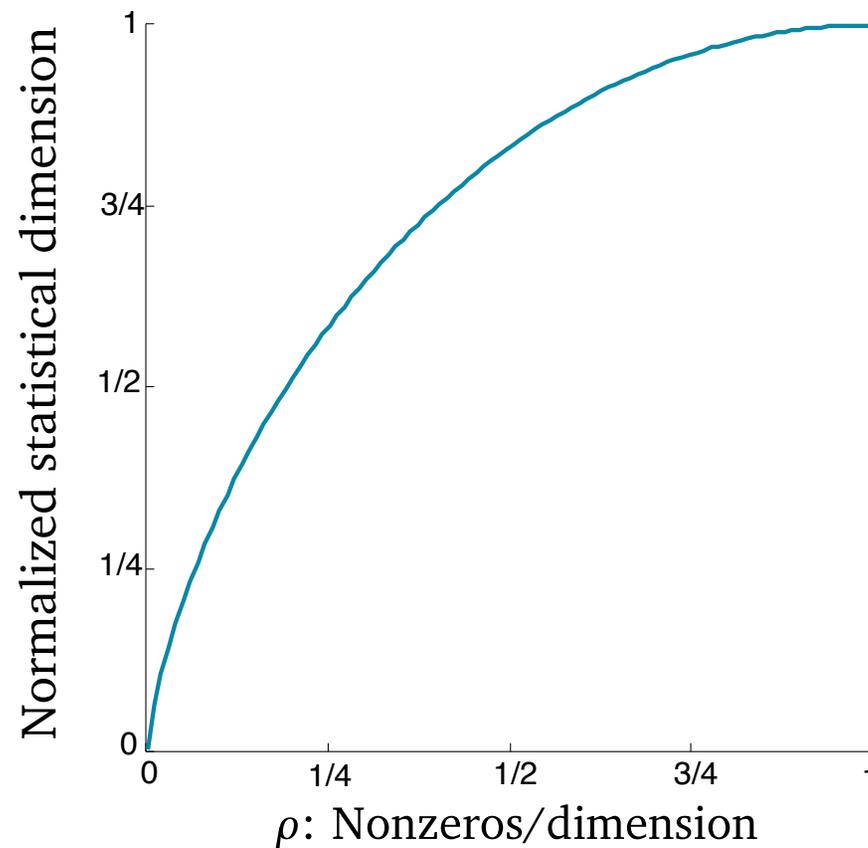
- Suppose $\mathbf{x}_\dagger \in \mathbb{R}^d$ has s nonzero entries
- Let $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ be a standard normal matrix
- Observe $\mathbf{z} = \mathbf{\Gamma} \mathbf{x}_\dagger$
- Find estimate $\hat{\mathbf{x}}$ by solving convex program

$$\text{minimize } \|\mathbf{x}\|_{\ell_1} \quad \text{subject to } \mathbf{\Gamma} \mathbf{x} = \mathbf{z}$$

- Hope:** $\hat{\mathbf{x}} = \mathbf{x}_\dagger$

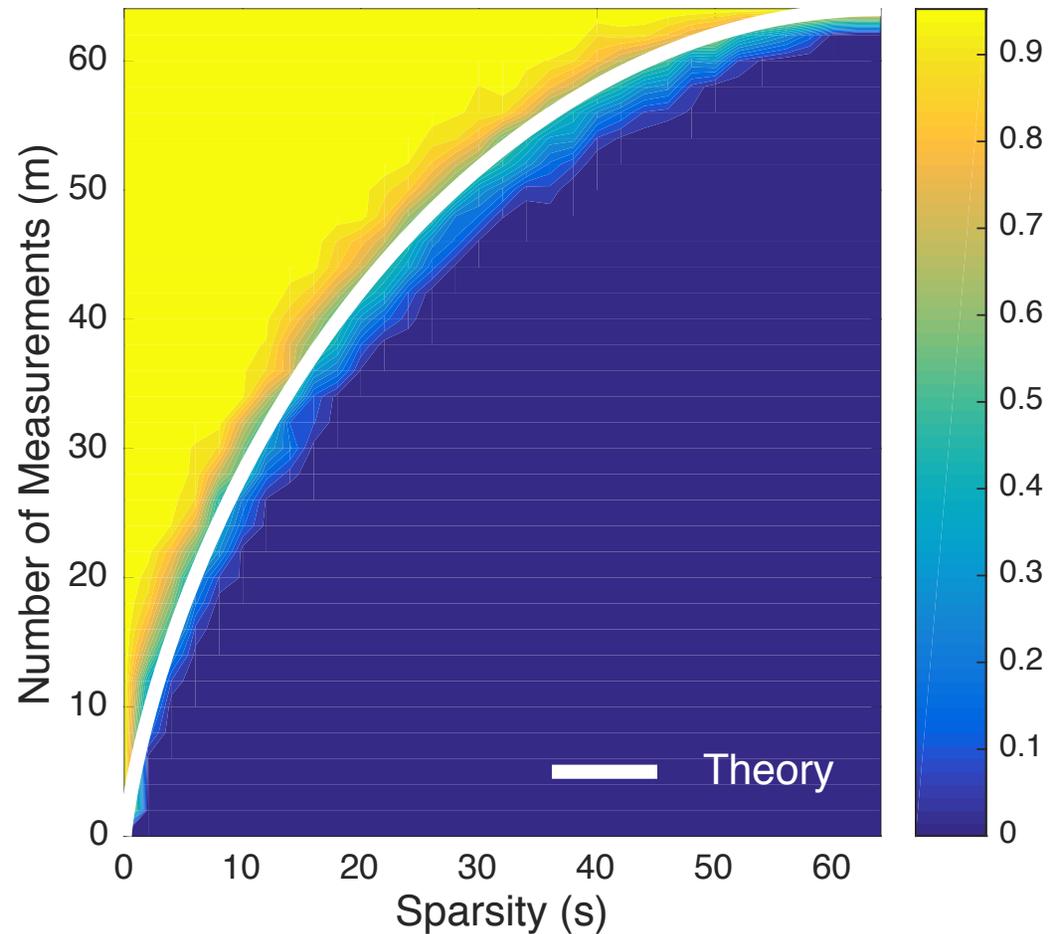
Sources: Donoho 2004, 2006; Candès & Tao 2006; Rudelson & Vershynin 2006; Donoho & Tanner 2009; Stojnic 2009, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Foygel & Mackey 2014; Goldstein et al. 2017.

ℓ_1 Statistical Dimension Curve

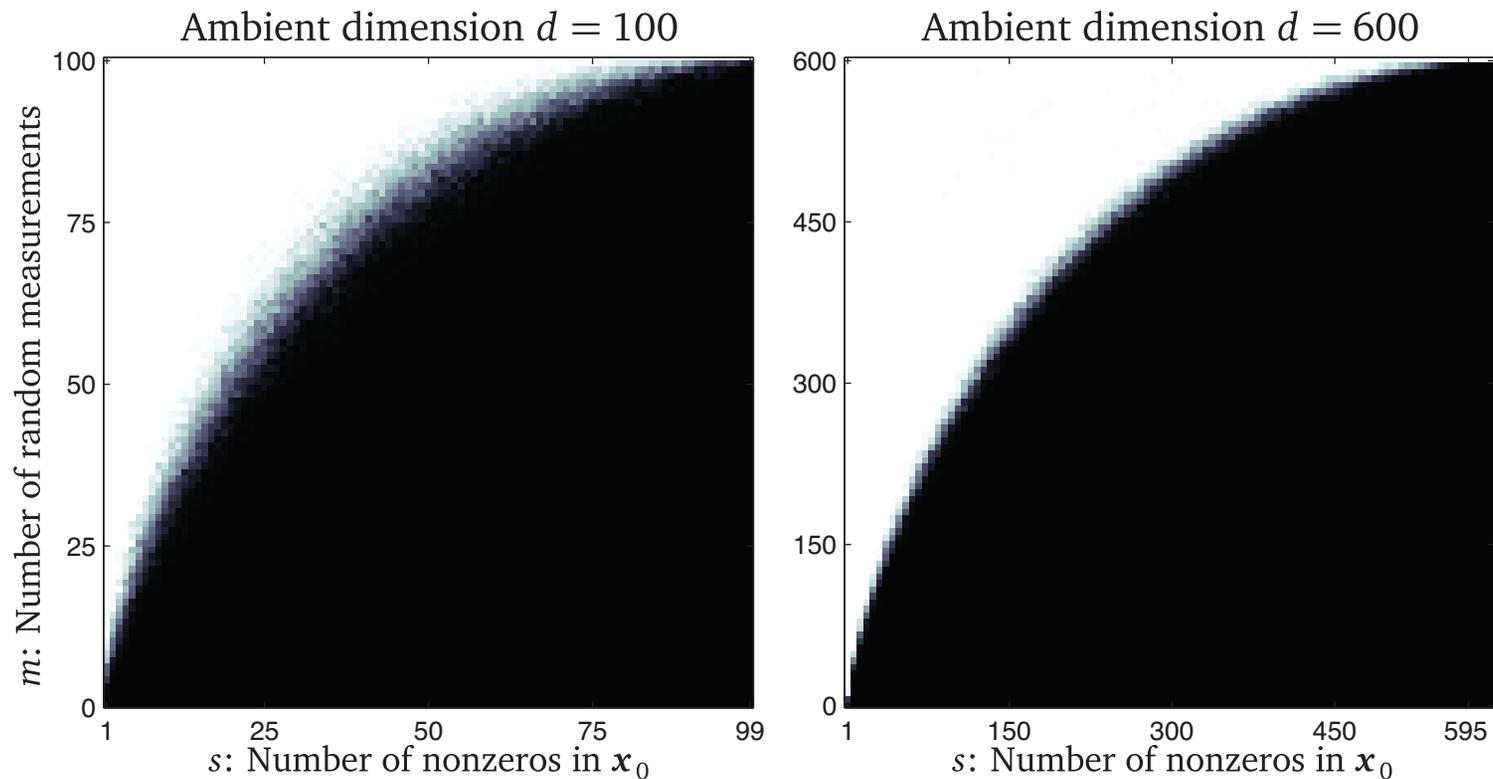


Sources: Affentranger & Schneider 1992; Betke & Henk 1993; Böröczky & Henk 1999; Donoho 2006; Donoho & Tanner 2009; Stojnic 2009, 2013; Chandrasekaran et al. 2012; McCoy & Tropp 2013; Amelunxen et al. 2014; Foygel & Mackey 2014.

Example: Performance of ℓ_1 Minimization



Example: Emergence of ℓ_1 Phase Transition



Example: S_1 Minimization

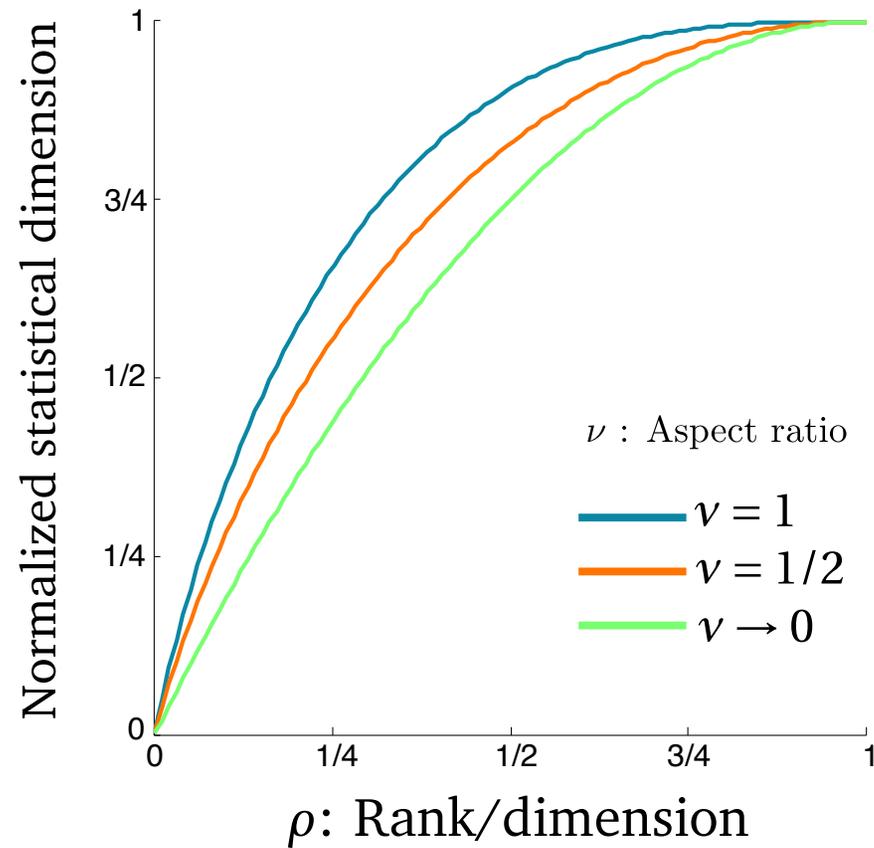
- Suppose $\mathbf{X}_\natural \in \mathbb{R}^{d_1 \times d_2}$ has rank r
- Let $\mathbf{\Gamma} \in \mathbb{R}^{m \times (d_1 \times d_2)}$ be a standard normal matrix
- Observe $\mathbf{z} = \mathbf{\Gamma}(\text{vec } \mathbf{X}_\natural)$
- Find estimate $\hat{\mathbf{X}}$ by solving convex program

$$\text{minimize } \|\mathbf{X}\|_{S_1} \quad \text{subject to } \mathbf{\Gamma}(\text{vec } \mathbf{X}) = \mathbf{z}$$

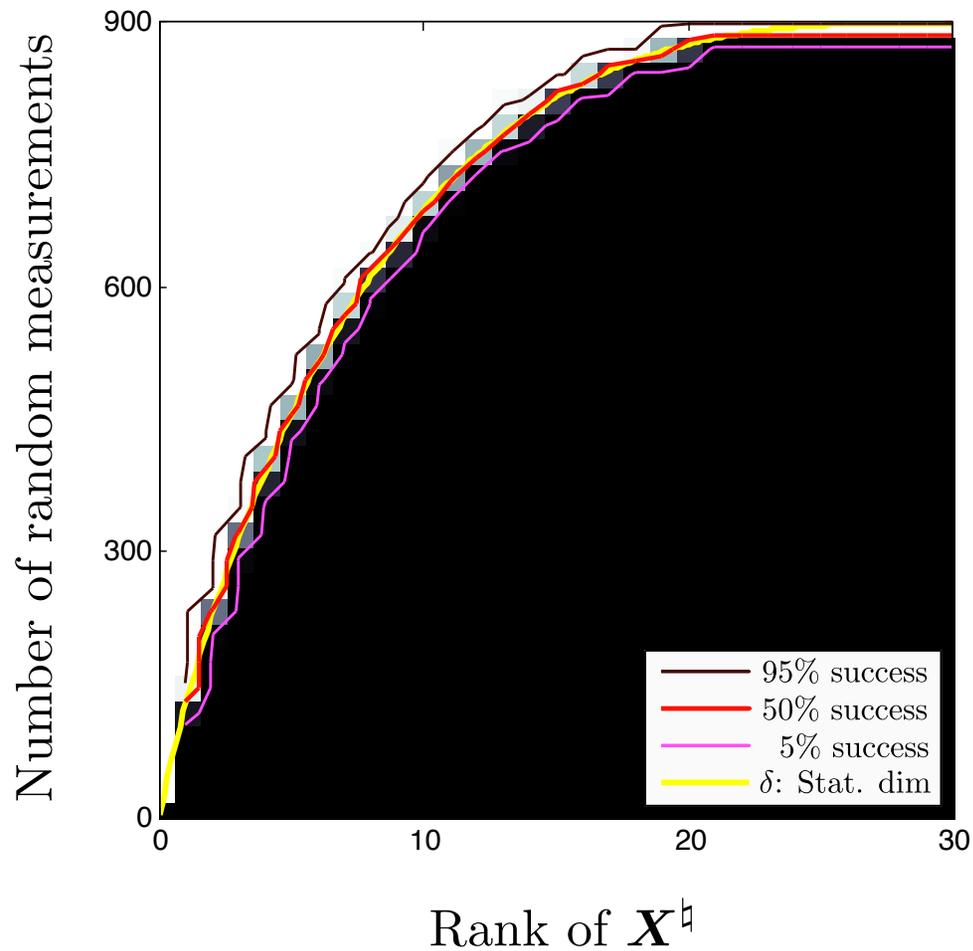
Hope: $\hat{\mathbf{X}} = \mathbf{X}_\natural$

Sources: Fazel 2002; Recht et al. 2010; Oymak et al. 2010, 2013; Chandrasekaran et al. 2012; Amelunxen et al. 2014; Thrampoulidis et al. 2014–2016; Tropp 2015; Goldstein et al. 2017.

S_1 Statistical Dimension Curve



Example: Performance of S_1 Minimization



Desserts + Digestifs

Complements

- 🦉 Universality
- 🦉 Signal recovery with noise
- 🦉 Sharp analysis for Gaussian signal recovery with Gaussian noise
- 🦉 Non-Gaussian measurements
- 🦉 Demixing
- 🦉 ...

Universality I

Theorem 41 (Oymak & Tropp 2015). *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_q \in \text{dom}(f)$
- $\Phi \in \mathbb{R}^{m \times d}$ has iid standardized, symmetric entries with 4+ moments
- Observe $\mathbf{z} = \Phi \mathbf{x}_q$
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \Phi \mathbf{x} = \mathbf{z}$$

Then

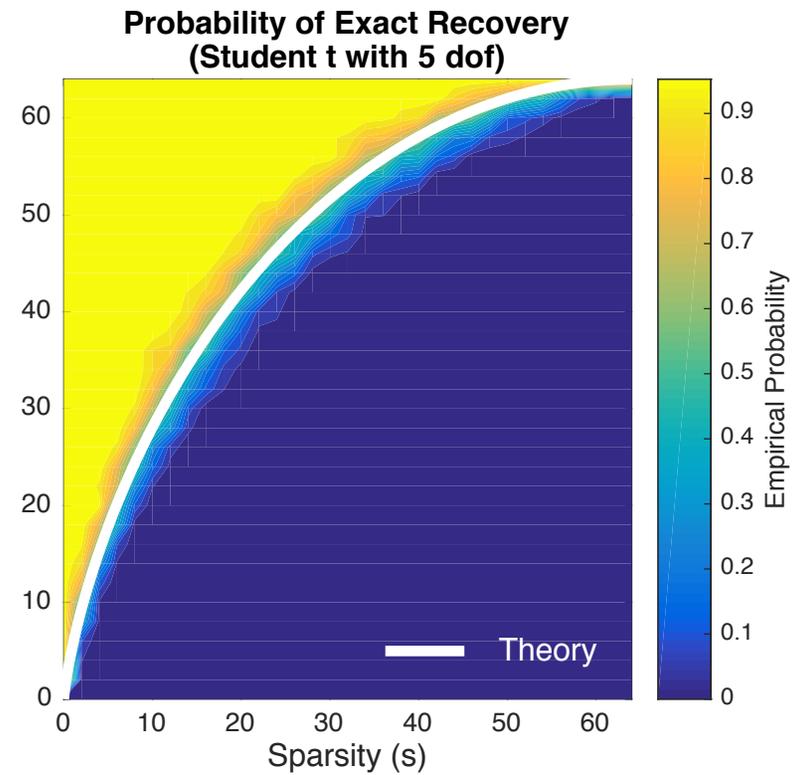
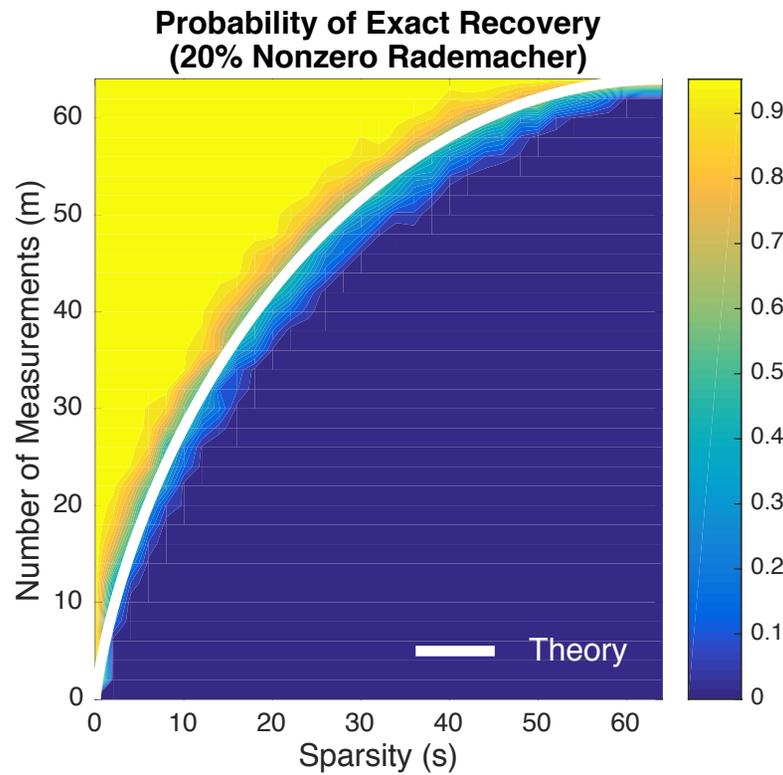
$m \leq \delta(K) - o(d)$ implies **success** with high probability

$m \geq \delta(K) + o(d)$ implies **failure** with high probability

where $K = \mathcal{D}(f, \mathbf{x}_q)$

Sources: Donoho & Tanner 2009; Bayati et al. 2015; Oymak & Tropp 2015.

Universality II



$$\text{minimize } \|\mathbf{x}\|_{\ell_1} \quad \text{subject to } \Phi \mathbf{x} = \Phi \mathbf{x}_s$$

Signal Recovery with Noise I

Theorem 42. *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_\dagger \in \text{dom}(f)$
- The matrix $\Phi \in \mathbb{R}^{m \times d}$
- Observe $\mathbf{z} = \Phi \mathbf{x}_\dagger + \mathbf{e}$ where $\|\mathbf{e}\| \leq \eta$
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \|\Phi \mathbf{x} - \mathbf{z}\| \leq \eta$$

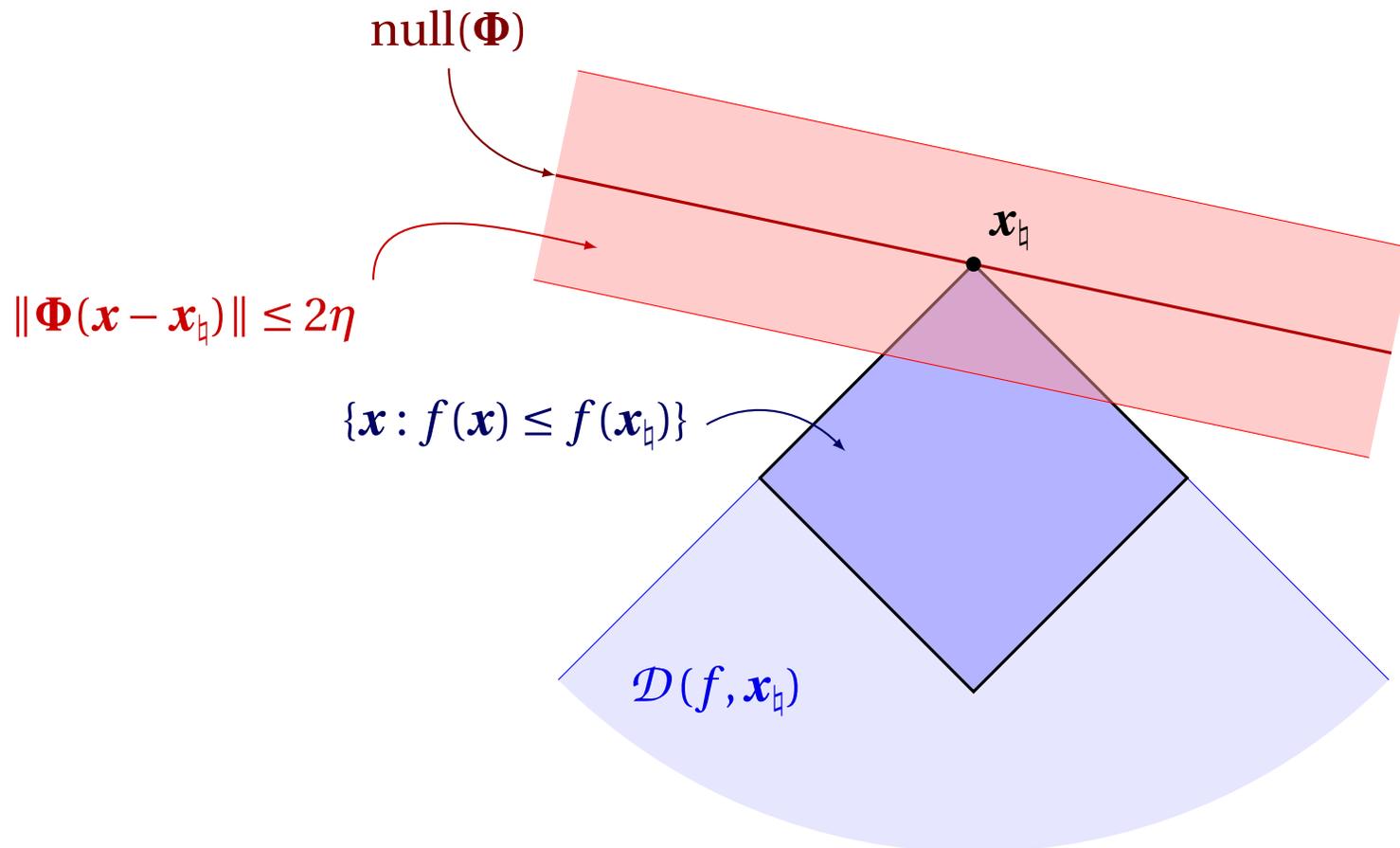
Then

$$\|\hat{\mathbf{x}} - \mathbf{x}_\dagger\| \leq \frac{2\eta}{\sigma_{\min}(\Phi; K)}$$

where $K = \mathcal{D}(f, \mathbf{x}_\dagger)$

Sources: Candès et al. 2006; Chandrasekaran et al. 2012; Tropp 2015.

Signal Recovery with Noise II



Signal Recovery with Noise III

Theorem 43 (Oymak et al. 2013). *Assume*

- $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is a convex function and $\mathbf{x}_\dagger \in \text{dom}(f)$
- Let $\mathbf{\Gamma} \in \mathbb{R}^{m \times d}$ is standard normal
- Observe $\mathbf{z} = \mathbf{\Phi} \mathbf{x}_\dagger + \eta \mathbf{g}$ where \mathbf{g} is standard normal
- The vector $\hat{\mathbf{x}}$ solves

$$\text{minimize } \|\mathbf{\Phi} \mathbf{x} - \mathbf{z}\|^2 \quad \text{subject to } f(\mathbf{x}) \leq f(\mathbf{x}_\dagger)$$

Then (roughly)

$$\sup_{\eta > 0} \frac{\mathbb{E} \|\hat{\mathbf{x}} - \mathbf{x}_\dagger\|^2}{\eta^2} = \lim_{\eta \downarrow 0} \frac{\mathbb{E} \|\hat{\mathbf{x}} - \mathbf{x}_\dagger\|^2}{\eta^2} = \frac{m}{m - \delta(K)}$$

where $K = \mathcal{D}(f, \mathbf{x}_\dagger)$

Sources: Oymak et al. 2013; Thrampoulidis et al. 2014–2016.

Non-Gaussian Measurements

Proposition 44 (Mendelson 2013). **Assume**

• The rows of $\Phi \in \mathbb{R}^{m \times d}$ are iid copies of $\varphi \in \mathbb{R}^d$

• K is a convex cone

• Define the small ball probability

$$Q = \inf_{\|u\| \leq 1, u \in K} \mathbb{P} \{ |\langle u, \varphi \rangle| \geq 1/6 \}$$

• For independent Rademacher variables $\{\varepsilon_i\}$, define the mean empirical width

$$W_m = \mathbb{E} \sup_{\|u\|=1, u \in K} \left\langle u, \frac{1}{\sqrt{m}} \sum_{i=1}^m \varepsilon_i \varphi_i \right\rangle$$

Then, for all $t > 0$,

$$\sigma_{\min}(\Phi; K) \geq \frac{1}{3} \sqrt{m} Q - 2W_m - t \quad \text{with prob.} \quad \geq 1 - e^{-18t^2}$$

Sources: Mendelson et al. 2013–2016; Tropp 2015.

Demixing I

Theorem 45 (Amelunxen et al. 2014). **Assume**

- $f, g : \mathbb{R}^d \rightarrow \mathbb{R}$ are convex functions
- $\mathbf{x}_\natural \in \text{dom}(f)$ and $\mathbf{y}_\natural \in \text{dom}(g)$
- Observe $\mathbf{z} = \mathbf{x}_\natural + \mathbf{Q}\mathbf{y}_\natural$ where $\mathbf{Q} \in \mathbb{R}^{d \times d}$ is random orthogonal
- The pair $(\hat{\mathbf{x}}, \hat{\mathbf{y}})$ solves

$$\text{minimize } f(\mathbf{x}) \quad \text{subject to } \mathbf{z} = \mathbf{x} + \mathbf{Q}\mathbf{y}, \quad g(\mathbf{y}) \leq g(\mathbf{y}_\natural)$$

Then

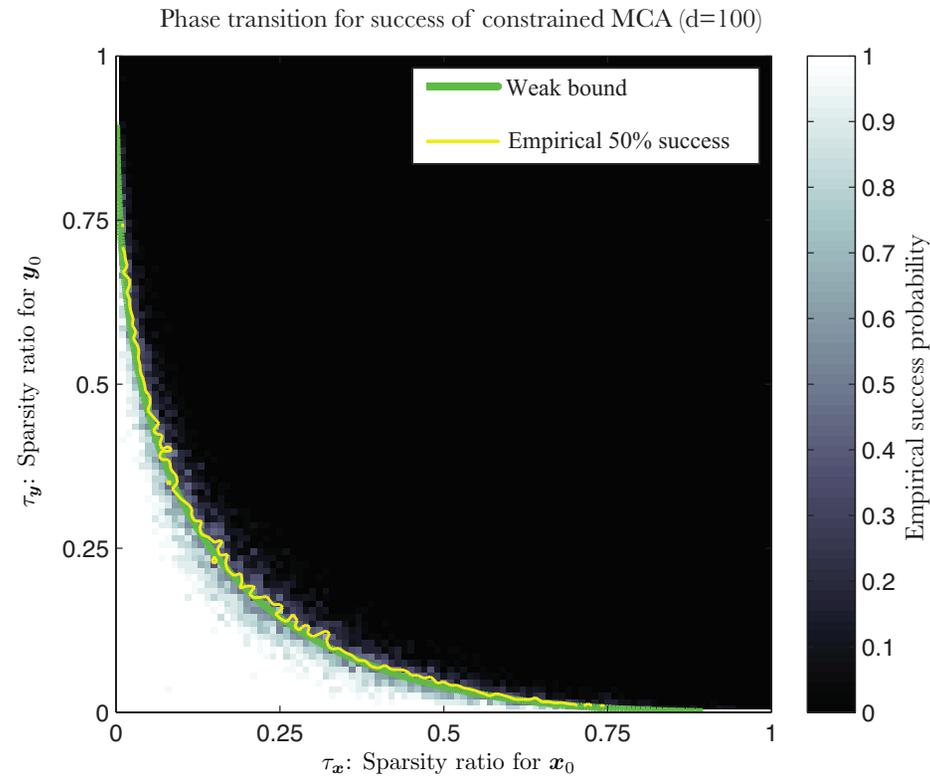
$$\delta(C) + \delta(K) \leq d - O(\sqrt{d}) \quad \text{implies} \quad (\hat{\mathbf{x}}, \hat{\mathbf{y}}) = (\mathbf{x}_\natural, \mathbf{y}_\natural) \quad \text{with high prob.}$$

$$\delta(C) + \delta(K) \geq d + O(\sqrt{d}) \quad \text{implies} \quad (\hat{\mathbf{x}}, \hat{\mathbf{y}}) \neq (\mathbf{x}_\natural, \mathbf{y}_\natural) \quad \text{with high prob.}$$

where $C = \mathcal{D}(f, \mathbf{x}_\natural)$ and $K = \mathcal{D}(g, \mathbf{y}_\natural)$

Sources: Amelunxen et al. 2014.

Demixing II



$$\text{minimize } \|\mathbf{x}\|_{\ell_1} \quad \text{subject to } \mathbf{z} = \mathbf{x} + \mathbf{Q}\mathbf{y}, \quad \|\mathbf{y}\|_{\ell_1} \leq \|\mathbf{y}_\dagger\|_{\ell_1}$$

Source: Starck et al. 2003; McCoy & Tropp 2013; McCoy et al. 2013; Amelunxen 2014.

Demixing III

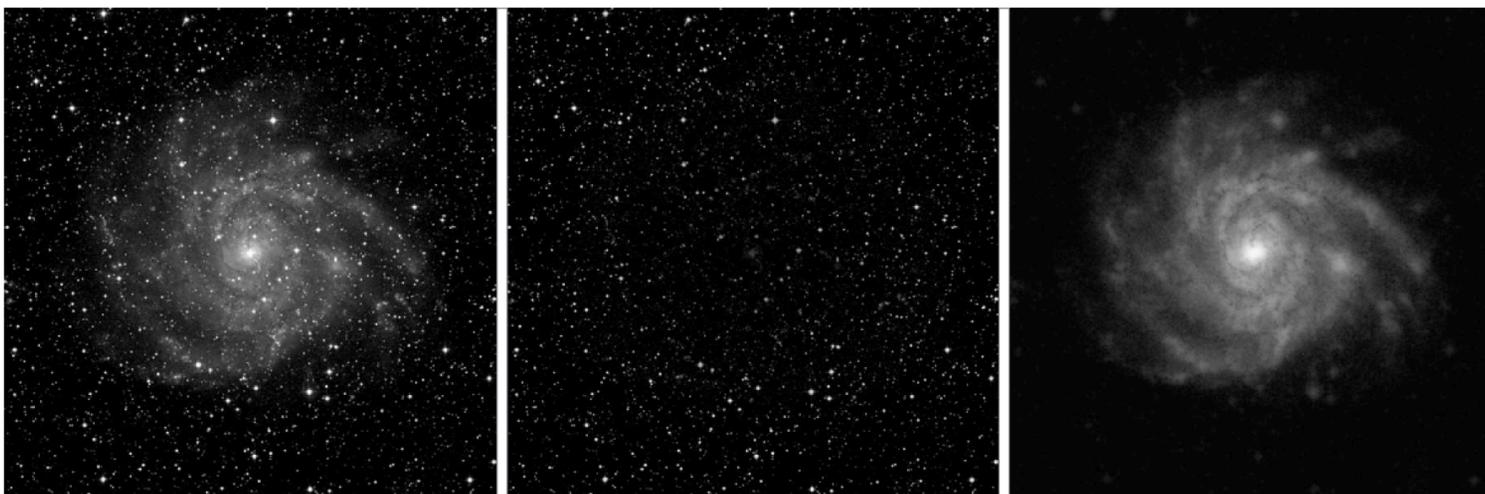


Image credit: NASA

Observation z_0

Sparse component x_0

DCT-sparse component y_0

Source: Starck et al. 2003; McCoy & Tropp 2013; McCoy et al. 2013; Amelunxen 2014.

To learn more...

E-mail: jtropp@cms.caltech.edu

Web: <http://users.cms.caltech.edu/~jtropp>

Some Sources:

- Ledoux & Talagrand, “Probability in Banach spaces,” Springer, 1991
- Chen et al., “Atomic decomposition by Basis Pursuit,” *SIAM Rev.*, 2001
- Tropp, “Topics in sparse approximation,” PhD Thesis, [UT-Austin](#), 2004
- Mallat & P  yre, “A wavelet tour of signal processing,” 3rd ed., Academic Press, 2009
- Jaggi, “Sparse convex optimization methods for machine learning,” PhD Thesis, ETH, 2011
- Chandrasekaran et al., “The convex geometry of linear inverse problems,” *Found. Comput. Math.*, 2012
- Foucart & Rauhut, “A mathematical introduction to compressive sampling,” Birkh  user, 2013
- McCoy, “A geometric analysis of convex demixing,” PhD Thesis, [Caltech](#), 2013
- McCoy & Tropp, “The achievable performance of convex demixing,” ACM Report 2017-02, [Caltech](#), 2013
- Amelunxen et al., “Living on the edge...,” *Inform. Inference*, 2014
- McCoy & Tropp, “From Steiner formulas for cones...,” *Discrete Comput. Geom.*, 2014
- Oymak, “Convex relaxation for low-dimensional representation...,” PhD Thesis, [Caltech](#), 2014
- Tropp, “Convex recovery of structured signals...,” in *Sampling Theory: A Renaissance*, Birkh  user, 2015
- Thrampoulidis, “Recovering structured signals in high dimensions...,” PhD Thesis, [Caltech](#), 2016
- Goldstein et al., “Gaussian phase transitions and conic intrinsic volumes...,” *Ann. Appl. Probab.*, 2017
- Amelunxen & Lotz, “Intrinsic volumes of polyhedral cones,” *Discrete Comput. Geom.*, 2017
- Oymak & Tropp, “Universality laws for randomized dimension reduction...,” arXiv:1511.09433, in revision