

Generic Chaining Meets (Non)convex Optimization

Mahdi Soltanolkotabi

June 23, 2016

Workshop on Chaining with Applications to Computer Science
Harvard University

Ming Hsieh Department of Electrical Engineering

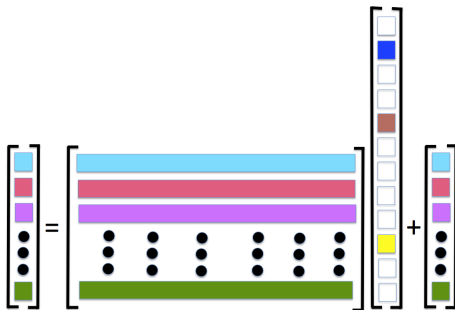


USC University of
Southern California

Collaborators: Samet Oymak and Ben Recht

Xiaodong Li and Emmanuel Candes

Linear inverse problems

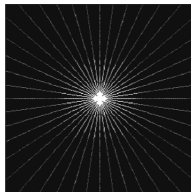


$$y = Ax + w$$

Compressed Sensing

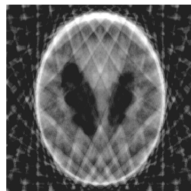


Fourier transform



highly subsampled

classical reconstruction



compressed sensing reconstruction

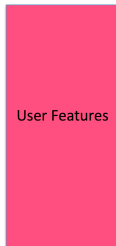


[Photo credit Candes, Romberg, Tao 2006]
linear measurements from a **structured signal**.

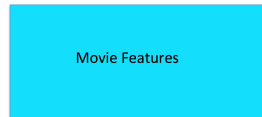
Recommender systems



=



User Features



Movie Features

The “power” of convex programming

Exciting research over the last decade demonstrating the effectiveness of convex programming/greedy algorithms.

The “power” of convex programming

Exciting research over the last decade demonstrating the effectiveness of convex programming/greedy algorithms.

Ideology

*“when life gives you lemons, **convexify**”*

The “power” of convex programming

Exciting research over the last decade demonstrating the effectiveness of convex programming/greedy algorithms.

Ideology

*“when life gives you lemons, **convexify**”*

- Sparse use ℓ_1 norm, Low-rank use nuclear norm, atomic norms, etc.

The “power” of convex programming

Exciting research over the last decade demonstrating the effectiveness of convex programming/greedy algorithms.

Ideology

*“when life gives you lemons, **convexify**”*

- Sparse use ℓ_1 norm, Low-rank use nuclear norm, atomic norms, etc.

convex algorithms are not perfect

convex algorithms are not perfect

- Computation and memory: convex programs maybe inefficient

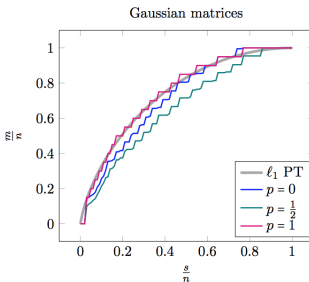


convex algorithms are not perfect

- Computation and memory: convex programs maybe inefficient

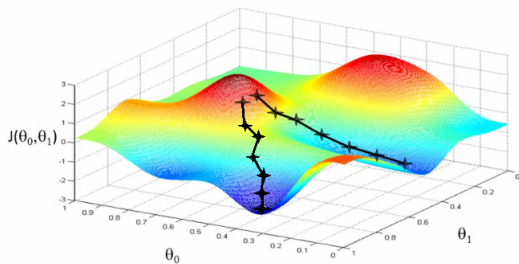


- Sometimes convex programs are inefficient in capturing the “structure” (usually require more samples)



Just follow the gradient?

In practice local search heuristics work really well..



Non-convex optimization is difficult!

- Non-convex optimization is quite tricky!

Non-convex optimization is difficult!

- Non-convex optimization is quite tricky!
-

$$f(\mathbf{x}) = \sum_{i,j=1}^n Q_{ij} x_i^2 x_j^2 \quad \nabla f(0) = 0 \quad \text{for all } Q$$

- Checking if 0 is a local minimum is NP-hard!

Going beyond worse case

Two stories with a common theme.

With randomized coefficients (e.g. functions of Gaussians) local search heuristics work.

- Story I: Solving quadratic equations (nonconvex objective)
- Story II: Linear inverse problems (nonconvex constraints)

Going beyond worse case

Two stories with a common theme.

With randomized coefficients (e.g. functions of Gaussians) local search heuristics work.

- Story I: Solving quadratic equations (nonconvex objective)
- Story II: Linear inverse problems (nonconvex constraints)

Challenge

Real data is not Gaussian ...

Going beyond worse case

Two stories with a common theme.

With randomized coefficients (e.g. functions of Gaussians) local search heuristics work.

- Story I: Solving quadratic equations (nonconvex objective)
- Story II: Linear inverse problems (nonconvex constraints)

Challenge

Real data is not Gaussian ...

Solution

Generic Chaining

*Story I:
Solving Quadratic Equations*

Solving quadratic equations

$$y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2 \quad r = 1, 2, \dots, m \quad \Leftrightarrow \quad \mathbf{y} = |\mathbf{A}\mathbf{x}|^2$$

Find a feasible point in the intersection of quadratic equations

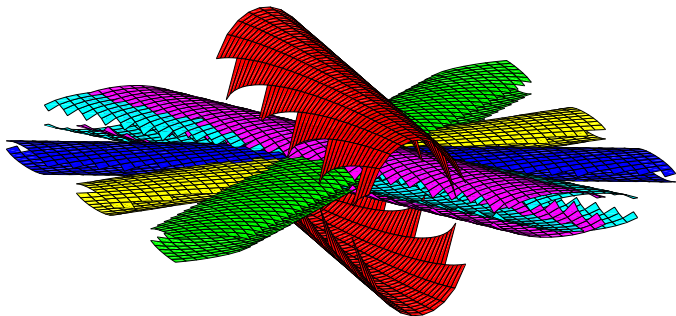
$$\mathbf{y}_r = \mathbf{x}^* \mathbf{A}_r \mathbf{x} \quad \text{for } r = 1, 2, \dots, m.$$

Solving quadratic equations

$$y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2 \quad r = 1, 2, \dots, m \quad \Leftrightarrow \quad \mathbf{y} = |\mathbf{A}\mathbf{x}|^2$$

Find a feasible point in the intersection of quadratic equations

$$\mathbf{y}_r = \mathbf{x}^* \mathbf{A}_r \mathbf{x} \quad \text{for } r = 1, 2, \dots, m.$$

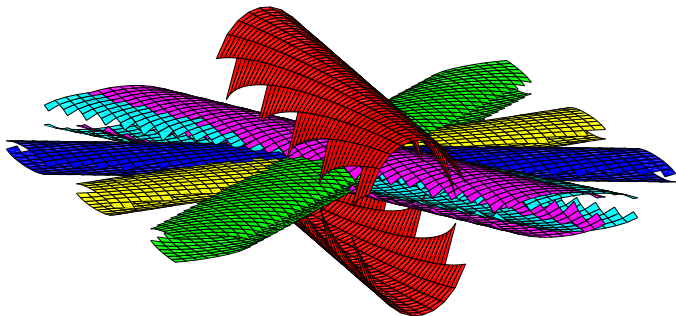


Solving quadratic equations

$$y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2 \quad r = 1, 2, \dots, m \quad \Leftrightarrow \quad \mathbf{y} = |\mathbf{A}\mathbf{x}|^2$$

Find a feasible point in the intersection of quadratic equations

$$\mathbf{y}_r = \mathbf{x}^* \mathbf{A}_r \mathbf{x} \quad \text{for } r = 1, 2, \dots, m.$$



One of the universal forms of combinatorial problems, NP-hard in general.

Missing phase problem

- Detectors only record intensities of diffracted rays (**magnitude measurements only!**)

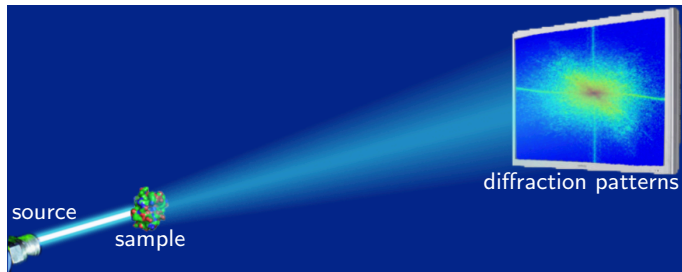


- Fraunhofer diffraction equation \Rightarrow optical field at the detector \approx Fourier transform

$$|\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-2\pi i(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Missing phase problem

- Detectors only record intensities of diffracted rays (**magnitude measurements only!**)



- Fraunhofer diffraction equation \Rightarrow optical field at the detector \approx Fourier transform

$$|\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-2\pi i(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Phase Retrieval Problem

How can we recover the phase (or equivalently signal $x(t_1, t_2)$) from $|\hat{x}(f_1, f_2)|^2$?

Phase retrieval (discrete 1D model)



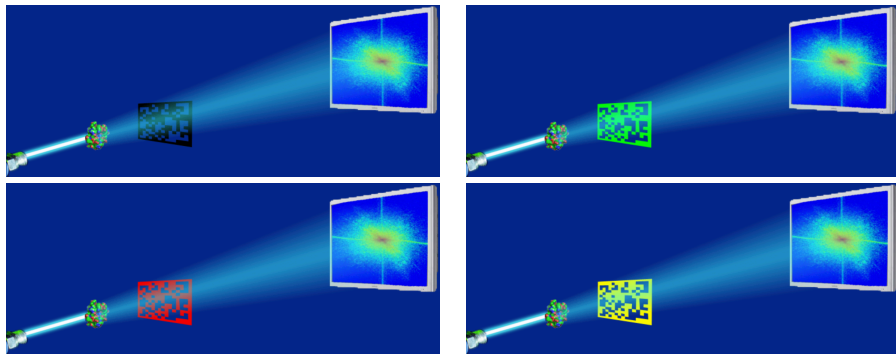
- Phaseless measurements about $\mathbf{x} \in \mathbb{C}^n$

$$|\mathbf{f}_k^* \mathbf{x}|^2 = \mathbf{y}_k \quad k \in \{1, 2, \dots, n\} = [n]$$

\mathbf{f}_k^* is k th row of the DFT matrix.

- Phase retrieval is impossible, inherent ambiguity.

Phase retrieval



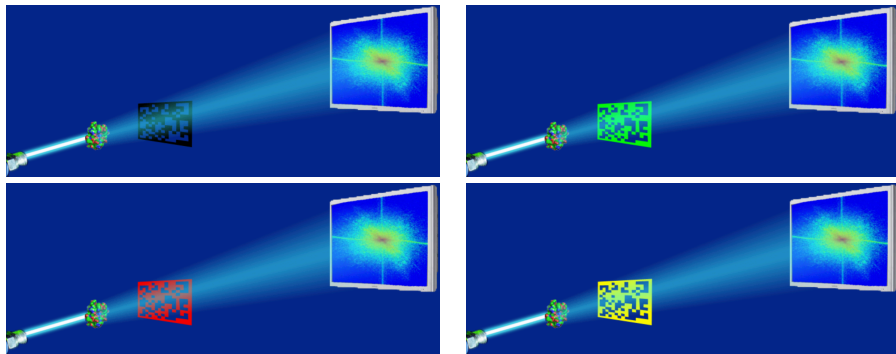
Measurements of the form

$$y_{k\ell} = |\mathbf{f}_k^* \mathbf{D}_\ell^* \mathbf{x}|^2 \quad k = 1, 2, \dots, n \quad \text{and} \quad \ell = 1, 2, \dots, L.$$

$$y_r = |\mathbf{a}_r^* \mathbf{x}|^2 \quad r = 1, 2, \dots, m.$$

Here $m = nL$ and $\mathbf{a}_r = \mathbf{D}_\ell \mathbf{f}_k$, r represents (k, ℓ)

Phase retrieval



Measurements of the form

$$y_{k\ell} = |\mathbf{f}_k^* \mathbf{D}_\ell^* \mathbf{x}|^2 \quad k = 1, 2, \dots, n \quad \text{and} \quad \ell = 1, 2, \dots, L.$$

$$y_r = |\mathbf{a}_r^* \mathbf{x}|^2 \quad r = 1, 2, \dots, m.$$

Here $m = nL$ and $\mathbf{a}_r = \mathbf{D}_\ell \mathbf{f}_k$, r represents (k, ℓ)

Phase Retrieval by non-convex optimization

Let $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m]$

$$\begin{aligned}\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) &:= \frac{1}{2m} \left\| \mathbf{y} - |\mathbf{A}^* \mathbf{z}|^2 \right\|_{\ell_2}^2 \\ &= \frac{1}{2m} \sum_{r=1}^m (\mathbf{y}_r - |\mathbf{a}_r^* \mathbf{z}|^2)^2\end{aligned}$$

For a vector $\mathbf{b} \in \mathbb{C}^n$ by $|b|^2$, I mean Matlab $\text{abs}(\mathbf{b}.^2)$.

- Pro: operates over vectors much less intensive!
- Con: Non-convex!

Wirtinger Flow (WF)

Algorithm 1 Wirtinger Flow (WF)

Input: Measurements y_r for $r = 1, 2, \dots, m$.

Initialization (WF-INIT):

Set \tilde{z}_0 to be the eigenvector corresponding to the largest eigenvalue of

$$\mathbf{Y} = \frac{1}{m} \sum_{r=1}^m y_r \mathbf{a}_r \mathbf{a}_r^*.$$

Set $\mathbf{z}_0 = \left(\sqrt{\frac{1}{m} \sum_{r=1}^m y_r} \right) \tilde{z}_0$.

Iterations:

for $\tau = 0$ **to** $t - 1$ **do**

Set

$$\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} - \frac{\mu_{\tau+1}}{\|\mathbf{z}_0\|_{\ell_2}^2} \left(\frac{1}{m} \sum_{r=1}^m (|\mathbf{a}_r^* \mathbf{z}|^2 - y_r) (\mathbf{a}_r \mathbf{a}_r^*) \mathbf{z} \right) := \mathbf{z}_{\tau} - \frac{\mu_{\tau+1}}{\|\mathbf{z}_0\|_{\ell_2}^2} \nabla f(\mathbf{z}_{\tau}).$$

end for

Output: $\hat{\mathbf{x}} = \mathbf{z}_t$.

Exact Phase Retrieval by WF (Gaussian Model)

For a vector $\mathbf{z} \in \mathbb{C}^n$

$$\text{dist}(\mathbf{z}, \mathbf{x}) = \min_{\phi \in [0, 2\pi]} \|\mathbf{z} - e^{i\phi} \mathbf{x}\|_{\ell_2}.$$

Theorem (Candes, Li, and Soltanolkotabi ('14), Soltanolkotabi ('14))

Assume $m \gtrsim n$. Using $0 \leq \mu \leq \mu_0/n$, with high probability

- Initialization:

$$\text{dist}(\mathbf{z}_0, \mathbf{x}) \leq \sqrt{\frac{5}{6}} \|\mathbf{x}\|_{\ell_2}.$$

- After t iterations:

$$\text{dist}(\mathbf{z}_t, \mathbf{x}) \leq e^{-c\mu t} \cdot \text{dist}(\mathbf{z}_0, \mathbf{x}) \leq \sqrt{\frac{5}{6}} e^{-c\mu t} \|\mathbf{x}\|_{\ell_2}.$$

[Chen and Candes 2015] also later established $m \gtrsim n$ via Truncated Wirtinger Flow

Computational complexity

- Initialization: essentially $L \log n$ FFTs $= o(nL \log^2 n)$.
- Iteration update: essentially $2L$ FFTs $= o(nL \log n)$.

Computational complexity

- Initialization: essentially $L \log n$ FFTs $= o(nL \log^2 n)$.
- Iteration update: essentially $2L$ FFTs $= o(nL \log n)$.

Computational Complexity

Near Optimal: Equivalent to a few thousand FFTs.

Computational complexity

- Initialization: essentially $L \log n$ FFTs $= o(nL \log^2 n)$.
- Iteration update: essentially $2L$ FFTs $= o(nL \log n)$.

Computational Complexity

Near Optimal: Equivalent to a few thousand FFTs.

Pros:

- No SVD's.
- No matrix inversions.
- Everything is in \mathbb{C}^n rather than the Lifted space.

A real experiment

Original Image

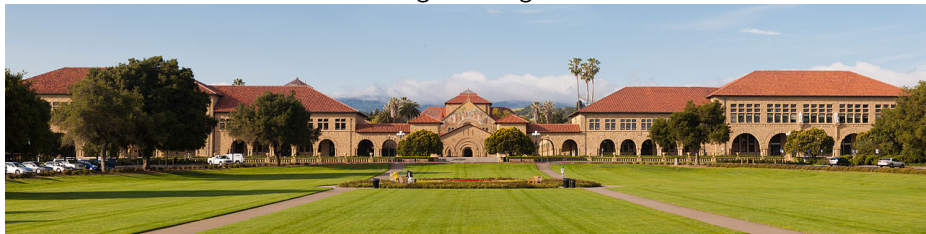


WF after 500 iterations: Time=**221.7697 sec**, Relerr= **2.5410×10^{-11}**



A real experiment

Original Image



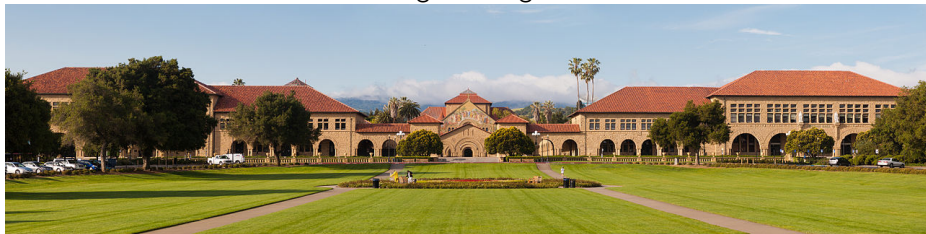
WF after 500 iterations: Time=**221.7697 sec**, Relerr= **2.5410×10^{-11}**



SDP based methods (PhaseLift, PhaseCut, ...) require xx^*

A real experiment

Original Image



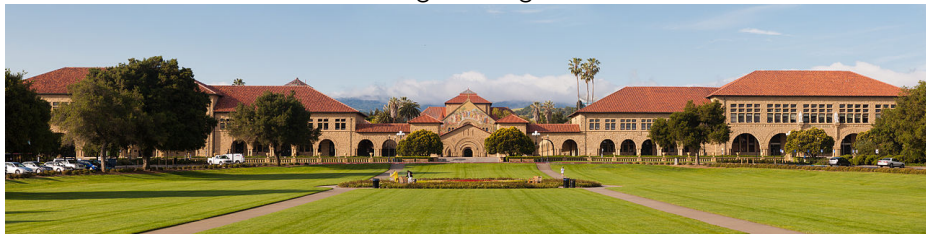
WF after 500 iterations: Time=**221.7697 sec**, Relerr= **2.5410×10^{-11}**



SDP based methods (PhaseLift, PhaseCut, ...) require xx^*
 $(320 \times 1280)^2 \times 8 \text{ Bytes} \approx 0.16777 \text{ Tera Bytes}$

A real experiment

Original Image



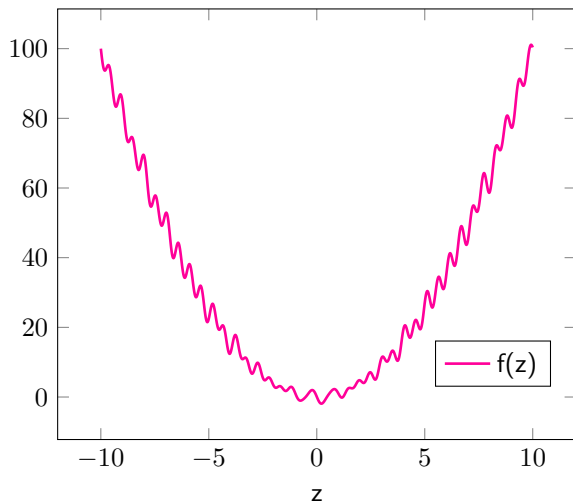
WF after 500 iterations: Time=**221.7697 sec**, Relerr= **2.5410×10^{-11}**



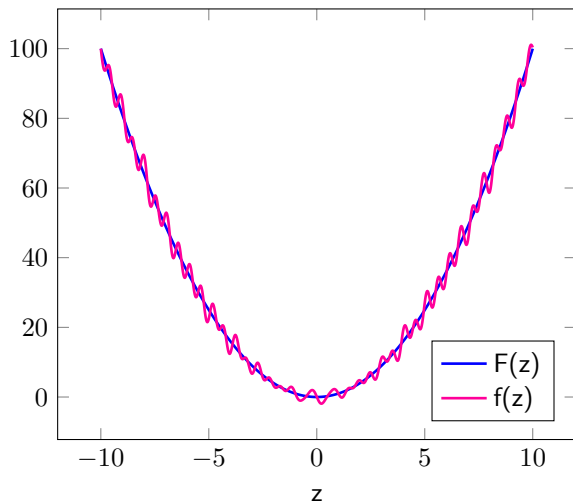
SDP based methods (PhaseLift, PhaseCut, ...) require xx^*
(320×1280)² \times 8 Bytes \approx 0.16777 Tera Bytes
does not even fit into memory!

Proof Sketch

Main Idea: Compare with an easy to analyze function



Main Idea: Compare with an easy to analyze function



For phase retrieval

$$f(\mathbf{z}) = \frac{1}{2m} \sum_{r=1}^m \left(y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2$$

and

$$F(\mathbf{z}) := \mathbb{E}[f(\mathbf{z})] = \mathbf{z}^* (\mathbf{I} - \mathbf{x}\mathbf{x}^*) \mathbf{z} + \frac{1}{2} \left(\|\mathbf{z}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2 \right)^2$$

For phase retrieval

$$f(\mathbf{z}) = \frac{1}{2m} \sum_{r=1}^m \left(y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2$$

and

$$F(\mathbf{z}) := \mathbb{E}[f(\mathbf{z})] = \mathbf{z}^* (\mathbf{I} - \mathbf{x}\mathbf{x}^*) \mathbf{z} + \frac{1}{2} \left(\|\mathbf{z}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2 \right)^2$$

Showing $f(\mathbf{z}) \approx F(\mathbf{z})$ not sufficient. Need to show this in terms of higher order derivatives.

Regularity Condition

Define

$$P := \{\mathbf{x}e^{i\phi} : \phi \in [0, 2\pi]\}.$$

and

$$E(\epsilon) := \{\mathbf{z} \in \mathbb{C}^n : \text{dist}(\mathbf{z}, P) \leq \epsilon\}.$$

Condition (Regularity Condition)

We say that the function f satisfies the regularity condition or $RC(\alpha, \beta, \epsilon)$ if for all vectors $\mathbf{z} \in E(\epsilon)$ we have

$$\text{Re} \left(\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x}e^{i\phi(\mathbf{z})} \rangle \right) \geq \frac{1}{\alpha} \text{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2.$$

Regularity condition leads to convergence

We will prove that if for all $\mathbf{z} \in E(\epsilon)$ we have

$$\operatorname{Re} \left(\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} \rangle \right) \geq \frac{1}{\alpha} \operatorname{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2.$$

then for all $0 < \mu \leq \frac{2}{\beta}$

$$\mathbf{z}_+ = \mathbf{z} - \mu \nabla f(\mathbf{z})$$

obeys

$$\operatorname{dist}^2(\mathbf{z}_+, \mathbf{x}) \leq \left(1 - \frac{2\mu}{\alpha} \right) \operatorname{dist}^2(\mathbf{z}, \mathbf{x}).$$

$$\begin{aligned} \|\mathbf{z}_+ - \mathbf{x} e^{i\phi(\mathbf{z}_+)}\|_{\ell_2}^2 &\leq \|\mathbf{z}_+ - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 = \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} - \mu \nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &= \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 - 2\mu \operatorname{Re} \left(\langle \nabla f(\mathbf{z}), (\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}) \rangle \right) + \mu^2 \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &\leq \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 - 2\mu \left(\frac{1}{\alpha} \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \right) + \mu^2 \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &= \left(1 - \frac{2\mu}{\alpha} \right) \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 + \mu \left(\mu - \frac{2}{\beta} \right) \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &\leq \left(1 - \frac{2\mu}{\alpha} \right) \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2, \end{aligned}$$

Regularity condition leads to convergence

We will prove that if for all $\mathbf{z} \in E(\epsilon)$ we have

$$\operatorname{Re} \left(\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} \rangle \right) \geq \frac{1}{\alpha} \operatorname{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2.$$

then for all $0 < \mu \leq \frac{2}{\beta}$

$$\mathbf{z}_+ = \mathbf{z} - \mu \nabla f(\mathbf{z})$$

obeys

$$\operatorname{dist}^2(\mathbf{z}_+, \mathbf{x}) \leq \left(1 - \frac{2\mu}{\alpha} \right) \operatorname{dist}^2(\mathbf{z}, \mathbf{x}).$$

$$\begin{aligned} \|\mathbf{z}_+ - \mathbf{x} e^{i\phi(\mathbf{z}_+)}\|_{\ell_2}^2 &\leq \|\mathbf{z}_+ - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 = \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} - \mu \nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &= \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 - 2\mu \operatorname{Re} \left(\langle \nabla f(\mathbf{z}), (\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}) \rangle \right) + \mu^2 \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &\leq \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 - 2\mu \left(\frac{1}{\alpha} \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \right) + \mu^2 \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &= \left(1 - \frac{2\mu}{\alpha} \right) \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2 + \mu \left(\mu - \frac{2}{\beta} \right) \|\nabla f(\mathbf{z})\|_{\ell_2}^2 \\ &\leq \left(1 - \frac{2\mu}{\alpha} \right) \|\mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})}\|_{\ell_2}^2, \end{aligned}$$

How do we prove the regularity condition?

$$\operatorname{Re} \left(\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} \rangle \right) \geq \frac{1}{\alpha} \operatorname{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2.$$

and

$$\nabla f(\mathbf{z}) = \frac{1}{m} \sum_{r=1}^m \left(|\mathbf{a}_r^* \mathbf{z}|^2 - y_r \right) (\mathbf{a}_r^* \mathbf{z}) \mathbf{a}_r$$

Condition (Local Curvature Condition)

$$\operatorname{Re} \left(\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} e^{i\phi(\mathbf{z})} \rangle \right) \geq \left(\frac{1}{\alpha} + \frac{(1-\delta)}{4} \right) \operatorname{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{10m} \sum_{r=1}^m \left| \mathbf{a}_r^* (\mathbf{z} - e^{i\phi(\mathbf{z})} \mathbf{x}) \right|^4$$

Condition (Local Smoothness Condition)

$$\|\nabla f(\mathbf{z})\|_{\ell_2}^2 \leq \beta \left(\frac{(1-\delta)}{4} \operatorname{dist}^2(\mathbf{z}, \mathbf{x}) + \frac{1}{10m} \sum_{r=1}^m \left| \mathbf{a}_r^* (\mathbf{z} - e^{i\phi(\mathbf{z})} \mathbf{x}) \right|^4 \right).$$

*Story II:
Linear inverse problems*

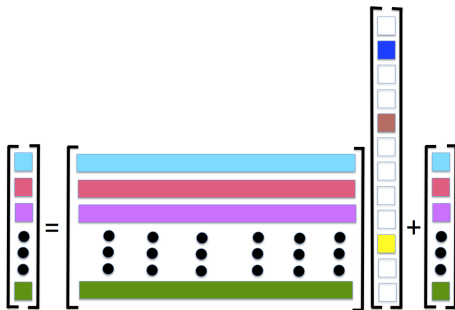
Structured signal recovery/Linear inverse problems

Optimization Problem

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{z}} \frac{1}{2} \sum_{i=1}^m (y_i - \langle \mathbf{a}_i, \mathbf{z} \rangle)^2 = \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_{\ell_2}^2$$

subject to $f(\mathbf{z}) \leq f(\mathbf{x})$

f is a function that enforces structure



$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$$

Minimal data/measurements for exact recovery?

(Review)

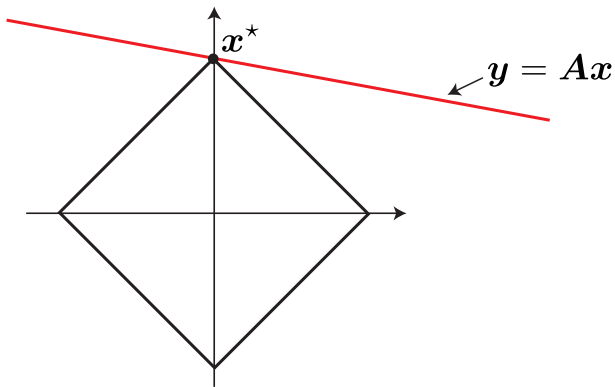
Minimal data/measurements for exact recovery?

$\mathbf{y} = \mathbf{A}\mathbf{x}$, $\mathbf{y} \in \mathbb{R}^m$, $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{x} \in \mathbb{R}^n$ with $m \ll n$.

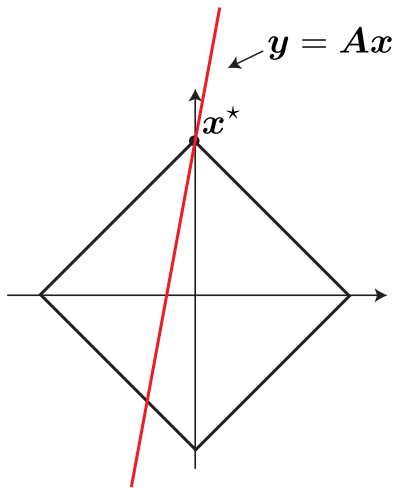
$$\hat{\mathbf{x}} = \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_{\ell_2}^2 \quad \text{subject to} \quad f(\mathbf{z}) \leq f(\mathbf{x}).$$

When is $\hat{\mathbf{x}} = \mathbf{x}$? m ?

Why does ℓ_1 work?

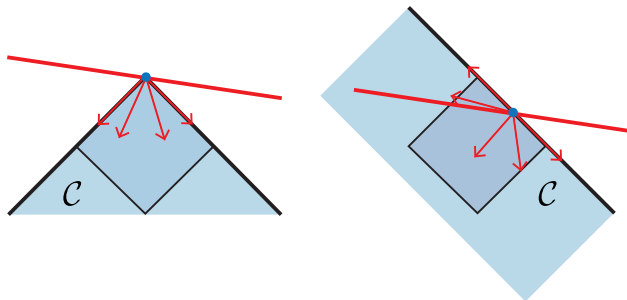


Why ℓ_1 may not always work



Geometry

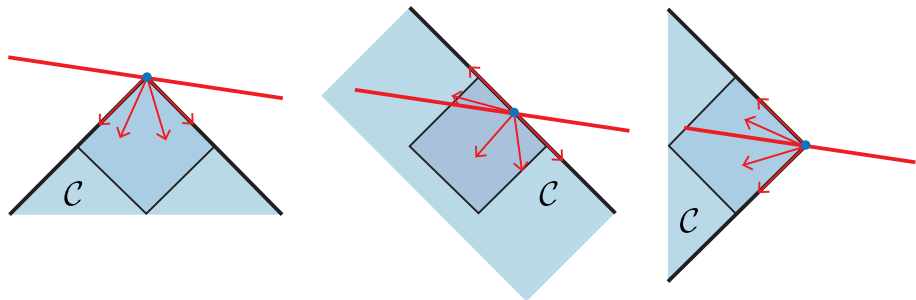
$$\mathcal{C} = \{h : \|x + th\| \leq \|x\| \text{ for some } t > 0\} \quad \text{cone of descent}$$



Exact recovery if $\mathcal{C} \cap \text{null}(A) = \{0\}$

Geometry

$$\mathcal{C} = \{h : \|x + th\| \leq \|x\| \text{ for some } t > 0\} \quad \text{cone of descent}$$



Exact recovery if $\mathcal{C} \cap \text{null}(A) = \{0\}$

Mean width

$$\omega(\mathcal{T}) = \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{x} \in \mathcal{T}} \mathbf{g}^T \mathbf{x}]$$

- finite set \mathcal{T}

Mean width

$$\omega(\mathcal{T}) = \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{x} \in \mathcal{T}} \mathbf{g}^T \mathbf{x}]$$

- finite set \mathcal{T}

$$\omega^2(\mathcal{T}) = 2 \log |\mathcal{T}|$$

Mean width

$$\omega(\mathcal{T}) = \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{x} \in \mathcal{T}} \mathbf{g}^T \mathbf{x}]$$

- finite set \mathcal{T}

$$\omega^2(\mathcal{T}) = 2 \log |\mathcal{T}|$$

- \mathcal{C} cone of decent of ℓ_1 norm at an s -sparse signal

Mean width

$$\omega(\mathcal{T}) = \mathbb{E}_{\mathbf{g}}[\sup_{\mathbf{x} \in \mathcal{T}} \mathbf{g}^T \mathbf{x}]$$

- finite set \mathcal{T}

$$\omega^2(\mathcal{T}) = 2 \log |\mathcal{T}|$$

- \mathcal{C} cone of decent of ℓ_1 norm at an s -sparse signal

$$\omega^2(\mathcal{C} \cap \mathbb{S}^{n-1}) \approx 2s \log(n/s)$$

Theorem (Chandrasekaran, Recht, Parrilo, and Willskey 2012)

For i.i.d. normal matrices as long as

$$m \geq m_0(f, \mathbf{x}) := \omega^2(\mathcal{C}_f(\mathbf{x})),$$

then

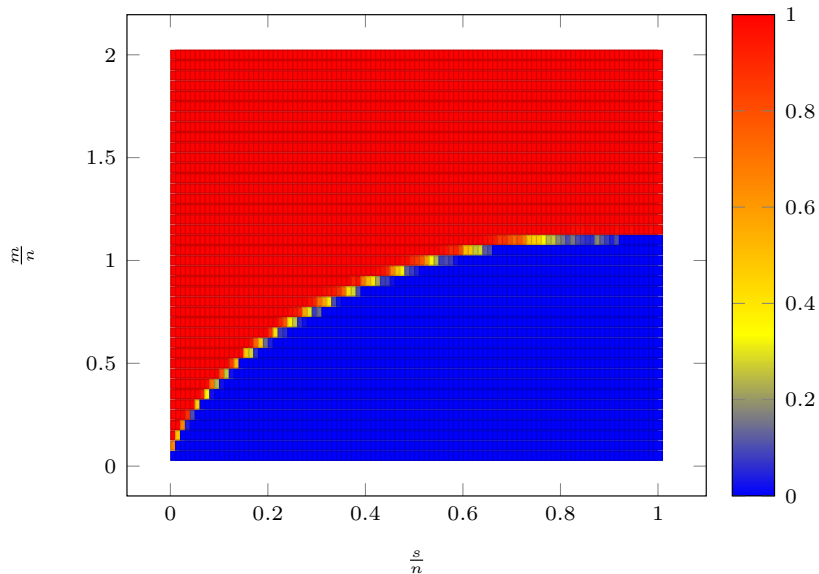
$$\hat{\mathbf{x}} = \mathbf{x}$$

holds with high probability.

- f is **convex** and \mathbf{A} Gaussian
(Chandrasekaran, Recht, Parrilo, Willskey 2012), (Amelunxen, Lotz, McCoy, Tropp 2013), (Stojnic 2009, 2013), Special cases [Donoho-Tanner 2009], many works by (Donoho, Maleki, Montanari 2009), (Bayati and Montanari 2011)
- sub-Gaussians: [Klartag and Mendelson 2005], [Mendelson, Pajor and Tomczak-Jaegermann 2007], [Dirksen 2015], [Bayati, Lelarge, Montanari 2014], [Oymak and Tropp 2015]
- non-Gaussians and non-i.i.d. not known for general structures
- Non-Gaussians, non-i.i.d. not known for the most part. Special structures: [Candes, Romberg, Tao 2004], [Rudelson and Vershynin 2007], [Krahmer, Mendelson, Rauhut 2012], [Bourgain, Dirksen, Nelson 2015].

Phase Transitions

Sparse recovery via ℓ_1 minimization

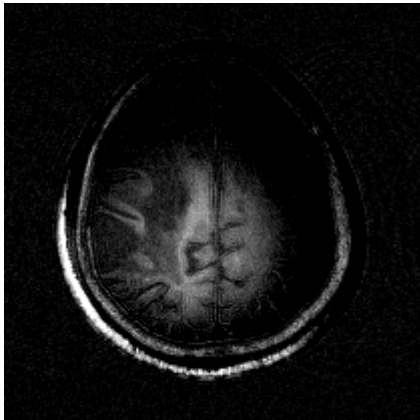


Going beyond convexity

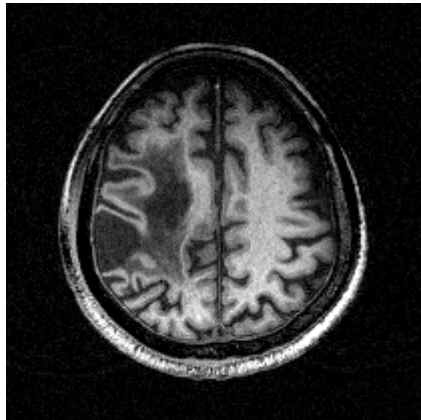
- A huge literature on convex relaxations/greedy methods over the last decade

Going beyond convexity

- A huge literature on convex relaxations/greedy methods over the last decade
- The function that best captures signal structure may not be convex...



real part



imaginary part

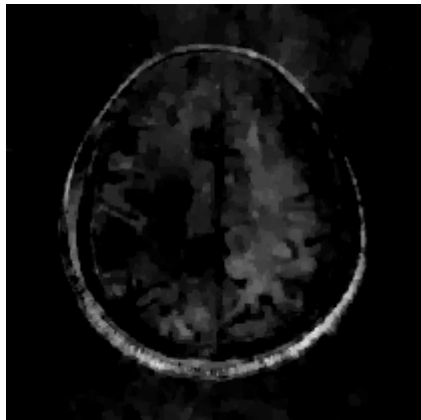
Data courtesy of USC biomedical imaging group

Accelerated MRI via total variation minimization

Undersample by a factor of 3 and use TV reconstruction

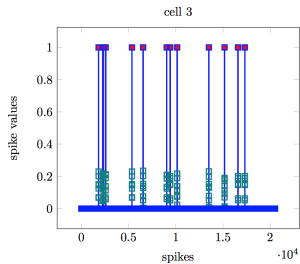
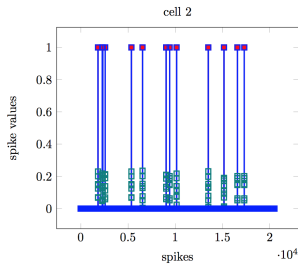
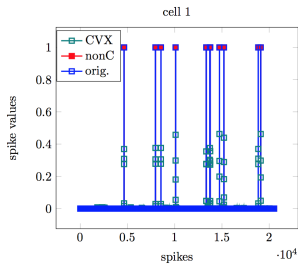


real part



imaginary part

This happens even for sparse signals



Algorithms

Projected Gradient Descent

$$\hat{\mathbf{x}} = \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_{\ell_2}^2 \quad \text{subject to} \quad f(\mathbf{z}) \leq f(\mathbf{x}).$$

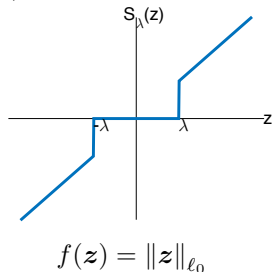
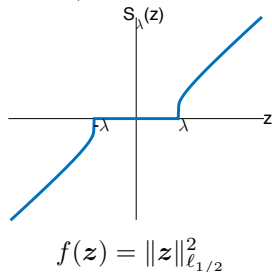
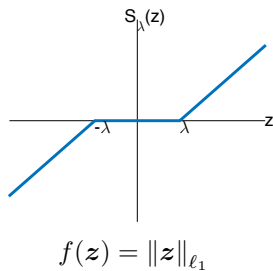
Projection: $\mathcal{P}_{\mathcal{K}}(\mathbf{z}), \mathcal{K} = \{\mathbf{z} \in \mathbb{R}^n : f(\mathbf{z}) \leq f(\mathbf{x})\}$

Projected Gradient Descent

$$\hat{\mathbf{x}} = \underset{\mathbf{z}}{\operatorname{argmin}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{z}\|_{\ell_2}^2 \quad \text{subject to} \quad f(\mathbf{z}) \leq f(\mathbf{x}).$$

Projection: $\mathcal{P}_{\mathcal{K}}(\mathbf{z})$, $\mathcal{K} = \{\mathbf{z} \in \mathbb{R}^n : f(\mathbf{z}) \leq f(\mathbf{x})\}$

- Start from $\mathbf{z}_0 = \mathbf{0}$.
- Apply the update $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$



Theory

A simple thought experiment

Reminder: $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$.

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

A simple thought experiment

Reminder: $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$.

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Forget about the projection: $\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau})$

A simple thought experiment

Reminder: $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$.

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Forget about the projection: $\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau})$

$$\mathbf{z}_{\tau+1} - \mathbf{x} = (\mathbf{I} - \mu\mathbf{A}^*\mathbf{A})(\mathbf{z}_{\tau} - \mathbf{x}) - \mu\mathbf{A}^*\mathbf{w}.$$

A simple thought experiment

Reminder: $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$.

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Forget about the projection: $\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau})$

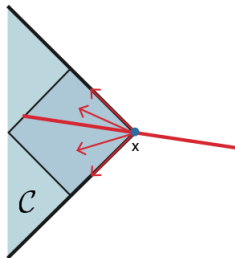
$$\mathbf{z}_{\tau+1} - \mathbf{x} = (\mathbf{I} - \mu \mathbf{A}^* \mathbf{A})(\mathbf{z}_{\tau} - \mathbf{x}) - \mu \mathbf{A}^* \mathbf{w}.$$

$$\|\mathbf{z}_{\tau+1} - \mathbf{x}\|_{\ell_2} \leq \|\mathbf{I} - \mu \mathbf{A}^* \mathbf{A}\| \|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} + \mu \|\mathbf{A}\| \|\mathbf{w}\|_{\ell_2}.$$

Geometry

Definition (Descent Cone)

$$\mathcal{C}_f(\mathbf{x}) = \{\mathbf{h} : f(\mathbf{x} + \tau\mathbf{h}) \leq f(\mathbf{x})\}.$$



A deterministic result

Algorithm:

- Start with $z_0 = \mathbf{0}$,
- Run $z_{\tau+1} = \mathcal{P}_{\mathcal{K}}(z_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}z_{\tau}))$

A deterministic result

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Soltanolkotabi, and Recht 2016)

For any \mathbf{A} , and any convex f

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq (\rho(\mu))^{\tau} \|\mathbf{x}\|_{\ell_2} + \frac{1 - (\rho(\mu))^{\tau}}{1 - \rho(\mu)} \xi_{\mu}(\mathbf{A}) \|\mathbf{w}\|_{\ell_2}.$$

- $\rho(\mu)$ is the convergence rate

$$\rho(\mu) := \rho(\mu, \mathbf{A}, f, \mathbf{x}) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathcal{B}^n} \mathbf{u}^* (\mathbf{I} - \mu \mathbf{A}^* \mathbf{A}) \mathbf{v},$$

- $\xi_{\mu}(\mathbf{A})$ is the noise amplification factor

$$\xi_{\mu}(\mathbf{A}) := \xi_{\mu}(\mathbf{A}, f, \mathbf{x}, \mathbf{w}) = \mu \cdot \sup_{\mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathcal{B}^n} \mathbf{v}^* \mathbf{A}^* \frac{\mathbf{w}}{\|\mathbf{w}\|_{\ell_2}}.$$

A deterministic result

Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Soltanolkotabi, and Recht 2016)

For any \mathbf{A} , and any convex f

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq (\rho(\mu))^{\tau} \|\mathbf{x}\|_{\ell_2} + \frac{1 - (\rho(\mu))^{\tau}}{1 - \rho(\mu)} \xi_{\mu}(\mathbf{A}) \|\mathbf{w}\|_{\ell_2}.$$

- $\rho(\mu)$ is the convergence rate

$$\rho(\mu) := \rho(\mu, \mathbf{A}, f, \mathbf{x}) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathcal{B}^n} \mathbf{u}^* (\mathbf{I} - \mu \mathbf{A}^* \mathbf{A}) \mathbf{v},$$

- $\xi_{\mu}(\mathbf{A})$ is the noise amplification factor

$$\xi_{\mu}(\mathbf{A}) := \xi_{\mu}(\mathbf{A}, f, \mathbf{x}, \mathbf{w}) = \mu \cdot \sup_{\mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathcal{B}^n} \mathbf{v}^* \mathbf{A}^* \frac{\mathbf{w}}{\|\mathbf{w}\|_{\ell_2}}.$$

- for **non-convex** f just replace ρ with 2ρ !

Geometric Convergence!

- For convex f traditional literature will tell you $1/\tau$ convergence!

Geometric Convergence!

- For convex f traditional literature will tell you $1/\tau$ convergence!
- For convex f we have geometric convergence if $\rho < 1$ even though the objective is not strongly convex!

Geometric Convergence!

- For convex f traditional literature will tell you $1/\tau$ convergence!
- For convex f we have geometric convergence if $\rho < 1$ even though the objective is not strongly convex!
- for non-convex f we have geometric convergence to the global optimum when $\rho < \frac{1}{2}$.

Gaussian measurements with different step size

Natural tradeoffs via step sizes:

- greedy: $\mu = 1/m$
- conservative: $\mu \approx \frac{1}{m+n}$

Gaussian measurements with greedy step size

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Gaussian measurements with greedy step size

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Soltanolkotabi, and Recht 2016)

For Gaussian \mathbf{A} , $\mu = 1/m$, with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(\sqrt{8\kappa_f^2 \frac{m_0}{m}} \right)^{\tau} \|\mathbf{z}_0 - \mathbf{x}\|_{\ell_2} + \kappa_f \sqrt{\frac{\pi}{2}} \frac{\sqrt{m_0}}{m} \|\mathbf{w}\|_{\ell_2}$$

- Sample complexity $m \geq 8\kappa_f^2 m_0$
- Linear rate
- $\kappa_f = 1$ for convex f and $\kappa_f = 2$ for non-convex f .
- Works for any function (including non-convex)!!

Gaussian measurements with conservative step size

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Gaussian measurements with conservative step size

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (2016)

For Gaussian \mathbf{A} and convex f , $\mu = \frac{0.99}{(\sqrt{m} + \sqrt{n})^2}$, with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(1 - \frac{0.3}{m+n} (\sqrt{m} - \sqrt{m_0})^2\right)^{\tau} \cdot \|\mathbf{x}\|_{\ell_2} + \frac{3.5}{(1 - \sqrt{\frac{m_0}{m}})^2} \frac{\sqrt{m_0}}{m} \|\mathbf{w}\|_{\ell_2}.$$

- Sample complexity $m \geq m_0$
- The rate is now geometric instead of linear

Some historical notes

- for convex and decomposable norms, up to constants (Agarwal, Negahban and Wainwright 2012), convergence upto “statistical accuracy” (Loh and Wainwright 2014)
- approximate message passing algorithm asymptotically achieves $\sqrt{\frac{m_0}{m}}$ for separable and pseudo Lipschitz f (Bayati and Montanari 11)

Nonlinear Observations

Nonlinear observe from a structured signal

$$\mathbf{y} = g(\mathbf{A}\mathbf{x}).$$

How can we recover \mathbf{x} from these measurements?

Nonlinear Observations

Nonlinear observe from a structured signal

$$\mathbf{y} = g(\mathbf{A}\mathbf{x}).$$

How can we recover \mathbf{x} from these measurements?

Pretend the model was linear and run the same algorithm

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Nonlinear Observations

Nonlinear observe from a structured signal

$$\mathbf{y} = g(\mathbf{Ax}).$$

How can we recover \mathbf{x} from these measurements?

Pretend the model was linear and run the same algorithm

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{Az}_{\tau}))$

Theorem (Oymak and Soltanolkotabi 2016)

For Gaussian \mathbf{A} and $\mu = 1/m$, with high probability

$$\|\mathbf{z}_{\tau} - \mu \mathbf{x}\|_{\ell_2} \leq \left(\sqrt{8\kappa_f^2 \frac{m_0}{m}} \right)^{\tau} \|\mathbf{z}_0 - \mu \mathbf{x}\|_{\ell_2} + 2\sigma \sqrt{\frac{m_0}{m}},$$

where $\mu := \mathbb{E}[wg(w)]$ and $\sigma^2 := \mathbb{E}[(g(w) - \mu w)^2]$ with $w \sim \mathcal{N}(0, 1)$.

Interpretation

$$g(\mathbf{Ax}) \approx \mu \mathbf{Ax} + \sigma w.$$

See [Plan and Vershynin 2014] for related results using convex programming.

Non-Gaussians

Subsampled Orthogonal with Random Sign (SORS)

Definition (Subsampled Orthogonal with Random Sign (SORS) matrices)

$\mathbf{F} \in \mathbb{R}^{n \times n}$ orthonormal matrix

$$\mathbf{F}^* \mathbf{F} = \mathbf{I} \quad \text{and} \quad \max_{i,j} |\mathbf{F}_{ij}| \leq \frac{\Delta}{\sqrt{n}}.$$

subsampled matrix $\mathbf{H} \in \mathbb{R}^{m \times n}$ with i.i.d. rows uniformly at random from rows of \mathbf{F} . $\mathbf{A} = \mathbf{H}\mathbf{D}$, with $\mathbf{D} \in \mathbb{R}^{n \times n}$ diagonal sign pattern.

SORS results

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu\mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

SORS results

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Recht and Soltanolkotabi 2016)

\mathbf{A} is a SORS matrix

$$m > c_{\Delta} \cdot (\log n)^4 \cdot m_0,$$

with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(c_{\Delta} \frac{m_0}{m} \log^4 n \right)^{\frac{\tau}{2}} \|\mathbf{x}\|_{\ell_2}.$$

SORS results

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Recht and Soltanolkotabi 2016)

\mathbf{A} is a SORS matrix

$$m > c_{\Delta} \cdot (\log n)^4 \cdot m_0,$$

with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(c_{\Delta} \frac{m_0}{m} \log^4 n \right)^{\frac{\tau}{2}} \|\mathbf{x}\|_{\ell_2}.$$

- Optimal (up to logs and constant) for “fast” matrices

SORS results

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Recht and Soltanolkotabi 2016)

\mathbf{A} is a SORS matrix

$$m > c_{\Delta} \cdot (\log n)^4 \cdot m_0,$$

with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(c_{\Delta} \frac{m_0}{m} \log^4 n\right)^{\frac{\tau}{2}} \|\mathbf{x}\|_{\ell_2}.$$

- Optimal (up to logs and constant) for “fast” matrices
- For sparse matrices see [Bourgain, Dirksen, Nelson 2014]

SORS results

$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{w}$ Algorithm:

- Start with $\mathbf{z}_0 = \mathbf{0}$,
- Run $\mathbf{z}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} + \mu \mathbf{A}^T(\mathbf{y} - \mathbf{A}\mathbf{z}_{\tau}))$

Theorem (Oymak, Recht and Soltanolkotabi 2016)

\mathbf{A} is a SORS matrix

$$m > c_{\Delta} \cdot (\log n)^4 \cdot m_0,$$

with high probability

$$\|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2} \leq \left(c_{\Delta} \frac{m_0}{m} \log^4 n \right)^{\frac{\tau}{2}} \|\mathbf{x}\|_{\ell_2}.$$

- Optimal (up to logs and constant) for “fast” matrices
- For sparse matrices see [Bourgain, Dirksen, Nelson 2014]
- Can we get down to $c_{\Delta}(\log n)(\log m_0)^2$? Maybe by [Haviv and Regev 15]?

Proof sketch

Gordon's lemma

- $\mathbf{C} \in \mathbb{R}^n$
- $\mathbf{A} \in \mathbb{R}^{m \times n}$, i.i.d. entries distributed as $\mathcal{N}(0, 1)$
- $b_m = \mathbb{E}[\|\mathbf{g}\|_{\ell_2}]$ with $\mathbf{g} \in \mathbb{R}^m$ distributed as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$

Gordon's lemma

- $\mathcal{C} \in \mathbb{R}^n$
- $\mathbf{A} \in \mathbb{R}^{m \times n}$, i.i.d. entries distributed as $\mathcal{N}(0, 1)$
- $b_m = \mathbb{E}[\|\mathbf{g}\|_{\ell_2}]$ with $\mathbf{g} \in \mathbb{R}^m$ distributed as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$

Then for all $\mathbf{u} \in \mathcal{C}$

- with prob. at least $1 - e^{-\frac{\eta^2}{2}}$

$$\frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \geq b_m - (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)$$

- with prob. at least $1 - e^{-\frac{\eta^2}{2}}$

$$\frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \leq b_m + (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)$$

Gordon's lemma

- $\mathcal{C} \in \mathbb{R}^n$
- $\mathbf{A} \in \mathbb{R}^{m \times n}$, i.i.d. entries distributed as $\mathcal{N}(0, 1)$
- $b_m = \mathbb{E}[\|\mathbf{g}\|_{\ell_2}]$ with $\mathbf{g} \in \mathbb{R}^m$ distributed as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$

Then for all $\mathbf{u} \in \mathcal{C}$

- with prob. at least $1 - e^{-\frac{\eta^2}{2}}$

$$\frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \geq b_m - (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)$$

- with prob. at least $1 - e^{-\frac{\eta^2}{2}}$

$$\frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \leq b_m + (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)$$

Thus for all $\mathbf{u} \in \mathcal{C}$

$$-2\frac{\omega + \eta}{b_m} \|\mathbf{u}\|_{\ell_2}^2 + \left(\frac{\omega + \eta}{b_m}\right)^2 \|\mathbf{u}\|_{\ell_2}^2 \leq \frac{1}{b_m^2} \|\mathbf{A}\mathbf{u}\|_{\ell_2}^2 - \|\mathbf{u}\|_{\ell_2}^2 \leq 2\frac{\omega + \eta}{b_m} \|\mathbf{u}\|_{\ell_2}^2 + \left(\frac{\omega + \eta}{b_m}\right)^2 \|\mathbf{u}\|_{\ell_2}^2$$

Bounding the convergence rate in the convex case

$$\rho\left(\frac{1}{m}\right) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathbb{S}^{n-1}} \mathbf{u}^T \left(\mathbf{I} - \frac{1}{m} \mathbf{A}^T \mathbf{A} \right) \mathbf{v}$$

Bounding the convergence rate in the convex case

$$\rho\left(\frac{1}{m}\right) = \sup_{\mathbf{u}, \mathbf{v} \in \mathcal{C}_f(\mathbf{x}) \cap \mathbb{S}^{n-1}} \mathbf{u}^T \left(\mathbf{I} - \frac{1}{m} \mathbf{A}^T \mathbf{A} \right) \mathbf{v}$$

Simple algebra

$$\begin{aligned} \mathbf{u}^T \left(\mathbf{I} - \frac{1}{m} \mathbf{A}^T \mathbf{A} \right) \mathbf{v} &= 2 \left(\left\| \frac{\mathbf{u} + \mathbf{v}}{2} \right\|_{\ell_2}^2 - \frac{1}{m} \left\| \mathbf{A} \left(\frac{\mathbf{u} + \mathbf{v}}{2} \right) \right\|_{\ell_2}^2 \right) \\ &\quad + \frac{1}{2} \left(\frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}^2}{m} - \|\mathbf{u}\|_{\ell_2}^2 \right) \\ &\quad + \frac{1}{2} \left(\frac{\|\mathbf{A}\mathbf{v}\|_{\ell_2}^2}{m} - \|\mathbf{v}\|_{\ell_2}^2 \right). \end{aligned}$$

Gordon for SORS matrices?

Definition (Subsampled Orthogonal with Random Sign (SORS) matrices)

$\mathbf{F} \in \mathbb{R}^{n \times n}$ orthonormal matrix

$$\mathbf{F}^* \mathbf{F} = \mathbf{I} \quad \text{and} \quad \max_{i,j} |\mathbf{F}_{ij}| \leq \frac{\Delta}{\sqrt{n}}.$$

subsampled matrix $\mathbf{H} \in \mathbb{R}^{m \times n}$ with i.i.d. rows uniformly at random from rows of \mathbf{F} . $\mathbf{A} = \mathbf{H}\mathbf{D}$, with $\mathbf{D} \in \mathbb{R}^{n \times n}$ diagonal sign pattern.

Gordon Holds for SORS matrices

Gordon for SORS matrices?

Definition (Subsampled Orthogonal with Random Sign (SORS) matrices)

$\mathbf{F} \in \mathbb{R}^{n \times n}$ orthonormal matrix

$$\mathbf{F}^* \mathbf{F} = \mathbf{I} \quad \text{and} \quad \max_{i,j} |\mathbf{F}_{ij}| \leq \frac{\Delta}{\sqrt{n}}.$$

subsampled matrix $\mathbf{H} \in \mathbb{R}^{m \times n}$ with i.i.d. rows uniformly at random from rows of \mathbf{F} . $\mathbf{A} = \mathbf{H}\mathbf{D}$, with $\mathbf{D} \in \mathbb{R}^{n \times n}$ diagonal sign pattern.

Gordon Holds for SORS matrices

Theorem (Oymak, Recht, and Soltanolkotabi 2016)

Assume $\mathbf{A} \in \mathbb{R}^{m \times n}$ is a SORS matrix. Then, with prob. at least $1 - 2e^{-\eta}$

$$\sup_{\mathbf{x} \in \mathcal{T}} \left| \|\mathbf{A}\mathbf{x}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2 \right| \leq \max(\delta, \delta^2) \cdot (\text{rad}(\mathcal{T}))^2,$$

as long as $m \geq C\Delta^2(1 + \eta)^2(\log n)^4 \frac{\max\left(1, \frac{\omega^2(\mathcal{T})}{(\text{rad}(\mathcal{T}))^2}\right)}{\delta^2}$

Open problem

- $\mathcal{C} \in \mathbb{R}^n$
- $\mathbf{A} \in \mathbb{R}^{m \times n}$, sub-sampled Fourier multiplied random sign.
- $b_m = \mathbb{E}[\|\mathbf{g}\|_{\ell_2}]$ with $\mathbf{g} \in \mathbb{R}^m$ distributed as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$

Open problem

- $\mathcal{C} \in \mathbb{R}^n$
- $\mathbf{A} \in \mathbb{R}^{m \times n}$, sub-sampled Fourier multiplied random sign.
- $b_m = \mathbb{E}[\|\mathbf{g}\|_{\ell_2}]$ with $\mathbf{g} \in \mathbb{R}^m$ distributed as $\mathcal{N}(\mathbf{0}, \mathbf{I}_m)$

Conjecture

- with prob. at least $1 - ???$

$$\inf_{\mathbf{u} \in \mathcal{C}} \frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \geq b_m - (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)$$

- with prob. at least $1 - ????$

$$\sup_{\mathbf{u} \in \mathcal{C}} \frac{\|\mathbf{A}\mathbf{u}\|_{\ell_2}}{\|\mathbf{u}\|_{\ell_2}} \leq (b_m + (\omega(\mathcal{C} \cap \mathbb{S}^{n-1}) + \eta)) \sqrt{\log n}????$$

This would settle universality conjecture in compressive sensing!

SORS matrices: from RIP to JL

- Restricted Isometry Property-RIP(s, δ)

$$\sup_{\mathbf{x}: \|\mathbf{x}\|_{\ell_0} \leq s} \left| \|\mathbf{F}\mathbf{x}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2 \right| \leq \max(\delta, \delta^2) \|\mathbf{x}\|_{\ell_2}^2$$

Theorem (Discrete JL embedding via RIP, Krahmer-Ward 2012)

Assume $\mathcal{T} \in \mathbb{R}^n$ finite points. Suppose

- $\mathbf{H} \in \mathbb{R}^{m \times n}$ obeys RIP(s, δ)
- $s \lesssim \log(|\mathcal{T}|)$ and $0 < \delta \leq \frac{\epsilon}{4}$

Then $\mathbf{A} = \mathbf{H}\mathbf{D}$ with \mathbf{D} diagonal sign pattern obeys

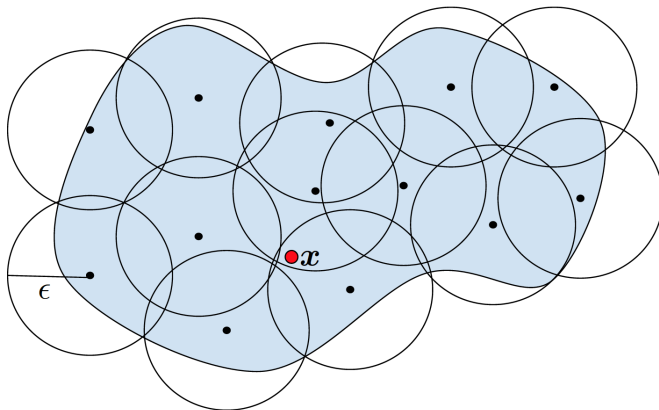
$$\sup_{\mathbf{x} \in \mathcal{T}} \left| \|\mathbf{A}\mathbf{x}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2 \right| \leq \max(\epsilon, \epsilon^2) \|\mathbf{x}\|_{\ell_2}^2,$$

with high probability.

From JL to Gordon?

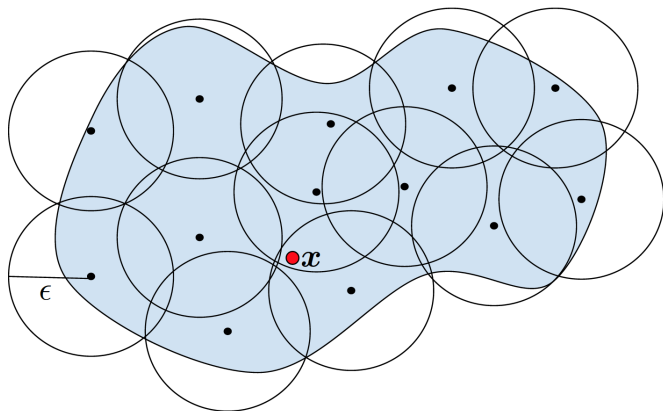
From JL to Gordon?

First attempt: covering



From JL to Gordon?

First attempt: covering

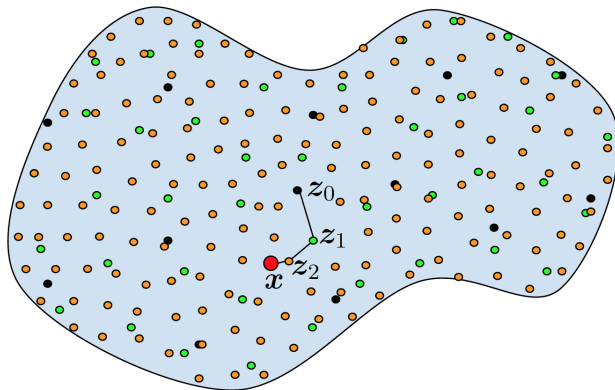


$$m \gtrsim \frac{\log |\mathcal{N}|}{\delta^2} \log^4 n \sim \frac{n}{m} \frac{\omega^2(\mathcal{T})}{\delta^4} \log^4 n \quad \Leftrightarrow \quad m \gtrsim \sqrt{n} \frac{\omega(\mathcal{T})}{\delta^2} \log^2 n,$$

From JL to Gordon?

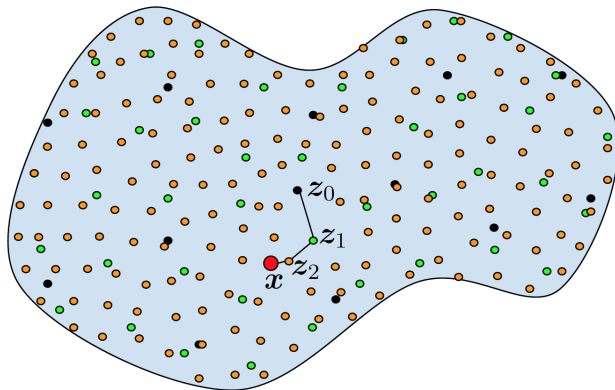
From JL to Gordon?

Second attempt: generic chaining



From JL to Gordon?

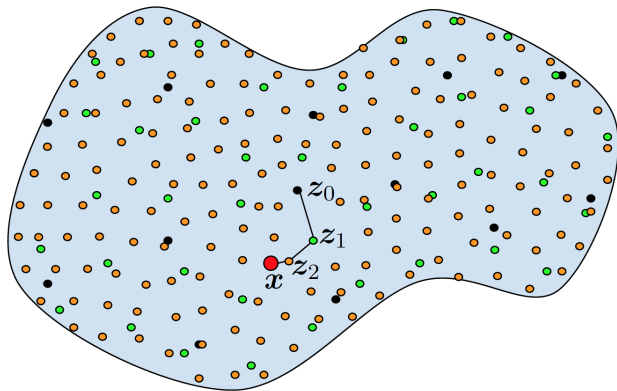
Second attempt: generic chaining



- Successive approximations of size $|\mathcal{T}_\ell| = 2^{2^\ell}$

From JL to Gordon?

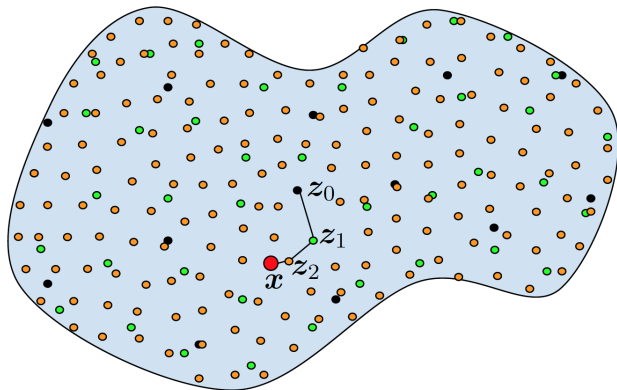
Second attempt: generic chaining



- Successive approximations of size $|\mathcal{T}_\ell| = 2^{2^\ell}$
- Different distortion levels $\delta_\ell = 2^{\ell/2} \frac{\delta}{\omega(\mathcal{T})}$

From JL to Gordon?

Second attempt: generic chaining



- Successive approximations of size $|\mathcal{T}_\ell| = 2^{2^\ell}$
- Different distortion levels $\delta_\ell = 2^{\ell/2} \frac{\delta}{\omega(\mathcal{T})}$

$$m \gtrsim \max_{\ell=1,2,\dots,L} \frac{\log |\mathcal{T}_\ell|}{\delta_\ell^2} \log^4 n = \max_{\ell=1,2,\dots,L} \frac{\log |2^{2^\ell}|}{2^\ell \frac{\delta^2}{\omega^2(\mathcal{T})}} \log^4 n = \frac{\omega^2(\mathcal{T})}{\delta^2} \sim \log^4 n \frac{\omega^2(\mathcal{T})}{\delta^2},$$

Is the over all distortion small enough?

After some work Discrete JL implies that with high probability

- For all $\mathbf{v} \in \mathcal{T}_{\ell-1} \cup \mathcal{T}_\ell \cup (\mathcal{T}_{\ell-1} - \mathcal{T}_\ell)$,

$$\|\mathbf{A}\mathbf{v}\|_{\ell_2} \leq (1 + \delta_\ell) \|\mathbf{v}\|_{\ell_2}.$$

- For all $\mathbf{v} \in \mathcal{T}_{\ell-1} \cup \mathcal{T}_\ell \cup (\mathcal{T}_{\ell-1} - \mathcal{T}_\ell)$,

$$|\|\mathbf{A}\mathbf{v}\|_{\ell_2}^2 - \|\mathbf{v}\|_{\ell_2}^2| \leq \max(\delta_\ell, \delta_\ell^2) \cdot \|\mathbf{v}\|_{\ell_2}^2.$$

- For all $\mathbf{u} \in \mathcal{T}_{\ell-1}$ and $\mathbf{v} \in \mathcal{T}_\ell - \{\mathbf{u}\} := \{\mathbf{y} - \mathbf{u} : \mathbf{y} \in \mathcal{T}_\ell\}$,

$$|\mathbf{u}^* \mathbf{A}^* \mathbf{A} \mathbf{v} - \mathbf{u}^* \mathbf{v}| \leq \max(\delta_\ell, \delta_\ell^2) \cdot \|\mathbf{u}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2}.$$

where

$$\delta_\ell = 2^{\ell/2} \frac{\delta}{\omega(\mathcal{T})}$$

Is the over all distortion small enough?

We are interested in bounding $|\|\mathbf{Ax}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2|$ for all $\mathbf{x} \in \mathcal{T}$.

Define $\tilde{L} = \max\left(0, \lfloor 2 \log_2\left(\frac{\omega(\mathcal{T})}{\delta}\right) \rfloor\right)$

$$\begin{aligned} |\|\mathbf{Ax}\|_{\ell_2}^2 - \|\mathbf{x}\|_{\ell_2}^2| &\leq \left| \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2}^2 - \|\mathbf{z}_{\tilde{L}}\|_{\ell_2}^2 \right| \\ &\quad + \left| \|\mathbf{Ax}\|_{\ell_2}^2 - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2}^2 \right| + \left| \|\mathbf{x}\|_{\ell_2}^2 - \|\mathbf{z}_{\tilde{L}}\|_{\ell_2}^2 \right| \\ &\leq \sum_{\ell=1}^{\tilde{L}} \left(\left| \|\mathbf{Az}_{\ell}\|_{\ell_2}^2 - \|\mathbf{z}_{\ell}\|_{\ell_2}^2 \right| - \left| \|\mathbf{Az}_{\ell-1}\|_{\ell_2}^2 - \|\mathbf{z}_{\ell-1}\|_{\ell_2}^2 \right| \right) \\ &\quad + \left| \|\mathbf{Ax}\|_{\ell_2}^2 - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2}^2 \right| + \left| \|\mathbf{x}\|_{\ell_2}^2 - \|\mathbf{z}_{\tilde{L}}\|_{\ell_2}^2 \right| \\ &\quad + \left| \|\mathbf{Az}_0\|_{\ell_2}^2 - \|\mathbf{z}_0\|_{\ell_2}^2 \right|. \end{aligned}$$

Is the over all distortion small enough? (First term)

For the first term

$$\left(\left| \|\mathbf{A}\mathbf{z}_\ell\|_{\ell_2}^2 - \|\mathbf{z}_\ell\|_{\ell_2}^2 \right| - \left| \|\mathbf{A}\mathbf{z}_{\ell-1}\|_{\ell_2}^2 - \|\mathbf{z}_{\ell-1}\|_{\ell_2}^2 \right| \right) \leq 10e_{\ell-1}\delta_\ell,$$

where $e_\ell = \text{dist}(\mathbf{x}, \mathcal{T}_\ell)$. Then

$$\begin{aligned} \sum_{\ell=1}^{\tilde{L}} \left(\left| \|\mathbf{A}\mathbf{z}_\ell\|_{\ell_2}^2 - \|\mathbf{z}_\ell\|_{\ell_2}^2 \right| - \left| \|\mathbf{A}\mathbf{z}_{\ell-1}\|_{\ell_2}^2 - \|\mathbf{z}_{\ell-1}\|_{\ell_2}^2 \right| \right) &\leq 10 \frac{\delta}{\omega(\mathcal{T})} \left(\sum_{\ell=1}^{\tilde{L}} 2^{\ell/2} e_{\ell-1} \right) \\ &= 10\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \left(\sum_{\ell=0}^{\tilde{L}-1} 2^{\ell/2} e_\ell \right) \\ &= 10\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \gamma_2(\mathcal{T}) \\ &\leq c\delta. \end{aligned}$$

Is the over all distortion small enough? (second term)

$$\begin{aligned} \left| \|A\mathbf{w}\|_{\ell_2} - \|A\mathbf{z}_{\bar{L}}\|_{\ell_2} \right| &= \left| \|A\mathbf{w}\|_{\ell_2} - \|A\mathbf{z}_L\|_{\ell_2} + \|A\mathbf{z}_L\|_{\ell_2} - \|A\mathbf{z}_{\bar{L}}\|_{\ell_2} \right| \\ &\leq \|A(\mathbf{w} - \mathbf{z}_L)\|_{\ell_2} + \|A(\mathbf{z}_L - \mathbf{z}_{\bar{L}})\|_{\ell_2} \\ &\leq \|A\| \|\mathbf{w} - \mathbf{z}_L\|_{\ell_2} + \left\| \sum_{\ell=\bar{L}+1}^L A(\mathbf{z}_\ell - \mathbf{z}_{\ell-1}) \right\|_{\ell_2} \\ &\leq \left(\frac{1}{4} 2^{\frac{L}{2}} \frac{\delta}{\omega(\mathcal{T})} + 1 \right) e_L + \sum_{\ell=\bar{L}+1}^L \|A(\mathbf{z}_\ell - \mathbf{z}_{\ell-1})\|_{\ell_2} \\ &\leq \left(\frac{1}{4} 2^{\frac{L}{2}} \frac{\delta}{\omega(\mathcal{T})} + 1 \right) e_L + \sum_{\ell=\bar{L}+1}^L \left(1 + 2^{\ell/2} \frac{\delta}{\omega(\mathcal{T})} \right) \|\mathbf{z}_\ell - \mathbf{z}_{\ell-1}\|_{\ell_2} \\ &\leq \frac{5}{4} 2^{L/2} \frac{\delta}{\omega(\mathcal{T})} e_L + \sum_{\ell=\bar{L}+1}^L 2^{\ell/2+1} \frac{\delta}{\omega(\mathcal{T})} \|\mathbf{z}_\ell - \mathbf{z}_{\ell-1}\|_{\ell_2} \\ &\leq \frac{5}{4} \frac{\delta}{\omega(\mathcal{T})} 2^{L/2} e_L + 4\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \sum_{\ell=\bar{L}+1}^L 2^{(\ell-1)/2} e_{\ell-1} \\ &\leq 4\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \left(\sum_{\ell=\bar{L}}^L 2^{\ell/2} e_\ell \right) \\ &\leq 4\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \gamma_2(\mathcal{T}). \end{aligned}$$

Is the over all distortion small enough? (second term cont.)

$$\begin{aligned} \left| \|\mathbf{Ax}\|_{\ell_2}^2 - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2}^2 \right| &\leq \left| \|\mathbf{Ax}\|_{\ell_2} - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2} \right| \left| \|\mathbf{Ax}\|_{\ell_2} + \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2} \right| \\ &\leq \left| \|\mathbf{Ax}\|_{\ell_2} - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2} \right|^2 + 2 \left| \|\mathbf{Ax}\|_{\ell_2} - \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2} \right| \|\mathbf{Az}_{\tilde{L}}\|_{\ell_2} \\ &\leq 32 \frac{\delta^2}{\omega^2(\mathcal{T})} \gamma_2^2(\mathcal{T}) + 16\sqrt{2} \frac{\delta}{\omega(\mathcal{T})} \gamma_2(\mathcal{T}) \\ &\leq c\delta. \end{aligned}$$

Is the over all distortion small enough? (third term)

Easy

$$\begin{aligned} | \|\mathbf{A}\mathbf{z}_0\|_{\ell_2}^2 - \|\mathbf{z}_0\|_{\ell_2}^2 | &\leq \max \left(\frac{\delta}{\omega(\mathcal{T})}, \left(\frac{\delta}{\omega(\mathcal{T})} \right)^2 \right) \|\mathbf{z}_0\|_{\ell_2}^2 \\ &\leq \max \left(\frac{\delta}{\omega(\mathcal{T})}, \left(\frac{\delta}{\omega(\mathcal{T})} \right)^2 \right). \end{aligned}$$

Is the failure probability small enough?

After some work Discrete JL implies that with high probability

- For all $\mathbf{v} \in \mathcal{T}_{\ell-1} \cup \mathcal{T}_\ell \cup (\mathcal{T}_{\ell-1} - \mathcal{T}_\ell)$,

$$\|\mathbf{A}\mathbf{v}\|_{\ell_2} \leq (1 + \delta_\ell) \|\mathbf{v}\|_{\ell_2}.$$

- For all $\mathbf{v} \in \mathcal{T}_{\ell-1} \cup \mathcal{T}_\ell \cup (\mathcal{T}_{\ell-1} - \mathcal{T}_\ell)$,

$$|\|\mathbf{A}\mathbf{v}\|_{\ell_2}^2 - \|\mathbf{v}\|_{\ell_2}^2| \leq \max(\delta_\ell, \delta_\ell^2) \cdot \|\mathbf{v}\|_{\ell_2}^2.$$

- For all $\mathbf{u} \in \mathcal{T}_{\ell-1}$ and $\mathbf{v} \in \mathcal{T}_\ell - \{\mathbf{u}\} := \{\mathbf{y} - \mathbf{u} : \mathbf{y} \in \mathcal{T}_\ell\}$,

$$|\mathbf{u}^* \mathbf{A}^* \mathbf{A} \mathbf{v} - \mathbf{u}^* \mathbf{v}| \leq \max(\delta_\ell, \delta_\ell^2) \cdot \|\mathbf{u}\|_{\ell_2} \|\mathbf{v}\|_{\ell_2}.$$

where

$$\delta_\ell = 2^{\ell/2} \frac{\delta}{\omega(\mathcal{T})}$$

The overall probability

$$\sum_{\ell=1}^L e^{-\ell(\eta+1)} \leq \sum_{\ell=1}^{\infty} e^{-\ell(\eta+1)} = \frac{e^{-(\eta+1)}}{1 - e^{-(\eta+1)}} \leq e^{-\eta},$$

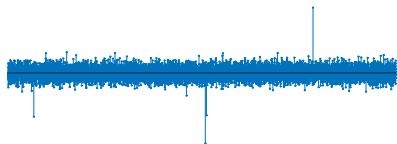
From de-noising to compressed sensing

Minimal number of data?

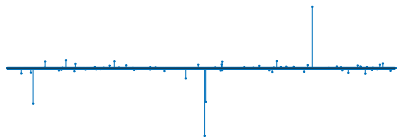
Minimal number of data?

Answer: Intimately related to de-noising capability of the function

Before Thresholding: $\mathbf{z}_\tau + \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{z}_\tau)$



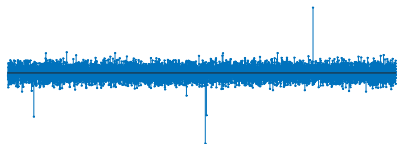
After Thresholding: $\mathcal{P}_\mathcal{K} (\mathbf{z}_\tau + \mathbf{A}^T (\mathbf{y} - \mathbf{A}\mathbf{z}_\tau))$



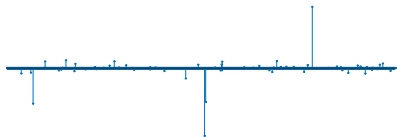
Minimal number of data?

Answer: Intimately related to de-noising capability of the function

Before Thresholding: $z_\tau + \mathbf{A}^T (\mathbf{y} - \mathbf{A}z_\tau)$



After Thresholding: $\mathcal{P}_\mathcal{K} (z_\tau + \mathbf{A}^T (\mathbf{y} - \mathbf{A}z_\tau))$



Conclusion

Intuitively *better de-noiser* should work with *less data*

Add Gaussian noise to parameter and then de-noise!

Add Gaussian noise to parameter and then de-noise!

Theorem (Oymak, Recht and Soltanolkotabi 2016)

For Gaussian \mathbf{A} , and convex \mathcal{K} minimal number of data is

$$\max_{\sigma} \frac{\mathbb{E} \|\mathcal{P}_{\mathcal{K}}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{x}\|_{\ell_2}^2}{\sigma^2} = m_0$$

For non-convex \mathcal{K}

$$\max_{\sigma} \frac{\mathbb{E} \|\mathcal{P}_{\mathcal{K}}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{x}\|_{\ell_2}^2}{\sigma^2} \leq 4m_0$$

$m_0 \propto \# \text{ params e.g. } (2s + 1) \log(n/s)$

Add Gaussian noise to parameter and then de-noise!

Theorem (Oymak, Recht and Soltanolkotabi 2016)

For Gaussian \mathbf{A} , and convex \mathcal{K} minimal number of data is

$$\max_{\sigma} \frac{\mathbb{E} \|\mathcal{P}_{\mathcal{K}}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{x}\|_{\ell_2}^2}{\sigma^2} = m_0$$

For non-convex \mathcal{K}

$$\max_{\sigma} \frac{\mathbb{E} \|\mathcal{P}_{\mathcal{K}}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{x}\|_{\ell_2}^2}{\sigma^2} \leq 4m_0$$

$m_0 \propto \#$ params e.g. $(2s + 1) \log(n/s)$ We establish a conjecture of Donoho, Johnston and Montanari 2011 up to a small constant.

Real experiment

Related AMP simulations
(Montanari IWT'12 plenary)
(Metzler-Maleki-Barinuk '14)

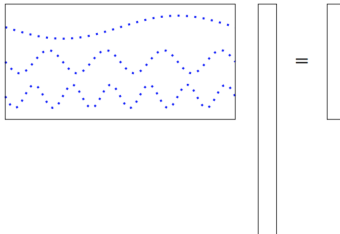
- Signal (1,000,080 pixels)



- Signal (1,000,080 pixels)

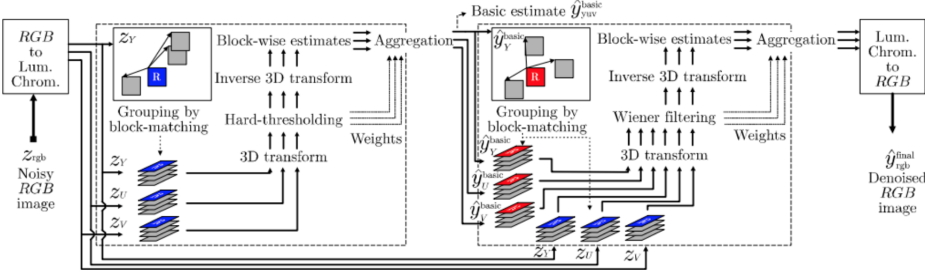


- Response (324,000 Measurements)



Good image denoiser

CBM3D



Is CBM3D a good denoiser

Signal

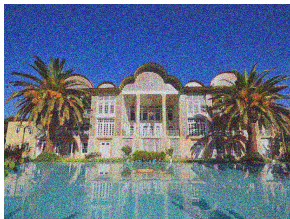


Is CBM3D a good denoiser

Signal



Signal+Gaussian noise

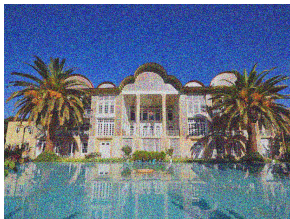


Is CBM3D a good denoiser

Signal



Signal+Gaussian noise

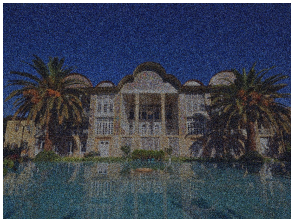


Result of CBMD3



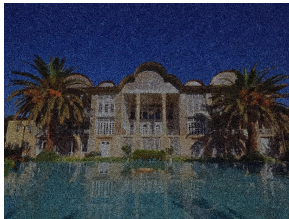
Iteration 1:

Before Thresholding:



Relerr=0.70687

After Thresholding:



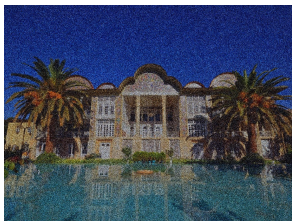
Relerr=0.57926

Original signal x :



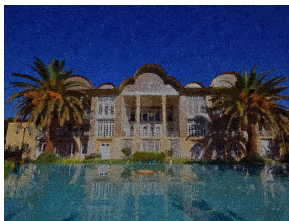
Iteration 2:

Before Thresholding:



Relerr=0.54722

After Thresholding:



Relerr=0.45153

Original signal x :



Iteration 5:

Before Thresholding:



Relerr=0.19012

After Thresholding:



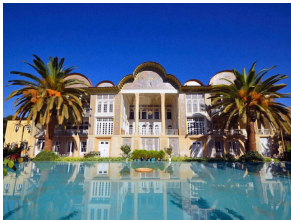
Relerr=0.15806

Original signal x :



Iteration 10:

Before Thresholding:



Relerr=0.062034

After Thresholding:



Relerr=0.081282

Original signal x :



Iteration 20:

Before Thresholding:



Relerr=0.04864

After Thresholding:



Relerr=0.061786

Original signal x :



Iteration 40:

Before Thresholding:

After Thresholding:

Original signal x :



Relerr=0.027613

Relerr=0.031075

Iteration 100:

Before Thresholding:

After Thresholding:

Original signal x :



Relerr=0.015708

Relerr=0.015732

Iteration 200:

Before Thresholding:

After Thresholding:

Original signal x :



Relerr=0.015603

Relerr=0.015618

Open Problem

Not very precise conjecture

Conjecture

Assume that for a mapping $\mathcal{S} : \mathbb{R}^m \rightarrow \mathbb{R}^m$ we have

$$\|\mathcal{S}(\mathbf{u}) - \mathcal{S}(\mathbf{v})\|_{\ell_2} \leq L \|\mathbf{u} - \mathbf{v}\|_{\ell_2}.$$

show that for all $\mathbf{z} \in \text{Range}(\mathcal{S})$

$$\left\| \mathcal{S} \left(\mathbf{x} + \left(\mathbf{I} - \frac{1}{m} \mathbf{A}^* \mathbf{A} \right) (\mathbf{z} - \mathbf{x}) \right) - \mathbf{x} \right\|_{\ell_2} \approx \left\| \mathcal{S} \left(\mathbf{x} + \frac{1}{\sqrt{m}} \| \mathbf{z} - \mathbf{x} \|_{\ell_2} \mathbf{g} \right) - \mathbf{x} \right\|_{\ell_2}$$

as long as

$$m \geq \max_{\sigma} \frac{\mathbb{E} \|\mathcal{S}(\mathbf{x} + \sigma \mathbf{z}) - \mathbf{x}\|_{\ell_2}^2}{\sigma^2}$$

Can prove it for a fixed \mathbf{z}

Summary

- provable (non)convex optimization with generic coefficients via local search
- the main challenge in nonconvex optimization is concentration
- for realistic data models need generic chaining

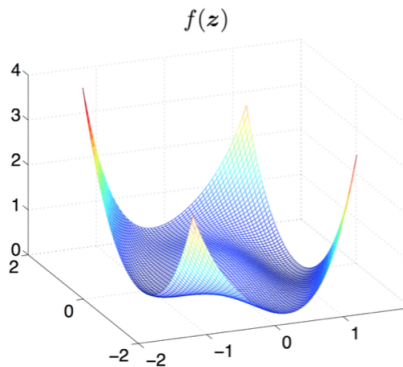
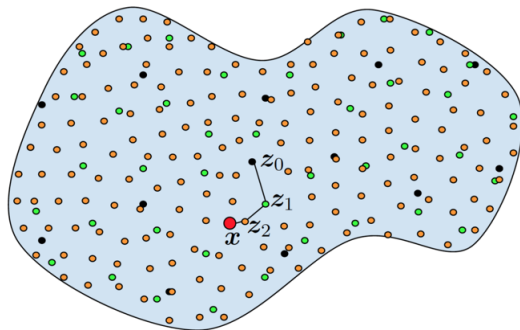
Summary

- provable (non)convex optimization with generic coefficients via local search
- the main challenge in nonconvex optimization is concentration
- for realistic data models need generic chaining

References:

- Phase retrieval
 - Phase retrieval via Wirtinger flow: Theory and algorithms E. J. Candes, X. Li, and M. Soltanolkotabi
 - Algorithms and theory for clustering and non-convex quadratic programming. M. Soltanolkotabi 2014.
 - Experimental robustness of Fourier Ptychography phase retrieval algorithms 2015 (in collaboration with the computational imaging lab at UC Berkeley)
- Low-rank matrix recovery
 - Low-rank Solutions of Linear Matrix Equations via Procrustes Flow. S. Tu, R. Boczar, M. Soltanolkotabi, and B. Recht 2015.
- Sharp time-data tradeoffs for (non)convex projected gradients
 - Sharp Time–Data Tradeoffs for Linear Inverse Problems. S. Oymak, B. Recht and M. Soltanolkotabi
 - Fast and reliable parameter estimation from nonlinear observations. S. Oymak and M. Soltanolkotabi.

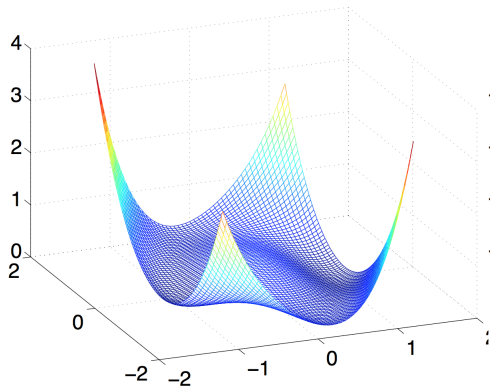
Generic chaining meets nonconvex optimization



Thanks!

Just follow the gradient?

$f(z)$



$-f(z)$

