

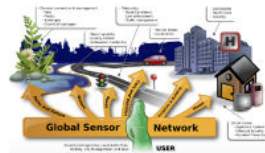
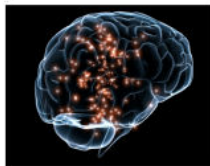
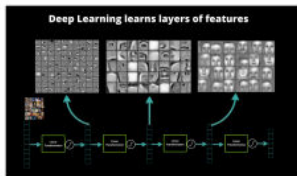
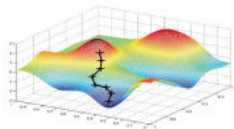
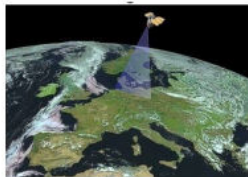
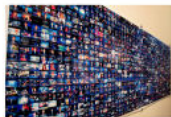
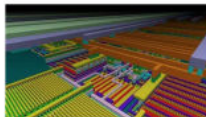
# Nonconvex Optimization for High-Dimensional Learning: From Phase Retrieval to Submodular Maximization

Mahdi Soltanolkotabi  
Department of Electrical Engineering



Interdisciplinary Seminar Series  
North Carolina State

# Nonconvexity is everywhere



# 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.

Idealogy

*“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  $\ell_1$  norm, Low-rank use nuclear norm, etc.

convex relaxations are not perfect

## convex relaxations are not perfect

- Computation and memory: convex programs maybe inefficient

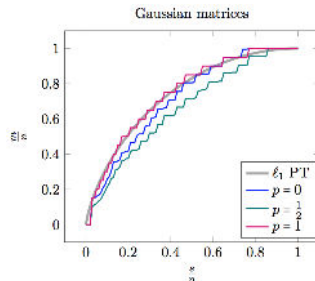


# convex relaxations are not perfect

- Computation and memory: convex programs maybe inefficient



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





Local search heuristics are often surprisingly effective...

# Local search heuristics are often surprisingly effective...

ANALYSE MATHÉMATIQUE. — *Méthode générale pour la résolution des systèmes d'équations simultanées*; par M. AUGUSTIN CAUCHY.

« Étant donné un système d'équations simultanées qu'il s'agit de résoudre, on commence ordinairement par les réduire à une seule, à l'aide d'éliminations successives, sauf à résoudre définitivement, s'il se peut, l'équation résultante. Mais il importe d'observer, 1<sup>o</sup> que, dans un grand nombre de cas, l'élimination ne peut s'effectuer en aucune manière; 2<sup>o</sup> que l'équation résultante est généralement très-compiquée, lors même que les équations données sont assez simples. Pour ces deux motifs, on conçoit qu'il serait très-utile de connaître une méthode générale qui pût servir à résoudre directement un système d'équations simultanées. Telle est celle que j'ai obtenue, et dont je vais dire ici quelques mots. Je me bornerai pour l'instant à indiquer les principes sur lesquels elle se fonde, me proposant de revenir avec plus de détails sur le même sujet, dans un prochain Mémoire.

\* Soit d'abord

$$u = f(x, y, z)$$

une fonction de plusieurs variables  $x, y, z, \dots$ , qui ne devienne jamais négative et qui reste continue, du moins entre certaines limites. Pour trouver les valeurs de  $x, y, z, \dots$ , qui vérifieront l'équation

$$(1) \quad u = 0,$$

il suffira de faire décroître indéfiniment la fonction  $u$ , jusqu'à ce qu'elle s'évanouisse. Or soient

$$x, y, z, \dots$$

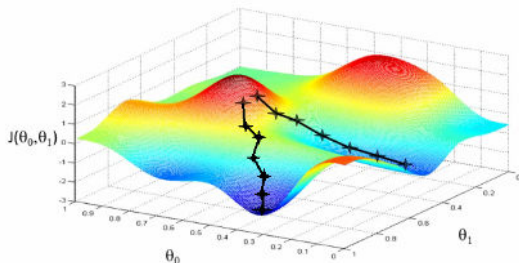
des valeurs particulières attribuées aux variables  $x, y, z, \dots$ ;  $u$  la valeur correspondante de  $u$ ;  $X, Y, Z, \dots$  les valeurs correspondantes de  $D_x u, D_y u, D_z u, \dots$ , et  $\alpha, \beta, \gamma, \dots$  des accroissements très-petits attribués aux valeurs particulières  $x, y, z, \dots$ . Quand on posera

$$x = x + \alpha, \quad y = y + \beta, \quad z = z + \gamma, \dots,$$

on aura sensiblement

$$(2) \quad u = f(x + \alpha, y + \beta, z + \gamma, \dots) = u + \alpha X + \beta Y + \gamma Z + \dots$$

# When should we just follow the gradient?



## Two stories with a common theme

- *Story I: Structured Signal Recovery from Quadratic Measurements*
- *Story II: Submodular Maximization*

## *Structured Signal Recovery from Quadratic Measurements*

*Specific example:*

*Specific example:*

*Sparse recovery from quadratic measurements*

# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements from an  $s$ -sparse signal

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



# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements from an  $s$ -sparse signal

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

Find an  $s$ -sparse signal from quadratic measurements

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

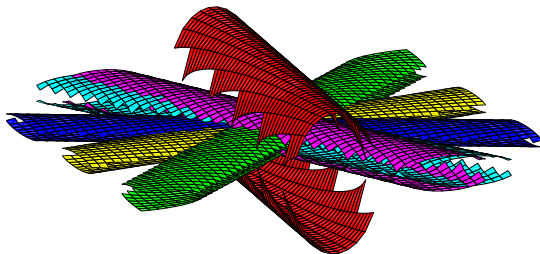
# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements from an  $s$ -sparse signal

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

Find an  $s$ -sparse signal from quadratic measurements

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

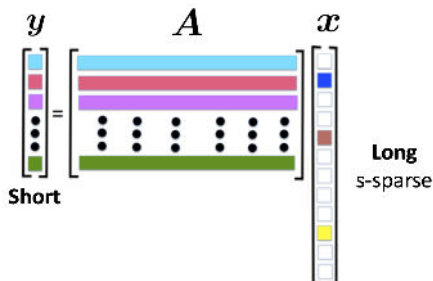


One of the universal forms of combinatorial problems, NP-hard in general.

# Sparse Signal Recovery from Linear Measurements

## Linear Measurements

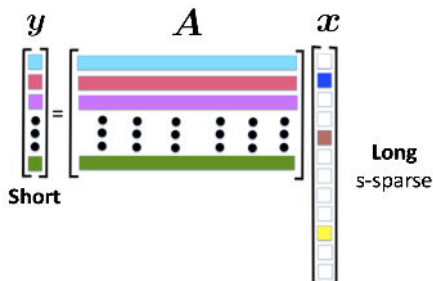
$$y_r = \langle \mathbf{a}_r, \mathbf{x} \rangle \quad r = 1, 2, \dots, m \quad \Leftrightarrow \quad \mathbf{y} = \mathbf{A}\mathbf{x}$$



# Sparse Signal Recovery from Linear Measurements

## Linear Measurements

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



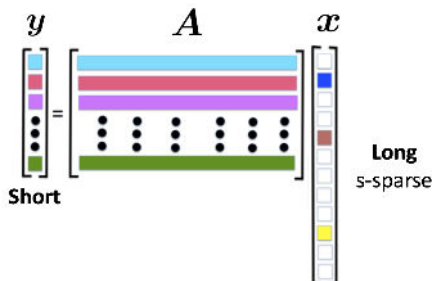
Sample complexity for **uniqueness**:

$$m \gtrsim s \log(n/s) \text{ generic measurements}$$

# Sparse Signal Recovery from Linear Measurements

## Linear Measurements

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

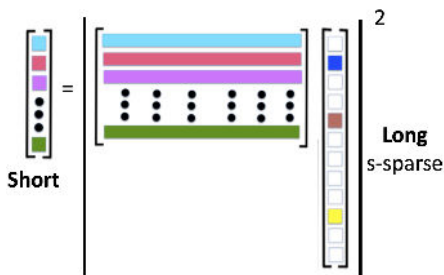


Sample complexity for **uniqueness**:  $m \gtrsim s \log(n/s)$  generic measurements  
Sample complexity of **convex relaxation**:  $m \gtrsim s \log(n/s)$  generic measurements

# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements

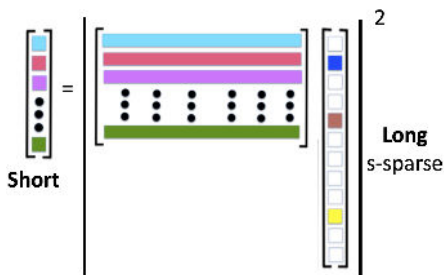
$$y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2 \quad r = 1, 2, \dots, m \quad \Leftrightarrow \quad \mathbf{y} = |\mathbf{A}\mathbf{x}|^2$$



# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements

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

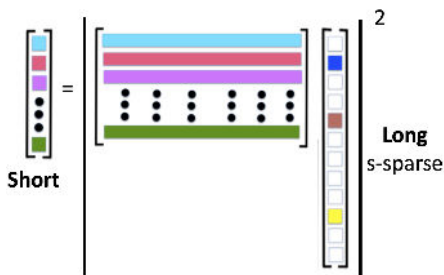


Sample complexity for **uniqueness**:  $m \gtrsim s \log(n/s)$  generic measurements

# Sparse Signal Recovery from Quadratic Measurements

Quadratic measurements

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



Sample complexity for **uniqueness**:  $m \gtrsim s \log(n/s)$  generic measurements

Sample complexity for **exact recovery**: ???????



*First attempt: Convex Optimization*

# Semidefinite Relaxation with Sparsity

$$\min \quad \|z\|_{\ell_0} \quad \text{subject to} \quad \mathbf{y}_r = |\mathbf{a}_r^* z|^2 = [\mathcal{A}(zz^*)]_r.$$

# Semidefinite Relaxation with Sparsity

$$\min \quad \|z\|_{\ell_0} \quad \text{subject to} \quad \mathbf{y}_r = |\mathbf{a}_r^* z|^2 = [\mathcal{A}(zz^*)]_r.$$

**Lifting:**  $\mathbf{Z} = zz^*$  Relax rank one constraint

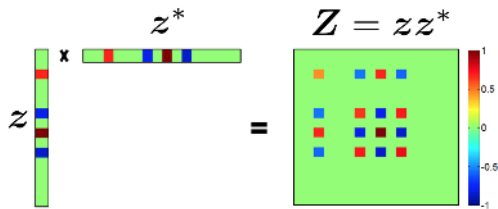
$$\min \quad \|z\|_{\ell_0} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

# Semidefinite Relaxation with Sparsity

$$\min \|z\|_{\ell_0} \quad \text{subject to} \quad \mathbf{y}_r = |\mathbf{a}_r^* z|^2 = [\mathcal{A}(zz^*)]_r.$$

**Lifting:**  $\mathbf{Z} = zz^*$  Relax rank one constraint

$$\min \|z\|_{\ell_0} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

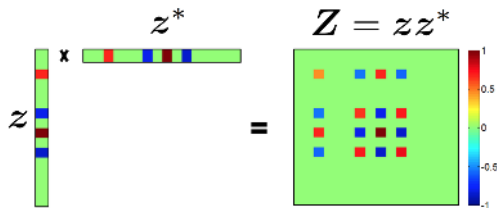


# Semidefinite Relaxation with Sparsity

$$\min \|z\|_{\ell_0} \quad \text{subject to} \quad y_r = |a_r^* z|^2 = [\mathcal{A}(zz^*)]_r.$$

**Lifting:**  $Z = zz^*$  Relax rank one constraint

$$\min \|z\|_{\ell_0} \quad \text{subject to} \quad y = \mathcal{A}(Z) \quad \text{and} \quad Z \succeq 0.$$



SDP relaxation

$$\min \|Z\|_{\ell_1} \quad \text{subject to} \quad y = \mathcal{A}(Z) \quad \text{and} \quad Z \succeq 0.$$

For Phase Retrieval [Shechtman et. al. 2011, Li and Voroninski 2013].

# Solving random quadratic equations

Given an  $s$ -sparse signal  $\mathbf{x} \in \mathbb{C}^n$ , measurements of the form

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

with  $\mathbf{a}_r$  i.i.d. complex random vector with each entry  $\sim \mathcal{CN}(0, 1)$ .

# Solving random quadratic equations

Given an  $s$ -sparse signal  $\mathbf{x} \in \mathbb{C}^n$ , measurements of the form

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

with  $\mathbf{a}_r$  i.i.d. complex random vector with each entry  $\sim \mathcal{CN}(0, 1)$ .

Theorem (Li and Voroninski (2013))

Using  $m \gtrsim s^2 \log n$  Gaussian measurements with high probability

$$\mathbf{x}\mathbf{x}^* = \arg \min \|\mathbf{Z}\|_{\ell_1} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

# Solving random quadratic equations

Given an  $s$ -sparse signal  $\mathbf{x} \in \mathbb{C}^n$ , measurements of the form

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

with  $\mathbf{a}_r$  i.i.d. complex random vector with each entry  $\sim \mathcal{CN}(0, 1)$ .

Theorem (Li and Voroninski (2013))

Using  $m \gtrsim s^2 \log n$  Gaussian measurements with high probability

$$\mathbf{x}\mathbf{x}^* = \arg \min \|\mathbf{Z}\|_{\ell_1} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

Maybe these results are not optimal...



# Solving random quadratic equations

Given an  $s$ -sparse signal  $\mathbf{x} \in \mathbb{C}^n$ , measurements of the form

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

with  $\mathbf{a}_r$  i.i.d. complex random vector with each entry  $\sim c\mathcal{N}(0, 1)$ .

Theorem (Li and Voroninski (2013))

Using  $m \gtrsim s^2 \log n$  Gaussian measurements with high probability

$$\mathbf{x}\mathbf{x}^* = \arg \min \|\mathbf{Z}\|_{\ell_1} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

Maybe these results are not optimal...

Theorem (Li and Voroninski (2013), [Oymak, Jalali, Fazel, Hassibi, Eldar (2014)])

With Gaussian measurements if

$$\mathbf{x}\mathbf{x}^* = \arg \min \|\mathbf{Z}\|_{\ell_1} \quad \text{subject to} \quad \mathbf{y} = \mathcal{A}(\mathbf{Z}) \quad \text{and} \quad \mathbf{Z} \succeq \mathbf{0}.$$

holds with high probability.

Then

$$m \gtrsim \frac{s^2}{\log^2 n}.$$

## Data Barriers...

$$m \gtrsim s \log(n/s) \quad \text{versus} \quad m \gtrsim \frac{s^2}{\log^2 n}$$

Uniqueness                      convex relaxation

# Data Barriers...

$$m \gtrsim s \log(n/s) \quad \text{versus} \quad m \gtrsim \frac{s^2}{\log^2 n}$$

Uniqueness

convex relaxation



## *Engineering Motivation*

# 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  $|\hat{x}(f_1, f_2)|$ ?*

# Phase retrieval (discrete 1D model)



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

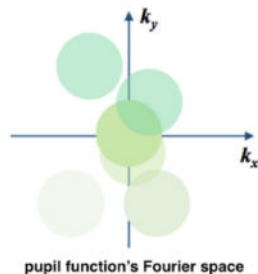
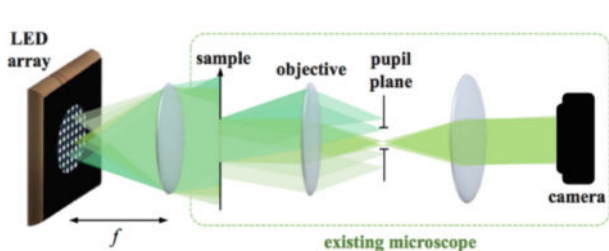
$$|\mathbf{f}_k^* \mathbf{x}|^2 = \mathbf{y}_k \quad k \in \{1, 2, \dots, n\} = [n]$$

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

- Phase retrieval is impossible, inherent ambiguity.

# Resolving ambiguity?

Solution: Create diversity



(stolen from the Waller Lab)

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2 \quad \text{where} \quad \mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \mathbf{A}_2 \\ \vdots \\ \mathbf{A}_L \end{bmatrix}$$

with  $\mathbf{A}_\ell \in \mathbb{C}^{n \times n}$



# Data Barriers...

$$m \gtrsim s \log(n/s) \quad \text{versus} \quad m \gtrsim \frac{s^2}{\log^2 n}$$

Uniqueness

convex relaxation



*Second attempt: Nonconvex Optimization*

# Solving quadratic equation by non-convex optimization (no constraints)

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

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \ell(\mathbf{y}_r, |\mathbf{a}_r^* \mathbf{z}|)$$

- 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

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

**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], [Wang and Giannakis], [Zhang and Liang 2016]  
established  $m \gtrsim n$  via variantes of Wirtinger Flow

Don't like initialization?

# Don't like initialization?

Theorem (Soltanolkotabi 2017)

*With  $m \gtrsim n \log n$  Gaussian measurements all local optima are global optima and cubic regularization converges to a global optima in  $\text{poly}(n)$  iterations.*

Earlier [Sun, Qu, Wright 2016]: All local optima are global optima with  $m \gtrsim n \log^3 n$  and trust region methods converge to a global optima in  $\text{poly}(n)$  iterations.

# Are local optima global optima?

Are saddles the only problem with nonconvexity?

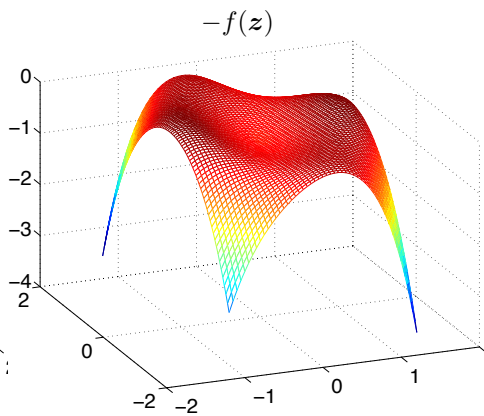
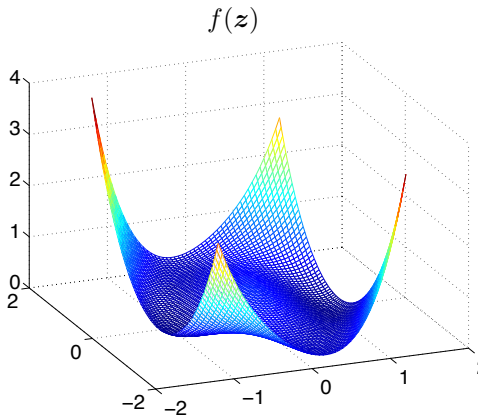


# Are local optima global optima?

Are saddles the only problem with nonconvexity?

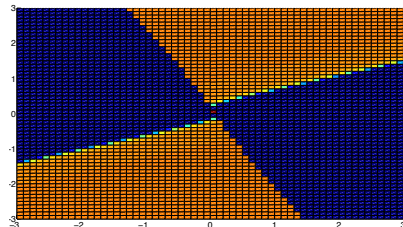
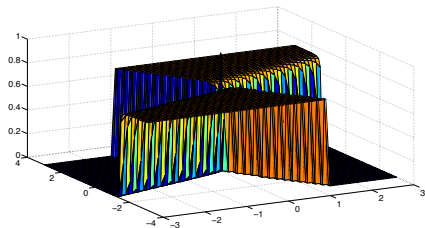
Example:  $\mathbf{x} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$ . Measurements  $y_r = |\mathbf{a}_r^* \mathbf{x}|^2$ ,  $r = 1, 2, \dots, m$ , with  $m = 4$ .

cost function:  $f(\mathbf{z}) = \frac{1}{4m} \sum_{r=1}^m (y_r - |\mathbf{a}_r^* \mathbf{x}|^2)^2$



# Which initial solutions work?

Run gradient descent ( $z_{\tau+1} = z_{\tau} - \mu \nabla f(z_{\tau})$ ) from different initial points.



## Solving quadratic equations via Projected Wirtinger Flow (PWF)

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

# Solving quadratic equations via Projected Wirtinger Flow (PWF)

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the gradient:

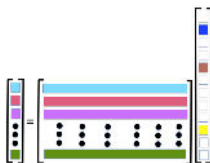
$$\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}} \left( \mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}) \right).$$

where

$$\mathcal{K} = \{ \mathbf{z} : \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}) \}$$

# What is the sample complexity of PWF?

Simpler question: Linear inverse problems



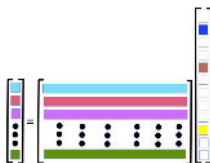
$y = Ax$ ,  $y \in \mathbb{R}^m$ ,  $A \in \mathbb{R}^{m \times n}$ , and  $x \in \mathbb{R}^n$  with  $m \ll n$ .

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

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

# What is the sample complexity of PWF?

Simpler question: Linear inverse problems



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

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

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

Theorem (Chandrasekaran, Recht, Parrilo, and Willskey 2012-Amelunxen, Lotz, McCoy, Tropp 2014)

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

$$m \approx m_0(\mathcal{R}, x),$$

*then with high probability  $\hat{x} = x$*

e.g. for an  $s$ -sparse signal  $m \geq 2s \log(n/s)$

## What is the sample complexity of PWF? (local)

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the gradient:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

Assume  $m \gtrsim m_0 \log n$ . Using  $0 \leq \mu \leq \mu_0/n$ , with high probability

Starting from any initial point  $\mathbf{z}_0$  obeying

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

we have

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

# What is the sample complexity of PWF? (local)

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the gradient:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

Assume  $m \gtrsim m_0 \log n$ . Using  $0 \leq \mu \leq \mu_0/n$ , with high probability

Starting from any initial point  $\mathbf{z}_0$  obeying

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

we have

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

- e.g. for sparsity  $m \gtrsim 2s \log(n/s) \log n$
- previous known result for local neighborhood via Thresholded WF  
 $m \gtrsim s^2 \log n$  [Cai, Li, Ma 2015]



# What is the sample complexity of PWF? (global)

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the gradient:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

*With  $m \gtrsim m_0 \log n$  Gaussian measurements all local optima are global optima and cubic regularization converges in  $\text{poly}(n)$  iterations.*

# What is the sample complexity of PWF? (global)

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m \left( y_r - |\mathbf{a}_r^* \mathbf{z}|^2 \right)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the gradient:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

*With  $m \gtrsim m_0 \log n$  Gaussian measurements all local optima are global optima and cubic regularization converges in  $\text{poly}(n)$  iterations.*

- e.g. for sparsity  $m \gtrsim 2s \log(n/s) \log n$

Removing logs and other things...

## Removing logs and other things...

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m (\sqrt{y_r} - |\mathbf{a}_r^* \mathbf{z}|)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the “gradient”:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

Assume  $m \gtrsim m_0$ . Using  $0 \leq \mu \leq \mu_0$ , with high probability  
Starting from any initial point  $\mathbf{z}_0$  obeying

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

we have

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

## Removing logs and other things...

Let  $\mathbf{a}_r \in \mathbb{R}^n$  be i.i.d.  $\mathcal{N}(\mathbf{0}, \mathbf{I})$  and  $y_r = |\langle \mathbf{a}_r, \mathbf{x} \rangle|^2$  for  $r = 1, 2, \dots, m$ .

$$\min_{\mathbf{z} \in \mathbb{C}^n} f(\mathbf{z}) := \frac{1}{2m} \sum_{r=1}^m (\sqrt{y_r} - |\mathbf{a}_r^* \mathbf{z}|)^2 \quad \text{subject to} \quad \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x}).$$

Follow the “gradient”:  $\mathbf{z}_{\tau+1} := \mathcal{P}_{\mathcal{K}}(\mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}))$  with  $\mathcal{K} = \{\mathbf{z} : \mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})\}$ .

Theorem (Soltanolkotabi 2017)

Assume  $m \gtrsim m_0$ . Using  $0 \leq \mu \leq \mu_0$ , with high probability  
Starting from any initial point  $\mathbf{z}_0$  obeying

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

we have

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

- This result also holds for nonconvex regularizers!

# Connecting sample complexity to mini-max denoising

Theorem (Soltanolkotabi 2016)

*For any set  $\mathcal{K}$ , as long as*

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

*PWF works.*

[Oymak and Hassibi 2014-Oymak, Recht, Soltanolkotabi 2016] + [Amelunxen, Lotz, McCoy, Tropp 2014] shows equivalence between min-max denoising and data complexity of linear inverse problems

# Theoretical implications

- signal with entries  $\pm 1$

# Theoretical implications

- signal with entries  $\pm 1$   
no problem best Gaussian denoiser is actually *tanh*



# Theoretical implications

- signal with entries  $\pm 1$   
no problem best Gaussian denoiser is actually *tanh*
- optimization over integers?

# Theoretical implications

- signal with entries  $\pm 1$   
no problem best Gaussian denoiser is actually *tanh*
- optimization over integers?  
no problem just threshold to the closest integer...
- many others

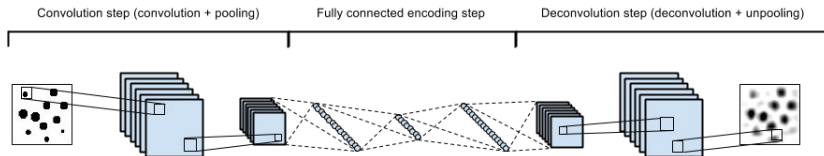
# Implications for imaging systems

What projection or non-linear shrinkage should you use?

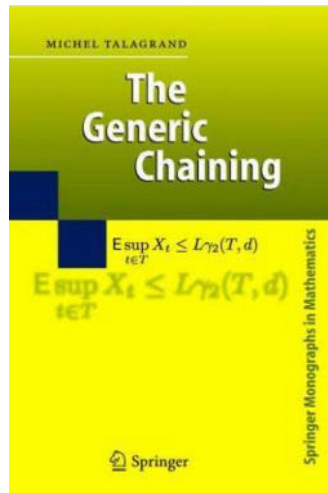
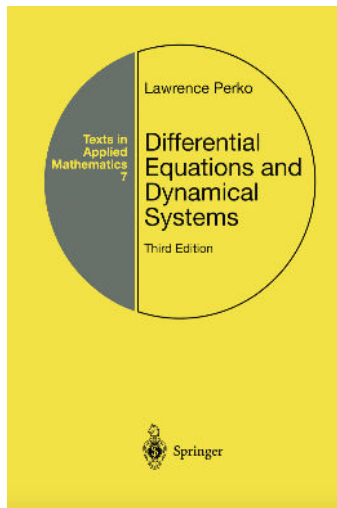
# Implications for imaging systems

What projection or non-linear shrinkage should you use?

We use GDS file from IBM add Gaussian noise and just learn the best denoiser...



# Tools



## Regularity condition?

$$\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \frac{1}{\alpha} \|\mathbf{z} - \mathbf{x}\|_{\ell_2}^2 + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2$$

## Regularity condition?

$$\langle \nabla f(\mathbf{z}), \mathbf{z} - \mathbf{x} \rangle \geq \frac{1}{\alpha} \|\mathbf{z} - \mathbf{x}\|_{\ell_2}^2 + \frac{1}{\beta} \|\nabla f(\mathbf{z})\|_{\ell_2}^2$$

Not really ...

# Proof Sketch

$$\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}).$$

Want to prove

$$\|\mathbf{z}_{\tau+1} - \mathbf{x}\|_{\ell_2} \leq \frac{1}{2} \|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2}$$



# Proof Sketch

$$\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}).$$

Want to prove

$$\|\mathbf{z}_{\tau+1} - \mathbf{x}\|_{\ell_2} \leq \frac{1}{2} \|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2}$$

Define the stochastic process

$$X_{\mathbf{u}, \mathbf{z}} = \frac{\mathbf{u}^T (\mathbf{z} - \mu \nabla f(\mathbf{z}))}{\|\mathbf{z} - \mathbf{x}\|_{\ell_2}}$$

# Proof Sketch

$$\mathbf{z}_{\tau+1} = \mathbf{z}_{\tau} - \mu_{\tau} \nabla f(\mathbf{z}_{\tau}).$$

Want to prove

$$\|\mathbf{z}_{\tau+1} - \mathbf{x}\|_{\ell_2} \leq \frac{1}{2} \|\mathbf{z}_{\tau} - \mathbf{x}\|_{\ell_2}$$

Define the stochastic process

$$X_{\mathbf{u}, \mathbf{z}} = \frac{\mathbf{u}^T (\mathbf{z} - \mu \nabla f(\mathbf{z}))}{\|\mathbf{z} - \mathbf{x}\|_{\ell_2}}$$

We prove that for all  $\mathbf{u} \in \mathbb{R}^n$  and  $\mathbf{z}$  obeying  $\mathcal{R}(\mathbf{z}) \leq \mathcal{R}(\mathbf{x})$

$$\sup_{\mathbf{u} \in \mathbb{S}^{n-1}, \mathbf{z} \in \mathcal{K}} X_{\mathbf{u}, \mathbf{z}} \leq \frac{1}{2}$$

## *Submodular Maximization*

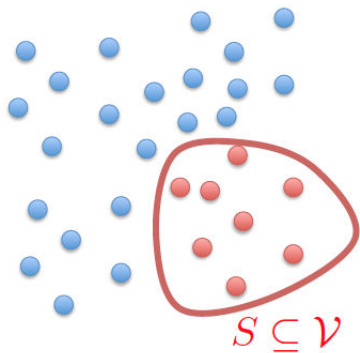
Collaborators: **Hamed Hassani** and Amin Karbasi

## *Submodular Maximization*

Collaborators: **Hamed Hassani** and Amin Karbasi

(Introductory figures/slides stolen from Stefanie Jegelka and Andreas Krause)

# Set Function Maximization



- ground set  $\mathcal{V}$
- (scoring) function

$$F : 2^{\mathcal{V}} \rightarrow \mathbb{R}_+$$

$$\max F(S)$$

# Maximizing monotone functions

$$\max_{S \subseteq \mathcal{V}} F(S) \quad \text{subject to} \quad |S| \leq k$$

# Maximizing monotone functions

$$\max_{S \subseteq V} F(S) \quad \text{subject to} \quad |S| \leq k$$

Greedy algorithm

- $S_0 =$
- for  $i = 0, 1, \dots, k - 1$

$$e^* = \arg \max_{e \in V/S_i} F(S_i \cup \{e\})$$

$$S_{i+1} = S_i \cup \{e^*\}$$

## Theory for greedy

$$\max_{S \subseteq \mathcal{V}} F(S) \quad \text{subject to} \quad |S| \leq k$$



# Theory for greedy

$$\max_{S \subseteq V} F(S) \quad \text{subject to} \quad |S| \leq k$$

Theorem (Nemhauser, Fisher, Wolsey '78)

*F monotone submodular. Then solution of greedy obeys*

$$F(\hat{S}) \geq \left(1 - \frac{1}{e}\right) F(S^*)$$

# Theory for greedy

$$\max_{S \subseteq V} F(S) \quad \text{subject to} \quad |S| \leq k$$

Theorem (Nemhauser, Fisher, Wolsey '78)

*F monotone submodular. Then solution of greedy obeys*

$$F(\hat{S}) \geq \left(1 - \frac{1}{e}\right) F(S^*)$$

No poly-time algorithm can do better than that!

# Why not just use greedy

- Many cases don't have exact function evaluations
- Greedy takes  $O(nk)$  time. What if  $n$  is large?
- What if the function is not submodular

# Making things continuous

sample item  $e$  with probability  $x_e$

$$f_M(x) = \mathbb{E}_{S \sim x} [F(S)]$$

$$= \sum_{S \subseteq \mathcal{V}} F(S) \prod_{e \in S} x_e \prod_{e \notin S} (1 - x_e)$$

	$x$	
$p(1) =$	0.5	✗
$p(2) =$	1.0	●
$p(3) =$	0.5	●
	0.2	✗
	0.2	✗

Basis for continuous greedy [Vondrak et. al.]

## Just follow the gradient

$$\mathbf{x}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{x}_{\tau} + \mu_{\tau} \nabla f_M(\mathbf{x}_{\tau}))$$

where

$$\mathcal{K} = \{\mathbf{z} \in \mathbb{R}_+^n : \sum_{i=1}^n z_i = k \quad 0 \leq z_i \leq 1\}$$

## How well does it work?

$$\max_{\mathcal{S} \subset \{1, 2, \dots, n\}} F(\mathcal{S}) = \log \det(\mathbf{I} + \mathbf{A}_{\mathcal{S}, \mathcal{S}}) \quad \text{subject to} \quad |\mathcal{S}| \leq k$$

## How well does it work?

$$\max_{\mathcal{S} \subset \{1, 2, \dots, n\}} F(\mathcal{S}) = \log \det(\mathbf{I} + \mathbf{A}_{\mathcal{S}, \mathcal{S}}) \quad \text{subject to} \quad |\mathcal{S}| \leq k$$

Greedy: 67.1    Gradient Descent: 74.81

# Stochastic Methods

Assume access to a stochastic oracle

$$\mathbb{E}[\mathbf{g}_t] = \nabla f_M(\mathbf{x}_t).$$

Run

$$\mathbf{x}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{x}_{\tau} + \mu_{\tau} \mathbf{g}_{\tau})$$

where

$$\mathcal{K} = \{\mathbf{z} \in \mathbb{R}_+^n : \sum_{i=1}^n z_i = k \quad 0 \leq z_i \leq 1\}$$



## Some theory

$$\mathbf{x}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{x}_{\tau} + \mu_{\tau} \mathbf{g}_{\tau})$$

### Theorem (Stochastic Gradient Method)

#### Assumptions

- $R^2 = \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\ell_2}^2$
- $f_M$  is  $L$ -smooth, monotone and multilinear extension of submodular
- stochastic oracle  $\mathbf{g}_t$  obeying

$$\mathbb{E}[\mathbf{g}_t] = \nabla f_M(\mathbf{x}_t) \quad \text{and} \quad \mathbb{E}[\|\mathbf{g}_t - \nabla f_M(\mathbf{x}_t)\|_{\ell_2}^2] \leq \sigma^2.$$

Run stochastic gradient updates with  $\mu_t = \frac{1}{L + \frac{\sigma}{R}\sqrt{t}}$ . Then,

$$\mathbb{E}[f_M(\mathbf{x}_T)] \geq \text{OPT} \left( \frac{1}{2} - \left( \frac{R^2 L}{T} + 2 \frac{R\sigma}{\sqrt{T}} \right) \right).$$

## Some theory

$$\mathbf{x}_{\tau+1} = \mathcal{P}_{\mathcal{K}}(\mathbf{x}_{\tau} + \mu_{\tau} g_{\tau})$$

### Theorem (Stochastic Gradient Method)

#### Assumptions

- $R^2 = \sup_{\mathbf{x}, \mathbf{y} \in \mathcal{K}} \frac{1}{2} \|\mathbf{x} - \mathbf{y}\|_{\ell_2}^2$
- $f_M$  is  $L$ -smooth, monotone and multilinear extension of submodular
- stochastic oracle  $\mathbf{g}_t$  obeying

$$\mathbb{E}[\mathbf{g}_t] = \nabla f_M(\mathbf{x}_t) \quad \text{and} \quad \mathbb{E}[\|\mathbf{g}_t - \nabla f_M(\mathbf{x}_t)\|_{\ell_2}^2] \leq \sigma^2.$$

Run stochastic gradient updates with  $\mu_t = \frac{1}{L + \frac{\sigma}{R}\sqrt{t}}$ . Then,

$$\mathbb{E}[f_M(\mathbf{x}_T)] \geq \text{OPT} \left( \frac{1}{2} - \left( \frac{R^2 L}{T} + 2 \frac{R\sigma}{\sqrt{T}} \right) \right).$$

- With Mirror descent can ensure  $L$  is constant
- Can get better approximation ratio starting from 0

# Conclusion

- Convex relaxations may be inefficient in terms of sample complexity
- discussed results towards breaking this barrier
- a lot of exciting barriers to think about e.g. planted clique
- interesting directions for bridging the gap between discrete and continuous optimization

# 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.
  - Structured signal recovery from quadratic measurements: breaking data barriers via nonconvex optimization, M. Soltanolkotabi, 2017
- 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, 2016.
- Submodular maximization
  - Stochastic gradient methods for submodular maximization H. Hassani, M. Soltanolkotabi, and A. Karbasi.

# Thanks!

When should we just follow the gradient?

