

The Power of Nonconvex Optimization in Solving Random Quadratic Systems of Equations

Yuxin Chen (Princeton)



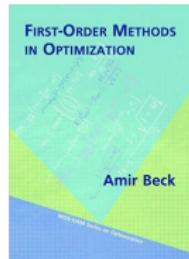
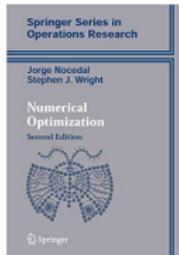
Emmanuel Candès (Stanford)



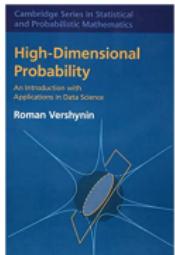
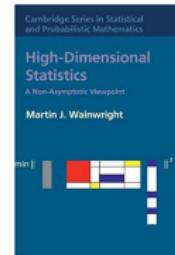
*Y. Chen, E. J. Candès, Communications on Pure and Applied Mathematics
vol. 70, no. 5, pp. 822-883, May 2017*

Agenda

1. The power of nonconvex optimization in solving random quadratic systems of equations (Aug. 28)
2. Random initialization and implicit regularization in nonconvex statistical estimation (Aug. 29)
3. The projected power method: an efficient nonconvex algorithm for joint discrete assignment from pairwise data (Sep. 3)
4. Spectral methods meets asymmetry: two recent stories (Sep. 4)
5. Inference and uncertainty quantification for noisy matrix completion (Sep. 5)



nonconvex optimization

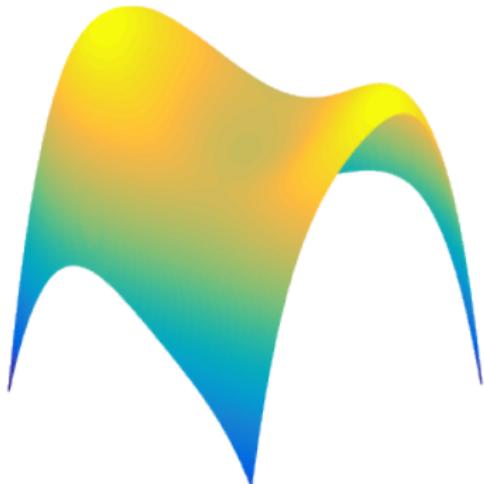


(high-dimensional) statistics

Nonconvex problems are everywhere

Maximum likelihood estimation is usually nonconvex

$$\begin{array}{lll} \text{maximize}_{\boldsymbol{x}} & \ell(\boldsymbol{x}; \text{data}) & \rightarrow \text{may be nonconcave} \\ \text{subj. to} & \boldsymbol{x} \in \mathcal{S} & \rightarrow \text{may be nonconvex} \end{array}$$

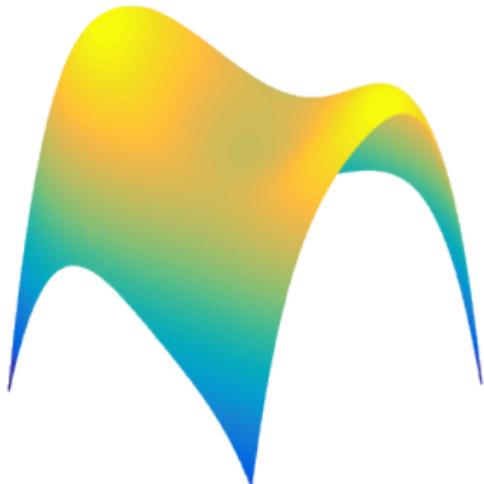


Nonconvex problems are everywhere

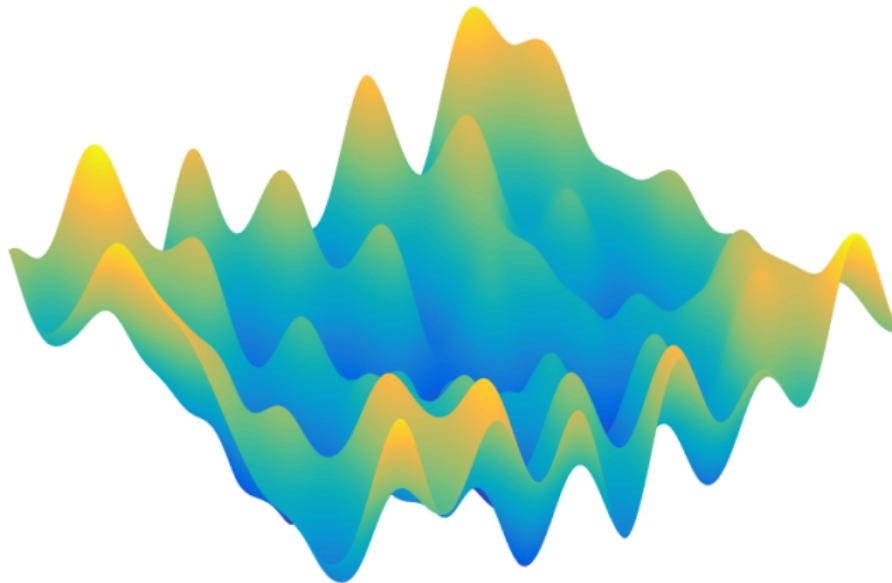
Maximum likelihood estimation is usually nonconvex

$$\begin{array}{ll} \text{maximize}_{\boldsymbol{x}} & \ell(\boldsymbol{x}; \text{data}) \rightarrow \text{may be nonconcave} \\ \text{subj. to} & \boldsymbol{x} \in \mathcal{S} \rightarrow \text{may be nonconvex} \end{array}$$

- low-rank matrix completion
- robust principal component analysis
- graph clustering
- dictionary learning
- blind deconvolution
- learning neural nets
- ...



Nonconvex optimization may be super scary



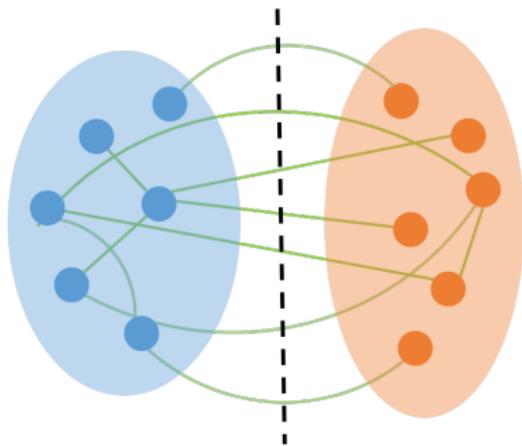
There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

Example: solving quadratic programs is hard

Finding maximum cut in a graph is

$$\begin{array}{ll}\text{maximize}_x & \boldsymbol{x}^\top \boldsymbol{W} \boldsymbol{x} \\ \text{subj. to} & x_i^2 = 1, \quad i = 1, \dots, n\end{array}$$



Example: solving quadratic programs is hard



"I can't find an efficient algorithm, but neither can all these people."

Fig credit: coding horror

\$1,000,000 question

One strategy: convex relaxation

Can relax into convex problems by

- finding convex surrogates (e.g. compressed sensing, matrix completion)
- lifting into higher dimensions (e.g. Max-Cut)

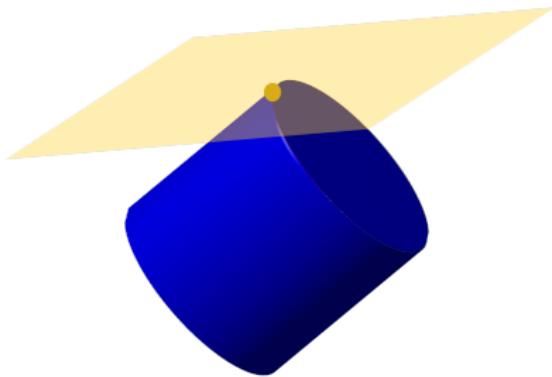
Example of convex surrogate: low-rank matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$$\text{minimize}_{\mathbf{M}} \text{ rank}(\mathbf{M}) \quad \text{subj. to data constraints}$$

↓ cvx surrogate

$$\text{minimize}_{\mathbf{M}} \text{ nuc-norm}(\mathbf{M}) \quad \text{subj. to data constraints}$$



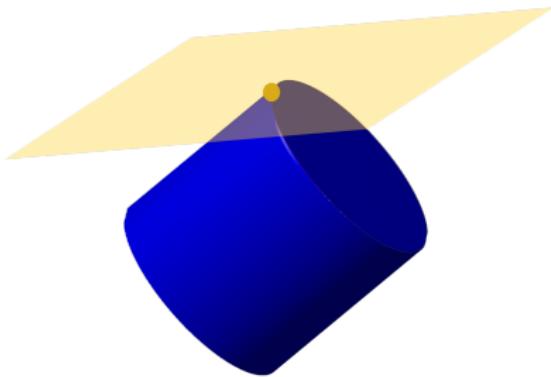
Example of convex surrogate: low-rank matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$$\text{minimize}_M \text{ rank}(M) \quad \text{subj. to data constraints}$$

↓ cvx surrogate

$$\text{minimize}_M \text{ nuc-norm}(M) \quad \text{subj. to data constraints}$$



Robust variation used everyday by Netflix

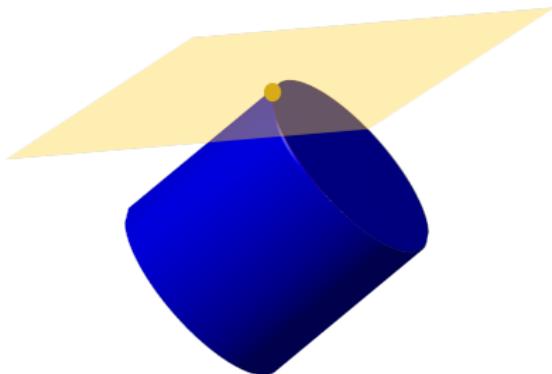
Example of convex surrogate: low-rank matrix completion

— Fazel '02, Recht, Parrilo, Fazel '10, Candès, Recht '09

$$\text{minimize}_M \text{ rank}(M) \quad \text{subj. to data constraints}$$

↓ cvx surrogate

$$\text{minimize}_M \text{ nuc-norm}(M) \quad \text{subj. to data constraints}$$

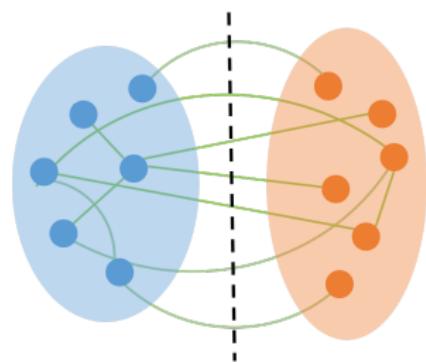


Robust variation used everyday by Netflix

Problem: operate in *full* matrix space even though X is low-rank

Example of lifting: Max-Cut

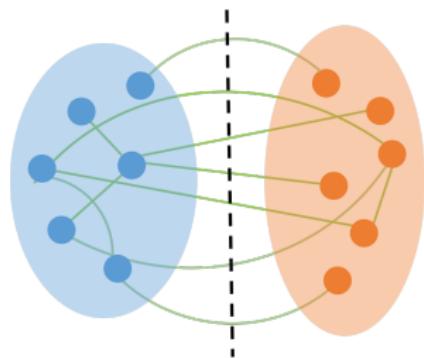
— Goemans, Williamson '95



$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

Example of lifting: Max-Cut

— Goemans, Williamson '95



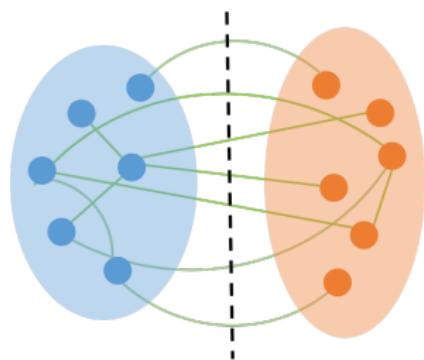
$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let X be xx^\top

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

Example of lifting: Max-Cut

— Goemans, Williamson '95



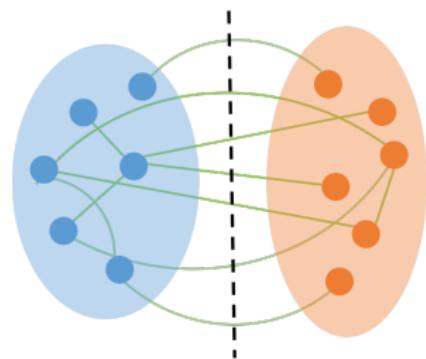
$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let X be xx^\top

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

Example of lifting: Max-Cut

— Goemans, Williamson '95



$$\begin{aligned} & \text{maximize}_x && x^\top W x \\ & \text{subj. to} && x_i^2 = 1, \quad i = 1, \dots, n \end{aligned}$$

↓ let X be xx^\top

$$\begin{aligned} & \text{maximize}_X && \langle X, W \rangle \\ & \text{subj. to} && X_{i,i} = 1, \quad i = 1, \dots, n \\ & && X \succeq 0 \\ & && \text{rank}(X) = 1 \end{aligned}$$

Problem: explosion in dimensions ($\mathbb{R}^n \rightarrow \mathbb{R}^{n \times n}$)

How about optimizing nonconvex problems directly without lifting?

A case study: solving random quadratic systems of equations

Solving quadratic systems of equations

$$\begin{array}{c} \mathbf{A} \\ \mathbf{x} \\ \mathbf{Ax} \\ \mathbf{y} = |\mathbf{Ax}|^2 \end{array} = \rightarrow$$

The diagram illustrates the computation of quadratic residuals. On the left, a matrix \mathbf{A} is shown as a 4x4 grid of red and white squares. To its right is a vector \mathbf{x} represented by a vertical column of four colored squares (blue, dark blue, light blue, blue). Below them is the product \mathbf{Ax} , shown as a vertical column of ten numbers: 1, -3, 2, -1, 4, 2, -2, -1, 3, 4. A large black arrow points from the \mathbf{Ax} column to the right, leading to the final result $\mathbf{y} = |\mathbf{Ax}|^2$, which is also a vertical column of ten numbers: 1, 9, 4, 1, 16, 4, 4, 1, 9, 16.

| |
|----|
| 1 |
| 9 |
| 4 |
| 1 |
| 16 |
| 4 |
| 4 |
| 1 |
| 9 |
| 16 |

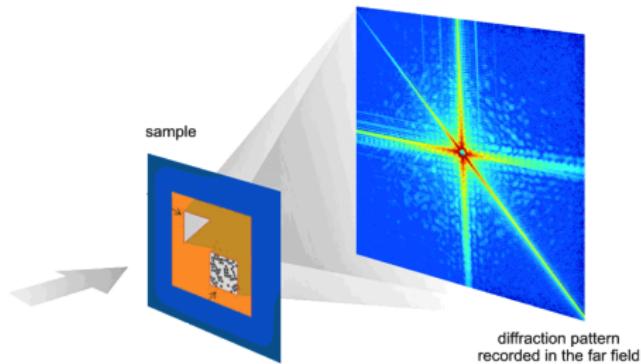
Solve for $\mathbf{x} \in \mathbb{C}^n$ in m quadratic equations

$$y_k \approx |\langle \mathbf{a}_k, \mathbf{x} \rangle|^2, \quad k = 1, \dots, m$$

Motivation: a missing phase problem in imaging science

Detectors record **intensities** of diffracted rays

- $x(t_1, t_2) \longrightarrow$ Fourier transform $\hat{x}(f_1, f_2)$

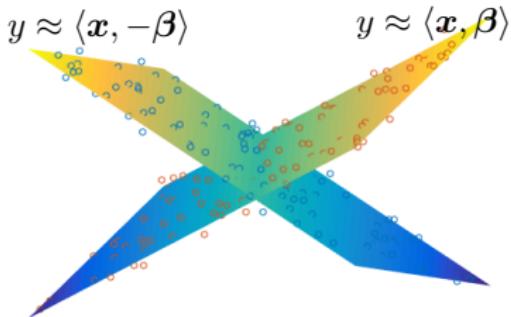


$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Phase retrieval: recover true signal $x(t_1, t_2)$ from intensity measurements

Motivation: latent variable models

Example: mixture of regression

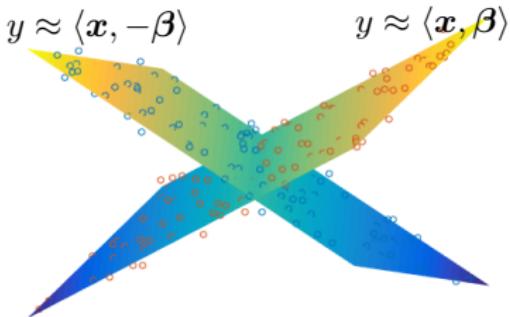


- Samples $\{(y_k, \mathbf{x}_k)\}$: drawn from one of two unknown regressors β and $-\beta$

$$y_k \approx \begin{cases} \langle \mathbf{x}_k, \beta \rangle, & \text{with prob. 0.5} \\ \langle \mathbf{x}_k, -\beta \rangle, & \text{else} \end{cases} \quad (\text{labels: latent variables})$$

Motivation: latent variable models

Example: mixture of regression



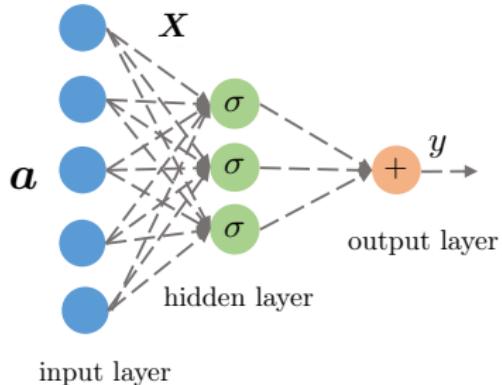
- Samples $\{(y_k, \mathbf{x}_k)\}$: drawn from one of two unknown regressors β and $-\beta$

$$y_k \approx \begin{cases} \langle \mathbf{x}_k, \beta \rangle, & \text{with prob. 0.5} \\ \langle \mathbf{x}_k, -\beta \rangle, & \text{else} \end{cases} \quad (\text{labels: latent variables})$$

- equivalent to observing $|y_k|^2 \approx |\langle \mathbf{x}_k, \beta \rangle|^2$
- Goal: estimate β

Motivation: learning neural nets with quadratic activation

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17



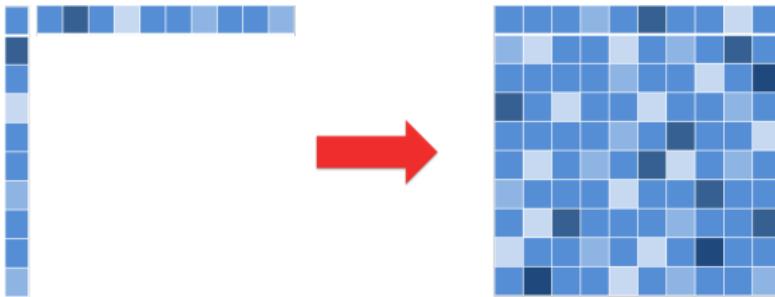
input features: \mathbf{a} ; weights: $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_r]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i) \stackrel{\sigma(z)=z^2}{=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i)^2$$

An equivalent view: low-rank factorization

Lifting: introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ to linearize constraints

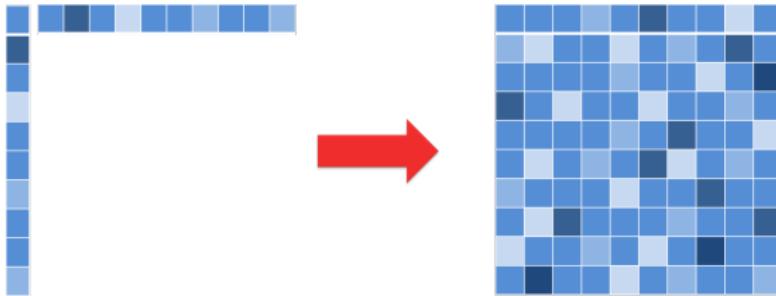
$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*) \mathbf{a}_k \implies y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k$$



An equivalent view: low-rank factorization

Lifting: introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ to linearize constraints

$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*) \mathbf{a}_k \implies y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k$$

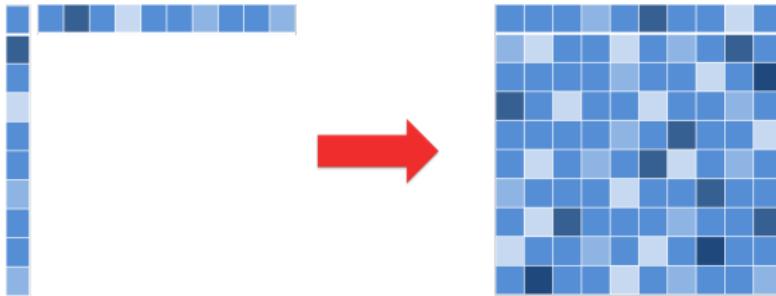


$$\begin{aligned} & \text{find} \quad \mathbf{X} \succeq 0 \\ & \text{s.t.} \quad y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m \\ & \quad \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

An equivalent view: low-rank factorization

Lifting: introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ to linearize constraints

$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*) \mathbf{a}_k \implies y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k$$

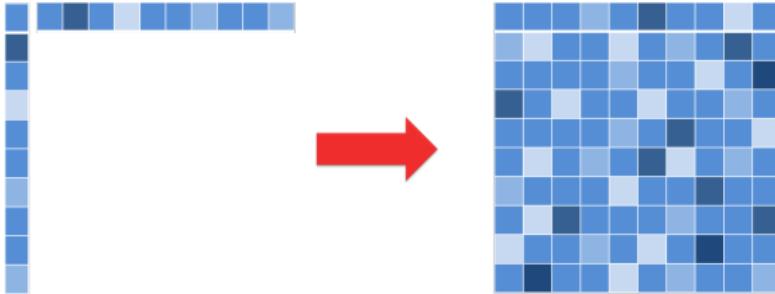


$$\begin{aligned} & \text{find} \quad \mathbf{X} \succeq 0 \\ & \text{s.t.} \quad y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m \\ & \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

An equivalent view: low-rank factorization

Lifting: introduce $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ to linearize constraints

$$y_k = |\mathbf{a}_k^* \mathbf{x}|^2 = \mathbf{a}_k^* (\mathbf{x}\mathbf{x}^*) \mathbf{a}_k \implies y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k$$



$$\begin{aligned} & \text{find} \quad \mathbf{X} \succeq 0 \\ & \text{s.t.} \quad y_k = \mathbf{a}_k^* \mathbf{X} \mathbf{a}_k, \quad k = 1, \dots, m \\ & \text{rank}(\mathbf{X}) = 1 \end{aligned}$$

Works well if $\{\mathbf{a}_k\}$ are random, but huge increase in dimensions

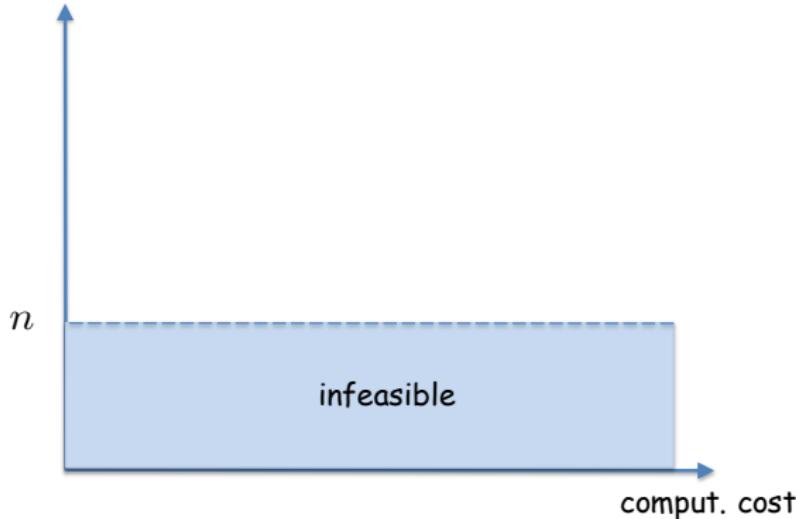
Prior art (before our work)

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$

sample complexity

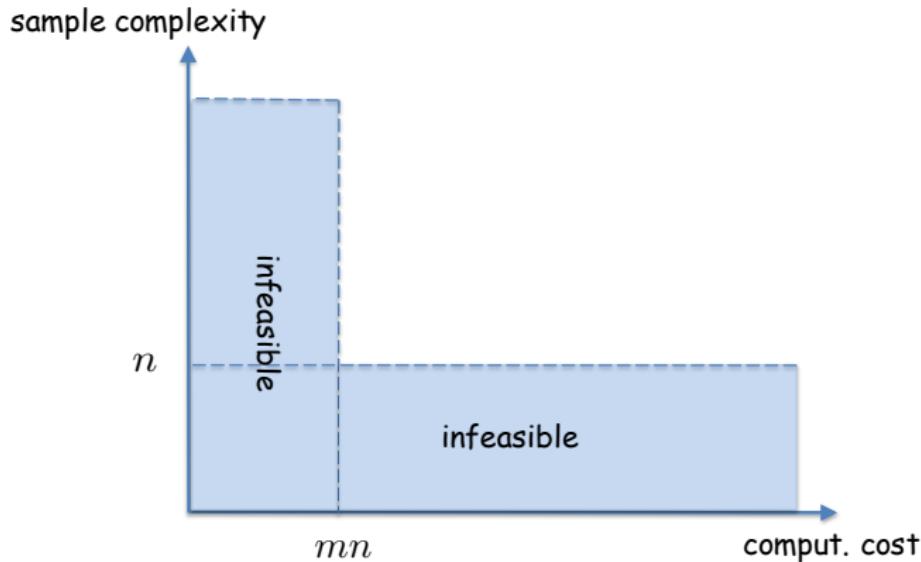


Prior art (before our work)

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$

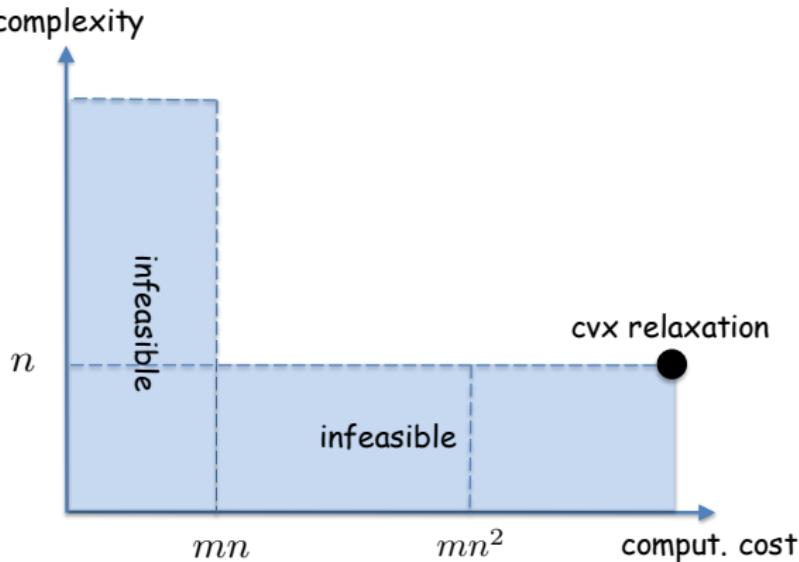


Prior art (before our work)

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$

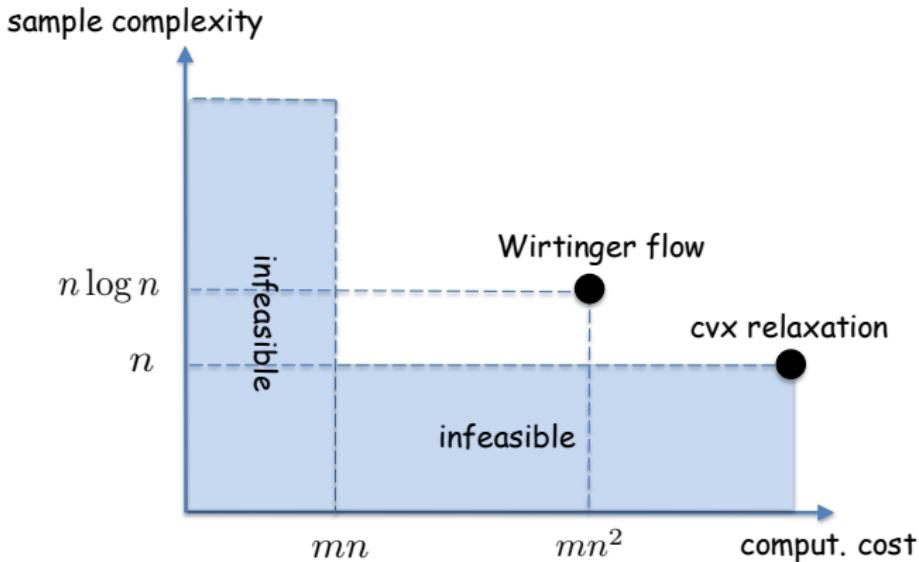


Prior art (before our work)

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$

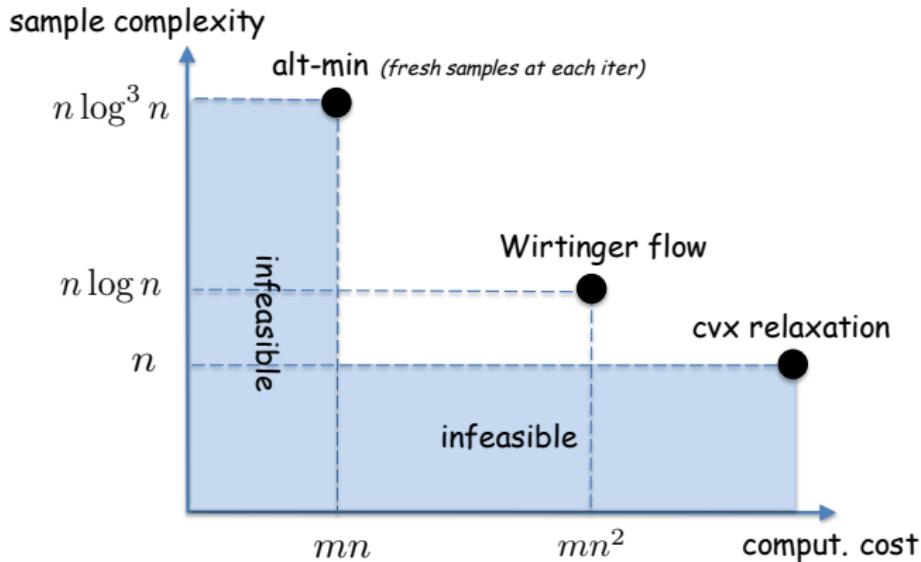


Prior art (before our work)

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$

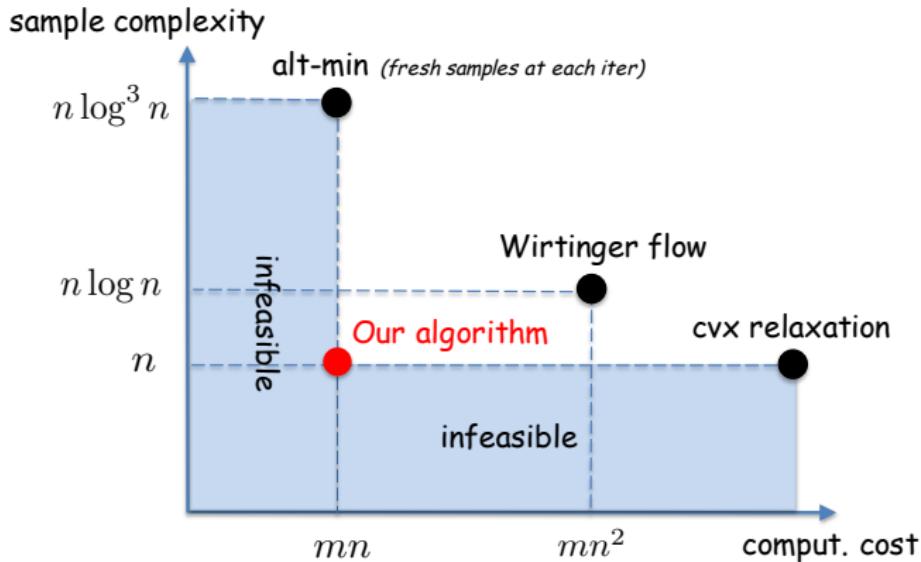


A glimpse of our results

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$



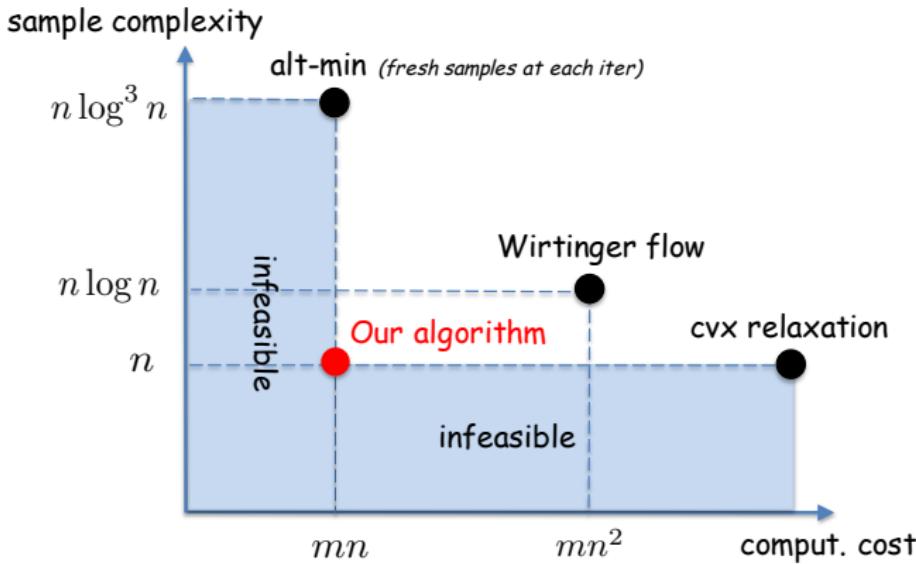
This work: random quadratic systems are solvable in linear time!

A glimpse of our results

n : # unknowns;

m : sample size (# eqns);

$$\mathbf{y} = |\mathbf{A}\mathbf{x}|^2, \mathbf{A} \in \mathbb{R}^{m \times n}$$



This work: random quadratic systems are solvable in linear time!

- ✓ minimal sample size
- ✓ optimal statistical accuracy

A first impulse: maximum likelihood estimate

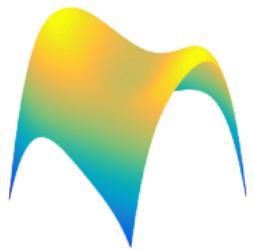
$$\text{maximize}_{\boldsymbol{z}} \quad \ell(\boldsymbol{z}) = \frac{1}{m} \sum_{k=1}^m \ell_k(\boldsymbol{z})$$

A first impulse: maximum likelihood estimate

$$\text{maximize}_{\mathbf{z}} \quad \ell(\mathbf{z}) = \frac{1}{m} \sum_{k=1}^m \ell_k(\mathbf{z})$$

- Gaussian data: $y_k \sim |\mathbf{a}_k^* \mathbf{x}|^2 + \mathcal{N}(0, \sigma^2)$

$$\ell_k(\mathbf{z}) = -(y_k - |\mathbf{a}_k^* \mathbf{z}|^2)^2$$

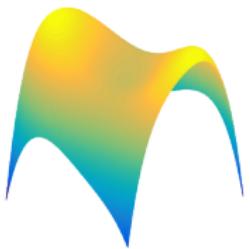


A first impulse: maximum likelihood estimate

$$\text{maximize}_{\mathbf{z}} \quad \ell(\mathbf{z}) = \frac{1}{m} \sum_{k=1}^m \ell_k(\mathbf{z})$$

- Gaussian data: $y_k \sim |\mathbf{a}_k^* \mathbf{x}|^2 + \mathcal{N}(0, \sigma^2)$

$$\ell_k(\mathbf{z}) = -(y_k - |\mathbf{a}_k^* \mathbf{z}|^2)^2$$



- Poisson data: $y_k \sim \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2)$

$$\ell_k(\mathbf{z}) = -|\mathbf{a}_k^* \mathbf{z}|^2 + y_k \log |\mathbf{a}_k^* \mathbf{z}|^2$$

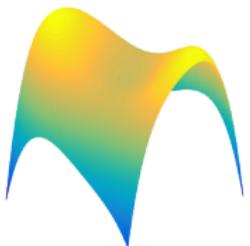


A first impulse: maximum likelihood estimate

$$\text{maximize}_{\mathbf{z}} \quad \ell(\mathbf{z}) = \frac{1}{m} \sum_{k=1}^m \ell_k(\mathbf{z})$$

- Gaussian data: $y_k \sim |\mathbf{a}_k^* \mathbf{x}|^2 + \mathcal{N}(0, \sigma^2)$

$$\ell_k(\mathbf{z}) = -(y_k - |\mathbf{a}_k^* \mathbf{z}|^2)^2$$



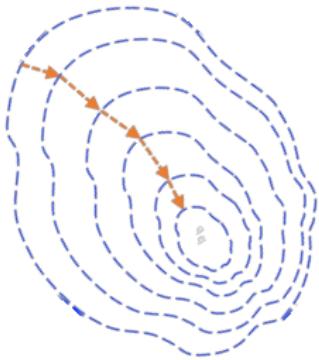
- Poisson data: $y_k \sim \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2)$

$$\ell_k(\mathbf{z}) = -|\mathbf{a}_k^* \mathbf{z}|^2 + y_k \log |\mathbf{a}_k^* \mathbf{z}|^2$$



Problem: $-\ell$ nonconvex, many local stationary points

Wirtinger flow: Candès, Li, Soltanolkotabi '14



- **Spectral initialization:** $z^0 \leftarrow$ leading eigenvector of

$$\frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^*$$

- **Iterative refinement:** for $t = 0, 1, \dots$

$$z^{t+1} = z^t + \mu_t \nabla \ell(z^t)$$

Already rich theory (see also Soltanolkotabi '14, Ma, Wang, Chi, Chen '17)

Interpretation of spectral initialization

Spectral initialization: $\mathbf{z}^0 \leftarrow$ leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^*$$

Interpretation of spectral initialization

Spectral initialization: $\mathbf{z}^0 \leftarrow$ leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^*$$

- Rationale: $\mathbb{E}[\mathbf{Y}] = \mathbf{I} + 2\mathbf{x}\mathbf{x}^*$ ($\|\mathbf{x}\|_2 = 1$) under Gaussian design

Interpretation of spectral initialization

Spectral initialization: $\mathbf{z}^0 \leftarrow$ leading eigenvector of

$$\mathbf{Y} := \frac{1}{m} \sum_{k=1}^m y_k \mathbf{a}_k \mathbf{a}_k^*$$

- Rationale: $\mathbb{E}[\mathbf{Y}] = \mathbf{I} + 2\mathbf{x}\mathbf{x}^*$ ($\|\mathbf{x}\|_2 = 1$) under Gaussian design
- Would succeed if $\mathbf{Y} \rightarrow \mathbb{E}[\mathbf{Y}]$

Empirical performance of initialization ($m = 12n$)

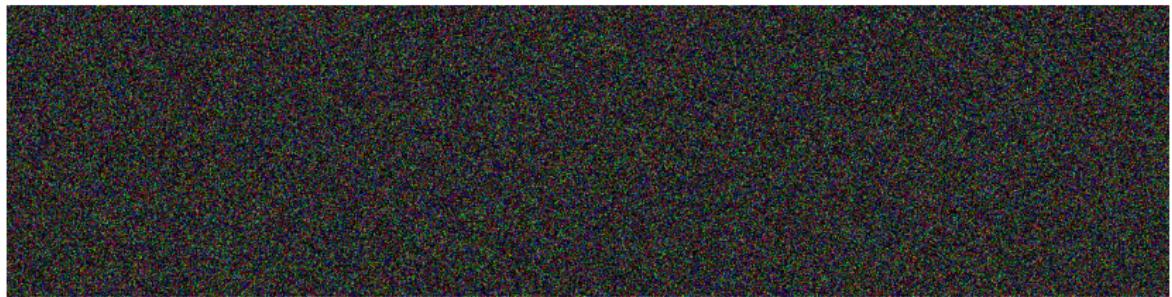


Ground truth $x \in \mathbb{R}^{409600}$

Empirical performance of initialization ($m = 12n$)



Ground truth $x \in \mathbb{R}^{409600}$



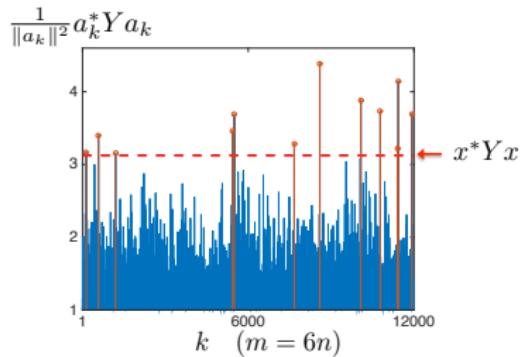
Spectral initialization

Improving initialization

$$Y = \frac{1}{m} \sum_k \underbrace{y_k a_k a_k^*}_{\text{heavy-tailed}} \quad \not\rightarrow \quad \mathbb{E}[Y] \quad \text{unless } m \gg n$$

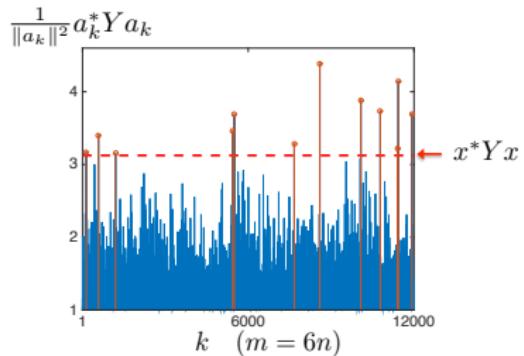
Improving initialization

$$Y = \frac{1}{m} \sum_k \underbrace{y_k a_k a_k^*}_{\text{heavy-tailed}} \quad \not\rightarrow \quad \mathbb{E}[Y] \quad \text{unless } m \gg n$$



Improving initialization

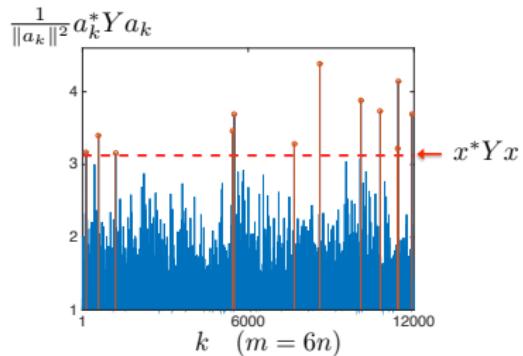
$$\mathbf{Y} = \frac{1}{m} \sum_k \underbrace{y_k \mathbf{a}_k \mathbf{a}_k^*}_{\text{heavy-tailed}} \quad \not\rightarrow \quad \mathbb{E}[\mathbf{Y}] \quad \text{unless } m \gg n$$



Problem large outliers $y_k = |\mathbf{a}_k^* \mathbf{x}|^2$ bear too much influence

Improving initialization

$$\mathbf{Y} = \frac{1}{m} \sum_k \underbrace{y_k \mathbf{a}_k \mathbf{a}_k^*}_{\text{heavy-tailed}} \quad \not\rightarrow \quad \mathbb{E}[\mathbf{Y}] \quad \text{unless } m \gg n$$

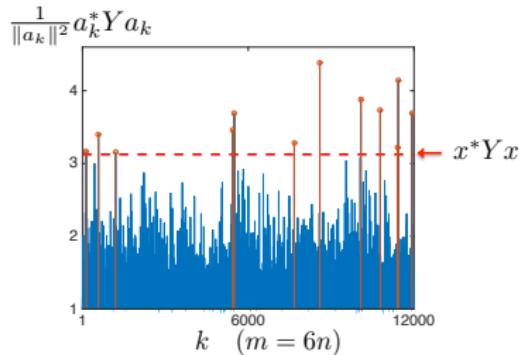


Problem large outliers $y_k = |\mathbf{a}_k^* \mathbf{x}|^2$ bear too much influence

Solution discard large samples and run PCA for $\frac{1}{m} \sum_k y_k \mathbf{a}_k \mathbf{a}_k^* \mathbf{1}_{\{|y_k| \lesssim \text{Avg}\{|y_l|\}\}}$

Improving initialization

$$\mathbf{Y} = \frac{1}{m} \sum_k \underbrace{y_k \mathbf{a}_k \mathbf{a}_k^*}_{\text{heavy-tailed}} \quad \not\rightarrow \quad \mathbb{E}[\mathbf{Y}] \quad \text{unless } m \gg n$$



Problem large outliers $y_k = |\mathbf{a}_k^* \mathbf{x}|^2$ bear too much influence

Solution discard large samples and run PCA for $\frac{1}{m} \sum_k y_k \mathbf{a}_k \mathbf{a}_k^* \mathbf{1}_{\{|y_k| \lesssim \text{Avg}\{|y_l|\}\}}$

— *improvable via more refined pre-processing*

(Wang, Giannakis, Eldar '16, Lu, Li '17, Mondelli, Montanari '17)

$$\frac{1}{m} \sum_k \rho(y_k) \mathbf{a}_k \mathbf{a}_k^* \quad \text{e.g. } \rho(y_k) = \max\{y_k, a\}$$

Empirical performance of initialization ($m = 12n$)



Ground truth $x \in \mathbb{R}^{409600}$



Regularized spectral initialization

Iterative refinement stage: search directions

Wirtinger flow: $\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{k=1}^m \underbrace{(y_k - |\mathbf{a}_k^\top \mathbf{z}^t|^2) \mathbf{a}_k \mathbf{a}_k^\top \mathbf{z}^t}_{= -\nabla \ell_k(\mathbf{z}^t)}$

Iterative refinement stage: search directions

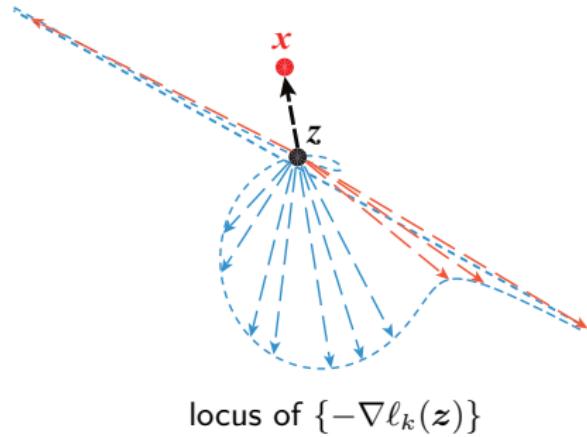
$$\text{Wirtinger flow: } \mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{k=1}^m \underbrace{(y_k - |\mathbf{a}_k^\top \mathbf{z}^t|^2) \mathbf{a}_k \mathbf{a}_k^\top \mathbf{z}^t}_{= -\nabla \ell_k(\mathbf{z}^t)}$$

Even in a local region around \mathbf{x} (e.g. $\{\mathbf{z} \mid \|\mathbf{z} - \mathbf{x}\|_2 \leq 0.1\|\mathbf{x}\|_2\}$):

- $f(\cdot)$ is NOT strongly convex unless $m \gg n$
- $f(\cdot)$ has huge smoothness parameter

Iterative refinement stage: search directions

Wirtinger flow: $\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{k=1}^m \underbrace{(y_k - |\mathbf{a}_k^\top \mathbf{z}^t|^2) \mathbf{a}_k \mathbf{a}_k^\top \mathbf{z}^t}_{= -\nabla \ell_k(\mathbf{z}^t)}$



Problem: descent direction has large variability

Our solution: variance reduction via proper trimming

More adaptive rule:

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{i=1}^m \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}^t|^2}{\mathbf{a}_i^\top \mathbf{z}^t} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}^t) \cap \mathcal{E}_2^i(\mathbf{z}^t)}$$

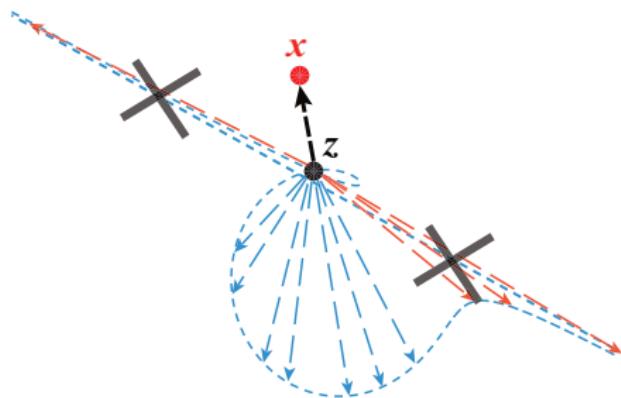
where $\mathcal{E}_1^i(\mathbf{z}) = \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \leq \alpha_z^{\text{ub}} \right\}$; $\mathcal{E}_2^i(\mathbf{z}) = \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\frac{\alpha_h}{m} \left\| \mathbf{y} - \mathcal{A}(\mathbf{z} \mathbf{z}^\top) \right\|_1 |\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \right\}$

Our solution: variance reduction via proper trimming

More adaptive rule:

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{i=1}^m \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}^t|^2}{\mathbf{a}_i^\top \mathbf{z}^t} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}^t) \cap \mathcal{E}_2^i(\mathbf{z}^t)}$$

where $\mathcal{E}_1^i(\mathbf{z}) = \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \leq \alpha_z^{\text{ub}} \right\}$; $\mathcal{E}_2^i(\mathbf{z}) = \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\frac{\alpha_h}{m} \left\| \mathbf{y} - \mathcal{A}(\mathbf{z} \mathbf{z}^\top) \right\|_1 |\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \right\}$

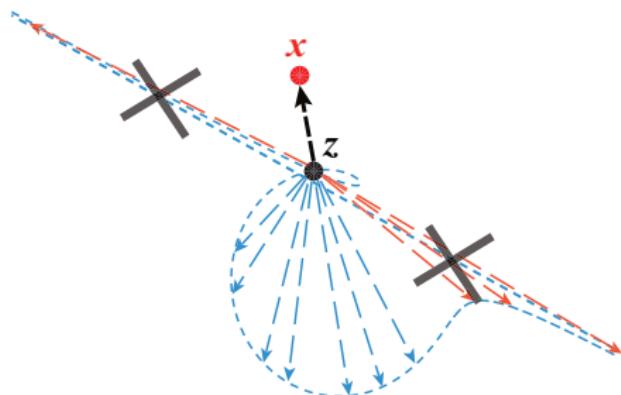


Our solution: variance reduction via proper trimming

More adaptive rule:

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{i=1}^m \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}^t|^2}{\mathbf{a}_i^\top \mathbf{z}^t} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}^t) \cap \mathcal{E}_2^i(\mathbf{z}^t)}$$

where $\mathcal{E}_1^i(\mathbf{z}) = \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \leq \alpha_z^{\text{ub}} \right\}$; $\mathcal{E}_2^i(\mathbf{z}) = \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\frac{\alpha_h}{m} \left\| \mathbf{y} - \mathcal{A}(\mathbf{z} \mathbf{z}^\top) \right\|_1 |\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \right\}$



informally, $\mathbf{z}^{t+1} = \mathbf{z}^t + \frac{\mu}{m} \sum_{k \in \mathcal{T}_t} \nabla \ell_k(\mathbf{z}^t)$

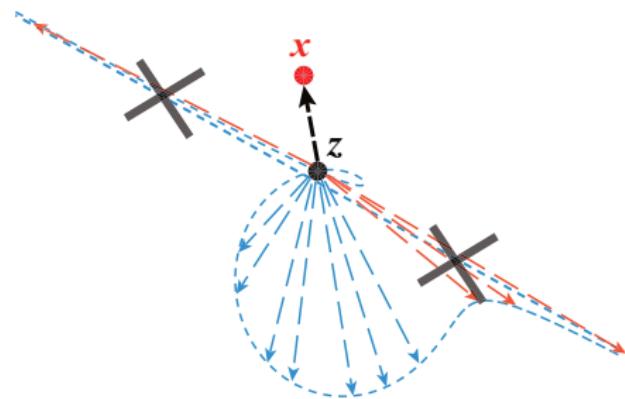
- \mathcal{T}_t trims away excessively large grad components

Our solution: variance reduction via proper trimming

More adaptive rule:

$$\mathbf{z}^{t+1} = \mathbf{z}^t - \frac{\mu_t}{m} \sum_{i=1}^m \frac{y_i - |\mathbf{a}_i^\top \mathbf{z}^t|^2}{\mathbf{a}_i^\top \mathbf{z}^t} \mathbf{a}_i \mathbf{1}_{\mathcal{E}_1^i(\mathbf{z}^t) \cap \mathcal{E}_2^i(\mathbf{z}^t)}$$

where $\mathcal{E}_1^i(\mathbf{z}) = \left\{ \alpha_z^{\text{lb}} \leq \frac{|\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \leq \alpha_z^{\text{ub}} \right\}$; $\mathcal{E}_2^i(\mathbf{z}) = \left\{ |y_i - |\mathbf{a}_i^\top \mathbf{z}|^2| \leq \frac{\frac{\alpha_h}{m} \left\| \mathbf{y} - \mathcal{A}(\mathbf{z} \mathbf{z}^\top) \right\|_1 |\mathbf{a}_i^\top \mathbf{z}|}{\|\mathbf{z}\|_2} \right\}$

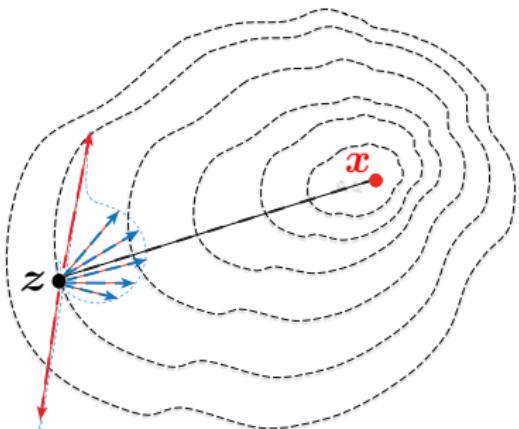


informally, $\mathbf{z}^{t+1} = \mathbf{z}^t + \frac{\mu}{m} \sum_{k \in \mathcal{T}_t} \nabla \ell_k(\mathbf{z}^t)$

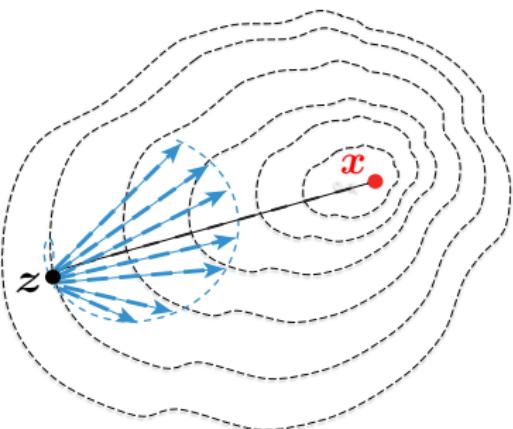
- \mathcal{T}_t trims away excessively large grad components

Slight bias + much reduced variance

Larger step size μ_t is feasible



without trimming: $\mu_t = O(1/n)$



with trimming: $\mu_t = O(1)$

With better-controlled descent directions, one proceeds far more aggressively

Summary: truncated Wirtinger flows (TWF)

1. **Regularized spectral initialization:** $\mathbf{z}^0 \leftarrow$ leading eigenvector of

$$\frac{1}{m} \sum_{k \in \mathcal{T}_0} y_k \mathbf{a}_k \mathbf{a}_k^*$$

2. **Regularized gradient descent**

$$\mathbf{z}^{t+1} = \mathbf{z}^t + \mu_t \underbrace{\frac{1}{m} \sum_{k \in \mathcal{T}_t} \nabla \ell_k(\mathbf{z}^t)}_{:= \nabla \ell^{\text{tr}}(\mathbf{z}^t)}$$

Key idea: adaptively discard high-leverage data

Performance guarantees of TWF (noiseless data)

$$\text{dist}(\mathbf{z}, \mathbf{x}) := \min\{\|\mathbf{z} \pm \mathbf{x}\|_2\}$$

Theorem (Chen & Candès '15). Under i.i.d. Gaussian design, TWF achieves

$$\text{dist}(\mathbf{z}^t, \mathbf{x}) \lesssim (1 - \rho)^t \|\mathbf{x}\|_2, \quad t = 0, 1, \dots$$

with high prob., provided that sample size $m \gtrsim n$. Here, $0 < \rho < 1$ is const.

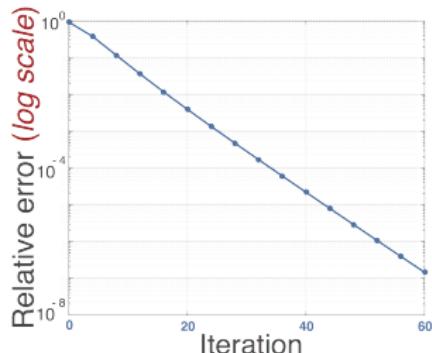
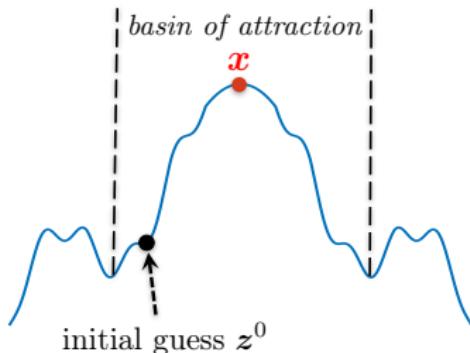
Performance guarantees of TWF (noiseless data)

$$\text{dist}(z, x) := \min\{\|z \pm x\|_2\}$$

Theorem (Chen & Candès '15). Under i.i.d. Gaussian design, TWF achieves

$$\text{dist}(z^t, x) \lesssim (1 - \rho)^t \|x\|_2, \quad t = 0, 1, \dots$$

with high prob., provided that sample size $m \gtrsim n$. Here, $0 < \rho < 1$ is const.



Computational complexity

$$\mathbf{A} := \{\mathbf{a}_k^*\}_{1 \leq k \leq m}$$

- **Initialization:** leading eigenvector \rightarrow a few applications of \mathbf{A} and \mathbf{A}^*

$$\sum_{k \in \mathcal{T}_0} y_k \mathbf{a}_k \mathbf{a}_k^* = \mathbf{A}^* \operatorname{diag}\{y_k \cdot \mathbf{1}_{k \in \mathcal{T}_0}\} \mathbf{A}$$

Computational complexity

$$\mathbf{A} := \{\mathbf{a}_k^*\}_{1 \leq k \leq m}$$

- **Initialization:** leading eigenvector \rightarrow a few applications of \mathbf{A} and \mathbf{A}^*

$$\sum_{k \in \mathcal{T}_0} y_k \mathbf{a}_k \mathbf{a}_k^* = \mathbf{A}^* \operatorname{diag}\{y_k \cdot 1_{k \in \mathcal{T}_0}\} \mathbf{A}$$

- **Iterations:** one application of \mathbf{A} and \mathbf{A}^* per iteration

$$\mathbf{z}^{t+1} = \mathbf{z}^t + \frac{\mu_t}{m} \nabla \ell_{\text{tr}}(\mathbf{z}^t) \quad -\nabla \ell_{\text{tr}}(\mathbf{z}^t) = \mathbf{A}^* \boldsymbol{\nu}$$
$$\boldsymbol{\nu} = 2 \frac{|\mathbf{A} \mathbf{z}^t|^2 - \mathbf{y}}{\mathbf{A} \mathbf{z}^t} \cdot 1_{\mathcal{T}}$$

Computational complexity

$$\mathbf{A} := \{\mathbf{a}_k^*\}_{1 \leq k \leq m}$$

- **Initialization:** leading eigenvector \rightarrow a few applications of \mathbf{A} and \mathbf{A}^*

$$\sum_{k \in \mathcal{T}_0} y_k \mathbf{a}_k \mathbf{a}_k^* = \mathbf{A}^* \operatorname{diag}\{y_k \cdot 1_{k \in \mathcal{T}_0}\} \mathbf{A}$$

- **Iterations:** one application of \mathbf{A} and \mathbf{A}^* per iteration

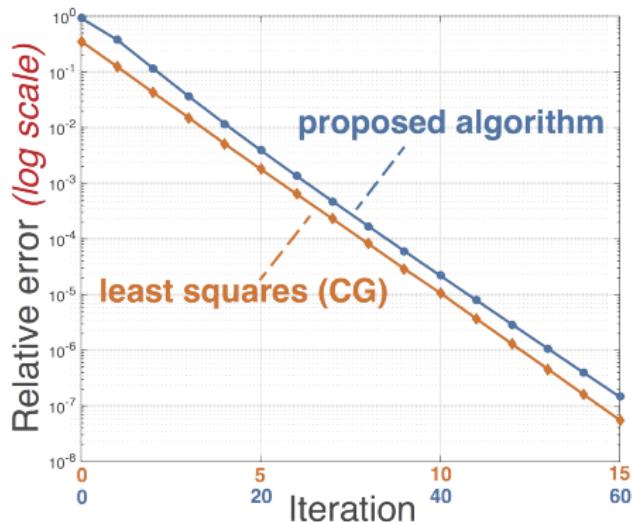
$$\begin{aligned} \mathbf{z}^{t+1} &= \mathbf{z}^t + \frac{\mu_t}{m} \nabla \ell_{\text{tr}}(\mathbf{z}^t) & -\nabla \ell_{\text{tr}}(\mathbf{z}^t) &= \mathbf{A}^* \boldsymbol{\nu} \\ & & \boldsymbol{\nu} &= 2 \frac{|\mathbf{A} \mathbf{z}^t|^2 - \mathbf{y}}{\mathbf{A} \mathbf{z}^t} \cdot 1_{\mathcal{T}} \end{aligned}$$

Approximate runtime: several tens of applications of \mathbf{A} and \mathbf{A}^*

Numerical surprise

- CG: solve $y = Ax$

- Our algorithm: solve $y = |Ax|^2$



For random quadratic systems ($m = 8n$)

comput. cost of our algo. \approx 4 \times comput. cost of least squares

Empirical performance



After regularized spectral initialization

Empirical performance



After regularized spectral initialization



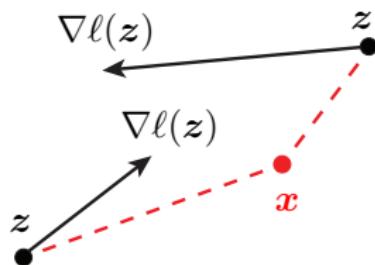
After 50 TWF iterations

Key convergence condition for gradient stage

If there are many samples:

$\forall z$ s.t. $\text{dist}(z, x) \leq \varepsilon \|x\|_2$:

$$\langle \nabla \ell(z), x - z \rangle \gtrsim \|z - x\|_2^2 + \|\nabla \ell(z)\|_2^2$$

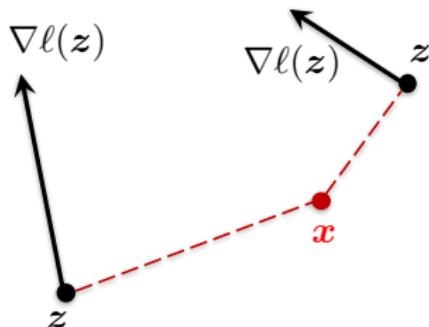


Key convergence condition for gradient stage

If there are **NOT** many samples, i.e. $m \asymp n$:

$\forall z$ s.t. $\text{dist}(z, x) \leq \varepsilon \|x\|_2$:

$$\langle \nabla \ell(z), x - z \rangle \gtrsim \|z - x\|_2^2 + \|\nabla \ell(z)\|_2^2$$

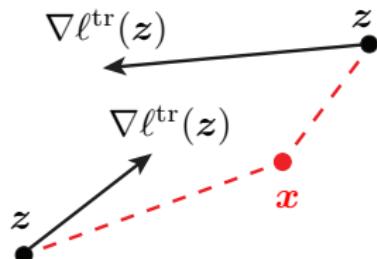
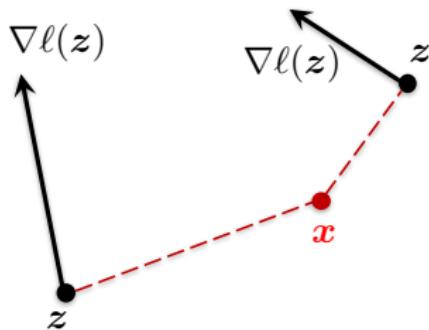


Key convergence condition for gradient stage

If there are **NOT** many samples, i.e. $m \asymp n$:

$\forall z$ s.t. $\text{dist}(z, x) \leq \varepsilon \|x\|_2$:

$$\langle \nabla \ell^{\text{tr}}(z), x - z \rangle \gtrsim \|z - x\|_2^2 + \|\nabla \ell^{\text{tr}}(z)\|_2^2$$



Stability under noisy data

- Noisy data: $y_k = |\mathbf{a}_k^* \mathbf{x}|^2 + \eta_k$
- Signal-to-noise ratio:

$$\text{SNR} := \frac{\sum_k |\mathbf{a}_k^* \mathbf{x}|^4}{\sum_k \eta_k^2} \approx \frac{3m \|\mathbf{x}\|^4}{\|\boldsymbol{\eta}\|^2}$$

- i.i.d. Gaussian design

Stability under noisy data

- Noisy data: $y_k = |\mathbf{a}_k^* \mathbf{x}|^2 + \eta_k$
- Signal-to-noise ratio:

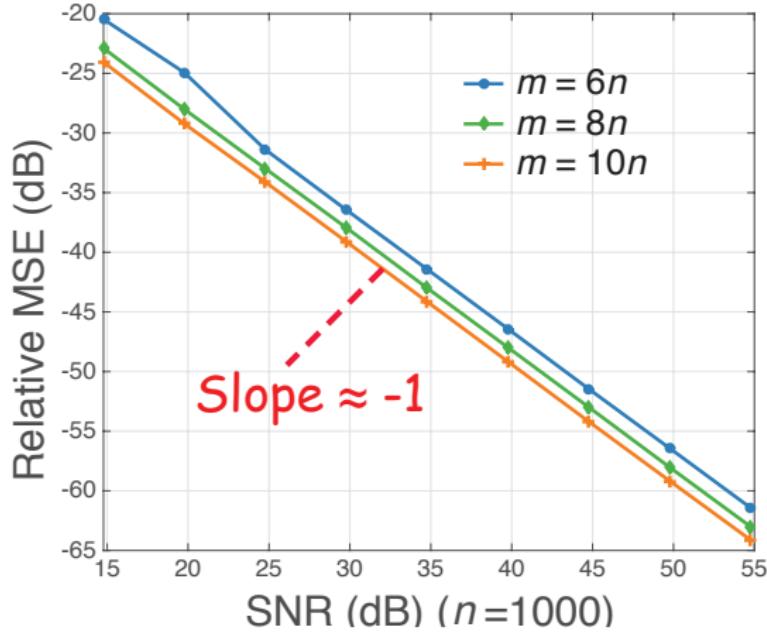
$$\text{SNR} := \frac{\sum_k |\mathbf{a}_k^* \mathbf{x}|^4}{\sum_k \eta_k^2} \approx \frac{3m \|\mathbf{x}\|^4}{\|\boldsymbol{\eta}\|^2}$$

- i.i.d. Gaussian design

Theorem (Soltanolkotabi) WF converges to MLE

Theorem (Chen, Candès) Relative error of TWF converges to $O(\frac{1}{\sqrt{\text{SNR}}})$

Relative MSE vs. SNR (Poisson data)



Empirical evidence: relative MSE scales inversely with SNR

This accuracy is nearly un-improvable (empirically)

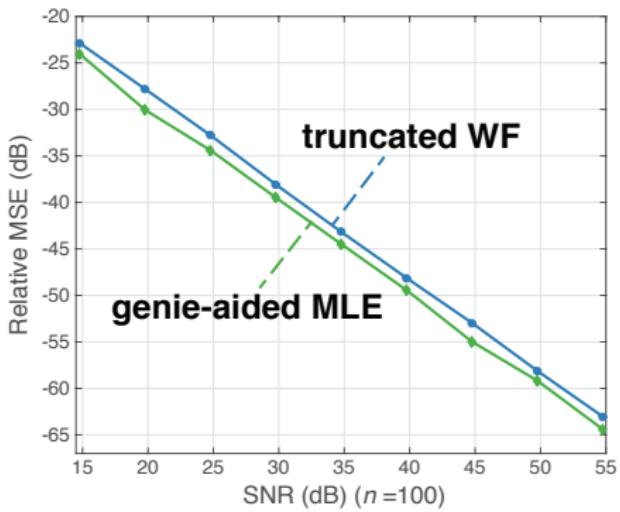
Comparison with genie-aided MLE (with sign info. revealed)

$$y_k \sim \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2) \quad \text{and} \quad \varepsilon_k = \text{sign}(\mathbf{a}_k^* \mathbf{x}) \quad (\text{revealed by a genie})$$

This accuracy is nearly un-improvable (empirically)

Comparison with genie-aided MLE (with sign info. revealed)

$$y_k \sim \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2) \quad \text{and} \quad \varepsilon_k = \text{sign}(\mathbf{a}_k^* \mathbf{x}) \quad (\text{revealed by a genie})$$



little empirical loss due to missing signs

This accuracy is nearly un-improvable (theoretically)

- Poisson data: $y_k \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2)$
- Signal-to-noise ratio:

$$\text{SNR} \approx \frac{\sum_k |\mathbf{a}_k^* \mathbf{x}|^4}{\sum_k \text{Var}(y_k)} \approx 3\|\mathbf{x}\|^2$$

This accuracy is nearly un-improvable (theoretically)

- Poisson data: $y_k \stackrel{\text{ind.}}{\sim} \text{Poisson}(|\mathbf{a}_k^* \mathbf{x}|^2)$
- Signal-to-noise ratio:

$$\text{SNR} \approx \frac{\sum_k |\mathbf{a}_k^* \mathbf{x}|^4}{\sum_k \text{Var}(y_k)} \approx 3\|\mathbf{x}\|^2$$

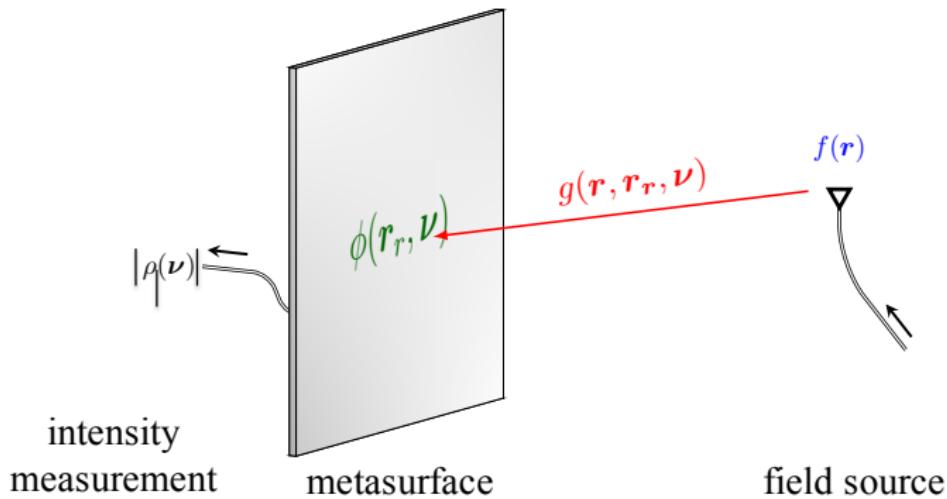
Theorem (Chen, Candès). Under i.i.d. Gaussian design, for any estimator $\hat{\mathbf{x}}$,

$$\inf_{\hat{\mathbf{x}}} \sup_{\mathbf{x}: \|\mathbf{x}\| \geq \log^{1.5} m} \frac{\mathbb{E} [\text{dist}(\hat{\mathbf{x}}, \mathbf{x}) \mid \{\mathbf{a}_k\}]}{\|\mathbf{x}\|} \gtrsim \frac{1}{\sqrt{\text{SNR}}},$$

provided that sample size $m \asymp n$.

Phaseless 3D computational imaging

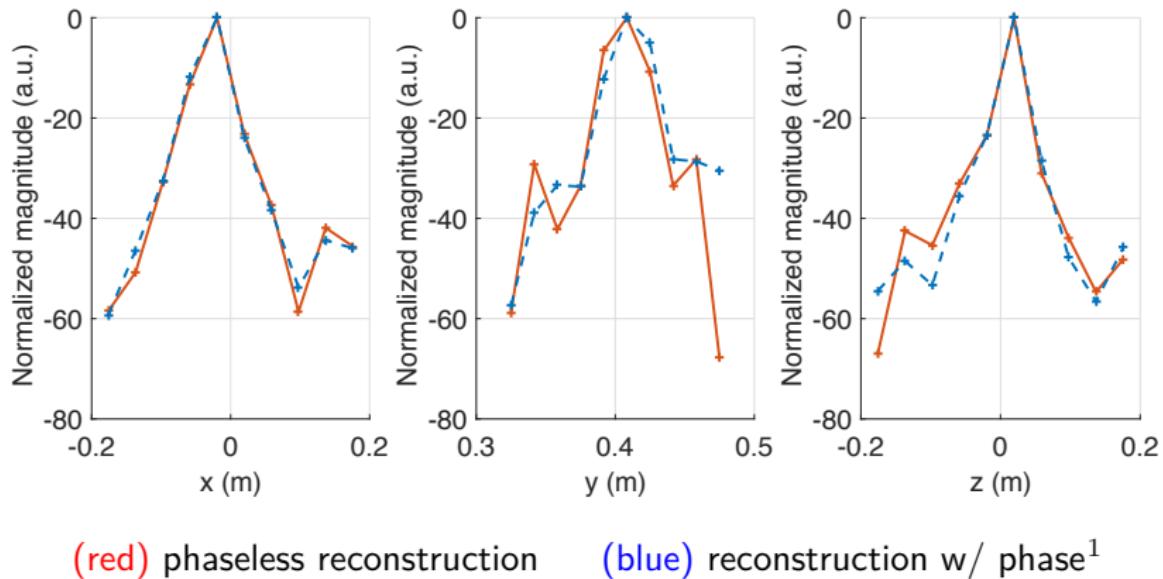
Fromenteze, Liu, Boyarsky, Gollub, & Smith '16



Measure intensities (with radiating metasurfaces) rather than complex signals
for sub-centimeter wavelengths

Phaseless 3D computational imaging

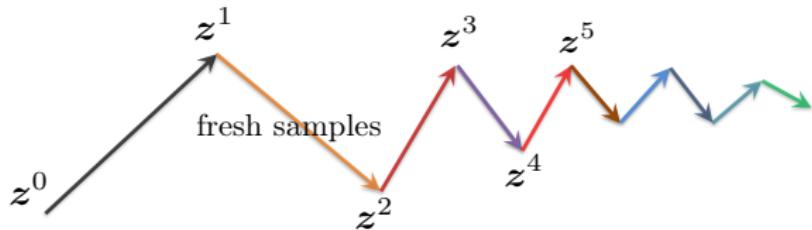
Fromenteze, Liu, Boyarsky, Gollub, & Smith '16



¹This demonstration is proposed in microwave range as proof of concept

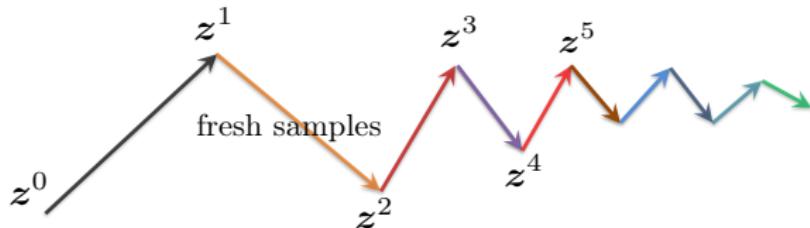
No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but much easier to analyze

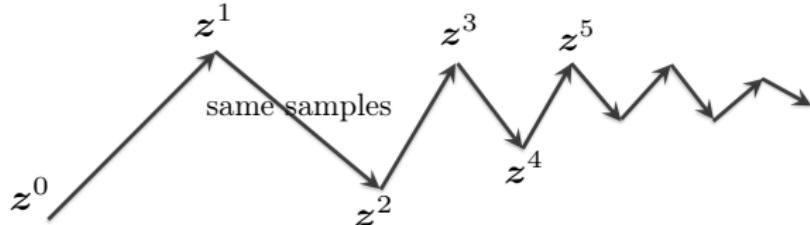


No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but much easier to analyze



- **Our works:** reuse all samples in all iterations



A small sample of more recent works

- other optimal algorithms
 - reshaped WF (Zhang et al.), truncated AF (Wang et al.), median-TWF (Zhang et al.)
 - alt-min w/o resampling (Waldspurger)
 - composite optimization (Duchi et al., Charisopoulos et al.)
 - approximate message passing (Ma et al.)
 - block coordinate descent (Barmherzig et al.)
 - PhaseMax (Goldstein et al., Bahmani et al., Salehi et al., Dhifallah et al., Hand et al.)
- stochastic algorithms (Kolte et al., Zhang et al., Lu et al., Tan et al., Jeong et al.)
- improved WF theory: iteration complexity $\rightarrow O(\log n \log \frac{1}{\varepsilon})$ (Ma et al.)
- improved initialization (Lu et al., Wang et al., Mondelli et al.)
- random initialization (Chen et al.)
- structured quadratic systems (Cai et al., Soltanolkotabi, Wang et al., Yang et al., Qu et al.)
- geometric analysis (Sun et al., Davis et al.)
- low-rank generalization (White et al., Li et al., Vaswani et al.)

Central message

- Simple nonconvex paradigms are surprisingly effective for computing MLE
- Importance of statistical thinking (initialization)

| | statistical accuracy | comput. cost |
|---------------------|----------------------|--------------|
| convex relaxation | | |
| nonconvex procedure | | |

- Y. Chen, E. Candès, "Solving random quadratic systems of equations is nearly as easy as solving linear systems," *Comm. Pure and Applied Math.*, 2017