

Solving Corrupted Quadratic Equations, Provably

Yuejie Chi

Electrical and Computer Engineering



THE OHIO STATE UNIVERSITY

University of Michigan

November 2016

Data science

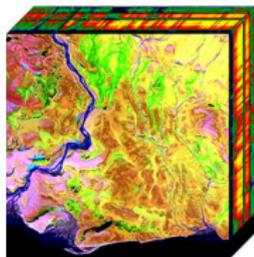
New imaging/sensing modalities allow us to probe the nature in unprecedented manners:



healthcare



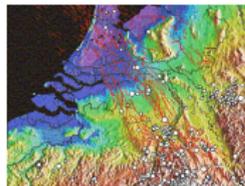
Radio astronomy



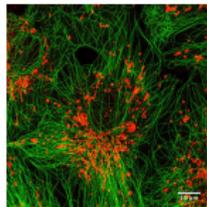
hyperspectral



Internet traffic



seismic imaging

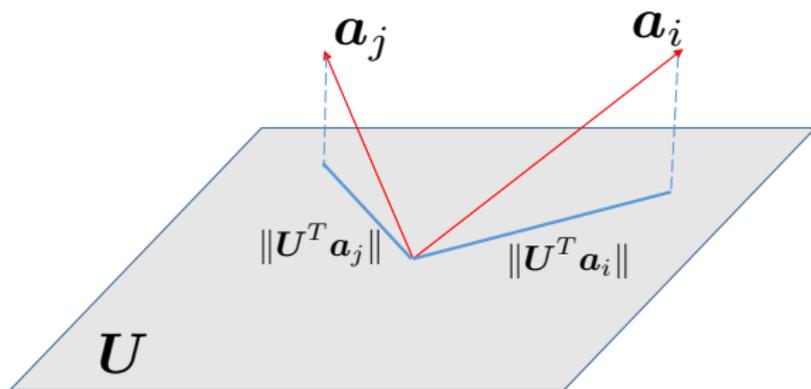


microscopy

but also with a lot of new (and exciting) challenges due to the unconventional manner these data are obtained.

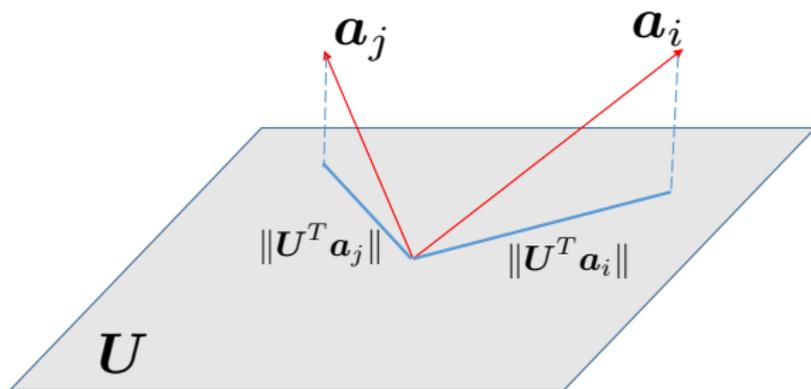
Subspace retrieval using intensity measurements only

- We wish to estimate a subspace $U \in \mathbb{R}^{n \times r}$ by interrogating it with vectors $\{\mathbf{a}_i\}_{i=1}^m$ and forming backprojections;



Subspace retrieval using intensity measurements only

- We wish to estimate a subspace $U \in \mathbb{R}^{n \times r}$ by interrogating it with vectors $\{\mathbf{a}_i\}_{i=1}^m$ and forming backprojections;



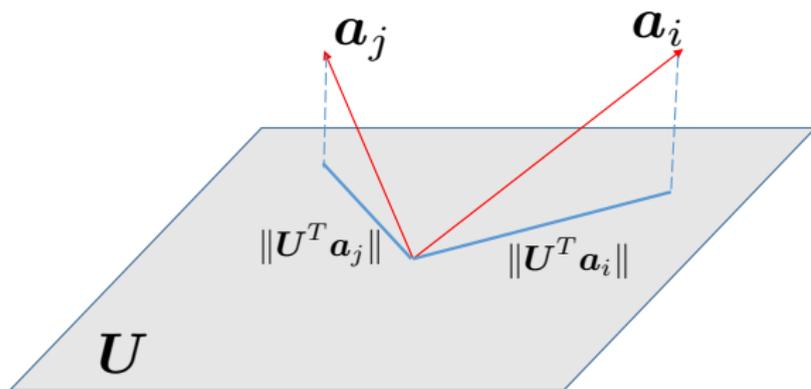
- We only observe the **intensity** of the backprojections, namely,

$$y_i = \|U^T \mathbf{a}_i\|_2^2 = \mathbf{a}_i^T (U U^T) \mathbf{a}_i, \quad i = 1, \dots, m.$$

They are **quadratic** with respect to U .

Subspace retrieval using intensity measurements only

- We wish to estimate a subspace $U \in \mathbb{R}^{n \times r}$ by interrogating it with vectors $\{\mathbf{a}_i\}_{i=1}^m$ and forming backprojections;



- We only observe the **intensity** of the backprojections, namely,

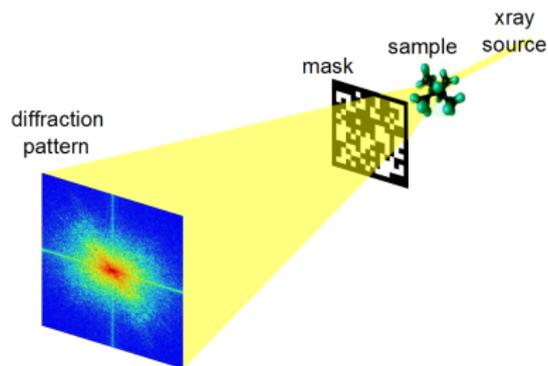
$$y_i = \|\mathbf{U}^T \mathbf{a}_i\|_2^2 = \mathbf{a}_i^T (\mathbf{U}\mathbf{U}^T) \mathbf{a}_i, \quad i = 1, \dots, m.$$

They are **quadratic** with respect to U .

- Intensity measurements are much easier to implement by an energy detector for high-frequency and wide-band (THz) applications.

Phase retrieval

How to recover structure of a sample from its diffraction pattern?



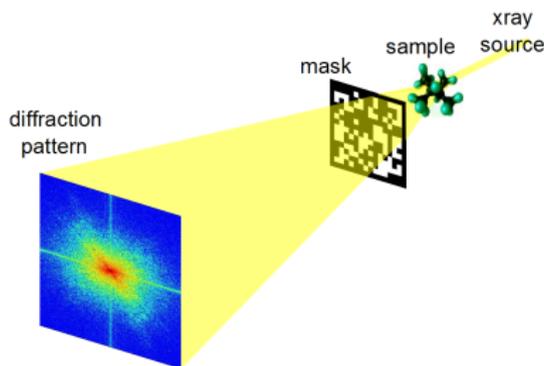
- In the important special case of $r = 1$, it becomes equivalent to **phase retrieval**^{*}, namely, recover $\mathbf{x} \in \mathbb{R}^n / \mathbb{C}^n$ from

$$y_i = |\mathcal{F}\{\mathbf{x}\}|^2, \quad \text{where } \mathcal{F} \text{ is Fourier transform,}$$

^{*}Image credit: E. J. Candès, Y. C. Eldar, T. Strohmer and V. Voroninski, "Phase retrieval via matrix completion," SIAM J. on Imaging Sciences.

Phase retrieval

How to recover structure of a sample from its diffraction pattern?



- In the important special case of $r = 1$, it becomes equivalent to **phase retrieval**^{*}, namely, recover $\mathbf{x} \in \mathbb{R}^n / \mathbb{C}^n$ from

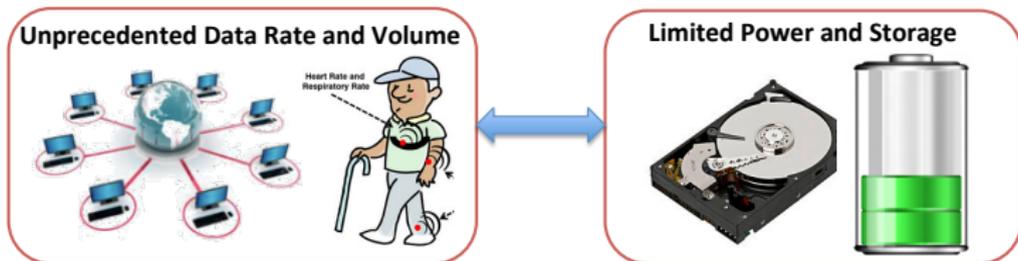
$$y_i = |\mathcal{F}\{\mathbf{x}\}|^2, \quad \text{where } \mathcal{F} \text{ is Fourier transform,}$$

This has wide applications in X-ray crystallography, electron microscopy and coherent diffractive imaging, and leads to winning of Nobel prize (e.g. discovery of double helix structure).

^{*}Image credit: E. J. Candès, Y. C. Eldar, T. Strohmer and V. Voroninski, "Phase retrieval via matrix completion," SIAM J. on Imaging Sciences.

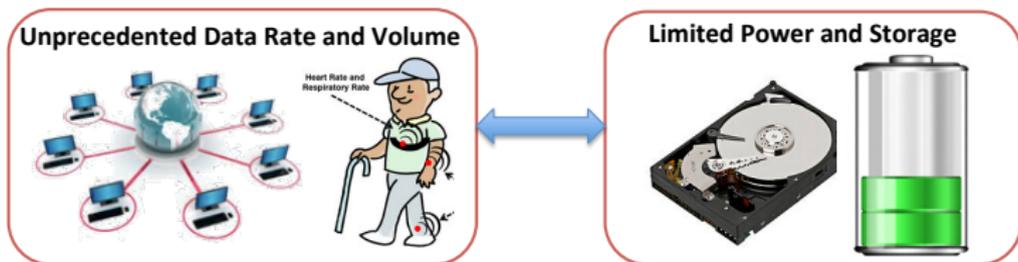
Covariance sketching for streaming data

Multivariate streaming data: a new data snapshot $\mathbf{x}_t \in \mathbb{C}^n / \mathbb{R}^n$ is generated by the sensor platform at each time t ;



Covariance sketching for streaming data

Multivariate streaming data: a new data snapshot $\mathbf{x}_t \in \mathbb{C}^n / \mathbb{R}^n$ is generated by the sensor platform at each time t ;



- **high-dimensional:** the number of variables, n , is large;
- **real-time:** data processed “on the fly”;
- **decentralized:** data collected at decentralized locations;
- **resource-constrained:** cannot store and transmit all data;

Covariance sketching

Observation: Fortunately, inference requires only statistics of the data stream, not the stream itself; we can “sketch” /compress the data at the hope of directly recovering its statistics!



Covariance sketching

Observation: Fortunately, inference requires only statistics of the data stream, not the stream itself; we can “sketch” /compress the data at the hope of directly recovering its statistics!

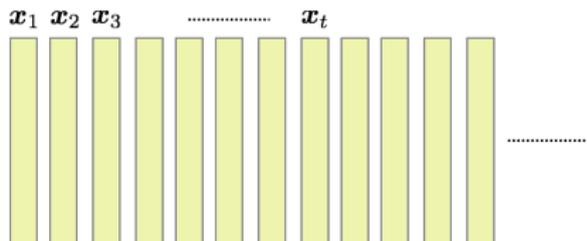


Approach: distributed data sketching and aggregation to recover the covariance structure or principal components.

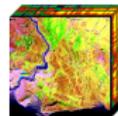
- access each data sample via quadratic (energy) sketches;
- aggregate the sketches into linear observations of the covariance matrix.

Quadratic sampling

How to sketch a high-dimensional data stream in order to recover its covariance matrix?



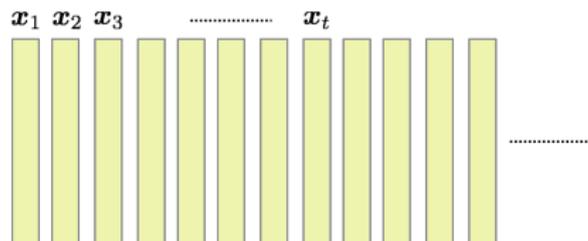
network traffic



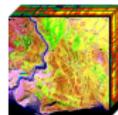
hyperspectral imagery

Quadratic sampling

How to sketch a high-dimensional data stream in order to recover its covariance matrix?



network traffic



hyperspectral imagery

- To meet resource constraints, we would like to sample **in a single pass on the fly**: a single *quadratic* sketch of x_t :

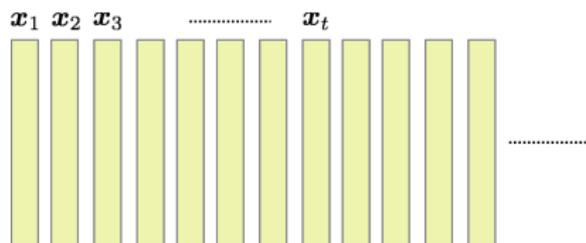
$$z_t = |\langle \mathbf{a}_t, \mathbf{x}_t \rangle|^2,$$

which reduces the dim. of each x_t to merely a scalar.

- sketching complexity is linear in length of x_t ;

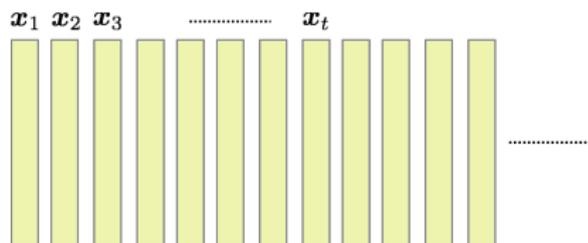
Quadratic sampling for covariance sketching

- Consider a data stream possibly distributively observed at m sensors, each with a sketching vector $\mathbf{a}_i \in \mathbb{R}^n$, $i = 1, \dots, m$:



Quadratic sampling for covariance sketching

- Consider a data stream possibly distributively observed at m sensors, each with a sketching vector $\mathbf{a}_i \in \mathbb{R}^n$, $i = 1, \dots, m$:



- Sketch a substream indexed by $\{\ell_t^i\}_{t=1}^T$ with $|\langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle|^2$ and compute the average:

$$y_{i,T} = \frac{1}{T} \sum_{t=1}^T \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2 = \mathbf{a}_i^T \left(\frac{1}{T} \sum_{t=1}^T \mathbf{x}_{\ell_t^i} \mathbf{x}_{\ell_t^i}^T \right) \mathbf{a}_i$$
$$\xrightarrow{T \rightarrow \infty} \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i,$$

where $\mathbf{X} = \mathbb{E}[\mathbf{x}\mathbf{x}^T]$ is the covariance matrix.

Low-rank covariance estimation

- More generally, quadratic samplers produce the following:

$$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i + \eta_i, \quad i = 1, \dots, m;$$

where η is an additive noise.

- **linear in the covariance matrix \mathbf{X} !**

Low-rank covariance estimation

- More generally, quadratic samplers produce the following:

$$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i + \eta_i, \quad i = 1, \dots, m;$$

where η is an additive noise.

- **linear in the covariance matrix \mathbf{X} !**
- **Low-rank covariance matrix:** Many high-dimensional data lie in a low-dimensional subspace, when a small number of components accounts for most of the variability in the data.

$$\mathbf{X} = \mathbf{U}\mathbf{U}^T = \begin{array}{c} \begin{array}{|c|} \hline u_1 \\ \hline \end{array} \dots \begin{array}{|c|} \hline u_r \\ \hline \end{array} \end{array} \begin{array}{c} \begin{array}{|c|} \hline u_1^T \\ \hline \end{array} \\ \vdots \\ \begin{array}{|c|} \hline u_r^T \\ \hline \end{array} \end{array}$$

- This yields the *subspace retrieval* problem.

Reconstruction?

Two sides of the same coin: We can recover

- either $\mathbf{X} = \mathbf{U}\mathbf{U}^T \in \mathbb{R}^{n \times n}$ (when r is possibly unknown) or
- the subspace $\mathbf{U} \in \mathbb{R}^{n \times r}$ (when r is known);

	\mathbf{X}	\mathbf{U}
measurements	$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i$	$y_i = \ \mathbf{U}^T \mathbf{a}_i\ _2^2$
loss	linear	quadratic
prior	\mathbf{X} is low-rank	-
dim. of unknowns	n^2	nr
optimization	convex	nonconvex

Reconstruction?

Two sides of the same coin: We can recover

- either $\mathbf{X} = \mathbf{U}\mathbf{U}^T \in \mathbb{R}^{n \times n}$ (when r is possibly unknown) or
- the subspace $\mathbf{U} \in \mathbb{R}^{n \times r}$ (when r is known);

	\mathbf{X}	\mathbf{U}
measurements	$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i$	$y_i = \ \mathbf{U}^T \mathbf{a}_i\ _2^2$
loss	linear	quadratic
prior	\mathbf{X} is low-rank	-
dim. of unknowns	n^2	nr
optimization	convex	nonconvex

We will discuss both convex (for reconstructing \mathbf{X}) and nonconvex methods (for reconstructing \mathbf{U}), possibly with additional corruptions in the measurements.

Low-rank covariance estimation via convex relaxation

- We would like to seek the covariance matrix satisfying the observations with the minimal rank:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \operatorname{rank}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

Low-rank covariance estimation via convex relaxation

- We would like to seek the covariance matrix satisfying the observations with the minimal rank:

$$\hat{\mathbf{X}} = \operatorname{argmin}_{\mathbf{M} \succeq 0} \operatorname{rank}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

- However this is non-convex and NP-hard. Therefore, we replace it by a convex relaxation which is the **trace minimization**, over all PSD matrices compatible with the measurements:

$$\hat{\mathbf{X}} = \operatorname{argmin}_{\mathbf{M} \succeq 0} \operatorname{Tr}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

Low-rank covariance estimation via convex relaxation

- We would like to seek the covariance matrix satisfying the observations with the minimal rank:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \operatorname{rank}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

- However this is non-convex and NP-hard. Therefore, we replace it by a convex relaxation which is the **trace minimization**, over all PSD matrices compatible with the measurements:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \operatorname{Tr}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

- Additionally, if \mathbf{X} is *Toeplitz*, solve:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0, \text{Toeplitz}}{\operatorname{argmin}} \operatorname{Tr}(\mathbf{M}) \quad \text{s.t.} \quad y_i = \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i, \quad i = 1, \dots, m.$$

Near-optimal recovery via convex programming

Theorem (Chen, C. and Goldsmith)

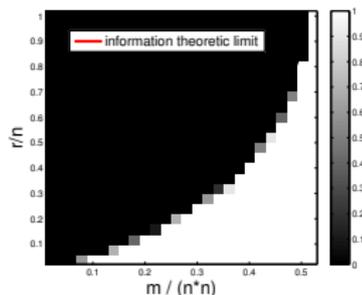
Assuming α_i 's are composed of i.i.d. Gaussian entries, with high probability, the solution \hat{X} exactly recovers all rank- r matrices X , provided that

$$m \gtrsim nr.$$

If there exists additional Toeplitz constraint, then similar guarantee holds provided

$$m \gtrsim r \text{polylog} n.$$

- **Exact recovery** with $m = O(nr)$;
- **Robust** against approximate low-rankness and bounded noise.
- **Under Toeplitz constraint:**

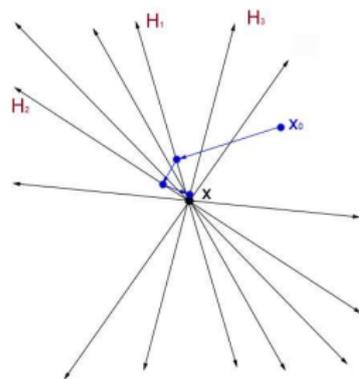
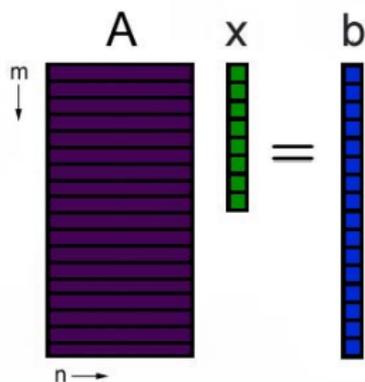


Kaczmarz method for solving quadratic equations

- Goal: reduce the memory and computational cost by directly estimating $\mathbf{U} \in \mathbb{R}^{n \times r}$.

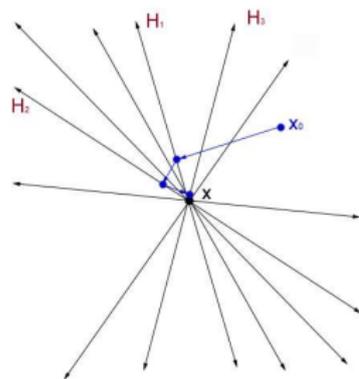
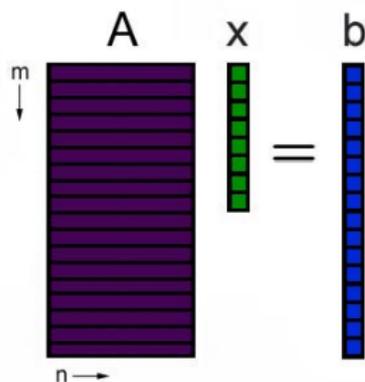
Kaczmarz method for solving quadratic equations

- Goal: reduce the memory and computational cost by directly estimating $U \in \mathbb{R}^{n \times r}$.
- The **Kaczmarz method** is a fast iterative algorithm for solving overdetermined linear system.



Kaczmarz method for solving quadratic equations

- Goal: reduce the memory and computational cost by directly estimating $U \in \mathbb{R}^{n \times r}$.
- The **Kaczmarz method** is a fast iterative algorithm for solving overdetermined linear system.



- Its randomized version [Strohmer and Vershynin] obtains linear rate in expectation.

Kaczmarz method for solving quadratic equations

- Extend Kaczmarz method by, at each iteration, project the current estimate to the closest signal that satisfies a (quadratic) constraint:[†]

$$\mathbf{U}_k = \underset{\mathbf{V}: \|\mathbf{V}^T \mathbf{a}_{\ell(k)}\|_2^2 = y_{\ell(k)}}{\operatorname{argmin}} \|\mathbf{U}_{k-1} - \mathbf{V}\|_{\mathbb{F}}^2,$$

[†]Y. Chi and Y. M. Lu, Kaczmarz Method for Solving Quadratic Equations, IEEE SPL 2016.

Kaczmarz method for solving quadratic equations

- Extend Kaczmarz method by, at each iteration, project the current estimate to the closest signal that satisfies a (quadratic) constraint:[†]

$$U_k = \underset{V: \|V^T \mathbf{a}_{\ell(k)}\|_2^2 = y_{\ell(k)}}{\operatorname{argmin}} \|U_{k-1} - V\|_F^2,$$

which can be solved in **closed form** via a **rank-one** update:

$$U_k = \left[\mathbf{I} - \left(\frac{\|U_{k-1}^T \mathbf{a}_{\ell(k)}\|_2 - \sqrt{y_{\ell(k)}}}{\|U_{k-1}^T \mathbf{a}_{\ell(k)}\|_2} \right) \frac{\mathbf{a}_{\ell(k)} \mathbf{a}_{\ell(k)}^T}{\|\mathbf{a}_{\ell(k)}\|_2^2} \right] U_{k-1}.$$

[†]Y. Chi and Y. M. Lu, Kaczmarz Method for Solving Quadratic Equations, IEEE SPL 2016.

Kaczmarz method for solving quadratic equations

- Extend Kaczmarz method by, at each iteration, project the current estimate to the closest signal that satisfies a (quadratic) constraint:[†]

$$U_k = \underset{\mathbf{V}: \|\mathbf{V}^T \mathbf{a}_{\ell(k)}\|_2^2 = y_{\ell(k)}}{\operatorname{argmin}} \|\mathbf{U}_{k-1} - \mathbf{V}\|_F^2,$$

which can be solved in **closed form** via a **rank-one** update:

$$U_k = \left[\mathbf{I} - \left(\frac{\|\mathbf{U}_{k-1}^T \mathbf{a}_{\ell(k)}\|_2 - \sqrt{y_{\ell(k)}}}{\|\mathbf{U}_{k-1}^T \mathbf{a}_{\ell(k)}\|_2} \right) \frac{\mathbf{a}_{\ell(k)} \mathbf{a}_{\ell(k)}^T}{\|\mathbf{a}_{\ell(k)}\|_2^2} \right] U_{k-1}.$$

- The solution is equivalent to

$$\min_{\mathbf{s}: \|\mathbf{s}\|_2=1} \underset{\mathbf{V}: \|\mathbf{V}^T \mathbf{a}_{\ell(k)}\|_2^2 = \mathbf{s}^T \sqrt{y_{\ell(k)}}}{\operatorname{argmin}} \|\mathbf{U}_{k-1} - \mathbf{V}\|_F^2$$

which corresponds to projecting the current estimate to the hyperplane with the phase that minimizes the projection.

[†]Y. Chi and Y. M. Lu, Kaczmarz Method for Solving Quadratic Equations, IEEE SPL 2016.

Performance Guarantee of Kaczmarz Method

Consider the phase retrieval case.

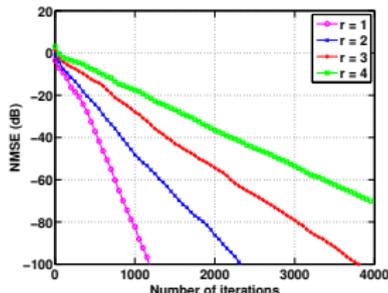
Theorem (Zhang, C., Liang)

Assume \mathbf{a}_i 's are generated with i.i.d. Gaussian entries, there exist some universal constants $\rho > 0$ such that if $m \gtrsim n$, then with high probability, randomized Kaczmarz update rule yields

$$\mathbb{E}_{i_t} \left[\text{dist}^2(\mathbf{z}^{(t+1)}, \mathbf{x}) \right] \leq \left(1 - \frac{\rho}{n} \right) \text{dist}^2(\mathbf{z}^{(t)}, \mathbf{x})$$

where $\mathbf{z}^{(0)}$ is initialized via the spectral method.

- This establishes linear convergence rate *in expectation*, despite the nonlinearity!
- We can obtain similar guarantees for the [block Kaczmarz](#) method which is further accelerated.



What about outliers?

- Outliers happen with
 - sensor failures, malicious attacks, ...
 - For covariance sketching, insufficiently aggregated sketches can be regarded as an outlier;

What about outliers?

- Outliers happen with
 - sensor failures, malicious attacks, ...
 - For covariance sketching, insufficiently aggregated sketches can be regarded as an outlier;
- We're interested when the measurements are corrupted by both *sparse outliers* and *bounded noise*:

$$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i + \eta_i + w_i, \quad i = 1, \dots, m,$$

where $\mathbf{X} = \mathbf{U}\mathbf{U}^T$, $\|\boldsymbol{\eta}\|_0 \leq sm$ and \mathbf{w} is a dense bounded noise.

What about outliers?

- Outliers happen with
 - sensor failures, malicious attacks, ...
 - For covariance sketching, insufficiently aggregated sketches can be regarded as an outlier;
- We're interested when the measurements are corrupted by both *sparse outliers* and *bounded noise*:

$$y_i = \mathbf{a}_i^T \mathbf{X} \mathbf{a}_i + \eta_i + w_i, \quad i = 1, \dots, m,$$

where $\mathbf{X} = \mathbf{U}\mathbf{U}^T$, $\|\boldsymbol{\eta}\|_0 \leq sm$ and \mathbf{w} is a dense bounded noise.

- **Goal:** develop algorithms that are *oblivious* to outliers, and statistically and computationally efficient.
 - small sample size: hopefully m is linear in n ;
 - large fraction of outliers: hopefully s is a small constant;
 - low computational complexity and easy to implement.

Outlier-robust recovery by convex programming

- To motivate, ideally one would like to look for low-rank matrices that maintain outlier sparsity:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \operatorname{cardinality}(\text{outliers}), \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{M}) = r$$

Outlier-robust recovery by convex programming

- To motivate, ideally one would like to look for low-rank matrices that maintain outlier sparsity:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \text{cardinality}(\text{outliers}), \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{M}) = r$$

- By *relaxing* the objective function to the ℓ_1 -norm minimization, and *dropping* the rank constraint, we propose to solve

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \sum_{i=1}^m |y_i - \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i|$$

Outlier-robust recovery by convex programming

- To motivate, ideally one would like to look for low-rank matrices that maintain outlier sparsity:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \text{cardinality}(\text{outliers}), \quad \text{s.t.} \quad \operatorname{rank}(\mathbf{M}) = r$$

- By *relaxing* the objective function to the ℓ_1 -norm minimization, and *dropping* the rank constraint, we propose to solve

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0}{\operatorname{argmin}} \sum_{i=1}^m |y_i - \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i|$$

- **Parameter-free** formulation without trace minimization or tuning parameters;
- No prior information is required for the matrix rank, corruption level or bounded noise level.

Performance guarantee of convex programming

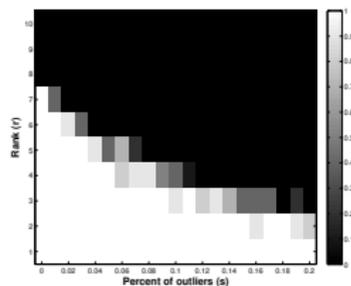
Theorem (Li, Sun and C., 2016)

Suppose that $\|\mathbf{w}\|_1 \leq \epsilon$. Assume the support of $\boldsymbol{\eta}$ is selected uniformly at random with the signs of $\boldsymbol{\eta}$ are generated from a symmetric Bernoulli distribution. Then as long as $m \gtrsim nr^2$, $s \lesssim 1/r$, the solution to the proposed algorithm satisfies

$$\|\hat{\mathbf{X}} - \mathbf{X}\|_{\text{F}} \lesssim \frac{r\epsilon}{m}$$

with high probability.

- Exact recovery when $\mathbf{w} = 0$ as long as $m \gtrsim nr^2$ and $s \lesssim 1/r$.
- When $r = 1$ recovers a previous result for the phase retrieval case[‡];
- RHS is phase transition for m vs r with 5% corruptions.



[‡]P. Hand, “Phaselift is robust to a constant fraction of arbitrary errors”.

Robust recovery of Toeplitz PSD Matrices

If \mathbf{X} is additionally Toeplitz, this can be incorporated:

$$\hat{\mathbf{X}} = \underset{\mathbf{M} \succeq 0, \text{ Toeplitz}}{\operatorname{argmin}} \sum_{i=1}^m |y_i - \mathbf{a}_i^T \mathbf{M} \mathbf{a}_i|.$$

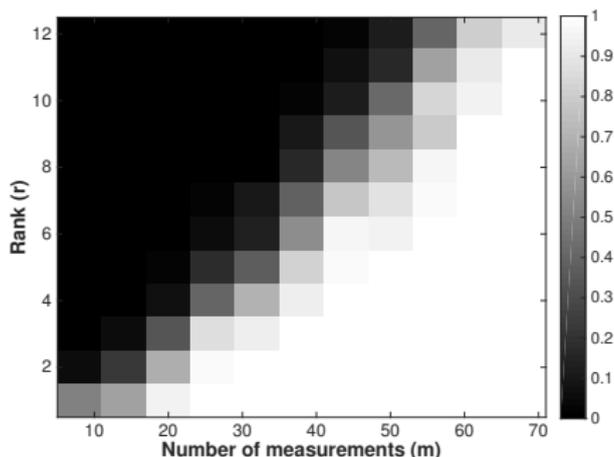
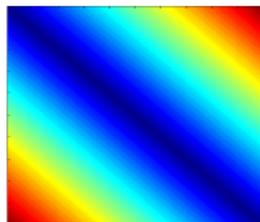


Figure : Phase transitions of low-rank Toeplitz PSD matrix recovery w.r.t. the number of measurements and the rank with 5% of measurements corrupted by standard Gaussian variables, when $n = 64$.

Non-convex approach based on factored model

Can we reduce the computational complexity?

- Recall $\mathbf{X} = \mathbf{U}\mathbf{U}^T$ where $\mathbf{U} \in \mathbb{R}^{n \times r}$, one can directly recover \mathbf{U} by attempting:

$$\hat{\mathbf{U}} = \operatorname{argmin}_{\mathbf{U} \in \mathbb{R}^{n \times r}} \ell(\mathbf{U}) := \operatorname{argmin}_{\mathbf{U} \in \mathbb{R}^{n \times r}} \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{U})$$

Non-convex approach based on factored model

Can we reduce the computational complexity?

- Recall $\mathbf{X} = \mathbf{U}\mathbf{U}^T$ where $\mathbf{U} \in \mathbb{R}^{n \times r}$, one can directly recover \mathbf{U} by attempting:

$$\hat{\mathbf{U}} = \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\operatorname{argmin}} \ell(\mathbf{U}) := \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{U})$$

for some **loss function** $\ell(y_i, \mathbf{U})$:

- quadratic loss of power: $\ell(\mathbf{U}; y_i) = \left(y_i - \|\mathbf{U}^T \mathbf{a}_i\|_2^2 \right)^2$
- quadratic loss of amplitude: $\ell(\mathbf{U}; y_i) = \left(\sqrt{y_i} - \|\mathbf{U}^T \mathbf{a}_i\|_2 \right)^2$
- Poisson loss: $\ell(\mathbf{U}; y_i) = \|\mathbf{U}^T \mathbf{a}_i\|_2^2 - y_i \log \|\mathbf{U}^T \mathbf{a}_i\|_2^2$

Non-convex approach based on factored model

Can we reduce the computational complexity?

- Recall $\mathbf{X} = \mathbf{U}\mathbf{U}^T$ where $\mathbf{U} \in \mathbb{R}^{n \times r}$, one can directly recover \mathbf{U} by attempting:

$$\hat{\mathbf{U}} = \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\operatorname{argmin}} \ell(\mathbf{U}) := \underset{\mathbf{U} \in \mathbb{R}^{n \times r}}{\operatorname{argmin}} \frac{1}{m} \sum_{i=1}^m \ell(y_i; \mathbf{U})$$

for some **loss function** $\ell(y_i, \mathbf{U})$:

- quadratic loss of power: $\ell(\mathbf{U}; y_i) = \left(y_i - \|\mathbf{U}^T \mathbf{a}_i\|_2^2 \right)^2$
 - quadratic loss of amplitude: $\ell(\mathbf{U}; y_i) = \left(\sqrt{y_i} - \|\mathbf{U}^T \mathbf{a}_i\|_2 \right)^2$
 - Poisson loss: $\ell(\mathbf{U}; y_i) = \|\mathbf{U}^T \mathbf{a}_i\|_2^2 - y_i \log \|\mathbf{U}^T \mathbf{a}_i\|_2^2$
- What are the challenges?
 - $\ell(\mathbf{U})$ can be non-convex and non-smooth.
 - With outliers, we want the loss to sum over only clean samples.

Non-convex phase retrieval

Exciting developments (without outliers) – all following the same recipe (for the phase retrieval or rank-1 case):

$$\hat{z} = \operatorname{argmin}_{z \in \mathbb{R}^n} \frac{1}{m} \sum_{i=1}^m \ell(y_i; z)$$

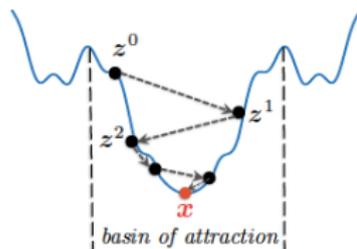
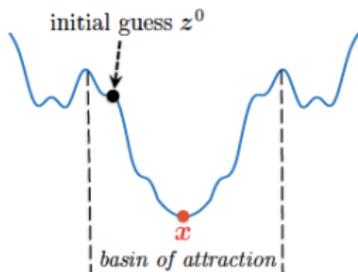
- Initialize $z^{(0)}$ via the (truncated) spectral method to land in the neighborhood of the ground truth;
- Iterative update using (truncated) gradient descent;



§

§Figure credit: Yuxin Chen.

Non-convex phase retrieval



Provable near-optimal performance for Gaussian measurement model:

- Statistically: $m = O(n)$ near-optimal sample complexity
- Computationally: linear convergence with near-linear run time

Non-convex phase retrieval



Provable near-optimal performance for Gaussian measurement model:

- Statistically: $m = O(n)$ near-optimal sample complexity
- Computationally: linear convergence with near-linear run time

Examples: Wirtinger Flow (WF) (Candès et.al. 2014), Truncated Wirtinger Flow (TWF) (Chen and Candès 2015), Reshaped Wirtinger Flow (Zhang and Liang 2016), Truncated Amplitude Flow (Wang, Giannakis and Eldar, 2016)

Non-convex phase retrieval with outliers

In the presence of *arbitrary outliers*, **existing approaches fail**:

- **Spectral initialization would fail**: the eigenvector of \mathbf{Y} can be arbitrarily perturbed

$$\underbrace{\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T}_{\text{WF}} \quad \text{or} \quad \underbrace{\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T \mathbb{1}_{\{|y_i| \leq \alpha_y \cdot \text{mean}(\{y_i\})\}}}_{\text{TWF}} \cdot$$

Non-convex phase retrieval with outliers

In the presence of *arbitrary outliers*, **existing approaches fail**:

- **Spectral initialization would fail**: the eigenvector of \mathbf{Y} can be arbitrarily perturbed

$$\underbrace{\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T}_{\text{WF}} \quad \text{or} \quad \underbrace{\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T \mathbb{1}_{\{|y_i| \leq \alpha_y \cdot \text{mean}(\{y_i\})\}}}_{\text{TWF}}.$$

- **Gradient descent would fail**: the search direction can be arbitrarily perturbed

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{\mu}{\|\mathbf{z}^{(0)}\|^2} \sum_{i \in \mathcal{T}_t} \nabla \ell(\mathbf{z}^{(t)}; y_i)$$

where $\mathcal{T}_t = \{1, \dots, m\}$ for WF and

$$\mathcal{T}_t = \left\{ i : |y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}|^2| \leq \alpha_h \cdot \text{mean}(\{|y_i - |\mathbf{a}_i^T \mathbf{z}^{(t)}|^2|\}) \right\} \quad \P$$

for TWF.

^{\P}with some details hiding

Robust phase retrieval via median-truncation

Need better strategy to eliminate outliers!

Key approach: “median-truncation”

- well-known in robust statistics to be outlier-resilient;
- little appearance in high-dimensional estimation;



Robust phase retrieval via median-truncation

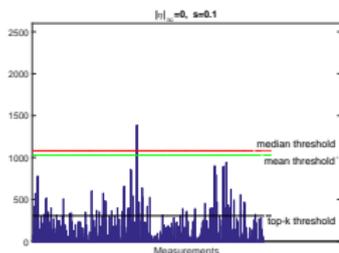
Need better strategy to eliminate outliers!

Key approach: “median-truncation”

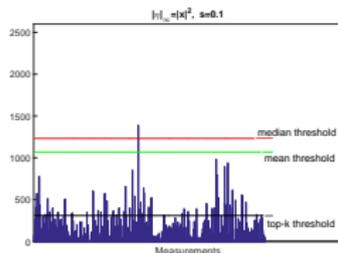
- well-known in robust statistics to be outlier-resilient;
- little appearance in high-dimensional estimation;



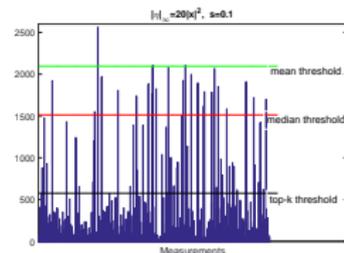
Median is more stable than mean and top-k truncation (which truncates a fixed amount of samples) for various levels of outliers.



no outliers



small outlier magnitudes



large outlier magnitudes

Median-Truncated Wirtinger Flow (median-TWF)

We adopt the Poisson loss function (other loss functions work too) and the Gaussian measurement model.

- **Median-truncated spectral initialization:** Set $\mathbf{z}^{(0)} := \lambda_0 \tilde{\mathbf{z}}$ where
 - Direction estimation: $\tilde{\mathbf{z}}$ is the leading eigenvector of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T \mathbb{1}_{\{|y_i| \leq 9/0.455 \cdot \text{median}(\{y_i\})\}}.$$

- Norm estimation: $\lambda_0 = \sqrt{\text{median}(\{y_i\})/0.455}$

$$y_i = |\mathbf{a}_i^T \mathbf{x}|^2 \sim \chi_1^2 \quad \text{and} \quad \mathbb{E}[\text{median}(\chi_1^2)] = 0.455$$

Median-Truncated Wirtinger Flow (median-TWF)

We adopt the Poisson loss function (other loss functions work too) and the Gaussian measurement model.

- **Median-truncated spectral initialization:** Set $\mathbf{z}^{(0)} := \lambda_0 \tilde{\mathbf{z}}$ where
 - Direction estimation: $\tilde{\mathbf{z}}$ is the leading eigenvector of

$$\mathbf{Y} = \frac{1}{m} \sum_{i=1}^m y_i \mathbf{a}_i \mathbf{a}_i^T \mathbb{1}_{\{|y_i| \leq 9/0.455 \cdot \text{median}(\{y_i\})\}}.$$

- Norm estimation: $\lambda_0 = \sqrt{\text{median}(\{y_i\})/0.455}$

$$y_i = |\mathbf{a}_i^T \mathbf{x}|^2 \sim \chi_1^2 \quad \text{and} \quad \mathbb{E}[\text{median}(\chi_1^2)] = 0.455$$

- As long as $m = O(n \log n)$ and $s = O(1)$, the initialization is provably close to the ground truth:

$$\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) \leq \frac{1}{10} \|\mathbf{x}\|,$$

where $\text{dist}(\mathbf{z}^{(0)}, \mathbf{x}) = \min\{\|\mathbf{z}^{(0)} + \mathbf{x}\|, \|\mathbf{z}^{(0)} - \mathbf{x}\|\}$.

Median-Truncated Wirtinger Flow (median-TWF)

- Median-truncated gradient descent:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{2\mu}{m} \underbrace{\sum_{i \in \mathcal{E}_1 \cap \mathcal{E}_2} \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|^2 - y_i}{\mathbf{a}_i^T \mathbf{z}^{(t)}} \mathbf{a}_i}_{\nabla \ell_{tr}(\mathbf{z})},$$

where

$$\mathcal{E}_1 = \left\{ i : 0.3 \leq \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|}{\|\mathbf{z}^{(t)}\|} \leq 5 \right\}, \mathcal{E}_2 = \left\{ i : r_i^{(t)} \leq 12 \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|}{\|\mathbf{z}^{(t)}\|} \cdot \text{median}(\{r_i^{(t)}\}) \right\},$$

$$\text{with } r_i^{(t)} = |y_i - (\mathbf{a}_i^T \mathbf{z}^{(t)})^2|.$$

Median-Truncated Wirtinger Flow (median-TWF)

- Median-truncated gradient descent:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} - \frac{2\mu}{m} \underbrace{\sum_{i \in \mathcal{E}_1 \cap \mathcal{E}_2} \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|^2 - y_i}{\mathbf{a}_i^T \mathbf{z}^{(t)}} \mathbf{a}_i}_{\nabla \ell_{tr}(\mathbf{z})},$$

where

$$\mathcal{E}_1 = \left\{ i : 0.3 \leq \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|}{\|\mathbf{z}^{(t)}\|} \leq 5 \right\}, \mathcal{E}_2 = \left\{ i : r_i^{(t)} \leq 12 \frac{|\mathbf{a}_i^T \mathbf{z}^{(t)}|}{\|\mathbf{z}^{(t)}\|} \cdot \text{median}(\{r_i^{(t)}\}) \right\},$$

with $r_i^{(t)} = |y_i - (\mathbf{a}_i^T \mathbf{z}^{(t)})^2|$.

- As long as $m = O(n \log n)$ and $s = O(1)$, $\nabla \ell_{tr}(\mathbf{z})$ satisfies the *Regularity Condition* $\text{RC}(\mu, \lambda)$ for all \mathbf{z} , $\mathbf{h} = \mathbf{z} - \mathbf{x}$:

$$-\left\langle \frac{1}{m} \nabla \ell_{tr}(\mathbf{z}), \mathbf{h} \right\rangle \geq \mu \left\| \frac{1}{m} \nabla \ell_{tr}(\mathbf{z}) \right\|^2 + \lambda \|\mathbf{h}\|^2, \quad \|\mathbf{h}\| \leq \frac{1}{10} \|\mathbf{z}\|.$$

which guarantees $\text{dist}(\mathbf{z}^{(t+1)}, \mathbf{x}) \leq (1 - \mu\lambda) \text{dist}(\mathbf{z}^{(t)}, \mathbf{x})$.

Performance guarantee of median-TWF

Theorem (Zhang, C. and Liang, 2016)

Assume $\|\mathbf{w}\|_\infty \leq c_1 \|\mathbf{x}\|^2$. Assume \mathbf{a}_i 's are generated with i.i.d. Gaussian entries. If $m \gtrsim n \log n$ and $s \lesssim s_0$, then with high probability, median-TWF yields

$$\text{dist}(\mathbf{z}^{(t)}, \mathbf{x}) \lesssim \frac{\|\mathbf{w}\|_\infty}{\|\mathbf{x}\|} + (1 - \rho)^t \|\mathbf{x}\|, \quad \forall t \in \mathbb{N}$$

simultaneously for all $\mathbf{x} \in \mathbb{R}^n \setminus \{\mathbf{0}\}$ for some $0 < \rho < 1$.

- **Exact recovery** when $\|\mathbf{w}\| = 0$ with slight more samples ($m = O(n \log n)$) but a constant fraction of outliers $s = O(1)$.
- **Stable recovery** with additional bounded noise;
- Resist outliers **obliviously**: no prior knowledge of outliers.
- **First** non-asymptotic robust recovery guarantee using median: much more involved due to the nonlinearity of median.

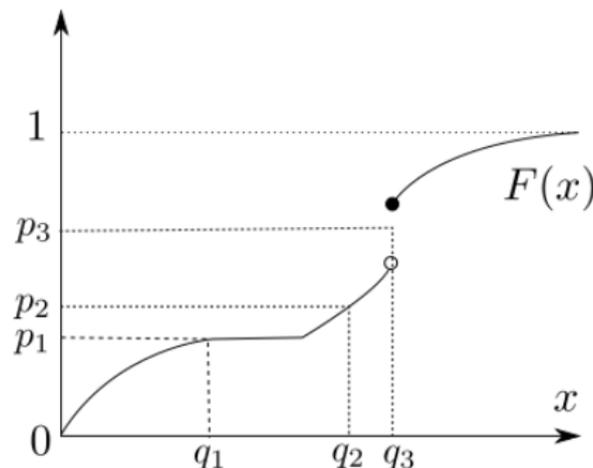
Proof sketch - preparation

Definition (Generalized quantile function)

Let $0 < p < 1$. If F is a CDF, the generalized quantile function is

$$F^{-1}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}.$$

Denote $\theta_p(F) := F^{-1}(p)$ and $\theta_p(\{X_i\}) := \theta_p(\hat{F})$, where \hat{F} is the empirical distribution of the samples $\{X_i\}_{i=1}^m$.



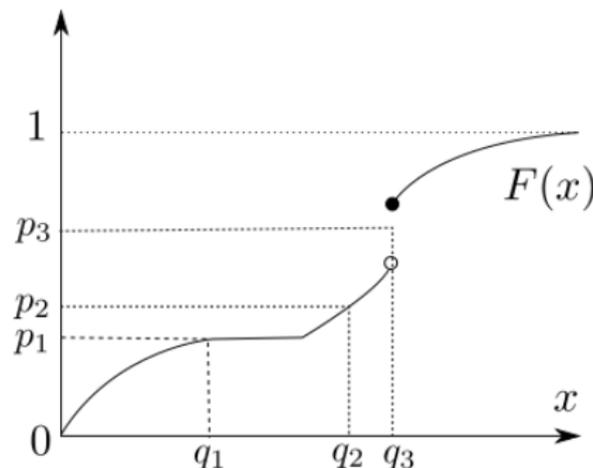
Proof sketch - preparation

Definition (Generalized quantile function)

Let $0 < p < 1$. If F is a CDF, the generalized quantile function is

$$F^{-1}(p) = \inf\{x \in \mathbb{R} : F(x) \geq p\}.$$

Denote $\theta_p(F) := F^{-1}(p)$ and $\theta_p(\{X_i\}) := \theta_p(\hat{F})$, where \hat{F} is the empirical distribution of the samples $\{X_i\}_{i=1}^m$.



Proof sketch

Lemma (Concentration of sample quantile)

Assume $\{X_i\}_{i=1}^m$ are i.i.d. drawn from some distribution F . Under some minor assumptions, w.h.p.

$$|\theta_p(\{X_i\}_{i=1}^m) - \theta_p(F)| < \epsilon$$

Lemma (Sandwich median by quantiles of clean samples)

Consider clean samples $\{\tilde{X}_i\}_{i=1}^m$ and contaminated samples $\{X_i\}_{i=1}^m$. Then

$$\theta_{\frac{1}{2}-s}(\{\tilde{X}_i\}) \leq \theta_{\frac{1}{2}}(\{X_i\}) \leq \theta_{\frac{1}{2}+s}(\{\tilde{X}_i\}).$$

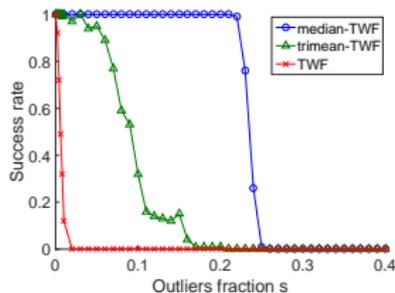
Lemma (Concentration of median)

If $m > c_0 n \log n$, then with probability at least $1 - c_1 \exp(-c_2 m)$, there exist constants β and β' such that

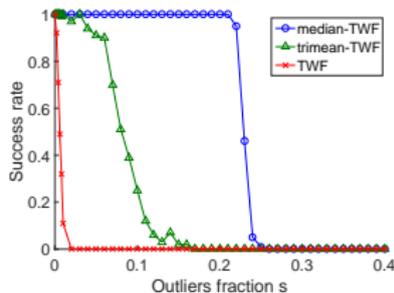
$$\beta \|z\| \|\mathbf{h}\| \leq \text{median}(\{|\mathbf{a}_i^T \mathbf{x}|^2 - |\mathbf{a}_i^T \mathbf{z}|^2\}_{i=1}^m) \leq \beta' \|z\| \|\mathbf{h}\|,$$

holds for all $\mathbf{z}, \mathbf{h} := \mathbf{z} - \mathbf{x}$ satisfying $\|\mathbf{h}\| < 1/11\|\mathbf{z}\|$.

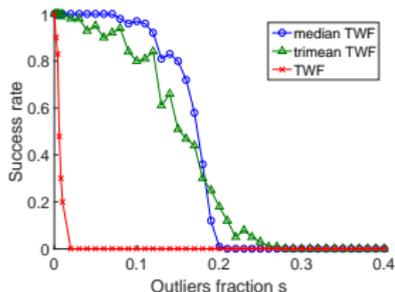
Numerical experiments with median-TWF



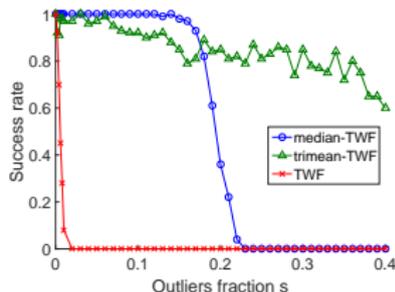
(a) $\|\eta\|_\infty = 0.1\|\mathbf{x}\|^2$



(b) $\|\eta\|_\infty = \|\mathbf{x}\|^2$



(c) $\|\eta\|_\infty = 10\|\mathbf{x}\|^2$



(d) $\|\eta\|_\infty = 100\|\mathbf{x}\|^2$

Figure : Success rate of **exact recovery** with outliers for **median-TWF**, **trimean-TWF**, and **TWF** at different levels of outlier magnitudes.

Numerical experiments with median-TWF

Recovery with both dense noise and sparse outliers:

- With outliers, median-TWF achieve better accuracy than TWF.
- Moreover, median-TWF with outliers achieves almost the same accuracy of TWF without outliers.

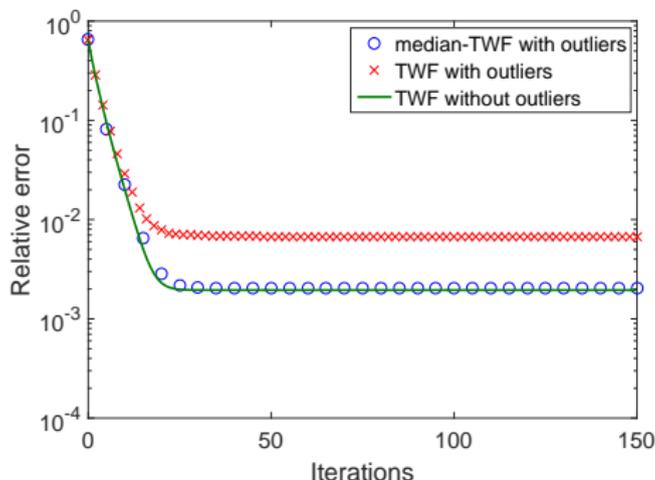


Figure : Relative error of median-TWF vs. TWF w.r.t. iteration when $s = 0.1$, $\|\mathbf{w}\|_{\infty} = 0.01\|\mathbf{x}\|^2$, and $\|\boldsymbol{\eta}\|_{\infty} = \|\mathbf{w}\|$.

Conclusions

We have discussed how to solve random quadratic systems of equations, possibly corrupted by a constant fraction of outliers, in a provable manner.

	X	U
measurements	$y_i = \mathbf{a}_i^T X \mathbf{a}_i$	$y_i = \ \mathbf{U}^T \mathbf{a}_i\ _2^2$
loss	linear/cvx	quadratic/ncvx
without outliers	Semidefinite Prog.	Kaczmarz/SGD
with outliers	Semidefinite Prog.	median-TWF

- **The class of convex methods** are based on convex relaxation for low-rank matrix completion and sparse recovery. It is easier to design but the computational cost is high;
- **The class of non-convex methods** are based on iterative updates with careful initializations. The computational cost is low but the design is a bit of an art.

References

1. Exact and Stable Covariance Estimation from Quadratic Sampling via Convex Programming, IEEE TIT 2015.
2. Low-Rank Positive Semidefinite Matrix Recovery from Corrupted Rank-One Measurements, IEEE TSP 2016.
3. Provable Non-convex Phase Retrieval with Outliers: Median Truncated Wirtinger Flow, ICML 2016.
4. Kaczmarz Method for Solving Quadratic Equations, IEEE SPL 2016.
5. Incremental Reshaped Wirtinger Flow and Its Connection to Kaczmarz Method, NIPS 2016 Workshop on Nonconvex Optimization.

<http://www.ece.osu.edu/~chi/>

Acknowledgement

- My collaborators: Yuxin Chen (Stanford), Andrea Goldsmith (Stanford), Yuanxin Li (OSU), Huishuai Zhang (Syracuse), Yingbin Liang (Syracuse) and Yue M. Lu (Harvard).



- Research supported by NSF, AFOSR and ONR.

