

Covariance Sketching via Quadratic Sampling

Yuejie Chi
Department of ECE and BMI
The Ohio State University

Tsinghua University
June 2015



THE OHIO STATE UNIVERSITY

Acknowledgement

- Thanks to my academic collaborators on some of the reported work:



Yuxin Chen



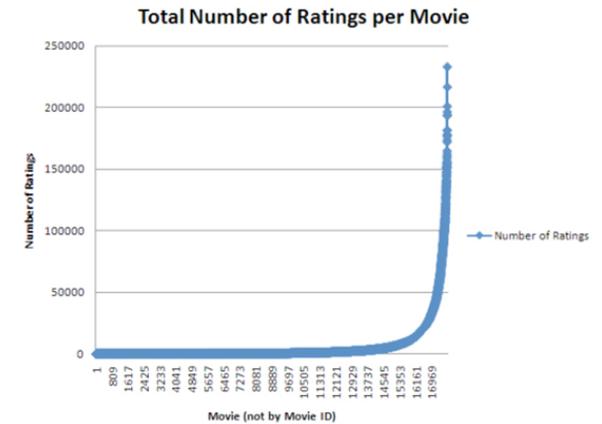
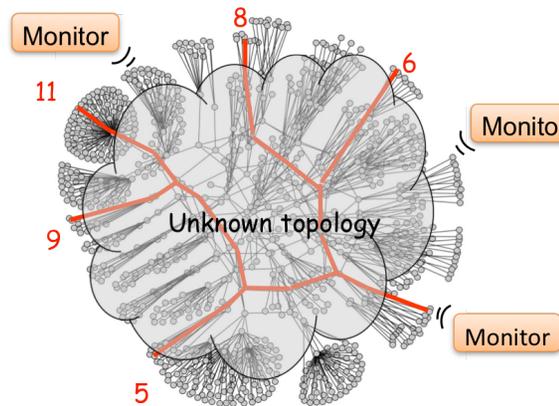
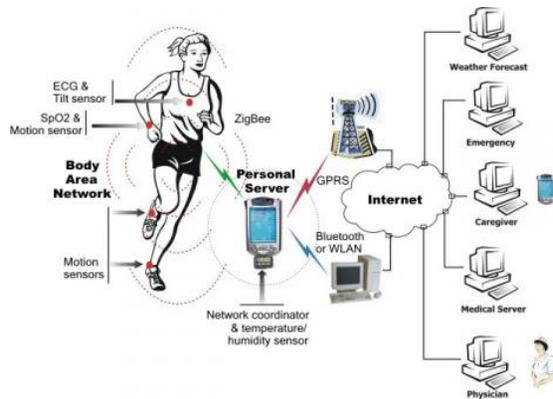
Andrea Goldsmith

- This work is supported by NSF and AFOSR.



High-dimensional Streaming Data

Each time snapshot a data vector $x_t \in \mathbb{R}^n$ is generated, with n large.



Wireless health monitoring:
transmission loss and
power-hungry sensors

Internet monitoring:
limited measurements and
massive data

Netflix Ratings:
incomplete (and skewed)
ratings

- Need to **learn and track** structures of minimally observed data streams.

Jovanov, Emil, et al. "A wireless body area network of intelligent motion sensors for computer assisted physical rehabilitation." Journal of NeuroEngineering and rehabilitation 2.1 (2005): 6.

Netflix statistics source: <http://www.hackedexistence.com/project-netflix.html>

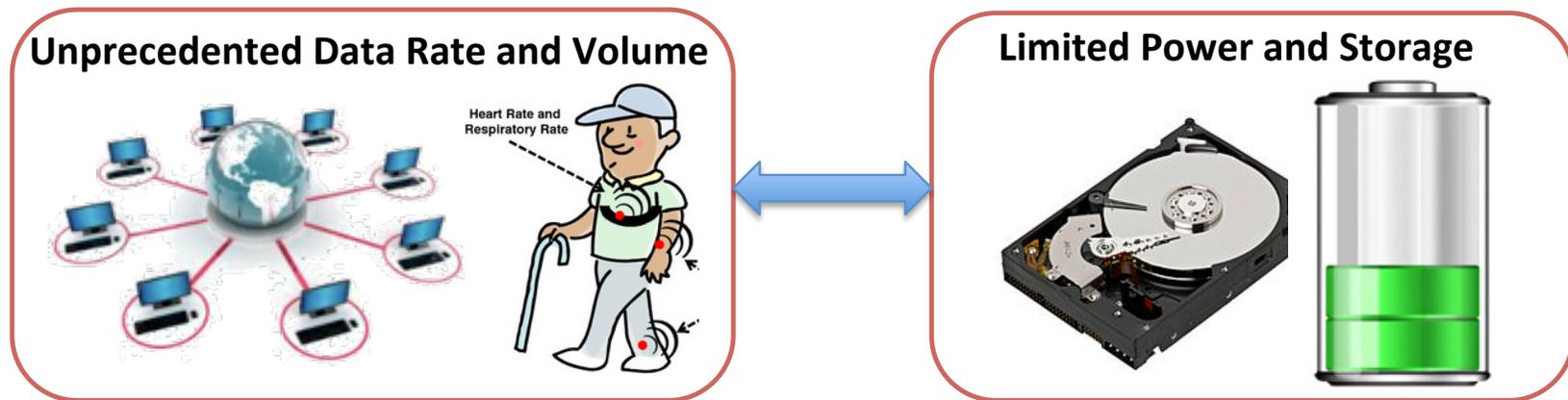
Challenges in Modern Data Acquisition

Data generation at unprecedented rate: data samples are

- not observable due to **privacy** or security constraints;
- **distributed** at multiple locations;
- online generated **on the fly** and can only access once.

Limited processing power at sensor platforms:

- **time-sensitive:** impossible to obtain a complete snapshot of the system;
- **storage-limited:** cannot store the whole data set;
- **power-hungry:** minimize the number of observations.



Covariance Sketching

Key Observation: the *covariance structure* can be recovered without measuring the whole data stream.

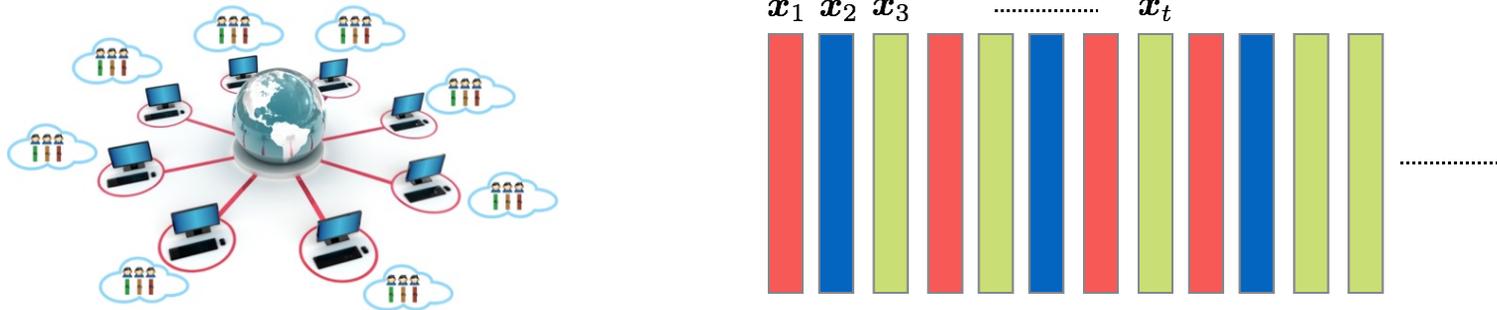


Approach: distributed data sketching and aggregation to recover the covariance structure or principal components.

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

Quadratic Sketching for Covariance Estimation

Consider a data stream possibly distributively observed at m sensors:



Quadratic Sketching: For each sensor $i = 1, \dots, m$:

- *randomly* select a sketching vector $\mathbf{a}_i \in \mathbb{R}^n$ with i.i.d. sub-Gaussian entries;
- Sketch an arbitrary substream indexed by $\{\ell_t^i\}_{t=1}^T$ with an energy measurement $\left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2$ and aggregate the average **energy** measurement:

$$y_{i,T} = \frac{1}{T} \sum_{t=1}^T \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2 = \frac{T-1}{T} y_{i,T-1} + \frac{1}{T} \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_T^i} \rangle \right|^2.$$

Quadratic Sketching for Covariance Estimation

- For each sketch:

$$y_{i,T} = \frac{1}{T} \sum_{t=1}^T \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2 = \frac{T-1}{T} y_{i,T-1} + \frac{1}{T} \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_T^i} \rangle \right|^2.$$

$$\frac{1}{T} \sum_{i=1}^T \left(\begin{array}{c|c|c} \mathbf{a}_i^T & \mathbf{x}_{\ell_t^i}^T & \\ \hline & \mathbf{x}_{\ell_t^i} & \\ \hline & & \mathbf{a}_i \end{array} \right) = \begin{array}{c|c} \mathbf{a}_i^T & \\ \hline & \left(\frac{1}{T} \sum_{i=1}^T \mathbf{x}_{\ell_t^i} \right) \\ \hline & \mathbf{a}_i \end{array} = \begin{array}{c|c} \mathbf{a}_i^T & \\ \hline & \mathbf{a}_i \end{array} \quad \boxed{\Sigma_T}$$

- As $T \rightarrow \infty$, $\Sigma_T \rightarrow \Sigma$,

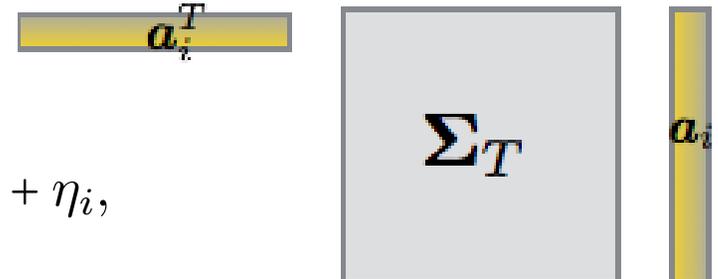
$$y_{i,T} \rightarrow \mathbf{a}_i^T \Sigma \mathbf{a}_i$$

yields a linear measurement of Σ !

- All sketches can be obtained in a fully distributed manner.

Covariance Estimation with Rank-One Measurements

- Quadratic Measurement Model:



$$y_{i,T} = \mathbf{a}_i^T \Sigma_T \mathbf{a}_i := \mathbf{a}_i^T \Sigma \mathbf{a}_i + \eta_i,$$

where $\eta_i = \mathbf{a}_i^T (\Sigma_T - \Sigma) \mathbf{a}_i$ is the additive noise.

- More generally, we assume the following measurement model:

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

or more concisely,

$$\mathbf{z} = \mathcal{A}(\Sigma) + \boldsymbol{\eta}.$$

- The measurements are **quadratic** in \mathbf{a}_i and **linear** in the rank-one matrix $\mathbf{a}_i \mathbf{a}_i^T$;

Sampling Model

We need some additional assumptions on the sampling model:

- **sub-Gaussian i.i.d. sketching vectors:** each a_i is i.i.d. copies of $a = [a_1, a_2, \dots, a_n]^T$ satisfying:

$$\mathbb{E}[a_i] = 0, \quad \mathbb{E}[a_i^2] = 1, \quad \text{and} \quad \mu_4 := \mathbb{E}a_i^4 > 1.$$

- **noise model:** deterministically bounded ℓ_1 -norm noise:

$$\|\boldsymbol{\eta}\|_1 \leq \epsilon$$

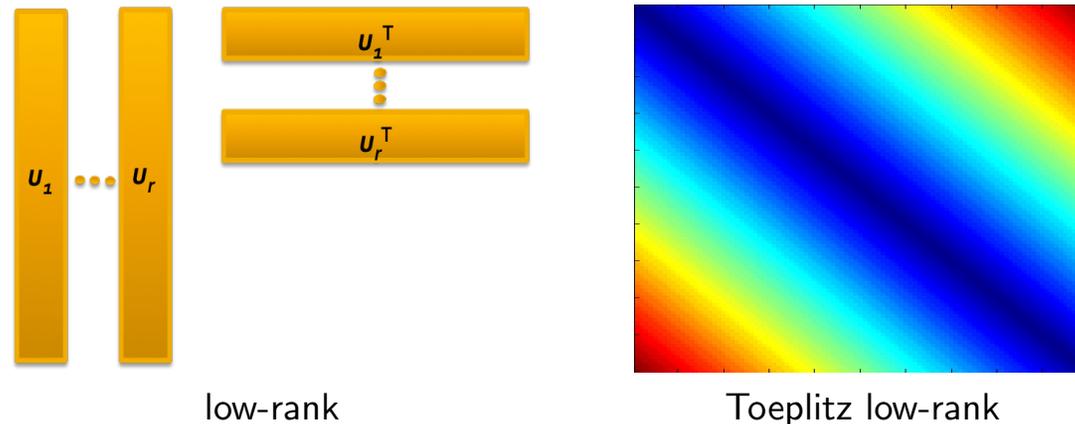
where ϵ is known *a priori*.

While we can solve for Σ via least-squares estimation if $m \geq n^2$ (the size of Σ), we can greatly reduce the number of m by **exploiting the low-dimensional structure of Σ** .

Geometry of Covariance Structures

Many high-dimensional data lie in a low-dimensional subspace, resulting in a low-rank covariance matrix:

- **Low-Rankness:** the covariance matrix is low-rank, which occurs when a small number of components accounts for most of the variability in the data.
- **Stationary Low-Rankness:** the covariance matrix is simultaneously Toeplitz and low rank, which has many applications in array signal processing.



Low-Rank Covariance Estimation via Convex Relaxation

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

$$\hat{\Sigma} = \underset{\Sigma \succeq 0}{\operatorname{argmin}} \operatorname{rank}(\Sigma) \quad \text{s.t.} \quad \|z - \mathcal{A}(\Sigma)\|_1 \leq \epsilon.$$

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

$$\hat{\Sigma} = \underset{\Sigma \succeq 0}{\operatorname{argmin}} \operatorname{Tr}(\Sigma) \quad \text{s.t.} \quad \|z - \mathcal{A}(\Sigma)\|_1 \leq \epsilon.$$

- Additionally, if Σ is Toeplitz, we add the additional structural constraint:

$$\hat{\Sigma} = \underset{\Sigma \succeq 0}{\operatorname{argmin}} \operatorname{Tr}(\Sigma) \quad \text{s.t.} \quad \|z - \mathcal{A}(\Sigma)\|_2 \leq \epsilon, \text{ and } \Sigma \text{ is Toeplitz.}$$

Our theoretical results considered ℓ_2 noise for Toeplitz covariance matrix recovery.

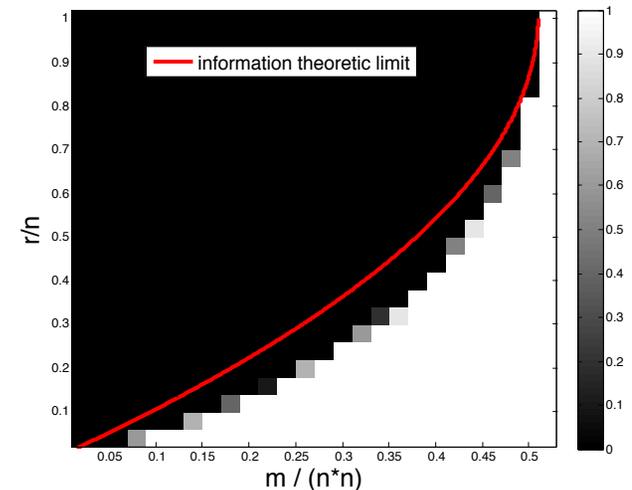
Near-Optimal Covariance Estimation

Theorem 1 (Chen, Chi and Goldsmith). *Consider the sub-Gaussian sampling model, then with probability exceeding $1 - \exp(-c_1 m)$, the solution $\hat{\Sigma}$ satisfies*

$$\|\hat{\Sigma} - \Sigma\|_F \leq \underbrace{C_1 \frac{\|\Sigma - \Sigma_r\|_*}{\sqrt{r}}}_{\text{due to imperfect structure}} + \underbrace{C_2 \frac{\epsilon}{m}}_{\text{due to noise}},$$

where Σ_r is the best rank- r approximation of Σ , provided that $m > c_0 nr$, where c_0, c_1, C_1 and C_2 are universal constants.

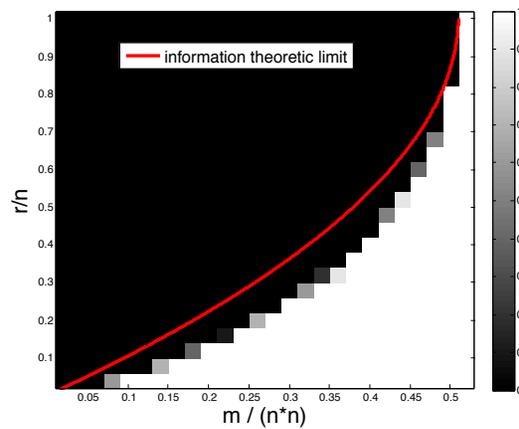
- **Exact** with $\Theta(nr)$ measurements;
- **Universal** for all low-rank matrices;
- **Robust** against approximate low-rankness and bounded noise;
- Results hold for *i.i.d. bilinear measurements* $a_i^T \Sigma b_i$ as well.



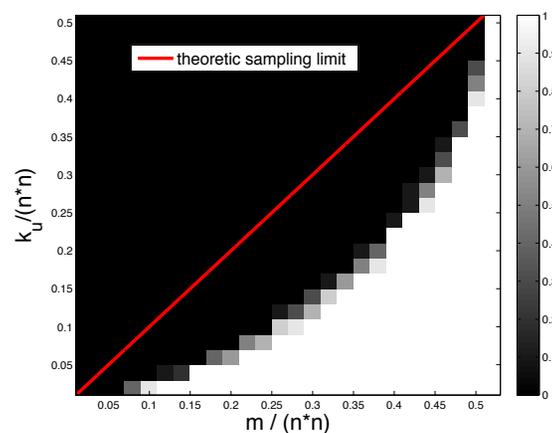
Phase Transition

Our covariance sketching scheme is

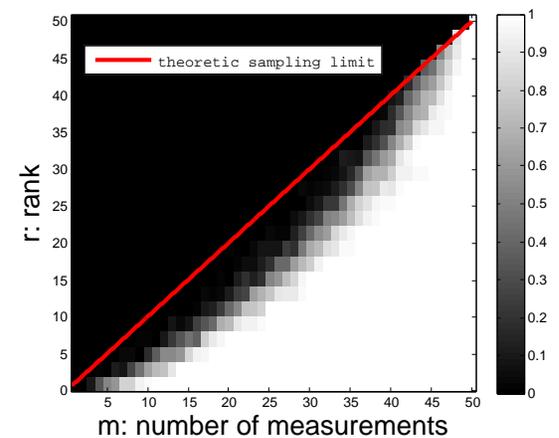
- universal
- works with sparse and Toeplitz low-rank covariance matrices as well;
- robust to noise and imperfect model assumptions



(a) Low-rank;



(b) Sparse;



(c) Toeplitz low-rank

Figure 1: Phase transition

Proof Ingredient: Mixed-Norm RIP

- **Restricted Isometry Property:** a powerful notion for compressed sensing

$$\forall \mathbf{X} \text{ in some class : } \quad \|\mathcal{B}(\mathbf{X})\|_2 \approx \|\mathbf{X}\|_F.$$

– Unfortunately, it does **NOT** hold for quadratic models.

- We proposed a **Mixed-norm Variant:** **RIP- ℓ_2/ℓ_1**

$$\forall \mathbf{X} \text{ in some class : } \quad \|\mathcal{B}(\mathbf{X})\|_1 \approx \|\mathbf{X}\|_F.$$

– does **NOT** hold for \mathcal{A} , but hold after \mathcal{A} is *debiased*:

$$\mathcal{B}_i(\Sigma) = \langle \Sigma, \mathbf{a}_{2i} \mathbf{a}_{2i}^T - \mathbf{a}_{2i+1} \mathbf{a}_{2i+1}^T \rangle$$

E. J. Candès, “The restricted isometry property and its implications for compressed sensing”. *Compte Rendus de l’Academie des Sciences, Paris, Serie I*, 346 589–592.

Comparisons with Related Work

- **low-rank/sparse matrix recovery** with linear measurements:

$$\mathbf{y} = \mathcal{A}(\Sigma)$$

- compressed sensing: $\mathcal{A}_i(\Sigma) = \text{Tr}(\mathbf{M}_i \Sigma)$
- matrix completion: $\mathcal{A}_i(\Sigma) = (\Sigma)_{i_1 i_2}$

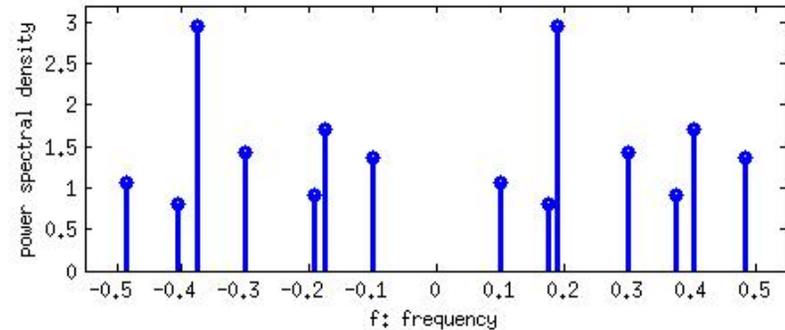
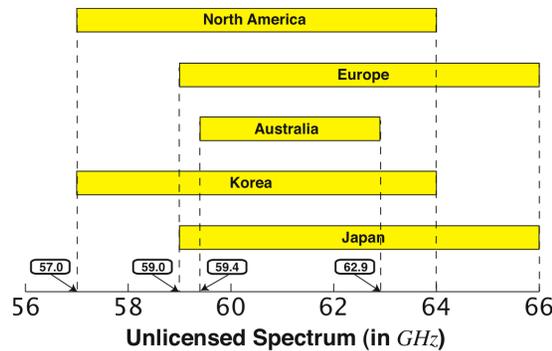
- **Sketching sparse matrices [Dasarathy et.al.]**:

$$\mathbf{Y} = \mathbf{A} \Sigma \mathbf{A}^T$$

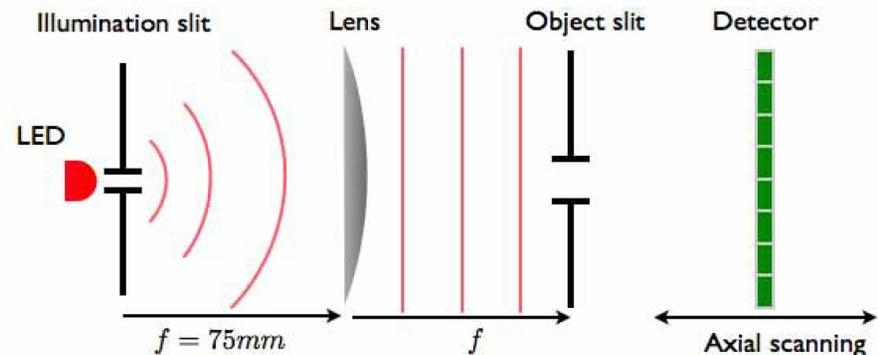
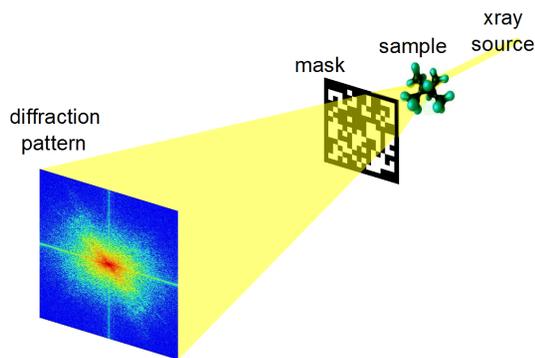
- $\mathbf{A} \in \mathbb{R}^{m \times n}$ with $m < n$, cannot estimate low-rank models;
 - no universal guarantees over sparse models;
- **Phaselift**: recover $\mathbf{x} \in \mathbb{C}^n$ from $\{|\langle \mathbf{a}_i, \mathbf{x} \rangle|^2\}_{i=1}^m$.
 - our algorithm is the same form of Phaselift when rank is one.
 - our algorithm extends the best performance guarantees of Phaselift to $O(n)$ *sub-Gaussian* measurements.

Other Applications of Quadratic Sensing

- **Energy measurements** are often more reliable with high-frequency applications for estimating *power spectral density*.



- **Quadratic measurements** arise in practical applications such as phase retrieval and phase space tomography.



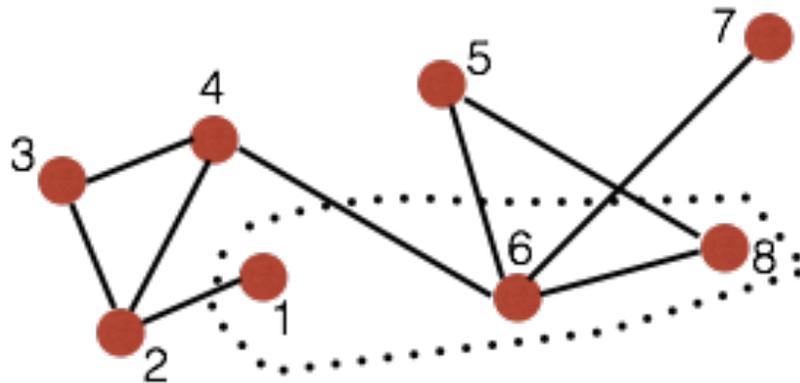
E. J. Candès, Y. C. Eldar, T. Strohmer and V. Voroninski, "Phase retrieval via matrix completion," SIAM J. on Imaging Sciences.

L. Tian, J. Lee, S. Oh, and G. Barbastathis, "Experimental compressive phase space tomography," Opt. Express.

Graph Sketching

Consider an undirected graph with bounded degree d and number of nodes n with adjacency matrix \mathbf{A} .

- Define the i th sketching vector $\mathbf{x}_i \in \{0, 1\}^n$ as composed of i.i.d. Bernoulli entries, then the sketch $y_i = \mathbf{x}_i^T \mathbf{A} \mathbf{x}_i$ amounts to counting twice the number of edges in random subgraphs.
- Example: $\mathbf{x} = [1, 0, 0, 0, 0, 1, 0, 1]^T$ whose support is $\mathcal{I} = \{1, 6, 8\}$, then $\mathbf{x}^T \mathbf{A} \mathbf{x} = 2$.



- Our results implies the graph can be perfectly reconstructed from $O(nd)$ quadratic sketches.

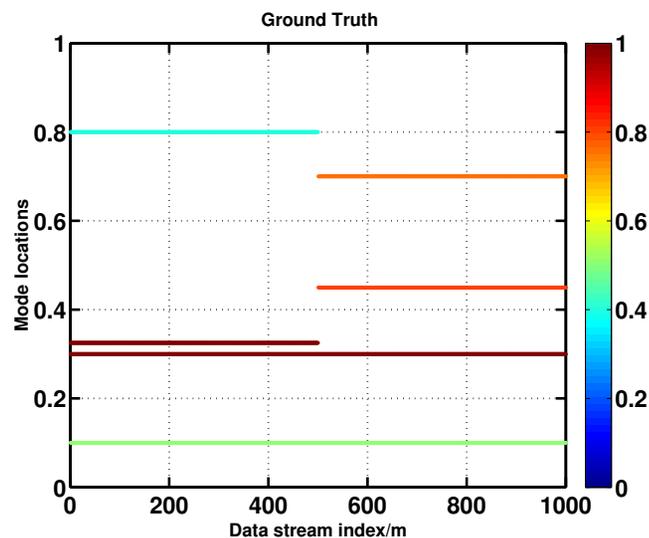
Covariance Tracking for DOA Estimation

The aggregation step in the sketching scheme can be easily implemented in an online manner to allow tracking.

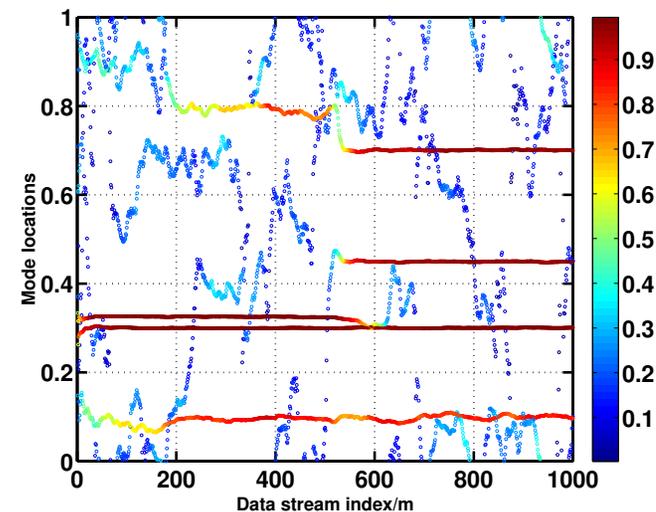
- **Aggregation:** The aggregates $y_{i,T}$ can be modified with a discounting factor λ for tracking:

$$y_{i,T} = \lambda y_{i,T-1} + |\langle \mathbf{a}_i, \mathbf{x}_T \rangle|^2$$

- **Estimation:** replace the trace minimization by a Projection onto Convex Sets (POCS) procedure ($n = 40, m = 600$).



(a) Ground Truth

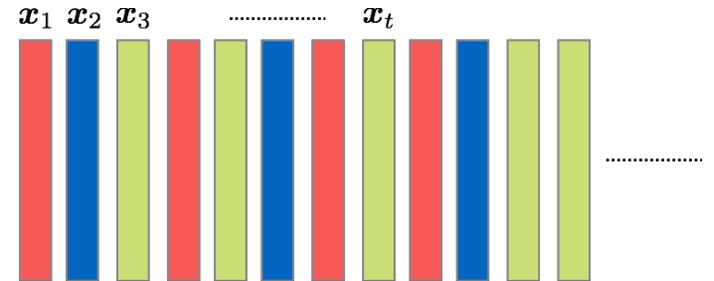


(b) Covariance Tracking

Can we estimate the covariance from Bits?

Assumption: Let $\Sigma = \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^H] = \mathbf{U} \mathbf{U}^H$ be a rank- r matrix with $\mathbf{U} \in \mathbb{C}^{n \times r}$.

One-Bit PCA: For each sensor $i = 1, \dots, m$:



- *randomly* select two sketching vectors $\mathbf{a}_i, \mathbf{b}_i \in \mathbb{C}^n$ with i.i.d. Gaussian entries;
- Sketch an arbitrary substream indexed by $\{\ell_t^i\}_{t=1}^T$ with two energy measurements $|\langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle|^2$, $|\langle \mathbf{b}_i, \mathbf{x}_{\ell_t^i} \rangle|^2$, and transmit a **binary bit** indicating the energy comparison outcome to the fusion center:

$$y_{i,T} = \text{sign} \left(\frac{1}{T} \sum_{t=1}^T |\langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle|^2 - \frac{1}{T} \sum_{t=1}^T |\langle \mathbf{b}_i, \mathbf{x}_{\ell_t^i} \rangle|^2 \right)$$

- **Estimation:** The fusion center recovers the principal components $\hat{\mathbf{U}} \in \mathbb{R}^{n \times r}$ by computing the top r eigenvectors of the surrogate matrix:

$$\mathbf{J}_m = \frac{1}{m} \sum_{i=1}^m y_{i,T} (\mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H).$$

Bit Comparisons are Robust

- With finite samples, the *numerical* energy difference measures the sample covariance Σ_T :

$$z_{i,T} = \frac{1}{T} \sum_{t=1}^T \left| \langle \mathbf{a}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2 - \frac{1}{T} \sum_{t=1}^T \left| \langle \mathbf{b}_i, \mathbf{x}_{\ell_t^i} \rangle \right|^2 = \langle \Sigma_T, \mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H \rangle.$$

The discrepancy $z_{i,T} - z_i = \langle \Sigma_T - \Sigma, \mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H \rangle \neq 0$.

- The *ordinal* energy difference measures the exact covariance Σ with high probability as soon as T is not too small:

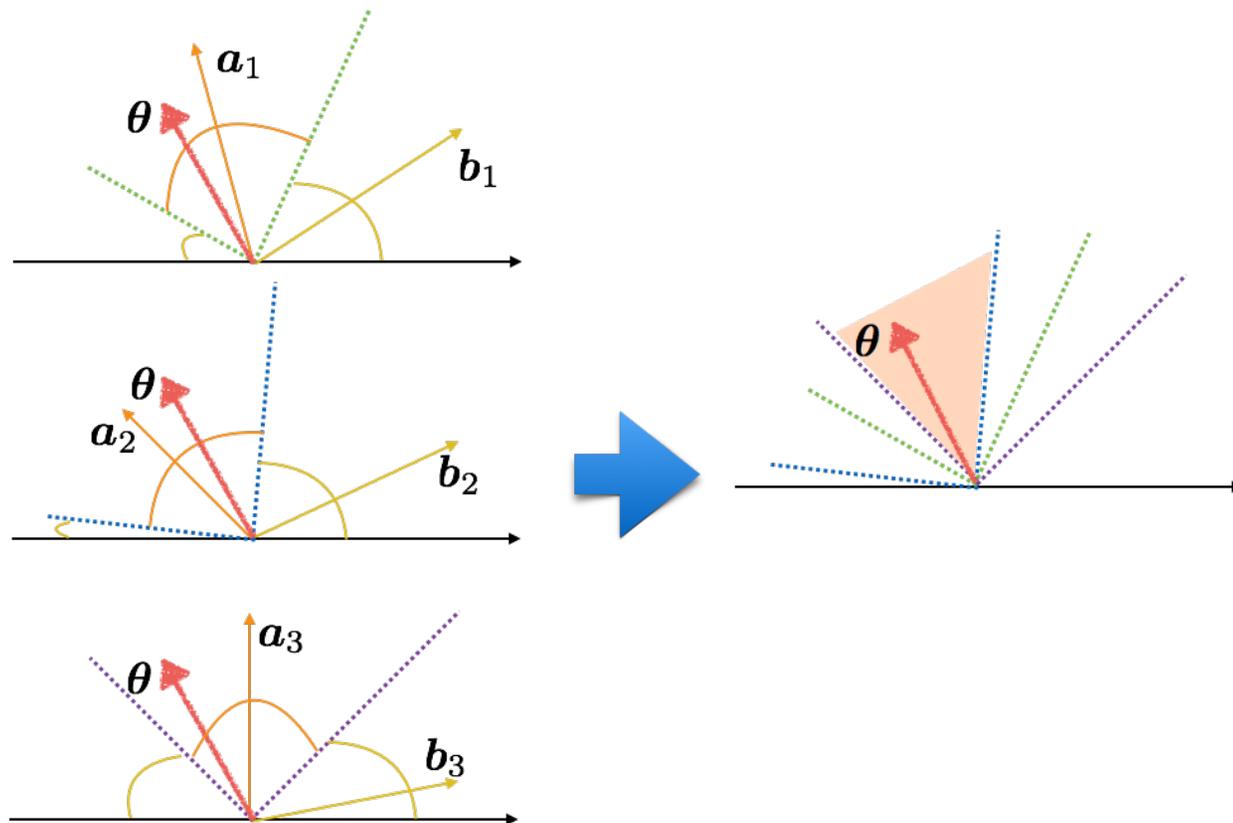
$$y_{i,T} = \text{sign} \left(\langle \Sigma_T, \mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H \rangle \right) = \text{sign} \left(\langle \Sigma, \mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H \rangle \right) = y_i.$$

Theorem 2 (Chi 2014). *Let $\mathbf{x}_t \sim \mathcal{N}(\mathbf{0}, \Sigma)$. Let $0 < \delta \leq 1$, then with probability at least $1 - \delta$ all bit measurements are exact, given that the number of samples observed by each sensor satisfies $T > c \text{Tr}(\Sigma) / \|\Sigma\|_F \log\left(\frac{m}{\delta}\right)$ for some sufficiently large constant c .*

One-Bit PCA: Why does it work?

Consider a rank-one example $\Sigma = \theta\theta^H$ with the eigenvector $\theta \in \mathbb{C}^2$:

- Each bit $y_i = \text{sign}(|\langle a_i, \theta \rangle|^2 - |\langle b_i, \theta \rangle|^2)$ selects the halfspace towards the direction with a smaller angle with either a_i or b_i .
- With enough bit measurements, we can trap the eigenvector θ accurately up to a sign difference.



One-Bit PCA: Performance Guarantee

Conditioned on that the bit measurements are exact, the principal subspace of \mathbf{J}_m agrees with \mathbf{U} with high probability given m is sufficiently large.

Theorem 3 (Chi 2014). Denote $\mathbf{U} \in \mathbb{C}^{n \times r}$ as the principal subspace of $\mathbf{\Sigma}$ and $\hat{\mathbf{U}}$ as the principal subspace of $\mathbf{J}_m = \frac{1}{m} \sum_{i=1}^m y_i (\mathbf{a}_i \mathbf{a}_i^H - \mathbf{b}_i \mathbf{b}_i^H)$. Let $0 < \delta < 1$, then with probability at least $1 - \delta$, there exists an $r \times r$ orthogonal matrix \mathbf{Q} such that

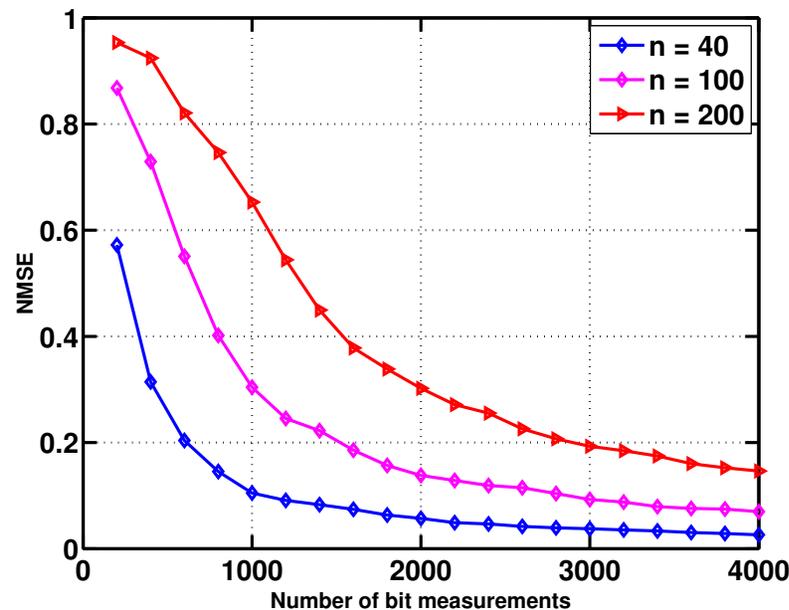
$$\|\hat{\mathbf{U}} - \mathbf{U}\mathbf{Q}\|_F \leq c_1 \sqrt{\frac{nr^2}{m} \log\left(\frac{2n}{\delta}\right)}$$

for all rank- r matrices $\mathbf{\Sigma}$, where c_1 is an absolute constant depending on r .

- The subspace estimate is accurate as soon as $m = \Theta(nr^2 \log n)$ which is near-optimal as the subspace requires at least n measurements.
- Not an exact recovery guarantee.

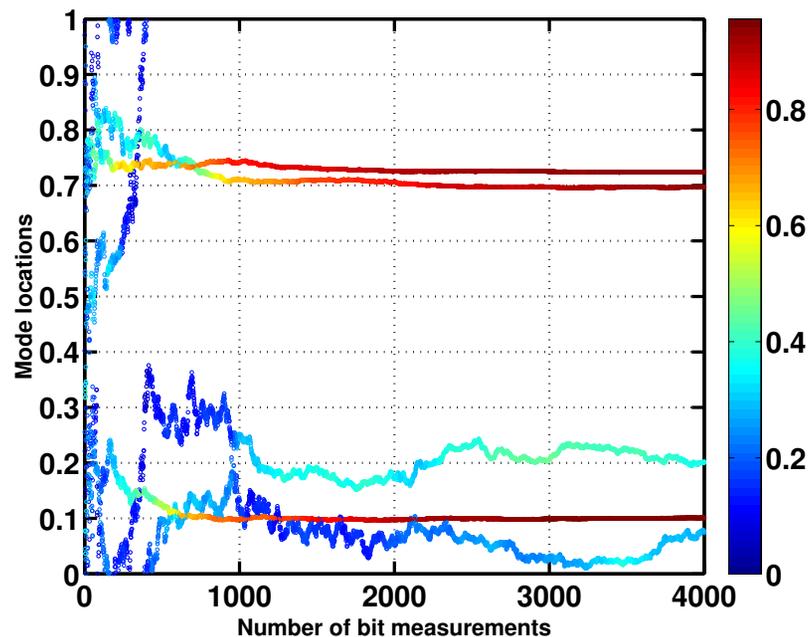
How many bits do we need?

- We generate the covariance matrix as $\Sigma = \mathbf{X}\mathbf{X}^T$, where $\mathbf{X} \in \mathbb{R}^{n \times 3}$ is composed of standard Gaussian entries. The sketching vectors a_i 's and b_i 's are also generated with standard Gaussian entries.
- The estimate $\hat{\mathbf{X}}$ is calculated via computing the top eigenvectors of \mathbf{J}_m .
- The error metric is calculated as $\|P_{\hat{\mathbf{X}}^\perp} \mathbf{X}\|_F^2 / \|\mathbf{X}\|_F^2$.



Online DOA Estimation with Bit Measurements

- The covariance matrix Σ is a low-rank Toeplitz PSD matrix with $n = 40$ and $r = 3$. The set of modes is $\mathcal{F} = [0.1, 0.7, 0.725]$ (notice the last two modes are separated by the Rayleigh limit $1/n$), and their variance is $\sigma^2 = 1$.
- The subspace is updated online as new bit measurements arrive sequentially; and ESPRIT is applied to estimate 5 modes using the subspace estimate.



1000 bits successfully distinguish two close modes separated by the Rayleigh limit.

Recap: Comparing the Two Schemes

Covariance Sketching with Real Measurements:

- *Near-Optimal Sample Complexity* for a variety of covariance structures;
- *Energy* measurements are easier to obtain;
- Requires a *noise estimate* to performance the algorithm with finite sample size and additive noise;
- The convex optimization might still be computationally expensive;

Covariance Sketching with One-bit Measurements:

- Communication overhead is minimized with bit measurements;
- Simple algorithm via computing the top eigenvectors;
- Robust measurement with respect to (possibly heterogeneous) noise.

Summary and Future Directions

Summary:

- Covariance estimation is possible without observing and reconstructing the whole data stream.
- The sensing and estimation procedure can be jointly designed to minimize complexity by leveraging the low-dimensional structures of data.
- Many potential applications in network traffic monitoring, video surveillance, and covariance estimation in privacy-aware and crowdsourcing environments.

Future Directions:

- Quality-Quantity-Computation Complexity Trade-offs for statistical inference.

List of Related Publications

- Y. Chen, Y. Chi and A. J. Goldsmith, “Universal and Robust Covariance Estimation via Convex Programming,” in International Symposium on Information Theory (ISIT), Honolulu, HI, Jun. 2014.
- Y. Chen, Y. Chi and A. J. Goldsmith, “Estimation of Simultaneously Structured Covariance Matrices from Quadratic Measurements,” in International Conference on Acoustics, Speech, and Signal Processing (ICASSP), Florence, Italy, May 2014.
- Y. Chen, Y. Chi, and A. J. Goldsmith, “Exact and Stable Covariance Estimation from Quadratic Sampling via Convex Programming,” *IEEE Trans. on Information Theory*, vol. 61, no. 7, pp. 4034-4059, 2015.
- Y. Chi, “One-Bit Principal Subspace Estimation”, IEEE Global Conference on Signal and Information Processing (GlobalSIP), Atlanta, GA, Dec. 2014.
- Y. Jiang and Y. Chi, “Covariance Tracking from Sketches of Rapid Data Streams”, in *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Brisbane, Australia, April 2015.
- Y. Chi, “Compressive Graph Clustering via Semidefinite Programming,” *International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, Brisbane, Australia, April 2015.