

Greedy and Randomized Projection Methods

Jamie Haddock

UCLA CAM Colloquium,
October 30, 2019

Computational and Applied Mathematics
UCLA



joint with Jesús A. De Loera, Deanna Needell, and Anna Ma

<https://arxiv.org/abs/1802.03126> (BIT Numerical Mathematics 2019)

<https://arxiv.org/abs/1605.01418> (SISC 2017)

BIG Data

The big data opportunity



DellWorld[™]15



The Economist

FEBRUARY 27TH - MARCH 5TH 2010

Economist.com

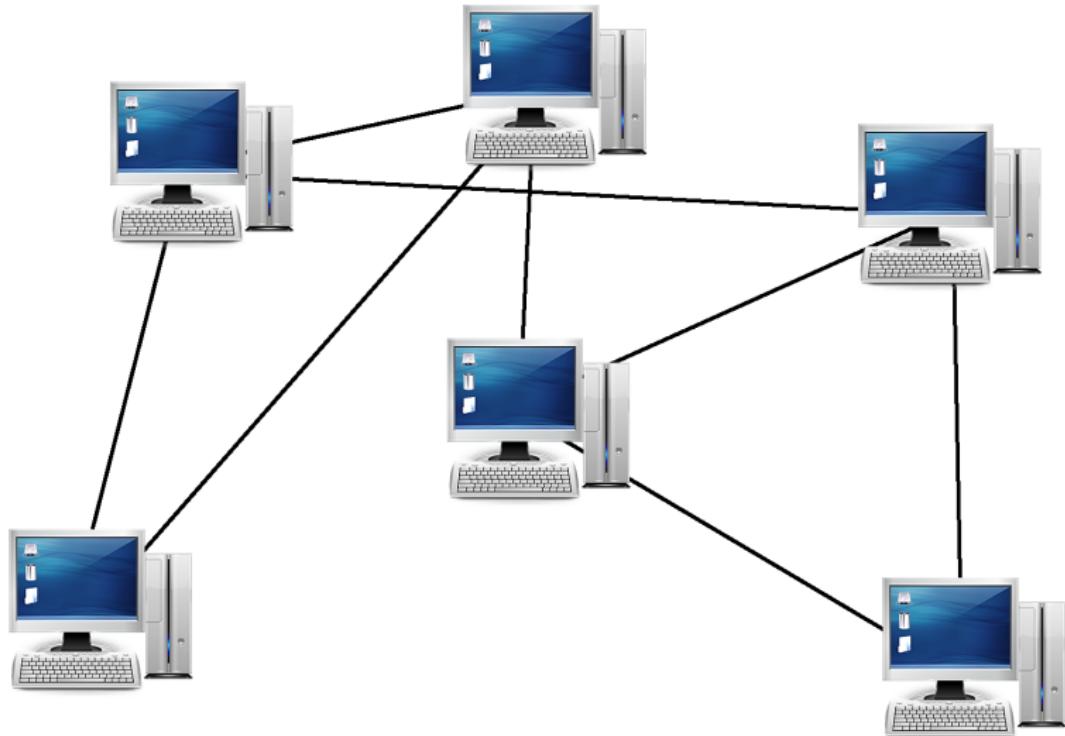
Obama the warrior
Misgoverning Argentina
The economic shift from West to East
Genetically modified crops blossom
The right to eat cats and dogs

The data deluge

AND HOW TO HANDLE IT: A 14-PAGE SPECIAL REPORT

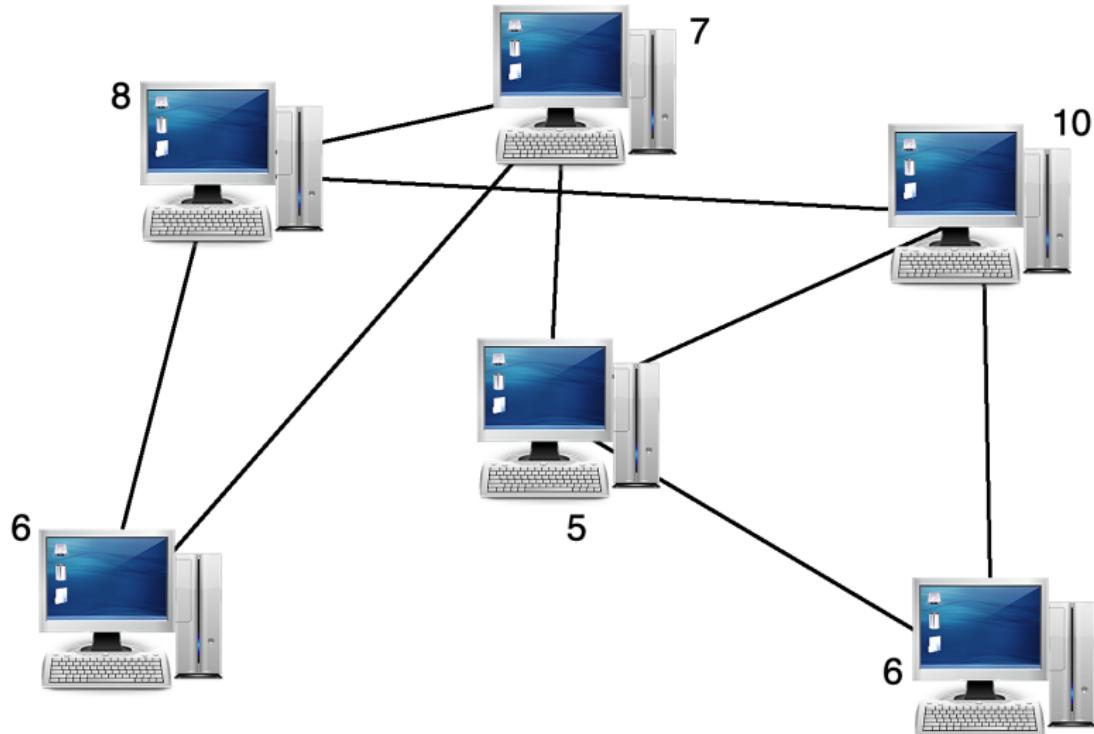
2

Motivation: Average Consensus



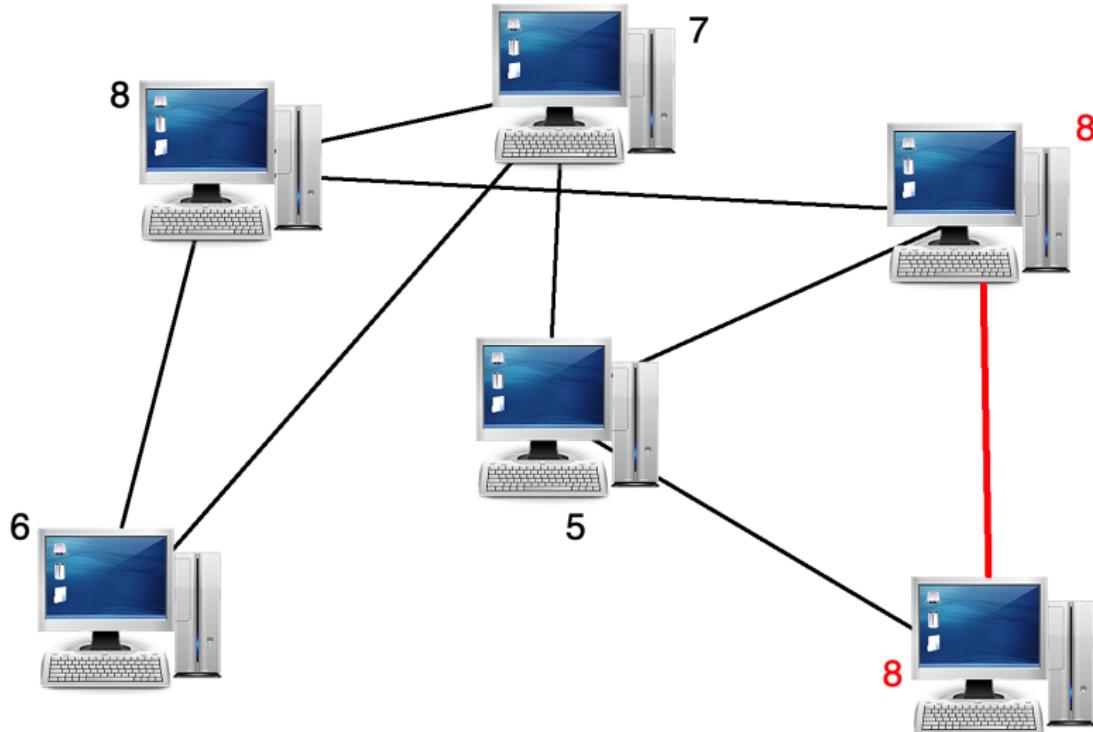
⁰ By Everaldo Coelho (YellowIcon); - All Crystal icons were posted by the author as LGPL on kde-look;; LGPL, <https://commons.wikimedia.org/w/index.php?curid=7288951>.

Motivation: Average Consensus



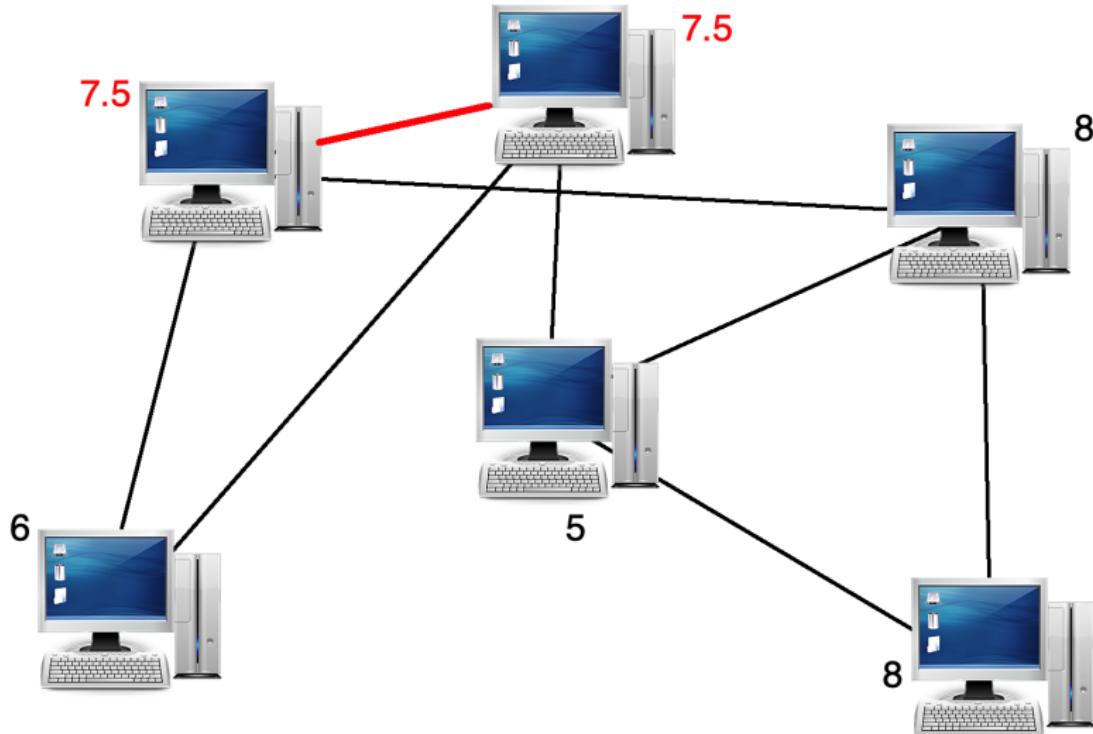
⁰ By Everaldo Coelho (YellowIcon); - All Crystal icons were posted by the author as LGPL on kde-look;; LGPL, <https://commons.wikimedia.org/w/index.php?curid=7288951>.

Motivation: Average Consensus



⁰ By Everaldo Coelho (YellowIcon); - All Crystal icons were posted by the author as LGPL on kde-look;; LGPL,
<https://commons.wikimedia.org/w/index.php?curid=7288951>.

Motivation: Average Consensus



0 By Everaldo Coelho (YellowIcon); - All Crystal icons were posted by the author as LGPL on kde-look;, LGPL,
<https://commons.wikimedia.org/w/index.php?curid=7288951>.

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_x \frac{1}{2} \|x - c\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} x = \mathbf{0}$$

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

Applications:

1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Applications:

▷ clock synchronization

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Applications:

- ▷ clock synchronization
- ▷ PageRank

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Applications:

- ▷ clock synchronization
- ▷ PageRank
- ▷ opinion formation

Motivation: Average Consensus

AC solution:

$$\operatorname{argmin}_{\mathbf{x}} \frac{1}{2} \|\mathbf{x} - \mathbf{c}\|^2 \text{ s.t. } \begin{bmatrix} 1 & 0 & -1 & 0 & \dots & 0 \\ 0 & 1 & -1 & 0 & \dots & 1 \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 1 & 0 & \dots & -1 \end{bmatrix} \mathbf{x} = \mathbf{0}$$

Gossip Method: Begin with $\mathbf{x}_0 = \mathbf{c} \in \mathbb{R}^{|\mathcal{V}|}$:

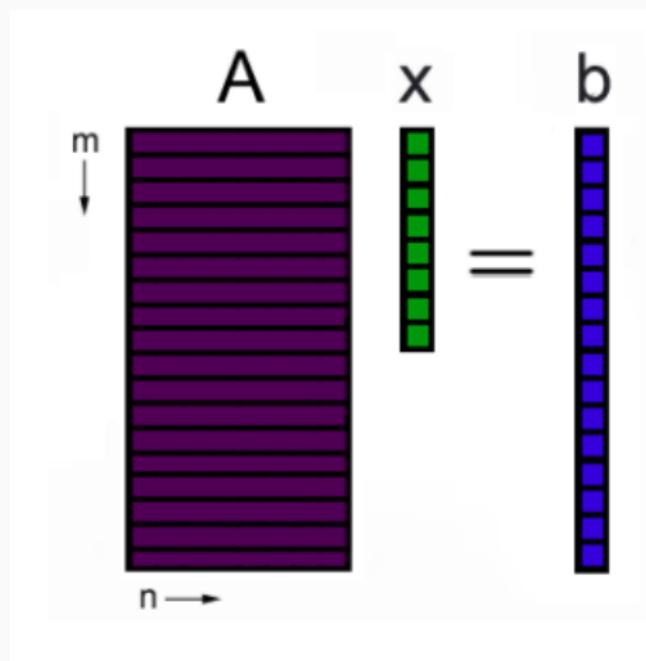
1. Choose $e_{ij} \in \mathcal{E}$.
2. Define $x_k^{(i)} = x_k^{(j)} = \frac{x_{k-1}^{(i)} + x_{k-1}^{(j)}}{2}$.
3. Repeat.

Applications:

- ▷ clock synchronization
- ▷ PageRank
- ▷ opinion formation
- ▷ blockchain technology

General Setup

We are interested in solving **highly overdetermined systems of equations (or inequalities)**, $Ax = b$ ($Ax \leq b$), where $A \in \mathbb{R}^{m \times n}$, $b \in \mathbb{R}^m$ and $m \gg n$. Rows are denoted a_i^T .



Iterative Projection Methods

If $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}\}$ is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method



Applications:

1. Tomography (Algebraic Reconstruction Technique)

Iterative Projection Methods

If $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}\}$ is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method
2. Motzkin's Method



Applications:

1. Tomography (Algebraic Reconstruction Technique)
2. Linear programming

Iterative Projection Methods

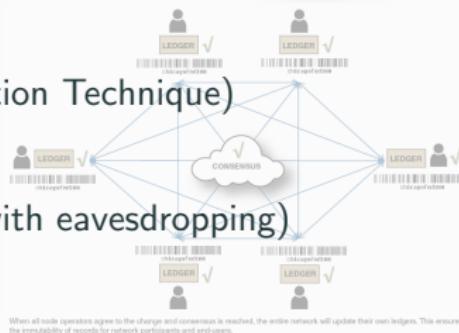
If $\{\mathbf{x} \in \mathbb{R}^n : A\mathbf{x} = \mathbf{b}\}$ is nonempty, these methods construct an **approximation** to a solution:

1. Randomized Kaczmarz Method
2. Motzkin's Method
3. Sampling Kaczmarz-Motzkin Methods (SKM)



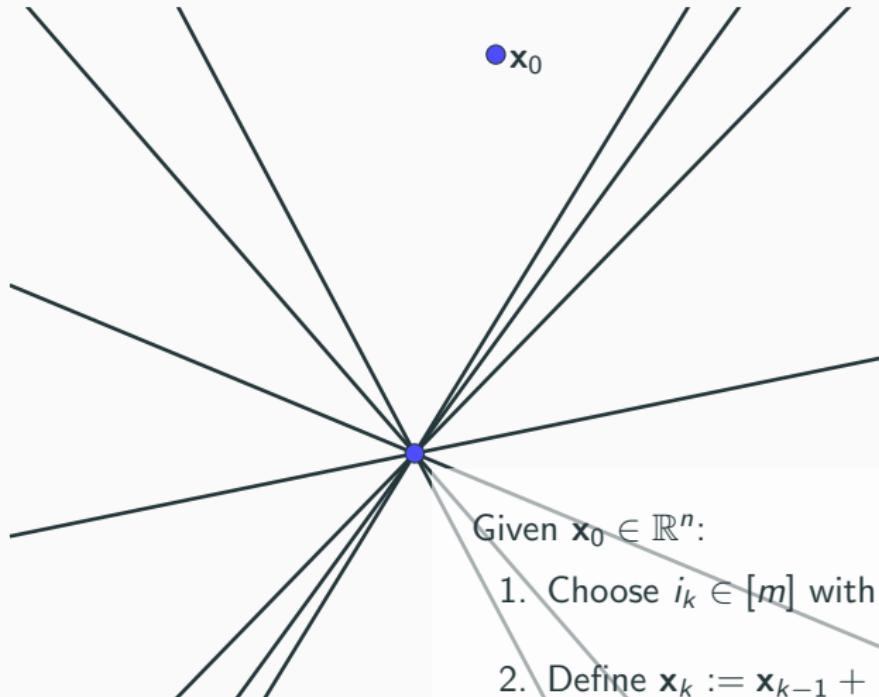
Applications:

1. Tomography (Algebraic Reconstruction Technique)
2. Linear programming
3. Average consensus (greedy gossip with eavesdropping)



When all node operators agree to the change and consensus is reached, the entire network will update their own ledgers. This ensures the immutability of records for network participants and end-users.

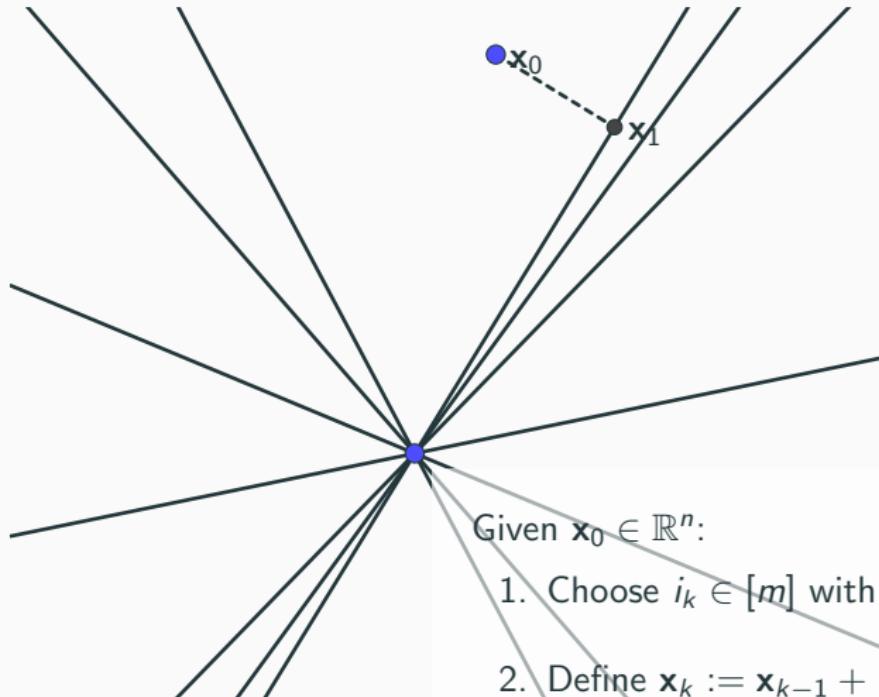
Kaczmarz Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $i_k \in [m]$ with probability $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$.
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
3. Repeat.

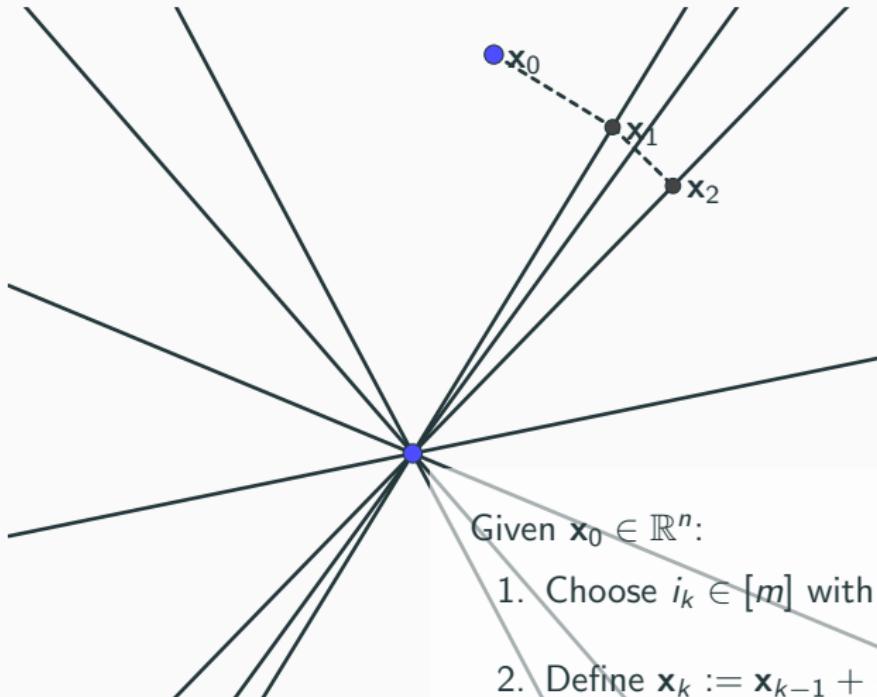
Kaczmarz Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $i_k \in [m]$ with probability $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$.
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
3. Repeat.

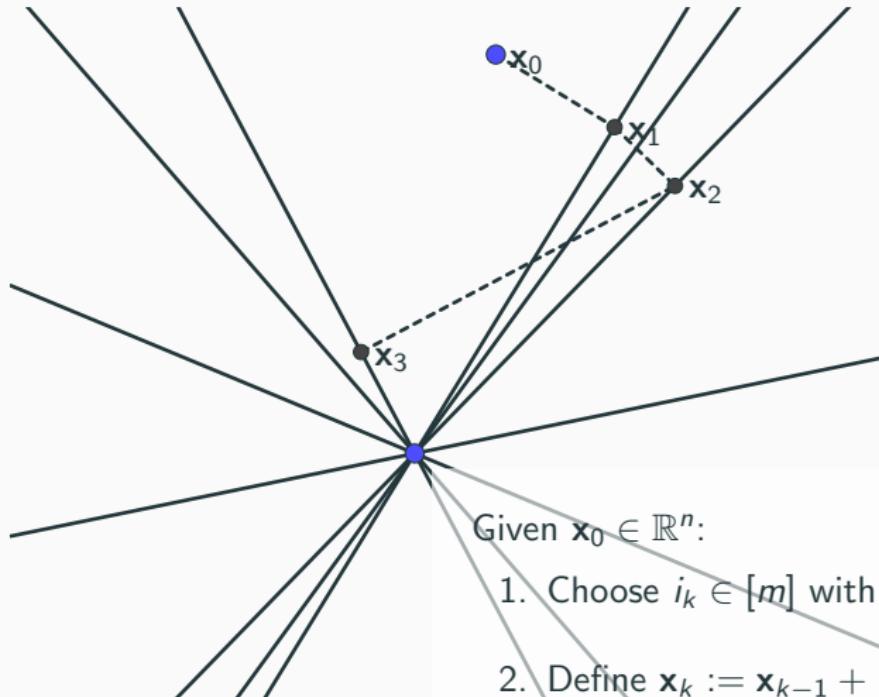
Kaczmarz Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $i_k \in [m]$ with probability $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$.
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
3. Repeat.

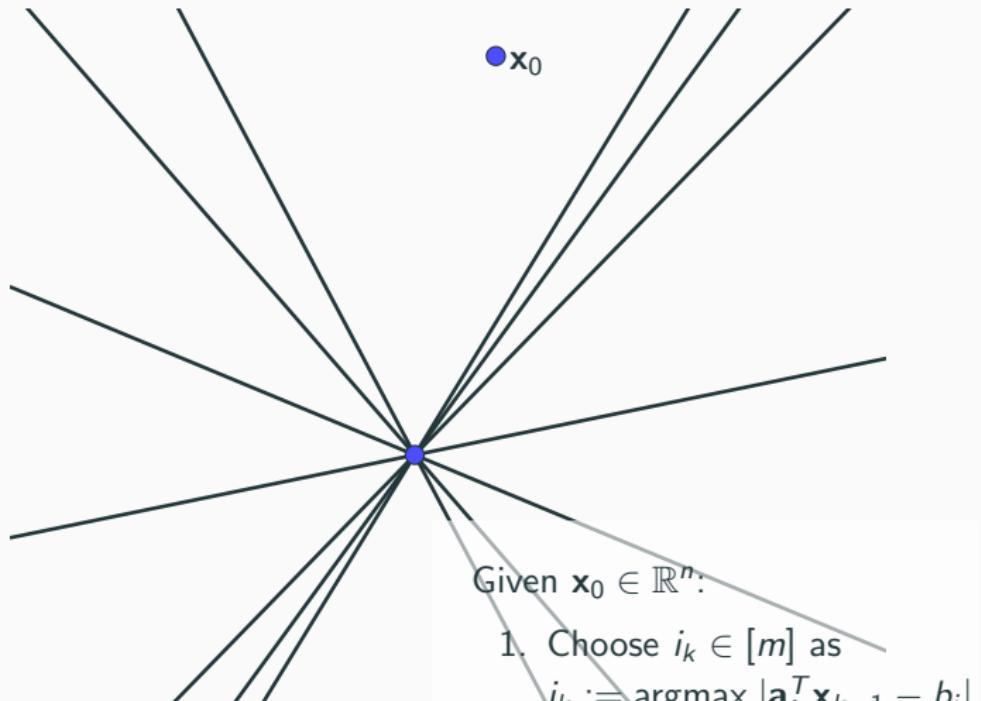
Kaczmarz Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

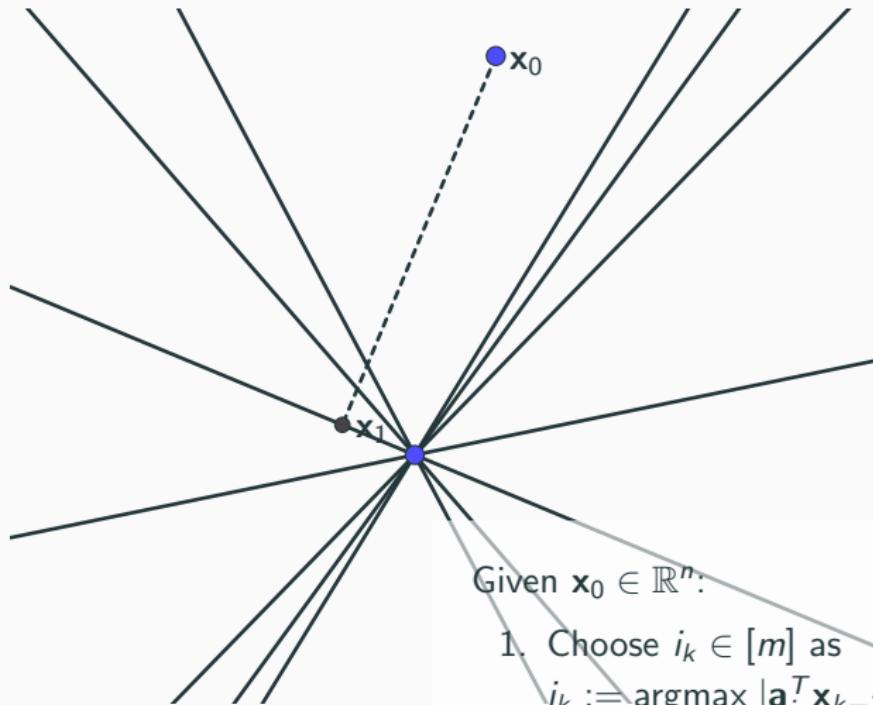
1. Choose $i_k \in [m]$ with probability $\frac{\|\mathbf{a}_{i_k}\|^2}{\|A\|_F^2}$.
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
3. Repeat.

Motzkin's Method



1. Choose $i_k \in [m]$ as
$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$
3. Repeat.

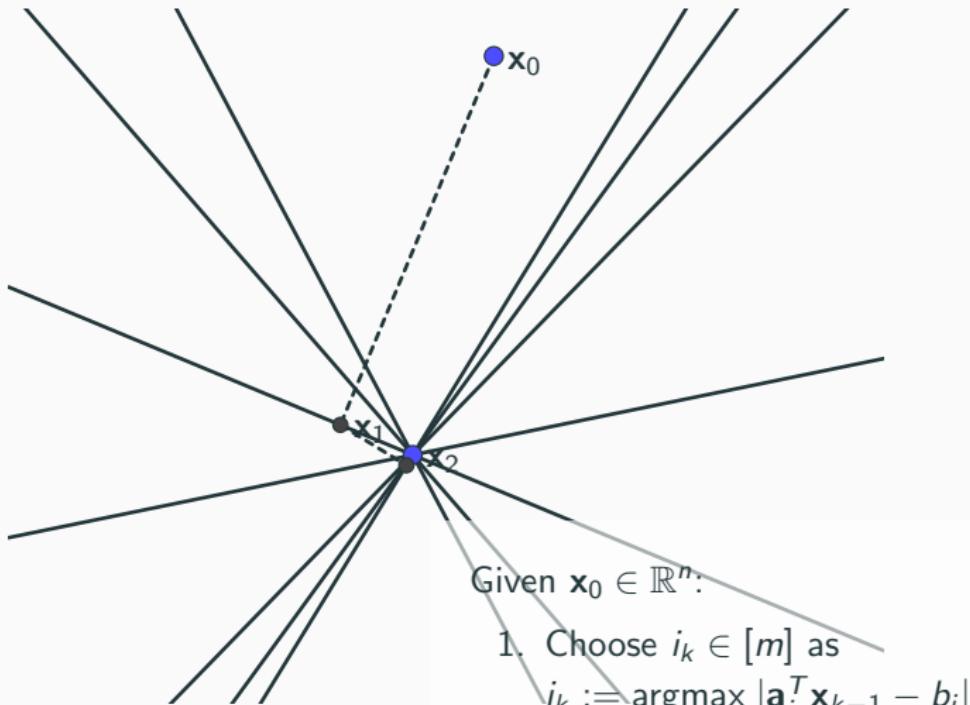
Motzkin's Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $i_k \in [m]$ as
$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$
3. Repeat.

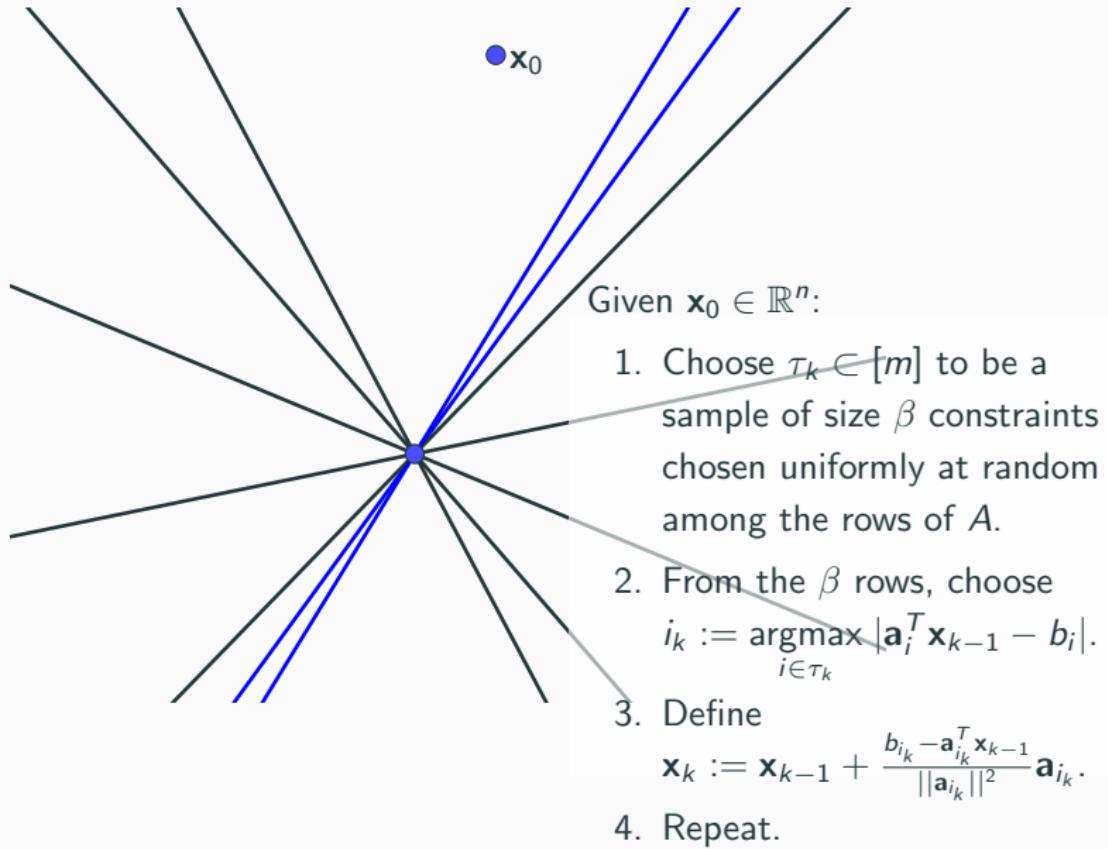
Motzkin's Method



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $i_k \in [m]$ as
$$i_k := \operatorname{argmax}_{i \in [m]} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|.$$
2. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$
3. Repeat.

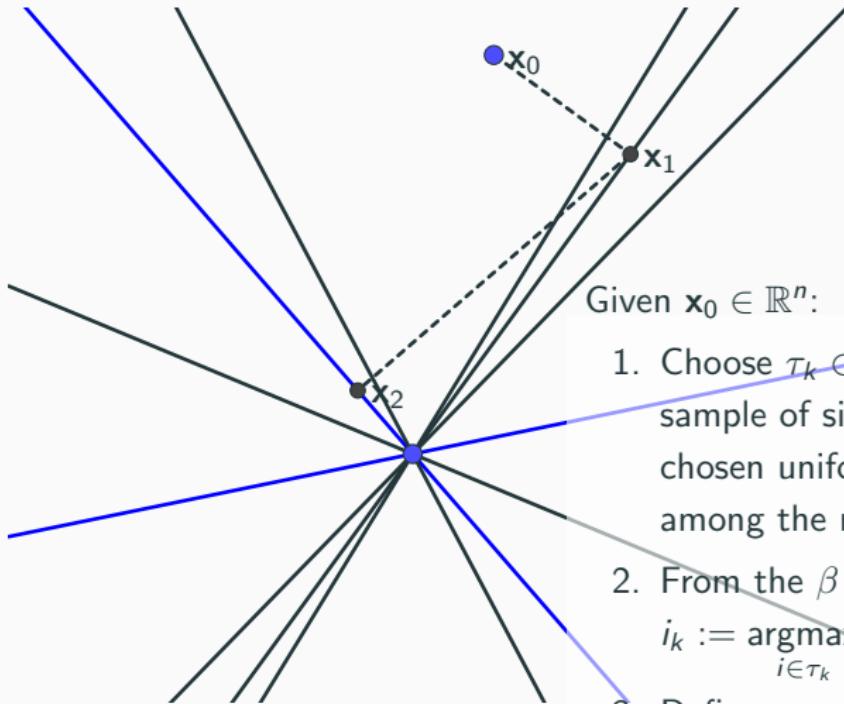
Our Hybrid Method (SKM)

- 
- Given $\mathbf{x}_0 \in \mathbb{R}^n$:
1. Choose $\tau_k \subset [m]$ to be a sample of size β constraints chosen uniformly at random among the rows of A .
 2. From the β rows, choose $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$.
 3. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
 4. Repeat.

Our Hybrid Method (SKM)

-
- Given $\mathbf{x}_0 \in \mathbb{R}^n$:
1. Choose $\tau_k \subset [m]$ to be a sample of size β constraints chosen uniformly at random among the rows of A .
 2. From the β rows, choose $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$.
 3. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
 4. Repeat.

Our Hybrid Method (SKM)



Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $\tau_k \subset [m]$ to be a sample of size β constraints chosen uniformly at random among the rows of A .
2. From the β rows, choose $i_k := \operatorname{argmax}_{i \in \tau_k} |\mathbf{a}_i^T \mathbf{x}_{k-1} - b_i|$.
3. Define
$$\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}.$$
4. Repeat.

Glimpse of HUGE Body of Literature

RK: [Strohmer-Vershynin '09], [Needell-Srebro-Ward '16]

Greedy: [Censor '81], [Nutini et al '16], [Bai-Wu '18], [Du-Gao '19]

Accel.: [Hanke-Niehamer '90], [Liu-Wright '16], [Morshed-Islam '19]

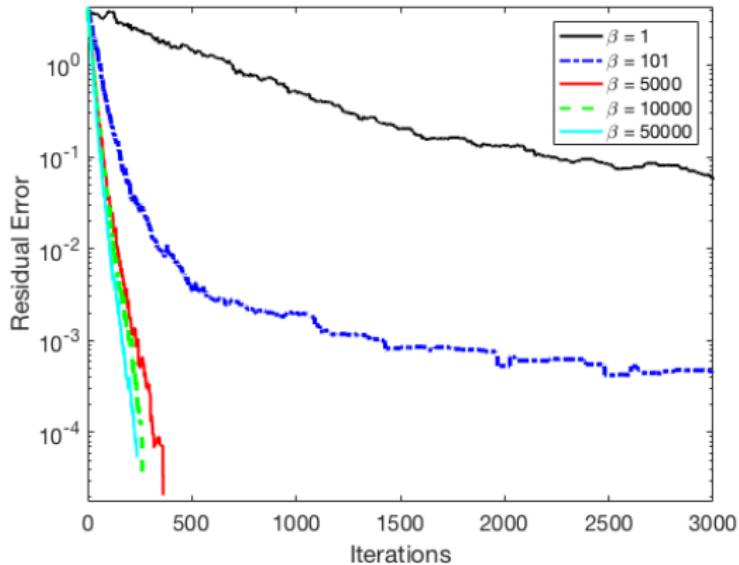
Block: [Popa et al '12], [Needell-Tropp '14], [Needell-Zhao-Zouzias '15],

Sketching: [Gower-Richtarik '15], [Needell-Rebrova '19]

Phase retrieval: [Tan-Vershynin '17], [Jeong-Güntürk '17]

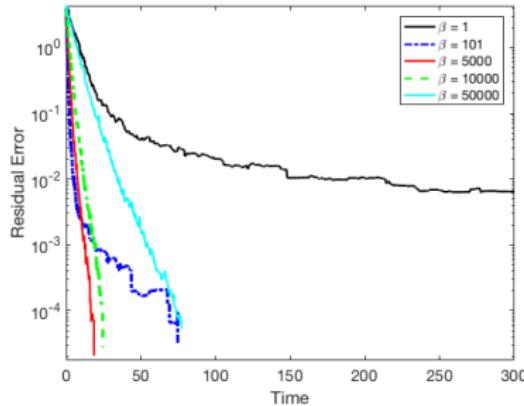
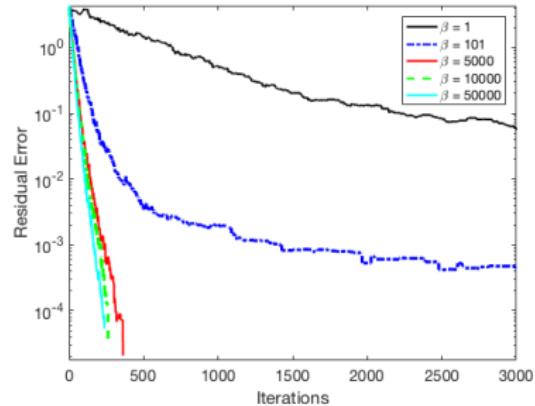
LP: [Motzkin-Schoenberg '54], [Agmon '54], [Goffin '80], [Chubanov '12]

Experimental Convergence



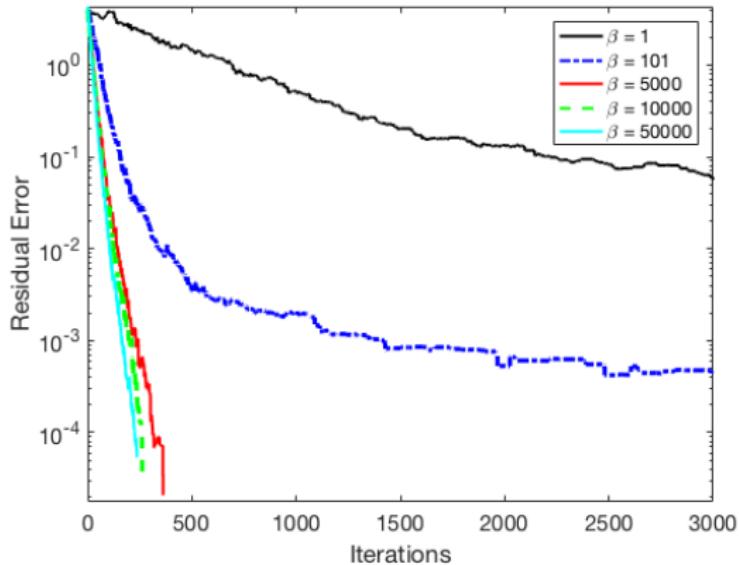
- ▷ β : sample size
- ▷ A is 50000×100 Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

Experimental Convergence



- ▷ β : sample size
- ▷ A is 50000×100 Gaussian matrix, consistent system
- ▷ ‘faster’ convergence for larger sample size

Experimental Convergence



- ▷ β : sample size
- ▷ A is 50000×100 Gaussian matrix, consistent system
- ▷ 'faster' convergence for larger sample size

Convergence Rates

Below are the convergence rates for the methods on a system, $\mathbf{A}\mathbf{x} = \mathbf{b}$, which is consistent with unique solution \mathbf{x} , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer - Vershynin '09):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

Convergence Rates

Below are the convergence rates for the methods on a system, $\mathbf{A}\mathbf{x} = \mathbf{b}$, which is consistent with unique solution \mathbf{x} , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer - Vershynin '09):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

Convergence Rates

Below are the convergence rates for the methods on a system, $\mathbf{A}\mathbf{x} = \mathbf{b}$, which is consistent with unique solution \mathbf{x} , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer - Vershynin '09):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ SKM (DeLoera - H. - Needell '17):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

Convergence Rates

Below are the convergence rates for the methods on a system, $\mathbf{A}\mathbf{x} = \mathbf{b}$, which is consistent with unique solution \mathbf{x} , whose rows have been normalized to have unit norm.

- ▷ RK (Strohmer - Vershynin '09):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

- ▷ MM (Agmon '54):

$$\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

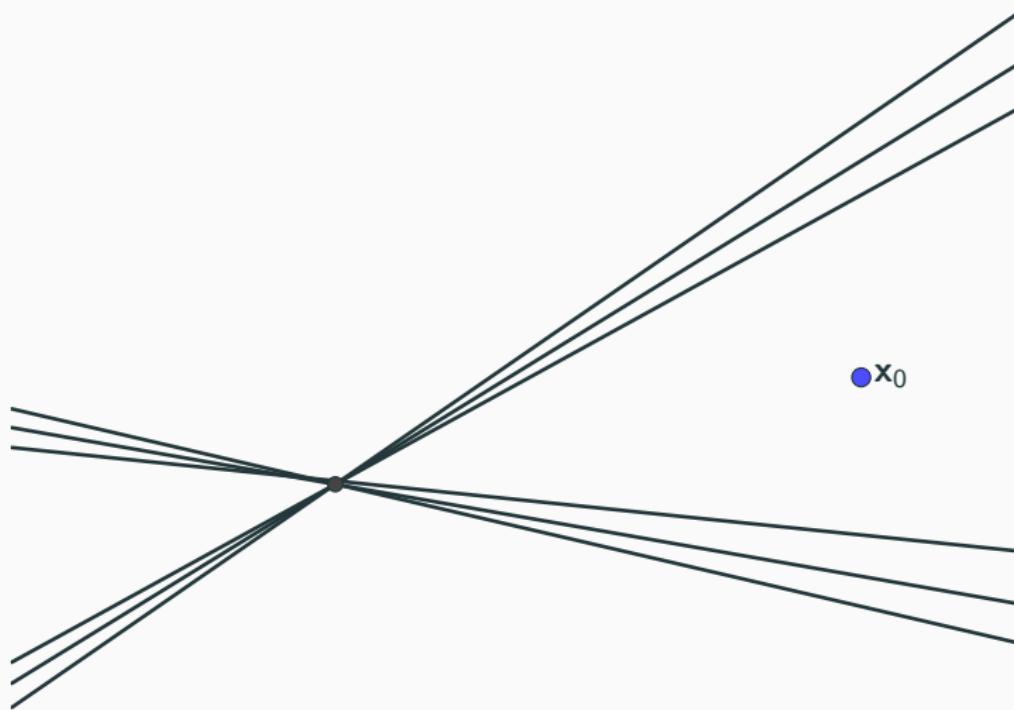
- ▷ SKM (DeLoera - H. - Needell '17):

$$\mathbb{E}\|\mathbf{x}_k - \mathbf{x}\|_2^2 \leq \left(1 - \frac{\sigma_{\min}^2(\mathbf{A})}{m}\right)^k \|\mathbf{x}_0 - \mathbf{x}\|_2^2$$

Why are these all the same?

A Pathological Example

Because.



Structure of the Residual

Several works have used sparsity of the residual to improve the convergence rate of greedy methods.

[De Loera, H., Needell '17], [Bai, Wu '18], [Du, Gao '19]

Structure of the Residual

Several works have used sparsity of the residual to improve the convergence rate of greedy methods.

[De Loera, H., Needell '17], [Bai, Wu '18], [Du, Gao '19]

However, not much sparsity can be expected in most cases. Instead, we'd like to use dynamic range of the residual to guarantee faster convergence.

$$\gamma_k := \frac{\|A\mathbf{x}_k - A\mathbf{x}\|^2}{\|A\mathbf{x}_k - A\mathbf{x}\|_\infty^2}$$

An Accelerated Convergence Rate

Theorem (H. - Needell '19)

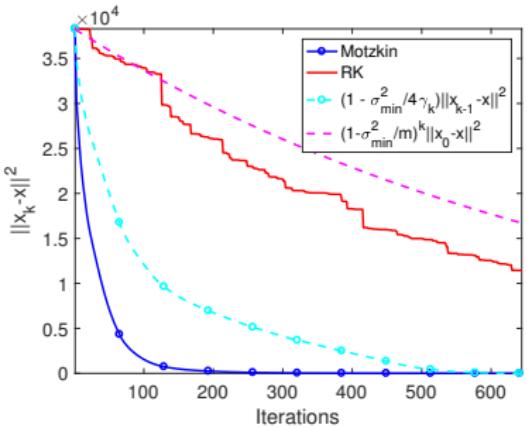
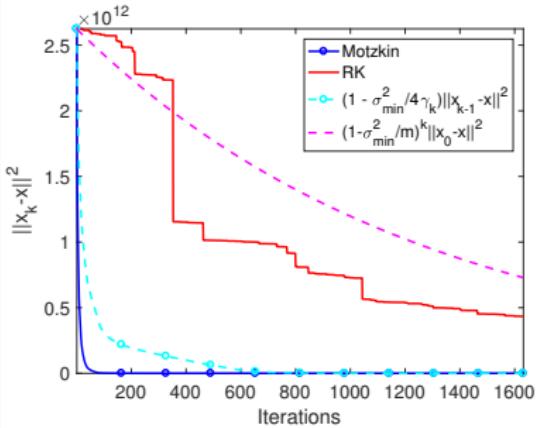
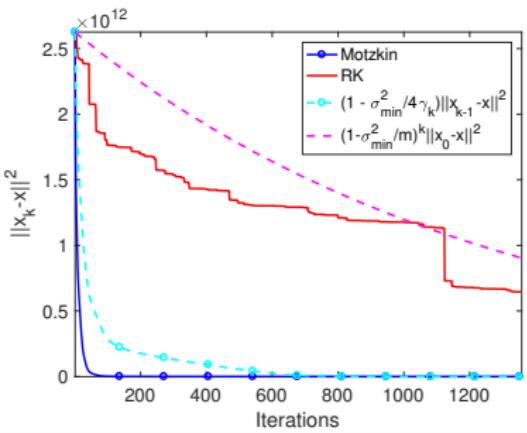
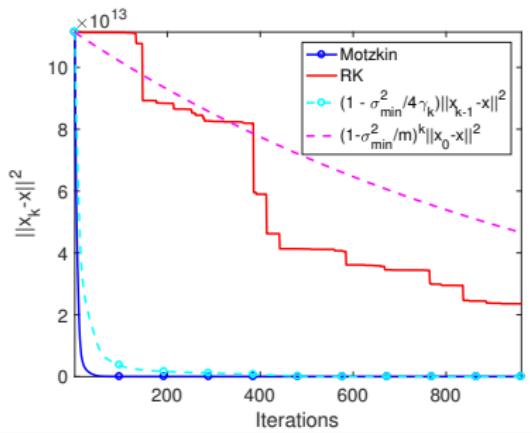
Let \mathbf{x} denote the solution of the consistent, normalized system $A\mathbf{x} = \mathbf{b}$. Motzkin's method exhibits the (possibly highly accelerated) convergence rate:

$$\|\mathbf{x}_T - \mathbf{x}\|^2 \leq \prod_{k=0}^{T-1} \left(1 - \frac{\sigma_{\min}^2(A)}{4\gamma_k}\right) \cdot \|\mathbf{x}_0 - \mathbf{x}\|^2$$

Here γ_k bounds the dynamic range of the k th residual, $\gamma_k := \frac{\|A\mathbf{x}_k - A\mathbf{x}\|^2}{\|A\mathbf{x}_k - A\mathbf{x}\|_\infty^2}$.

- ▷ improvement over previous result when $4\gamma_k < m$

Netlib LP Systems



Extending to SKM

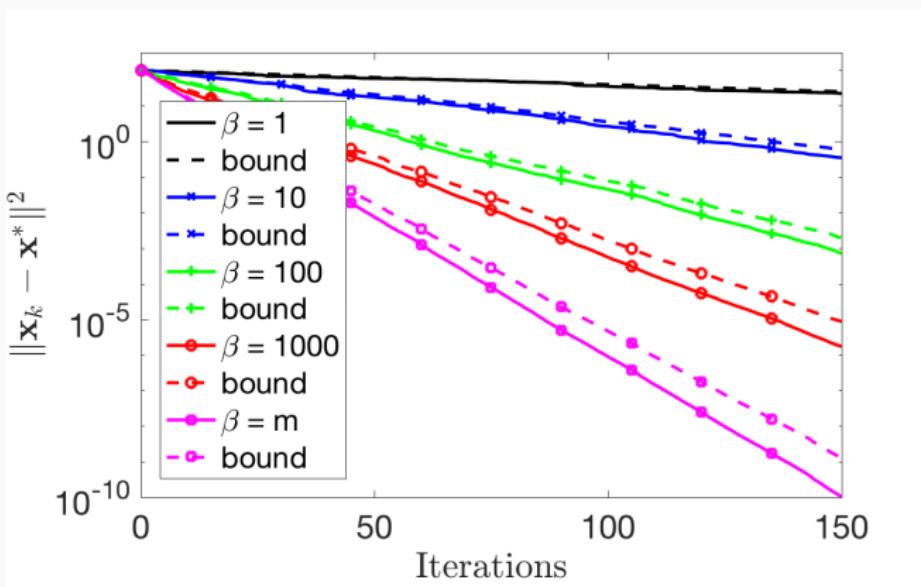
Corollary (H. - Ma 2019+)

Let A be normalized so $\|\mathbf{a}_i\|_2 = 1$ for all rows $i = 1, \dots, m$. If the system $A\mathbf{x} = \mathbf{b}$ is consistent with the unique solution \mathbf{x}^* then the SKM method converges at least linearly in expectation and the rate depends on the dynamic range of the random sample of rows of A , τ_j . Precisely, in the $j + 1$ st iteration of SKM, we have

$$\mathbb{E}_{\tau_j} \|\mathbf{x}_{j+1} - \mathbf{x}^*\|_2^2 \leq \left(1 - \frac{\beta \sigma_{\min}^2(A)}{\gamma_j m}\right) \|\mathbf{x}_j - \mathbf{x}^*\|_2^2$$

where $\gamma_j = \frac{\sum_{\tau_j \in \binom{[m]}{\beta}} \|A_{\tau_j} \mathbf{x}_j - \mathbf{b}_{\tau_j}\|_2^2}{\sum_{\tau_j \in \binom{[m]}{\beta}} \|A_{\tau_j} \mathbf{x}_j - \mathbf{b}_{\tau_j}\|_\infty^2}$.

Extending to SKM



- ▷ A is 50000×100 Gaussian matrix, consistent system
- ▷ bound uses dynamic range of sample of β rows

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$



$$A\mathbf{x}_k - \mathbf{b} = \mathbf{e}_i$$

$$A\mathbf{x}_k - \mathbf{b} = c\mathbf{1}$$

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$



$$A\mathbf{x}_k - \mathbf{b} = \mathbf{e}_i$$

"Best" case



$$A\mathbf{x}_k - \mathbf{b} = c\mathbf{1}$$

"Worst" case

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$

↗ ↘

$$\begin{array}{ll} A\mathbf{x}_k - \mathbf{b} = \mathbf{e}_i & A\mathbf{x}_k - \mathbf{b} = c\mathbf{1} \\ \text{"Best" case} & \text{"Worst" case} \end{array}$$

	Best Case	Worst Case	Previous Best	Previous Worst
MM	$1 - \sigma_{\min}^2(A)$		$1 - \frac{\sigma_{\min}^2(A)}{4}$	
SKM	$1 - \frac{\beta\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{m}$		$1 - \frac{\sigma_{\min}^2(A)}{m}$
RK	$1 - \frac{\sigma_{\min}^2(A)}{m}$		$1 - \frac{\sigma_{\min}^2(A)}{m}$	

Table 1: Contraction terms α such that $\mathbb{E}_{\tau_k} \|\mathbf{e}_k\|^2 \leq \alpha \|\mathbf{e}_{k-1}\|^2$.

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$

\nearrow \nwarrow

$$\begin{array}{ll} A\mathbf{x}_k - \mathbf{b} = \mathbf{e}_i & A\mathbf{x}_k - \mathbf{b} = c\mathbf{1} \\ \text{"Best" case} & \text{"Worst" case} \end{array}$$

	Best Case	Worst Case	Previous Best	Previous Worst
MM	$1 - \sigma_{\min}^2(A)$		$1 - \frac{\sigma_{\min}^2(A)}{4}$	
SKM	$1 - \frac{\beta\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{m}$
RK	$1 - \frac{\sigma_{\min}^2(A)}{m}$			

Table 1: Contraction terms α such that $\mathbb{E}_{\tau_k} \|\mathbf{e}_k\|^2 \leq \alpha \|\mathbf{e}_{k-1}\|^2$.

Nervous?

What can we say about γ_j ?

$$1 \leq \gamma_j \leq \beta$$

\nearrow \nwarrow

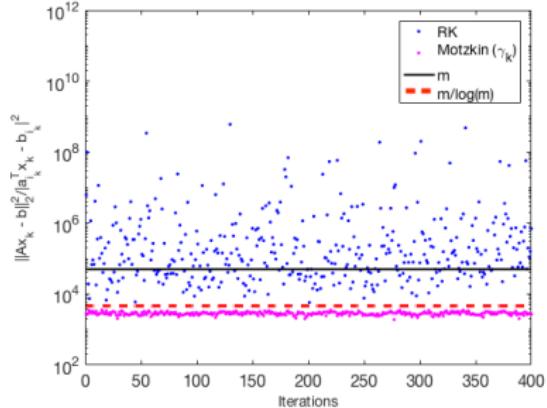
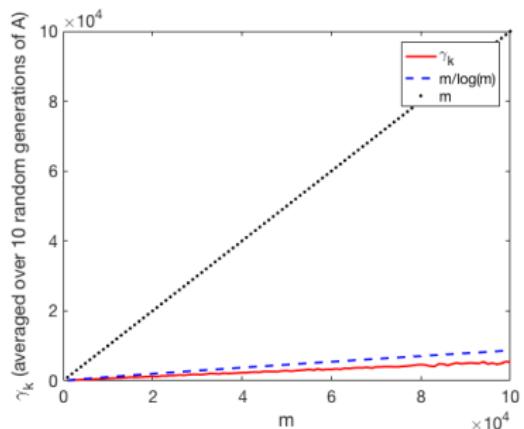
$$\begin{array}{ll} A\mathbf{x}_k - \mathbf{b} = \mathbf{e}_i & A\mathbf{x}_k - \mathbf{b} = c\mathbf{1} \\ \text{"Best" case} & \text{"Worst" case} \end{array}$$

	Best Case	Worst Case	Previous Best	Previous Worst
MM	$1 - \sigma_{\min}^2(A)$		$1 - \frac{\sigma_{\min}^2(A)}{4}$	
SKM	$1 - \frac{\beta\sigma_{\min}^2(A)}{m}$	$1 - \frac{\sigma_{\min}^2(A)}{m}$		$1 - \frac{\sigma_{\min}^2(A)}{m}$
RK	$1 - \frac{\sigma_{\min}^2(A)}{m}$		$1 - \frac{\sigma_{\min}^2(A)}{m}$	

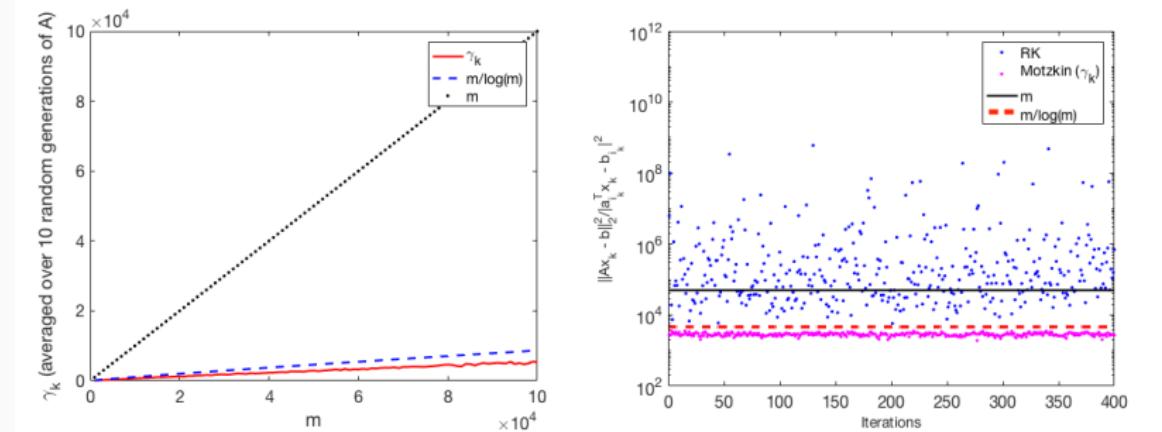
Table 1: Contraction terms α such that $\mathbb{E}_{\tau_k} \|\mathbf{e}_k\|^2 \leq \alpha \|\mathbf{e}_{k-1}\|^2$.

Nervous? $\gamma_k \geq \frac{\beta}{m} \sigma_{\min}^2(A)$ when A is row-normalized

γ_k : Gaussian systems

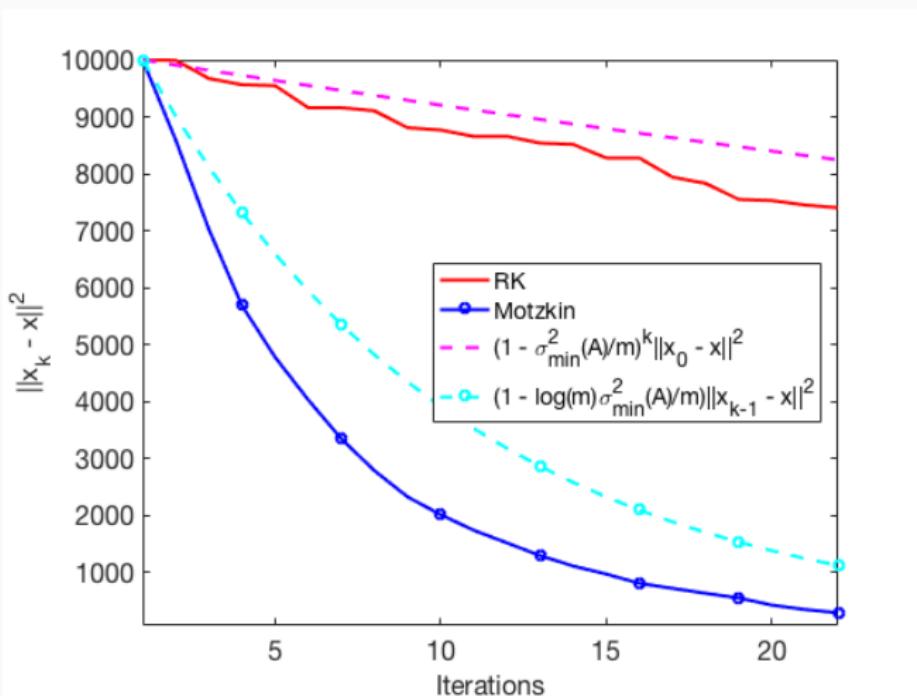


γ_k : Gaussian systems



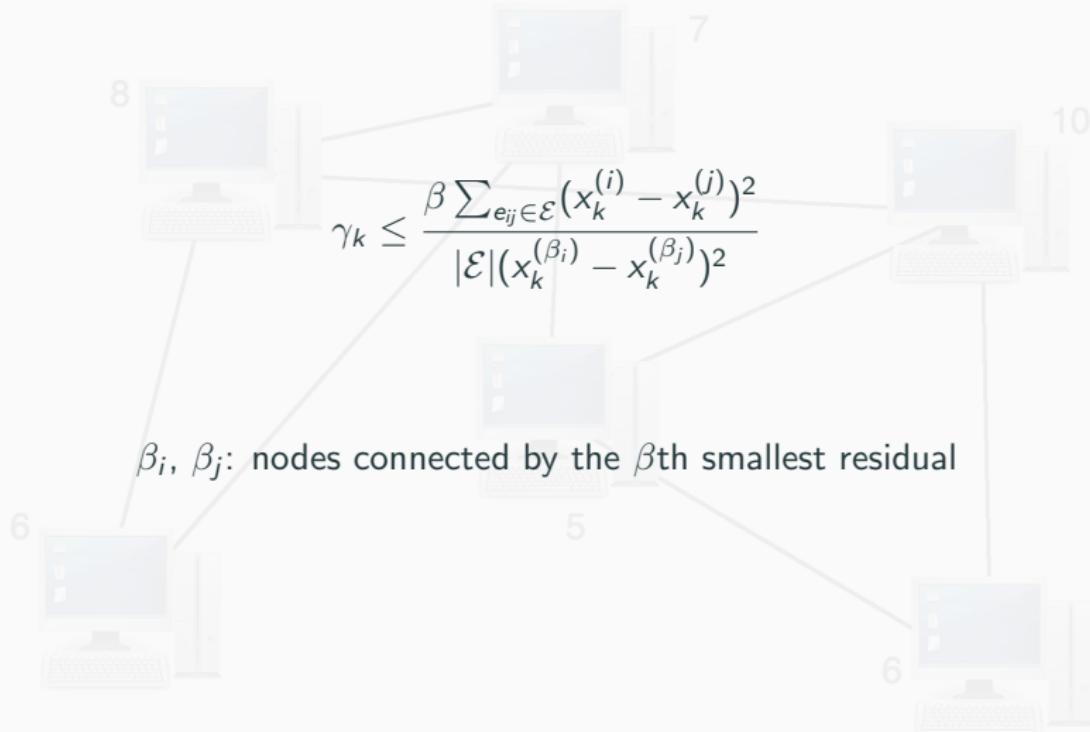
$$\gamma_k \lesssim \frac{n\beta}{\log \beta}$$

Gaussian Convergence



- A is 50000×100 Gaussian matrix, consistent system

γ_k : AC systems



Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need
 - a sampling distribution that depends upon $\|\mathbf{a}_i\|^2$

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need
 - a sampling distribution that depends upon $\|\mathbf{a}_i\|^2$
 - a sampling distribution that depends upon \mathbf{x}_k

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need
 - a sampling distribution that depends upon $\|\mathbf{a}_i\|^2$
 - a sampling distribution that depends upon \mathbf{x}_k

Generalized SKM sampling distribution: $p_{\mathbf{x}} : \binom{[m]}{\beta_k} \rightarrow [0, 1)$

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need
 - a sampling distribution that depends upon $\|\mathbf{a}_i\|^2$
 - a sampling distribution that depends upon \mathbf{x}_k

Generalized SKM sampling distribution: $p_{\mathbf{x}} : \binom{[m]}{\beta_k} \rightarrow [0, 1)$

- $t(\tau, \mathbf{x}) = \operatorname{argmax}_{t \in \tau} (\mathbf{a}_t \mathbf{x} - b_t)^2$

greedy subresidual choice

Generalizing the Convergence Result

- ▷ can immediately generalize to varying β (SKM with β_k)
- ▷ to generalize to non-normalized A , we need
 - a sampling distribution that depends upon $\|\mathbf{a}_i\|^2$
 - a sampling distribution that depends upon \mathbf{x}_k

Generalized SKM sampling distribution: $p_{\mathbf{x}} : \binom{[m]}{\beta_k} \rightarrow [0, 1)$

- $t(\tau, \mathbf{x}) = \operatorname{argmax}_{t \in \tau} (\mathbf{a}_t \mathbf{x} - b_t)^2$
greedy subresidual choice
- $p_{\mathbf{x}}(\tau_k) = \frac{\|\mathbf{a}_{t(\tau_k, \mathbf{x})}\|^2}{\sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x})}\|^2}$
proportional to norm of selected row

Generalized SKM

Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $\tau_k \in \binom{[m]}{\beta_k}$ according to $p_{\mathbf{x}_{k-1}}$.
2. Choose $i_k := t(\tau_k, \mathbf{x}_{k-1})$.
3. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{||\mathbf{a}_{i_k}||^2} \mathbf{a}_{i_k}$.
4. Repeat.

Generalized SKM

Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $\tau_k \in \binom{[m]}{\beta_k}$ according to $p_{\mathbf{x}_{k-1}}$.
2. Choose $i_k := t(\tau_k, \mathbf{x}_{k-1})$.
3. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{||\mathbf{a}_{i_k}||^2} \mathbf{a}_{i_k}$.
4. Repeat.

▷ If $\beta_k = 1$, this is the distribution in [Strohmer, Vershynin '09].

Generalized SKM

Given $\mathbf{x}_0 \in \mathbb{R}^n$:

1. Choose $\tau_k \in \binom{[m]}{\beta_k}$ according to $p_{\mathbf{x}_{k-1}}$.
2. Choose $i_k := t(\tau_k, \mathbf{x}_{k-1})$.
3. Define $\mathbf{x}_k := \mathbf{x}_{k-1} + \frac{b_{i_k} - \mathbf{a}_{i_k}^T \mathbf{x}_{k-1}}{\|\mathbf{a}_{i_k}\|^2} \mathbf{a}_{i_k}$.
4. Repeat.

- ▷ If $\beta_k = 1$, this is the distribution in [Strohmer, Vershynin '09].
- ▷ If $\|\mathbf{a}_i\|^2 = 1$, this is the uniform distribution over $\binom{[m]}{\beta_k}$.

Generalized Result

Theorem (H. - Ma 2019+)

Let \mathbf{x}^* denote the unique solution to the system of equations $A\mathbf{x} = \mathbf{b}$.

Then generalized SKM converges at least linearly in expectation and the bound on the rate depends on the dynamic range, γ_k of the random sample of β_k rows of A , τ_k . Precisely, in the k th iteration of generalized SKM, we have

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2.$$

Generalized Result

Theorem (H. - Ma 2019+)

Let \mathbf{x}^* denote the unique solution to the system of equations $A\mathbf{x} = \mathbf{b}$.

Then generalized SKM converges at least linearly in expectation and the bound on the rate depends on the dynamic range, γ_k of the random sample of β_k rows of A , τ_k . Precisely, in the k th iteration of generalized SKM, we have

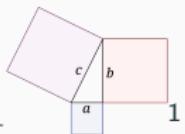
$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2.$$

▷ If all rows of A have the same norm, then

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \sigma_{\min}^2(A)}{\gamma_k \|A\|_F^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2.$$

Sketch of Proof

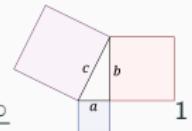
$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 = \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2}$$



¹ Originally created by en:User:Michael Hardy, then scaled, with colour and labels being added by en:User:Wapcaplet, transformed in svg format by fr:Utilisateur:Steff, changed colors and font by de:Leo2004. (<https://commons.wikimedia.org/wiki/File:Pythagorean.svg>)

Sketch of Proof

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 = \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2}$$

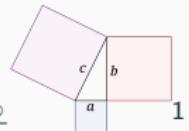


$$\mathbb{E}_{\tau_k} \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} = \sum_{\tau \in \binom{[m]}{\beta_k}} \frac{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \frac{\|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2}{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}$$

¹ Originally created by en:User:Michael Hardy, then scaled, with colour and labels being added by en:User:Wapcaplet, transformed in svg format by fr:Utilisateur:Steff, changed colors and font by de:Leo2004. (<https://commons.wikimedia.org/wiki/File:Pythagorean.svg>)

Sketch of Proof

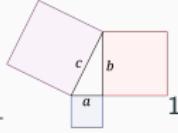
$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 = \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2}$$



$$\begin{aligned} \mathbb{E}_{\tau_k} \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} &= \sum_{\tau \in \binom{[m]}{\beta_k}} \frac{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \frac{\|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2}{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2} \\ &= \frac{1}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \sum_{\tau \in \binom{[m]}{\beta_k}} \|A_\tau \mathbf{x}_{k-1} - \mathbf{b}_\tau\|_\infty^2 \end{aligned}$$

¹ Originally created by en:User:Michael Hardy, then scaled, with colour and labels being added by en:User:Wapcaplet, transformed in svg format by fr:Utilisateur:Steff, changed colors and font by de:Leo2004. (<https://commons.wikimedia.org/wiki/File:Pythagorean.svg>)

Sketch of Proof

$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 = \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2 - \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2}$$


$$\begin{aligned} \mathbb{E}_{\tau_k} \frac{\|A_{\tau_k} \mathbf{x}_{k-1} - \mathbf{b}_{\tau_k}\|_\infty^2}{\|\mathbf{a}_{t(\tau_k, \mathbf{x}_{k-1})}\|^2} &= \sum_{\tau \in \binom{[m]}{\beta_k}} \frac{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \frac{\|A_{\tau} \mathbf{x}_{k-1} - \mathbf{b}_{\tau}\|_\infty^2}{\|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2} \\ &= \frac{1}{\sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \sum_{\tau \in \binom{[m]}{\beta_k}} \|A_{\tau} \mathbf{x}_{k-1} - \mathbf{b}_{\tau}\|_\infty^2 \\ &= \frac{\binom{m}{\beta_k} \beta_k}{\gamma_k m \sum_{\pi \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\pi, \mathbf{x}_{k-1})}\|^2} \|A \mathbf{x}_{k-1} - \mathbf{b}\|^2 \end{aligned}$$

¹ Originally created by en:User:Michael Hardy, then scaled, with colour and labels being added by en:User:Wapcaplet, transformed in svg format by fr:Utilisateur:Steff, changed colors and font by de:Leo2004. (<https://commons.wikimedia.org/wiki/File:Pythagorean.svg>)

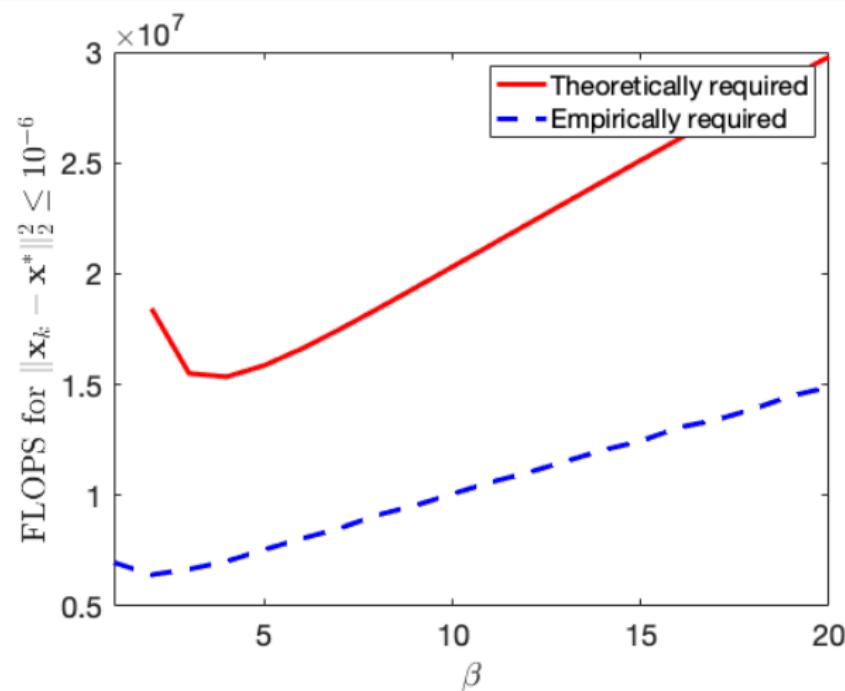
Now can we determine the optimal β ?

Now can we determine the optimal β ?

Roughly, if we know the value of γ_j , we can (just) do it.

Now can we determine the optimal β ?

Roughly, if we know the value of γ_j , we can (just) do it.



Conclusions and Future Work

Conclusions and Future Work

- ▷ presented a convergence rate for SKM which generalizes (and improves) results for RK and Motzkin

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

Conclusions and Future Work

- ▷ presented a convergence rate for SKM which generalizes (and improves) results for RK and Motzkin

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

- ▷ proved a result which “easily” yields results which illustrate SKM improvement for specific types of measurement matrices

Conclusions and Future Work

- ▷ presented a convergence rate for SKM which generalizes (and improves) results for RK and Motzkin

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

- ▷ proved a result which “easily” yields results which illustrate SKM improvement for specific types of measurement matrices
- ▷ specialized our result to Gaussian matrices and AC systems

Conclusions and Future Work

- ▷ presented a convergence rate for SKM which generalizes (and improves) results for RK and Motzkin

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

- ▷ proved a result which “easily” yields results which illustrate SKM improvement for specific types of measurement matrices
- ▷ specialized our result to Gaussian matrices and AC systems
- ▷ identify useful bounds on γ_k for other useful systems

Conclusions and Future Work

- ▷ presented a convergence rate for SKM which generalizes (and improves) results for RK and Motzkin

$$\mathbb{E}_{\tau_k} \|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\beta_k \binom{m}{\beta_k} \sigma_{\min}^2(A)}{\gamma_k m \sum_{\tau \in \binom{[m]}{\beta_k}} \|\mathbf{a}_{t(\tau, \mathbf{x}_{k-1})}\|^2}\right) \|\mathbf{x}_{k-1} - \mathbf{x}^*\|^2$$

- ▷ proved a result which “easily” yields results which illustrate SKM improvement for specific types of measurement matrices
- ▷ specialized our result to Gaussian matrices and AC systems
- ▷ identify useful bounds on γ_k for other useful systems
- ▷ identify optimal β of systems for which γ_k is known

Questions?

- [1] J. A. De Loera, J. Haddock, and D. Needell. **A sampling Kaczmarz-Motzkin algorithm for linear feasibility.** SIAM Journal on Scientific Computing, 39(5):S66–S87, 2017.
- [2] J. Haddock and D. Needell. **On Motzkins method for inconsistent linear systems.** BIT Numerical Mathematics, 59(2):387–401, 2019.
- [3] Nicolas Loizou and Peter Richtárik. **Revisiting randomized gossip algorithms: General framework, convergence rates and novel block and accelerated protocols.** arXiv preprint arXiv:1905.08645, 2019.
- [4] T. S. Motzkin and I. J. Schoenberg. **The relaxation method for linear inequalities.** Canadian J. Math., 6:393–404, 1954.
- [5] T. Strohmer and R. Vershynin. **A randomized Kaczmarz algorithm with exponential convergence.** J. Fourier Anal. Appl., 15:262–278, 2009.

Block Kaczmarz

