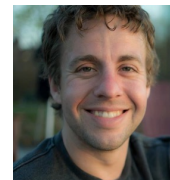
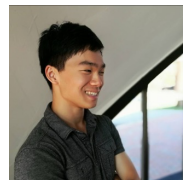


Short-Flat Decompositions and Faster Algorithms for Linear Inverse Problems

Kevin Tian (UT Austin)

Simons Institute Optimization and Algorithm Design Workshop

Based on joint work with:



Jonathan Kelner (MIT), Jerry Li (MSR), Allen Liu (MIT), Aaron Sidford (Stanford)

Roadmap

- Overview
 - Sparse recovery: SOTA and what's new
 - Matrix completion: SOTA and what's new
- Sparse recovery
 - Short-flat decompositions
 - Projected gradient descent
- Matrix completion

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Underconstrained regime: $n \ll d$

Clearly impossible in the worst case. Need to assume more!

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

x^* is s -sparse

Standard
assumption:

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

x^* is s -sparse

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Assume: \mathbf{A} satisfies RNP (no sparse vectors in kernel)

Algo:
$$\min_{\mathbf{A}x=b} \|x\|_1$$

Polynomial time algorithms
(convex programming)

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

x^* is s -sparse

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Assume: \mathbf{A} satisfies RNP (no sparse vectors in kernel)

Algo: $\min_{\mathbf{A}x=b} \|x\|_1$

Polynomial time algorithms
(convex programming)

Upshot: very flexible + general!

- Extends to noisy settings
- Essentially minimal assumptions
- ...potentially expensive in high-dim.

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

x^* is s -sparse

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Assume: \mathbf{A} satisfies RNP

(\mathbf{A} is near-isometry
on sparse vectors)

\mathbf{A} satisfies RIP

Algo:

$$\min_{\mathbf{A}x=b} \|x\|_1$$

Polynomial time algorithms
(convex programming)

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

for all $\|v\|_0 = O(s)$

Nearly-linear time algorithms

Sparse recovery

$$\mathbf{A}x^* = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

x^* is s -sparse

$$\begin{pmatrix} - & a_1 & - \\ - & a_2 & - \\ & \vdots & \\ - & a_n & - \end{pmatrix} \begin{pmatrix} [x^*]_1 \\ [x^*]_2 \\ \vdots \\ [x^*]_d \end{pmatrix} = \begin{pmatrix} [b]_1 \\ [b]_2 \\ \vdots \\ [b]_n \end{pmatrix}$$

Assume: \mathbf{A} satisfies RNP

(\mathbf{A} is near-isometry
on sparse vectors)

\mathbf{A} satisfies RIP

Algo:

$$\min_{\mathbf{A}x=b} \|x\|_1$$

Greedy: Pursuit, OMP

Non-convex: IHT, CoSaMP

Convex: Projected GD

Polynomial time algorithms
(convex programming)

Nearly-linear time algorithms

Sparse recovery

- Theory: both work under standard generative models
- Practice: fast methods much more brittle [Davenport, Needell, Wakin '13], [Jain, Tewari, Kar '14], [Polania, Carrillo, Blanco-Velasco, Barner '14], [Zhang, Wei, Wei, Li, Liu, Liu '16], ...

What's going on?

Assume: \mathbf{A} satisfies RNP

Algo:
$$\min_{\mathbf{A}x=b} \|\mathbf{x}\|_1$$

Polynomial time algorithms
(convex programming)

\mathbf{A} satisfies RIP

Greedy: Pursuit, OMP
Non-convex: IHT, CoSaMP
Convex: Projected GD

Nearly-linear time algorithms

Sparse recovery

Theory vs. practice: *what's going on?*

Not broken by semi-random adversary!



Easily broken by semi-random adversary!

Assume: \mathbf{A} satisfies RNP

Algo:
$$\min_{\mathbf{A}x=b} \|x\|_1$$

Polynomial time algorithms
(convex programming)

\mathbf{A} satisfies RIP

Greedy: Pursuit, OMP
Non-convex: IHT, CoSaMP
Convex: Projected GD

Nearly-linear time algorithms

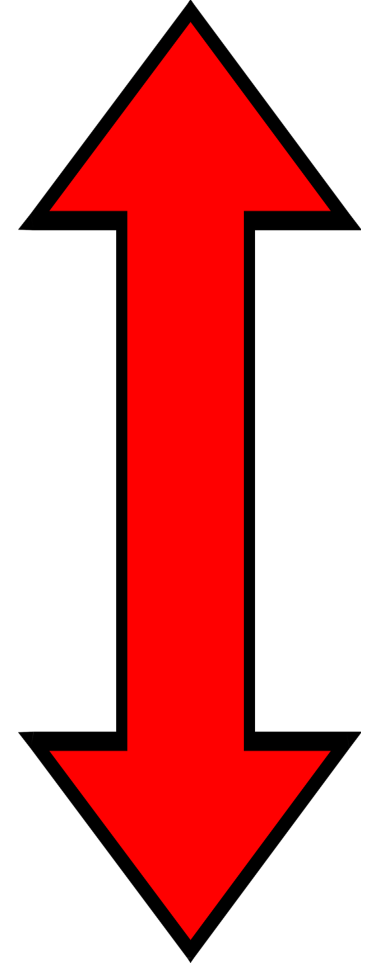
Semi-random models

“Fully random”

Easy! Polynomial-time (very fast?)

“Worst-case”

Hard! (NP-hard, info-impossible?)



Semi-random models

Philosophy

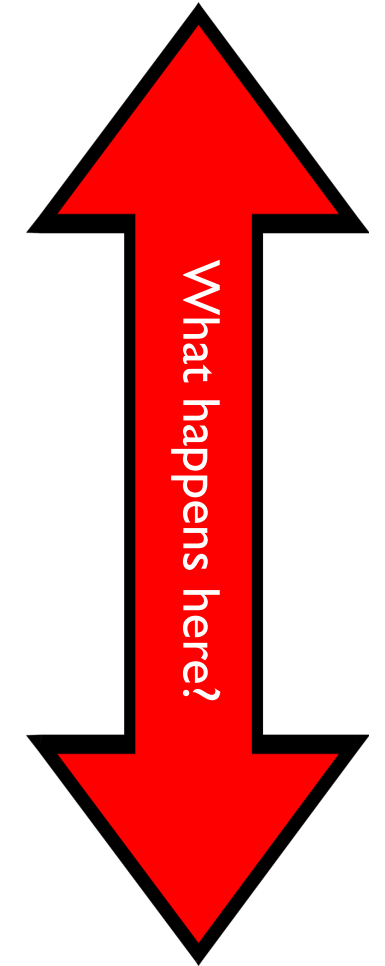
- “Beyond best-case analysis”
- Main q: design algorithms which are *robust* to input assumption violations?

If everything works when life is easy, choose the algorithm that is most robust to assumptions.

“Fully random”

“Worst-case”

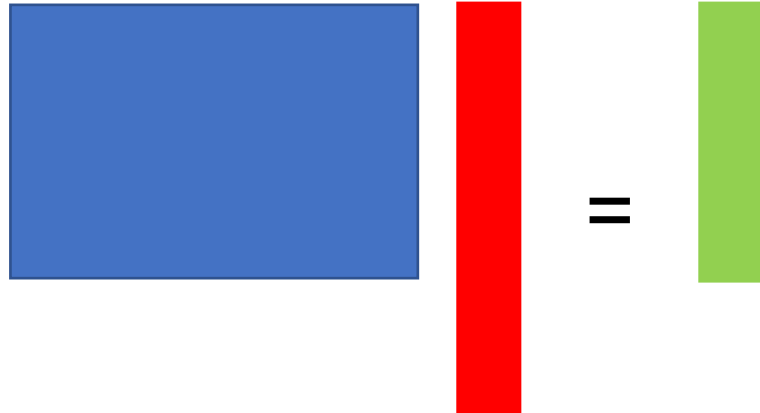
Easy! Polynomial-time (very fast?)



Hard! (NP-hard, info-impossible?)

Semi-random sparse recovery

$$\mathbf{A}x^* = b$$
$$\mathbf{A} \in \mathbb{R}^{n \times d}$$



Nearly-linear time algos:
assume restricted
isometry property (RIP)

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

for all $\|v\|_0 = O(s)$

Semi-random sparse recovery

$$\mathbf{A}x^* = b$$
$$\mathbf{A} \in \mathbb{R}^{n \times d}$$



=



Basic semi-random adversary:

1. Take RIP matrix \mathbf{G}
2. Augment with additional “consistent” measurements
3. Shuffle matrix, present \mathbf{A}

RIP: for all $\|v\|_0 = O(s)$

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

Semi-random sparse recovery

$$\min_{\mathbf{A}x=b} \|x\|_1$$

(polynomial time)

$$\mathbf{A}x^* = b$$
$$\mathbf{A} \in \mathbb{R}^{n \times d}$$



=



Basic semi-random adversary:

1. Take RIP matrix \mathbf{G}
2. Augment with additional “consistent” measurements
3. Shuffle matrix, present \mathbf{A}

RIP: for all $\|v\|_0 = O(s)$

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

Semi-random sparse recovery

$$\min_{\mathbf{A}x=b} \|x\|_1$$

(polynomial time)

$$\mathbf{A}x^* = b$$
$$\mathbf{A} \in \mathbb{R}^{n \times d}$$



=



Basic semi-random adversary:

1. Take RIP matrix \mathbf{G}
2. Augment with additional “consistent” measurements
3. Shuffle matrix, present \mathbf{A}

Fast algorithms?

1. Many greedy/non-convex iterative methods immediately fail (explicit counterexamples)
2. Convex iterative methods' analyses depend on *restricted conditioning*, easy to break

RIP: for all $\|v\|_0 = O(s)$

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

Our basic result

$$\mathbf{A}x^* = b$$
$$\mathbf{A} \in \mathbb{R}^{n \times d}$$



=



- pRIP adversary:
1. Take RIP matrix \mathbf{G}
 2. Augment with additional “consistent” measurements
 3. Shuffle matrix, present \mathbf{A}

Theorem [Kelner, Li, Liu, Sidford, Tian '23]:

Can solve linear systems in
entrywise-bounded* pRIP \mathbf{A} in time

$$\tilde{O}(nd)$$

*satisfied by standard RIP constructions, e.g.
Gaussian, subsampled Fourier/Hadamard matrices

Our general result

$$\|x - x^*\|_2^2 = O\left(\frac{1}{m} \|\xi_{\text{top } m}\|_2^2\right)$$

$$\mathbf{A}x^* + \xi = b$$

$$\mathbf{A} \in \mathbb{R}^{n \times d}$$

(Weighted RIP)



\approx

wRIP (>pRIP) adversary:

1. Exists diagonal reweighting \mathbf{W} such that $\mathbf{A}^T \mathbf{W} \mathbf{A}$ is RIP and \mathbf{A} is entrywise bounded
2. We define $m := \frac{\|\mathbf{w}\|_1}{\|\mathbf{w}\|_\infty}$

Theorem [Kelner, Li, Liu, Sidford, Tian '23]:

Can solve noisy linear systems in entrywise-bounded wRIP \mathbf{A} optimally in time

$$\tilde{O}\left(d \cdot \frac{ns}{m}\right)$$

In pRIP model:

- When all of \mathbf{A} is RIP and $n = m \gg s$, sublinear
- When \mathbf{A} contains minimum $m \approx s$, linear

Matrix “sparse recovery”

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Matrix “sparse recovery”

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumption: \mathbf{X}^* is rank- r

Matrix “sparse recovery”

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard
assumption:

\mathbf{X}^* is rank- r

Poster

Robust Matrix Sensing in the Semi-Random Model

Xing Gao · Yu Cheng

Great Hall & Hall B1+B2 #1720

Theorem 1.1 (Semi-Random Matrix Sensing (Informal)). *Given a set of wRIP sensing matrices $\{A_i\}_{i=1}^n$ and observation vector $b = \mathcal{A}[X^*]$, we can compute X such that $\|X - X^*\|_F \leq \epsilon$ with high probability in time $\tilde{O}(nd^{\omega+1})$, where n is the number of sensing matrices and $O(d^\omega)$ represents the matrix multiplication time for $X \in \mathbb{R}^{d \times d}$.*

Ask Xing and Yu @ NeurIPS '23!

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumption: \mathbf{X}^* is rank- r

$$\mathbf{A}_i = e_{j_i} e_{k_i}^\top \quad \forall i \in [n]$$

2						
		3				
					10	
		-7				
				2		

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumption: \mathbf{X}^* is rank- r

“RIP”-type assumption impossible:
dodge observations with single spike

$$\mathbf{A}_i = e_{j_i} e_{k_i}^\top \quad \forall i \in [n]$$

2						
		3				
					10	
		-7				
				2		

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumptions:

$$\mathbf{X}^* = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^\top$$

$\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$ are “spread”

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumptions:

$$\mathbf{X}^* = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$$

$\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$ are “spread”

$$\max_{i \in [d]} \|\mathbf{U}^\top e_i\|_2 \lesssim \sqrt{\frac{r}{d}}$$

“Incoherence”

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard assumptions:

$$\mathbf{X}^* = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$$

$\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$ are “spread”

Polynomial time:
[Recht '11]

$$\min_{\mathbf{X}^* \text{ matches obs}} \|\mathbf{X}^*\|_1$$

$\approx dr$ samples

Matrix completion

$$\langle \mathbf{A}_i, \mathbf{X}^* \rangle = b_i \quad \forall i \in [n]$$

$\gtrsim dr$

$$\{\mathbf{A}_i\}_{i \in [n]}, \mathbf{X}^* \in \mathbb{R}^{d \times d}$$

Standard
assumptions:

$$\mathbf{X}^* = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^\top$$

$\mathbf{U}, \mathbf{V} \in \mathbb{R}^{d \times r}$ are “spread”

Polynomial time:
[Recht '11]

$$\min_{\mathbf{X}^* \text{ matches obs}} \|\mathbf{X}^*\|_1$$

$\approx dr$ samples

Near-linear time:
[Jain-Netrapalli '15]

PGD + clipping

$$\approx dr^7 \text{ time}$$
$$\approx dr^5 \text{ samples}$$

Matrix completion

Open questions:

1. Improved “fast” rates?
2. Beyond incoherence?

Polynomial time:
[Recht '11]

$$\min_{\mathbf{X}^* \text{ matches obs}} \|\mathbf{X}^*\|_1$$

$\approx dr$ samples

Near-linear time:
[Jain-Netrapalli '15]

PGD + clipping

$$\approx dr^7 \text{ time}$$
$$\approx dr^5 \text{ samples}$$

Matrix completion

Open questions:

1. Improved “fast” rates?
2. Beyond incoherence?
3. Noise-robustness?

Observe: $\mathbf{M}^* + \mathbf{N}$, $\|\mathbf{N}\|_F \leq \Delta$

Recovery: $\|\mathbf{M} - \mathbf{M}^*\|_F \leq \sqrt{d}\Delta$

...SOTA even for polynomial time!
[Candes-Plan '10]

Polynomial time: $\min \|\mathbf{X}^*\|_1$
[Recht '11] \mathbf{X}^* matches obs
 $\approx dr$ samples

Near-linear time: PGD + clipping
[Jain-Netrapalli '15]
 $\approx dr^7$ time
 $\approx dr^5$ samples

Matrix completion

Theorem, Part I [Kelner, Li, Liu, Sidford, Tian '23]:

From rank- r $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ + ($\|\mathbf{N}\|_F \leq \Delta$),
can give $\mathbf{M} \in \mathbb{R}^{d \times d}$, $S \subseteq [d]$ with:

$$\|[\mathbf{M} - \mathbf{M}^*]_{S \times S}\|_F = O(\Delta), \quad |S| \geq 0.99d$$

Polynomial time: $\min \|\mathbf{X}^*\|_1$
[Recht '11] \mathbf{X}^* matches obs
 $\approx dr$ samples

Near-linear time: PGD + clipping
[Jain-Netrapalli '15]
 $\approx dr^7$ time
 $\approx dr^5$ samples

Matrix completion

Theorem, Part I [Kelner, Li, Liu, Sidford, Tian '23]:

From rank- r $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ + ($\|\mathbf{N}\|_F \leq \Delta$),
can give $\mathbf{M} \in \mathbb{R}^{d \times d}$, $S \subseteq [d]$ with:

$$\|[\mathbf{M} - \mathbf{M}^*]_{S \times S}\|_F = O(\Delta), \quad |S| \geq 0.99d$$

...using $\approx dr^{1+o(1)}$ samples

$\approx dr^{2+o(1)}$ time

Polynomial time: $\min \|\mathbf{X}^*\|_1$
[Recht '11] \mathbf{X}^* matches obs
 $\approx dr$ samples

Near-linear time: PGD + clipping
[Jain-Netrapalli '15]
 $\approx dr^7$ time
 $\approx dr^5$ samples

Matrix completion

Theorem, Part IIA [Kelner, Li, Liu, Sidford, Tian '23]:

From rank- r , “regular” $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ +
($\|\mathbf{N}\|_F \leq \Delta$), can give $\mathbf{M} \in \mathbb{R}^{d \times d}$ with:

$$\|\mathbf{M} - \mathbf{M}^*\|_F = O(r^{1.5} \Delta)$$

...using $\approx dr^{1+o(1)}$ samples

$\approx dr^{2+o(1)}$ time

Polynomial time: $\min \|\mathbf{X}^*\|_1$
[Recht '11] \mathbf{X}^* matches obs
 $\approx dr$ samples

Near-linear time: PGD + clipping
[Jain-Netrapalli '15]
 $\approx dr^7$ time
 $\approx dr^5$ samples

Matrix completion

Theorem, Part IIB [Kelner, Li, Liu, Sidford, Tian '23]:

From rank- r , “incoherent” $\mathbf{M}^* \in \mathbb{R}^{d \times d}$ +
($\|\mathbf{N}\|_F \leq \Delta$), can give $\mathbf{M} \in \mathbb{R}^{d \times d}$ with:

$$\|\mathbf{M} - \mathbf{M}^*\|_F = O(r^{1.5} \Delta)$$

...using $\approx dr^{2+o(1)}$ samples

$\approx dr^{3+o(1)}$ time

Polynomial time: $\min \|\mathbf{X}^*\|_1$
[Recht '11] \mathbf{X}^* matches obs
 $\approx dr$ samples

Near-linear time: PGD + clipping
[Jain-Netrapalli '15]
 $\approx dr^7$ time
 $\approx dr^5$ samples

Roadmap

- Overview
 - Sparse recovery: SOTA and what's new
 - Matrix completion: SOTA and what's new
- Sparse recovery
 - Short-flat decompositions
 - Projected gradient descent
- Matrix completion

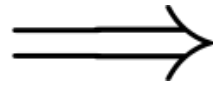
Optimization for sparse recovery?

$$\text{RIP: for all } \|v\|_0 = O(s)$$
$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

“Restricted well-conditioning”:
Well-conditioned restricted to some set

Optimization for sparse recovery?

$$\text{RIP: for all } \|v\|_0 = O(s)$$
$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$



(folklore, proof
via shelling)

$$\text{RIP+: for all NS } v$$

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Optimization for sparse recovery?

A first attempt

- Maintain $x - x^*$ is NS
- ??????
- Profit

RIP+: for all NS v

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Optimization for sparse recovery?

A first attempt

- Maintain $x - x^*$ is NS
- ??????
- Profit

Can maintain (?)
via ℓ_1 projection

RIP+: for all NS v

$$\Omega(\|v\|_2) = \|\mathbf{A}v\|_2 = O(\|v\|_2)$$

Question:

How to reason about effect of projection?

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Key geometric insight

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

Key geometric insight

$$\left(\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top \right) v = A^\top A$$

Random Gaussian matrix

$$\{a_i\}_{i \in [n]} \sim_{\text{i.i.d.}} \mathcal{N}(0, \mathbf{I})$$

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

Key geometric insight

$$\left(\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top \right) v$$

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

Write $a_i = \xi_i v + a_i^\perp$
 $\xi_i \sim \mathcal{N}(0, 1), a_i^\perp \sim \mathcal{N}(0, \mathbf{I} - vv^\top)$

Key geometric insight

$$\left(\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top \right) v$$

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

Write $a_i = \xi_i v + a_i^\perp$

$$a_i a_i^\top v = \xi_i^2 v + \xi_i a_i^\perp$$

Key geometric insight

$$\left(\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top \right) v$$

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

Write $a_i = \xi_i v + a_i^\perp$

$$a_i a_i^\top v = \xi_i^2 v + \xi_i a_i^\perp$$

$$\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top v \approx v + \frac{1}{n} \sum_{i \in [n]} \xi_i a_i^\perp$$

Key geometric insight

$$\left(\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top \right) v$$

Write $a_i = \xi_i v + a_i^\perp$
 $a_i a_i^\top v = \xi_i^2 v + \xi_i a_i^\perp$

Lemma (informal): If you hit a unit v with a random Gaussian matrix, it is “flat” in all directions except v

$$\frac{1}{n} \sum_{i \in [n]} a_i a_i^\top v \approx v + \frac{1}{n} \sum_{i \in [n]} \xi_i a_i^\perp$$

(essentially random)

“flat” := ℓ_∞ bounded

Short-flat decompositions

Lemma (formal): let \mathbf{A} be RIP with parameter s . For all NS unit v ,

$$\mathbf{A}^\top \mathbf{A}v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

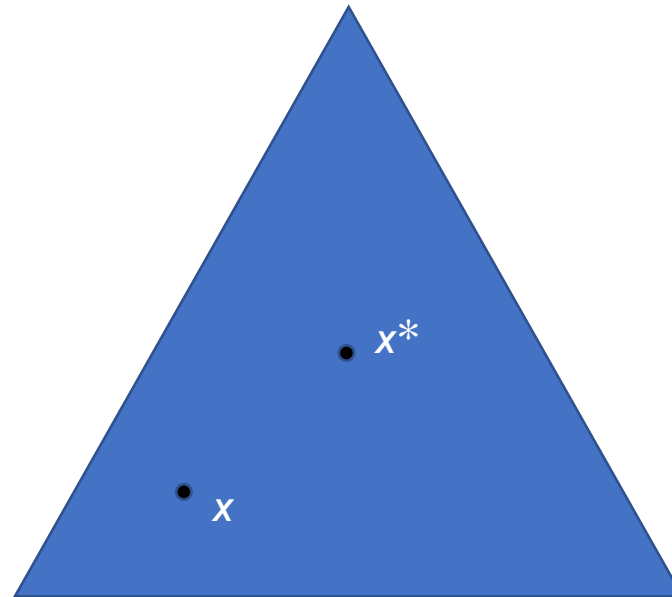
Why does PGD work?

Lemma: If you hit an NS unit v with $\mathbf{A}^T \mathbf{A}$ where \mathbf{A} is RIP, the result has a short-flat decomposition.

Let $v := x - x^*$

Suppose:

- $\|v\|_2 \leq 1$
- $\|v\|_1 \leq \sqrt{s}$



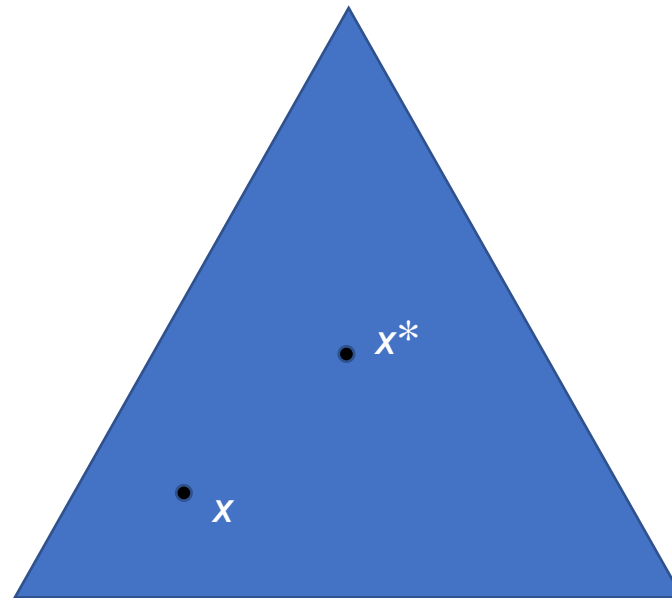
Why does PGD work?

Lemma: If you hit an NS unit v with $\mathbf{A}^T \mathbf{A}$ where \mathbf{A} is RIP, the result has a short-flat decomposition.

Let $v := x - x^*$

Suppose:

- $\|v\|_2 \leq 1$
- $\|v\|_1 \leq \sqrt{s}$



Case 1: $\|v\|_2 \leq \frac{1}{2}$

Halve our radius ☺

Case 2: $\|v\|_2 \geq \frac{1}{2}$

Use $\mathbf{A}^T(\mathbf{A}x - b) = \mathbf{A}^T \mathbf{A}v$ as descent direction

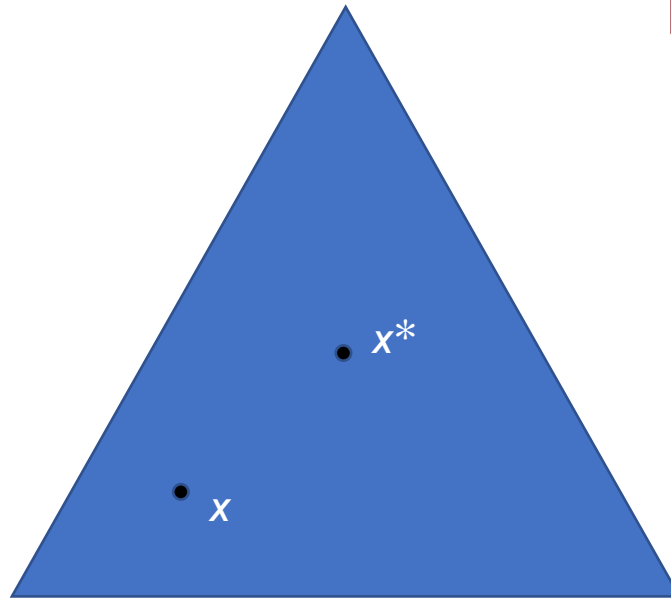
Why does PGD work?

Lemma: If you hit an NS unit v with $\mathbf{A}^\top \mathbf{A}$ where \mathbf{A} is RIP, the result has a short-flat decomposition.

Let $v := x - x^*$

Suppose:

- $\|v\|_2 \leq 1$
- $\|v\|_1 \leq \sqrt{s}$



Case 2 is good idea by Lemma:

$$\mathbf{A}^\top \mathbf{A} v \approx v + e$$

“flat” := ℓ_∞ bounded

Filtered by PGD against ℓ_1
ball + Hölder’s inequality

Case 1: $\|v\|_2 \leq \frac{1}{2}$

Halve our radius ☺

Case 2: $\|v\|_2 \geq \frac{1}{2}$

Use $\mathbf{A}^\top (\mathbf{A}x - b) = \mathbf{A}^\top \mathbf{A} v$ as
descent direction

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
- This set contains x^* by Cauchy-Schwarz

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Restricted W-C: for all NS v ,

$$\frac{1}{n} \sum_{i \in [n]} \langle a_i, v \rangle^2 = [\Omega(1), O(1)]$$

Short-flat: for all NS unit v ,

$$\mathbf{A}^\top \mathbf{A} v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
- $x \leftarrow x_{\text{in}}$
- For 10 iterations:
 - If v is not numerically sparse, we're done
 - If it is numerically sparse, we can PGD

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Restricted W-C: for all NS v ,

$$\frac{1}{n} \sum_{i \in [n]} \langle a_i, v \rangle^2 = [\Omega(1), O(1)]$$

Short-flat: for all NS unit v ,

$$\mathbf{A}^\top \mathbf{A} v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
- $x \leftarrow x_{\text{in}}$
- For 10 iterations:
 - $\Delta = \mathbf{A}x - b = \mathbf{A}v$ for $v = x - x^*$
 - If $\frac{1}{n} \sum_{1 \leq i \leq n} \Delta_i^2 \geq \Omega(1)$ and $\frac{1}{n} \mathbf{A}^T \Delta$ has a short-flat decomposition:
 - $x \leftarrow \operatorname{argmin}_{x' \in \mathcal{X}} \|x' - (x - \eta \mathbf{A}^T \Delta)\|_2$
 - **Constant progress in distance to x^***

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Restricted W-C: for all NS v ,

$$\frac{1}{n} \sum_{i \in [n]} \langle a_i, v \rangle^2 = [\Omega(1), O(1)]$$

Short-flat: for all NS unit v ,

$$\mathbf{A}^T \mathbf{A}v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
- $x \leftarrow x_{\text{in}}$
- For 10 iterations:
 - $\Delta = \mathbf{A}x - b = \mathbf{A}v$ for $v = x - x^*$
 - If $\frac{1}{n} \sum_{1 \leq i \leq n} \Delta_i^2 \geq \Omega(1)$ and $\frac{1}{n} \mathbf{A}^T \Delta$ has a short-flat decomposition:
 - $x \leftarrow \operatorname{argmin}_{x' \in \mathcal{X}} \|x' - (x - \eta \mathbf{A}^T \Delta)\|_2$
 - **Constant progress in distance to x^***
 - Else:
 - Break
 - **Not numerically sparse, radius loose**

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Restricted W-C: for all NS v ,

$$\frac{1}{n} \sum_{i \in [n]} \langle a_i, v \rangle^2 = [\Omega(1), O(1)]$$

Short-flat: for all NS unit v ,

$$\mathbf{A}^T \mathbf{A}v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
- $x \leftarrow x_{\text{in}}$
- For 10 iterations:
 - $\Delta = \mathbf{A}x - b = \mathbf{A}v$ for $v = x - x^*$
 - If $\frac{1}{n} \sum_{1 \leq i \leq n} \Delta_i^2 \geq \Omega(1)$ and $\frac{1}{n} \mathbf{A}^T \Delta$ has a short-flat decomposition:
 - $x \leftarrow \operatorname{argmin}_{x' \in \mathcal{X}} \|x' - (x - \eta \mathbf{A}^T \Delta)\|_2$
 - **Constant progress in distance to x^***
 - Else:
 - Break
- Return x truncated to s largest coordinates

v is numerically sparse (NS) if

$$\frac{\|v\|_1^2}{\|v\|_2^2} = O(s)$$

Restricted W-C: for all NS v ,

$$\frac{1}{n} \sum_{i \in [n]} \langle a_i, v \rangle^2 = [\Omega(1), O(1)]$$

Short-flat: for all NS unit v ,

$$\mathbf{A}^T \mathbf{A}v = p_v + e_v$$

$$\|p_v\|_2 = O(1), \|e_v\|_\infty = O\left(\frac{1}{\sqrt{s}}\right)$$

Algorithm sketch

Input: s -sparse x_{in} , $\|x_{\text{in}} - x^*\|_2 \leq R$

Output: s -sparse x_{out} , $\|x_{\text{out}} - x^*\|_2 \leq \frac{R}{2}$

- $\mathcal{X} := \{x \mid \|x_{\text{in}} - x\|_1 = O(\sqrt{s})R\}$
 - $x \leftarrow x_{\text{in}}$
 - For 10 iterations:
 - $\Delta = \mathbf{A}x - b = \mathbf{A}v$ for $v = x - x^*$
 - If $\frac{1}{n} \sum_{1 \leq i \leq n} \Delta_i^2 \geq \Omega(1)$ and $\frac{1}{n} \mathbf{A}^T \Delta$ has a short-flat decomposition: ←----- *Makes sense even in semi-random case!*
 - $x \leftarrow \operatorname{argmin}_{x' \in \mathcal{X}} \|x' - (x - \eta \mathbf{A}^T \Delta)\|_2$
 - **Constant progress in distance to x^***
 - Else:
 - Break
 - Return x truncated to s largest coordinates
- We find planted solution in near-linear time.

Analysis sketch

$$\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2 \geq 2\eta \langle \underbrace{g_t}_{:= \mathbf{A}^\top \mathbf{A}(x_t - x^*)}, x_{t+1} - x^* \rangle - \|x_t - x_{t+1}\|_2^2$$

Analysis sketch

$$\begin{aligned}\|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2 &\geq 2\eta \langle \underbrace{g_t}_{:= \mathbf{A}^\top \mathbf{A}(x_t - x^*)}, x_{t+1} - x^* \rangle - \|x_t - x_{t+1}\|_2^2 \\ &\geq 2\eta \langle g_t, x_t - x^* \rangle - 2\eta \langle e_t, x_t - x_{t+1} \rangle \\ &\quad - 2\eta \langle p_t, x_t - x_{t+1} \rangle - \|x_t - x_{t+1}\|_2^2\end{aligned}$$

Analysis sketch

$$\begin{aligned} \|x_t - x^*\|_2^2 - \|x_{t+1} - x^*\|_2^2 &\geq 2\eta \langle \underbrace{g_t}_{:= \mathbf{A}^\top \mathbf{A}(x_t - x^*)}, x_{t+1} - x^* \rangle - \|x_t - x_{t+1}\|_2^2 \\ &\geq 2\eta \langle g_t, x_t - x^* \rangle - 2\eta \langle e_t, x_t - x_{t+1} \rangle \\ &\quad \text{big (restricted W-C)} \qquad \text{small (flatness + Hölder)} \\ &\quad - 2\eta \langle p_t, x_t - x_{t+1} \rangle - \|x_t - x_{t+1}\|_2^2 \\ &\quad \text{small (shortness + Young)} \end{aligned}$$

Roadmap

- Overview
 - Sparse recovery: SOTA and what's new
 - Matrix completion: SOTA and what's new
- Sparse recovery
 - Short-flat decompositions
 - Projected gradient descent
- **Matrix completion**

Matrix short-flat decomposition?

$$\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega}$$

“gradient” (i.e. scaled residuals)

Matrix short-flat decomposition?

$$\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega} = \underbrace{\mathbf{M} - \mathbf{M}^*}_{\mathbf{P}} + \underbrace{\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega} - (\mathbf{M} - \mathbf{M}^*)}_{\mathbf{E}}$$

...hopefully flat
(opnorm bounded)?

Matrix short-flat decomposition?

$$\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega} - (\mathbf{M} - \mathbf{M}^*)$$

Matrix Bernstein controls opnorm via...

- “Prob. I bound”: entrywise small
- “Variance bound”: row-column norms small

Matrix short-flat decomposition?

$$\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega} - (\mathbf{M} - \mathbf{M}^*)$$

Matrix Bernstein controls opnorm via...

- “Prob. 1 bound”: entrywise small
- “Variance bound”: row-column norms small

Not true in general, but OK if we drop 1% of rows/cols.

Matrix short-flat decomposition?

$$\frac{1}{p} [\mathbf{M} - \mathbf{M}^*]_{\Omega} - (\mathbf{M} - \mathbf{M}^*)$$

Matrix Bernstein controls opnorm via...

- “Prob. 1 bound”: entrywise small
- “Variance bound”: row-column norms small

Not true in general, but OK if we drop 1% of rows/cols.

...recovering dropped rows/cols is most of the work...

...also need to maintain iterates are low-rank...

What else?

1. General framework for semi-random inverse problems?
 - Similar “fast algo/robust algo” gaps for other problems
 - Fine-grained guarantees?
2. Harder adversaries?
 - How far can we push definition of “bad” observations?
 - Weaker types of hidden structure?

Thank you!

Contact

kjtian.github.io

kjtian@cs.utexas.edu



Semi-Random Sparse Recovery in
Nearly-Linear Time



Matrix Completion in
Almost-Verification Time