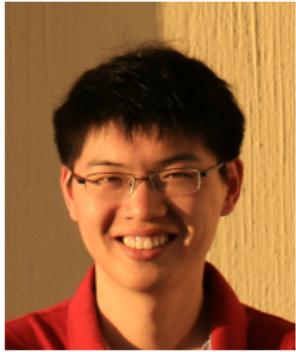


Random Initialization in Nonconvex Phase Retrieval

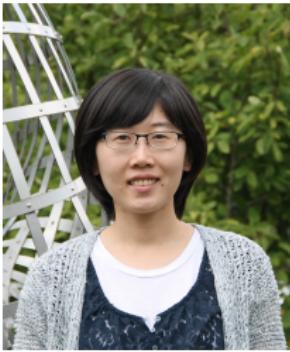


Yuxin Chen

Electrical Engineering, Princeton University



Cong Ma
Princeton ORFE



Yuejie Chi
CMU ECE

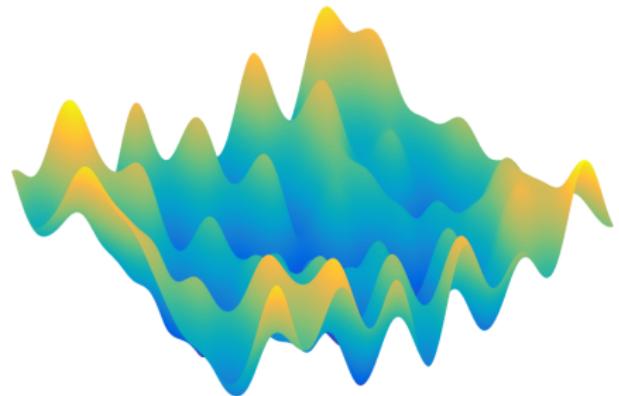


Jianqing Fan
Princeton ORFE

Nonconvex problems are everywhere

Empirical risk minimization is usually nonconvex

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}; \boldsymbol{y})$$

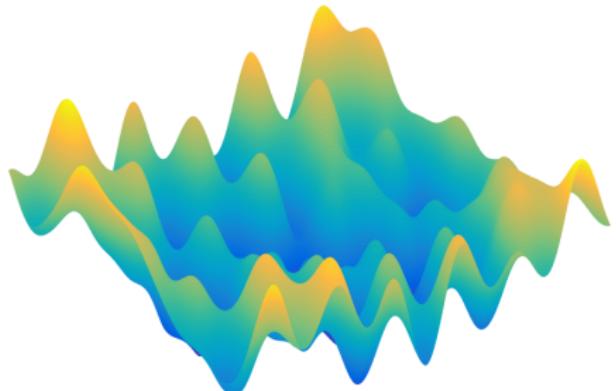


Nonconvex problems are everywhere

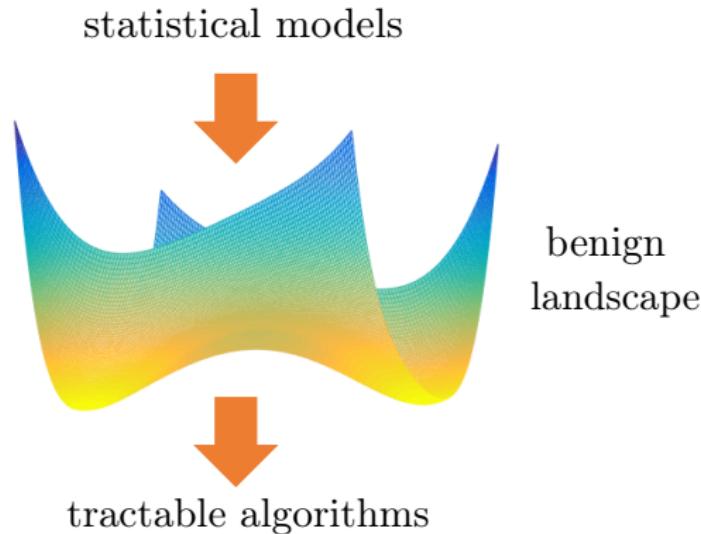
Empirical risk minimization is usually nonconvex

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}; \boldsymbol{y})$$

- low-rank matrix completion
- blind deconvolution
- dictionary learning
- mixture models
- deep learning
- ...

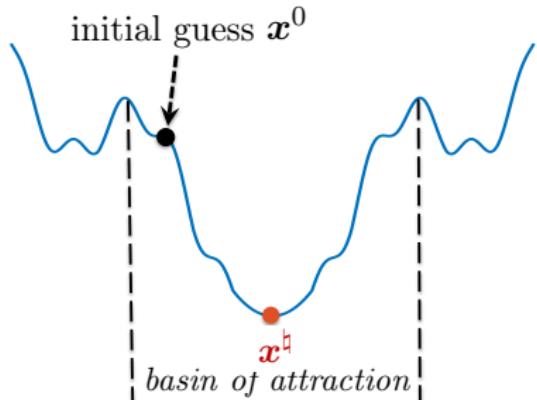


Statistical models come to rescue



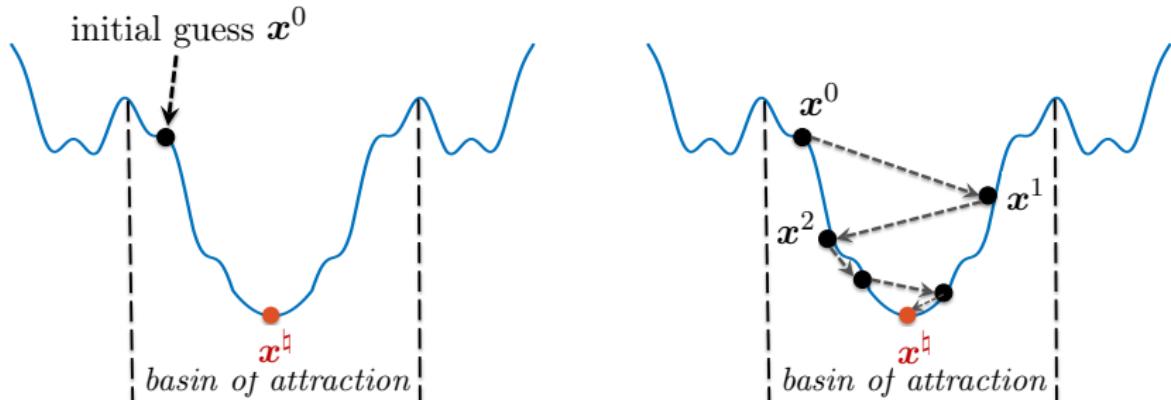
When data are generated by certain statistical models, problems are often much nicer than worst-case instances

A popular two-stage approach



1. initialize within local basin sufficiently close to x^*
(restricted) strongly convex; no local mins

A popular two-stage approach



1. initialize within local basin sufficiently close to x^\natural
(restricted) strongly convex; no local mins
2. iterative refinement

A highly incomplete list of two-stage methods

phase retrieval:

- Netrapalli, Jain, Sanghavi '13
- Candès, Li, Soltanolkotabi '14
- Chen, Candès '15
- Cai, Li, Ma '15
- Wang, Giannakis, Eldar '16
- Zhang, Zhou, Liang, Chi '16
- Kolte, Ozgur '16
- Zhang, Chi, Liang '16
- Soltanolkotabi '17
- Vaswani, Nayer, Eldar '16
- Chi, Lu '16
- Wang, Zhang, Giannakis, Akcakaya, Chen '16
- Tan, Vershynin '17
- Ma, Wang, Chi, Chen '17
- Duchi, Ruan '17
- Jeong, Gunturk '17
- Yang, Yang, Fang, Zhao, Wang, Neykov '17
- Qu, Zhang, Wright '17
- Goldstein, Studer '16
- Bahmani, Romberg '16
- Hand, Voroninski '16
- Wang, Giannakis, Saad, Chen '17
- Barmherzig, Sun '17
- ...

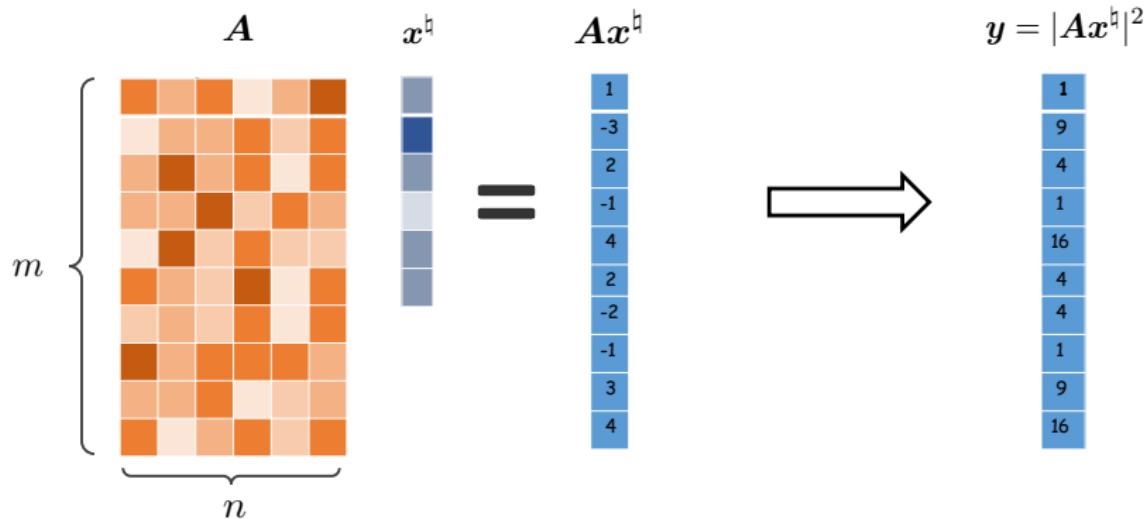
other problems:

- Keshavan, Montanari, Oh '09
- Sun, Luo '14
- Chen, Wainwright '15
- Tu, Boczar, Simchowitz, Soltanolkotabi, Recht '15
- Zheng, Lafferty '15
- Balakrishnan, Wainwright, Yu '14
- Chen, Suh '15
- Chen, Candès '16
- Li, Ling, Strohmer, Wei '16
- Yi, Park, Chen, Caramanis '16
- Jin, Kakade, Netrapalli '16
- Huang, Kakade, Kong, Valiant '16
- Ling, Strohmer '17
- Li, Ma, Chen, Chi '18
- Aghasi, Ahmed, Hand '17
- Lee, Tian, Romberg '17
- Li, Chi, Zhang, Liang '17
- Cai, Wang, Wei '17
- Abbe, Bandeira, Hall '14
- Chen, Kamath, Suh, Tse '16
- Zhang, Zhou '17
- Boumal '16
- Zhong, Boumal '17
- ...

Is carefully-designed initialization necessary for fast convergence?

A case study: phase retrieval

Phase retrieval / solving quadratic systems of equations



Recover $x^h \in \mathbb{R}^n$ from m random quadratic measurements

$$y_k = |\mathbf{a}_k^\top x^h|^2, \quad k = 1, \dots, m$$

assume w.l.o.g. $\|x^h\|_2 = 1$

A natural least squares formulation

$$\underset{\boldsymbol{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$

A natural least squares formulation

$$\text{minimize}_{\boldsymbol{x} \in \mathbb{R}^n} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$

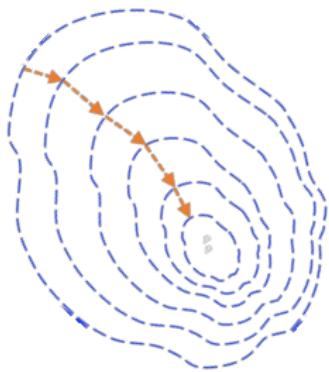
- **issue:** $f(\cdot)$ is highly nonconvex
→ *computationally challenging!*

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

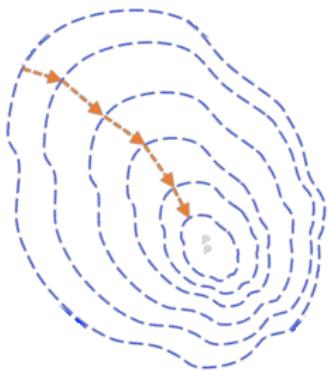
$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$



- **spectral initialization:** $\boldsymbol{x}^0 \leftarrow$ leading eigenvector of certain data matrix

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

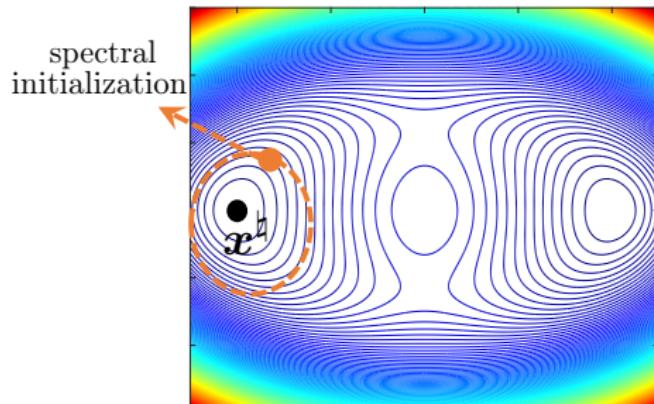
$$\text{minimize}_{\boldsymbol{x}} \quad f(\boldsymbol{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\boldsymbol{a}_k^\top \boldsymbol{x})^2 - y_k \right]^2$$



- **spectral initialization:** $\boldsymbol{x}^0 \leftarrow$ leading eigenvector of certain data matrix
- **gradient descent:**

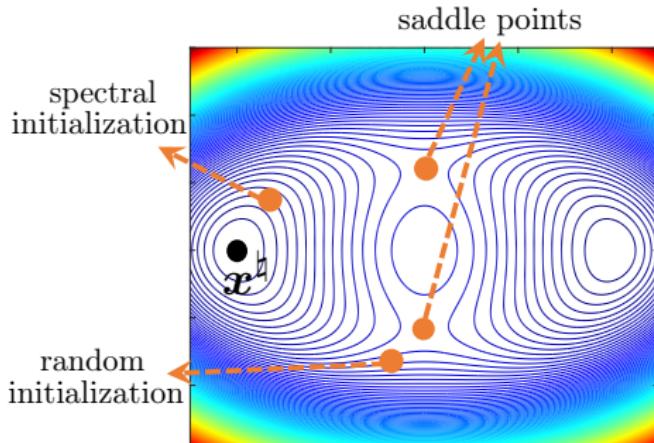
$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \eta_t \nabla f(\boldsymbol{x}^t), \quad t = 0, 1, \dots$$

Initialization



- spectral initialization gets us to (restricted) strongly cvx region

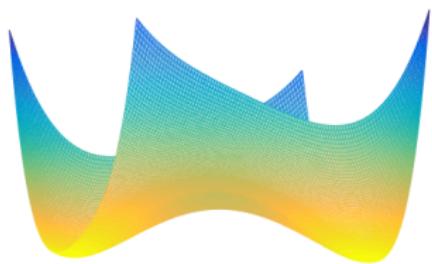
Initialization



- spectral initialization gets us to (restricted) strongly cvx region

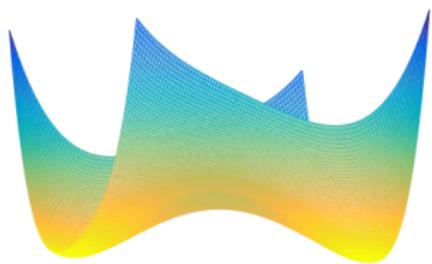
Can we initialize GD randomly, which is **simpler** and **model-agnostic**?

What does prior theory say?



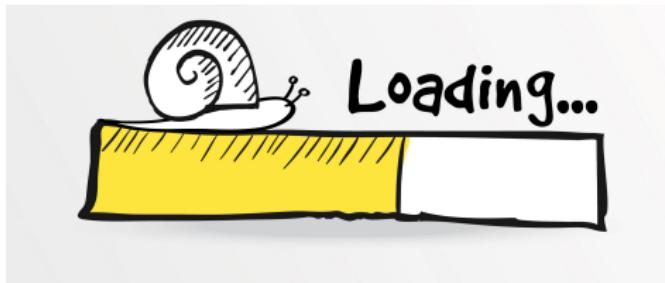
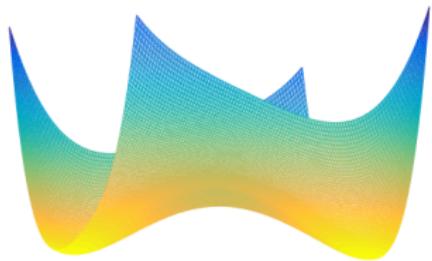
- **landscape:** no spurious local mins (Sun et al. '16)

What does prior theory say?



- **landscape:** no spurious local mins (Sun et al. '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

What does prior theory say?

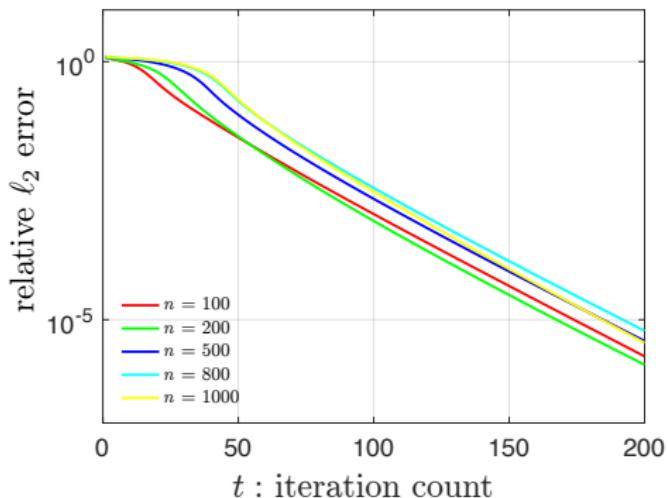


- **landscape:** no spurious local mins (Sun et al. '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

“almost surely” might mean “takes forever”

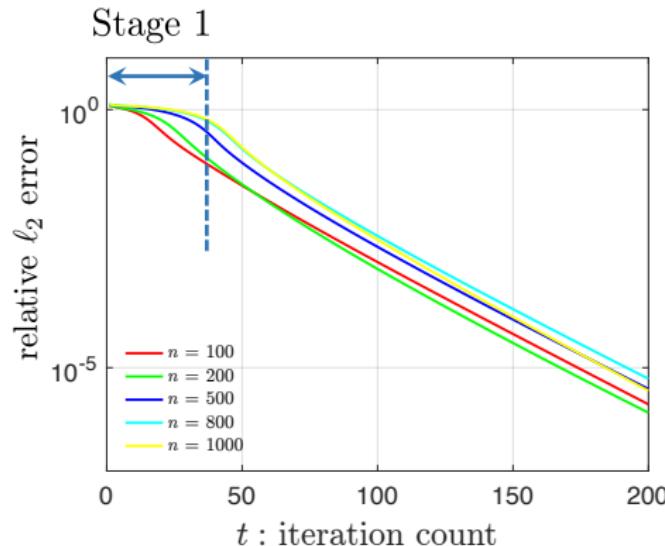
Numerical efficiency of randomly initialized GD

$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Numerical efficiency of randomly initialized GD

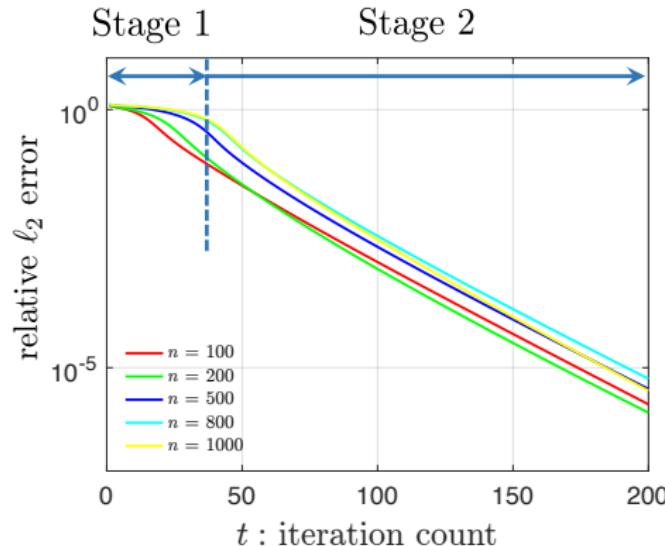
$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **a few iterations**

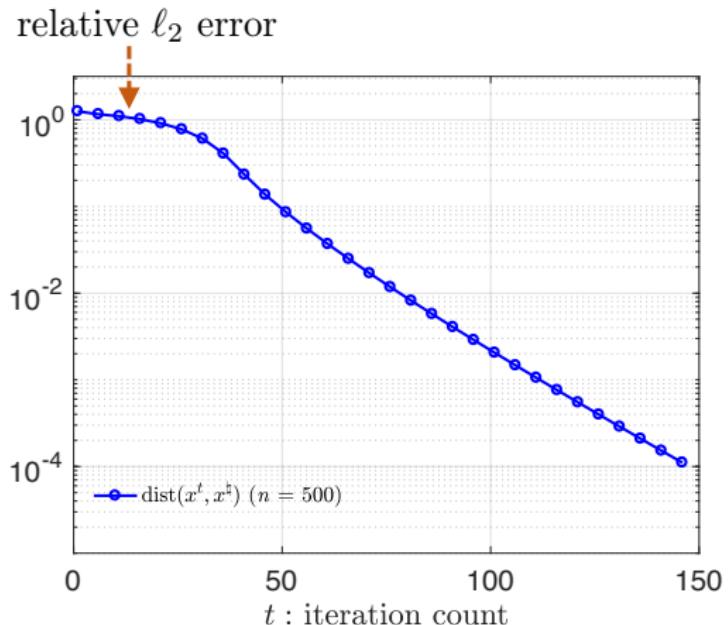
Numerical efficiency of randomly initialized GD

$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$

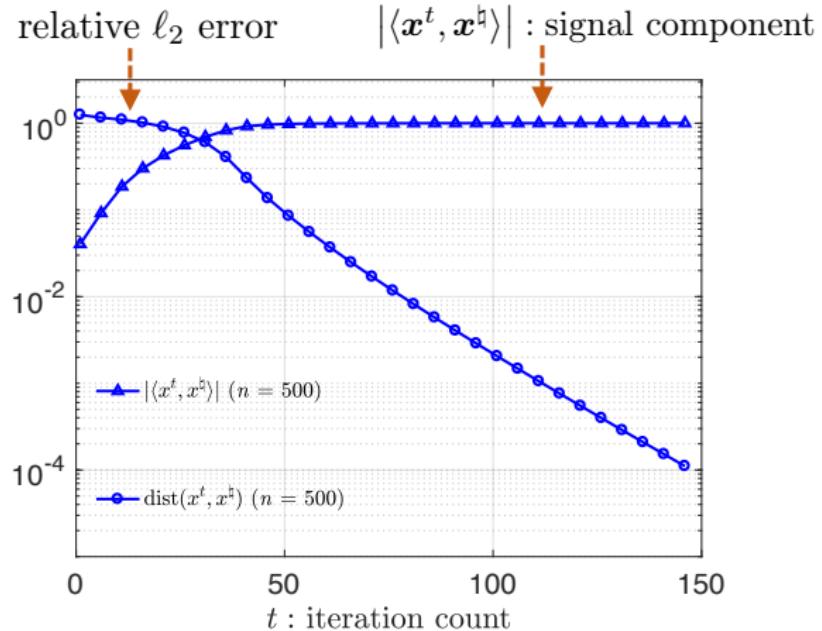


Randomly initialized GD enters local basin within **a few iterations**

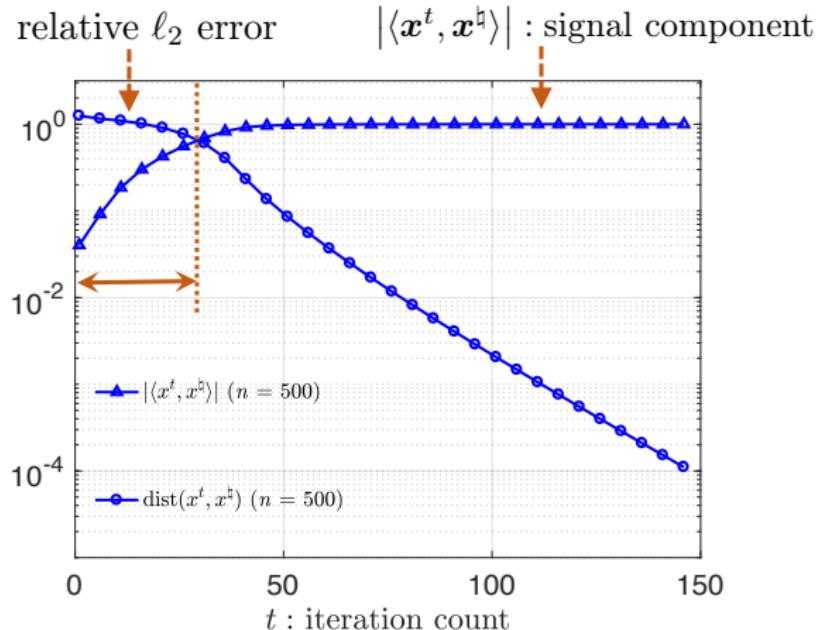
Exponential growth of signal strength in Stage 1



Exponential growth of signal strength in Stage 1



Exponential growth of signal strength in Stage 1



Numerically, a few iterations suffice for entering local region

Our theory

These numerical findings can be formalized when $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$:

Our theory

These numerical findings can be formalized when $a_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, I_n)$:

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

Theorem 1 (Chen, Chi, Fan, Ma '18)

Under i.i.d. Gaussian design, GD with $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} I_n)$ achieves

Our theory

These numerical findings can be formalized when $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$:

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

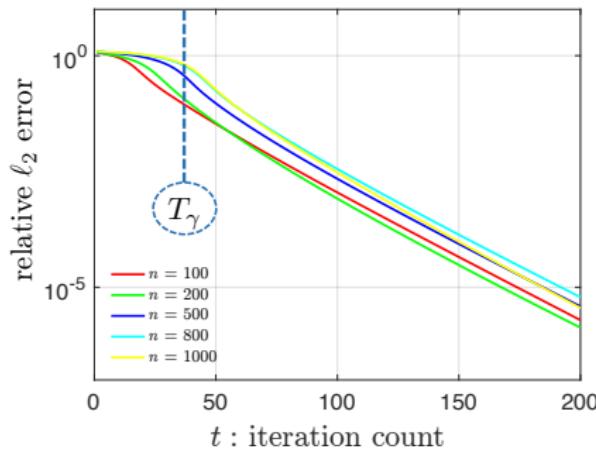
Theorem 1 (Chen, Chi, Fan, Ma '18)

Under i.i.d. Gaussian design, GD with $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$ achieves

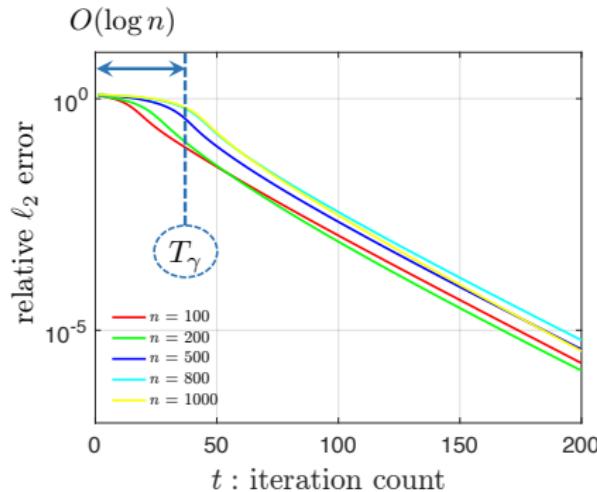
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\natural\|_2, \quad t \geq T_\gamma$$

for $T_\gamma \lesssim \log n$ and some constants $\gamma, \rho > 0$, provided that step size $\eta \asymp 1$ and sample size $m \gtrsim n \text{polylog } m$

Our theory

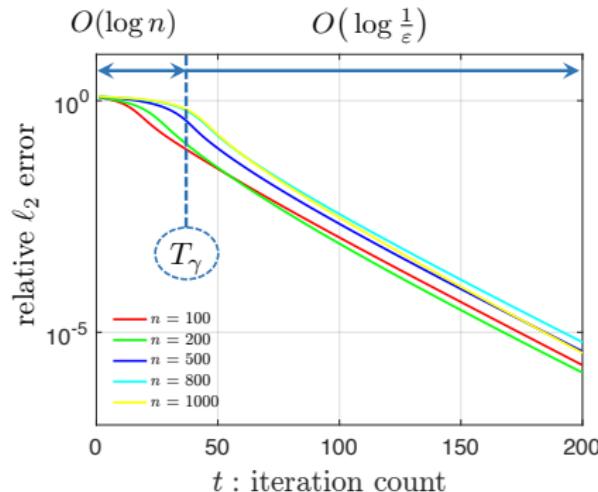


Our theory



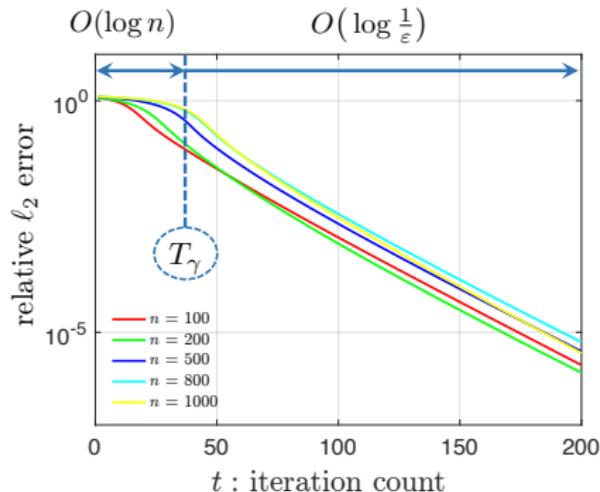
- Stage 1: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \underbrace{\gamma}_{\text{small}}$

Our theory



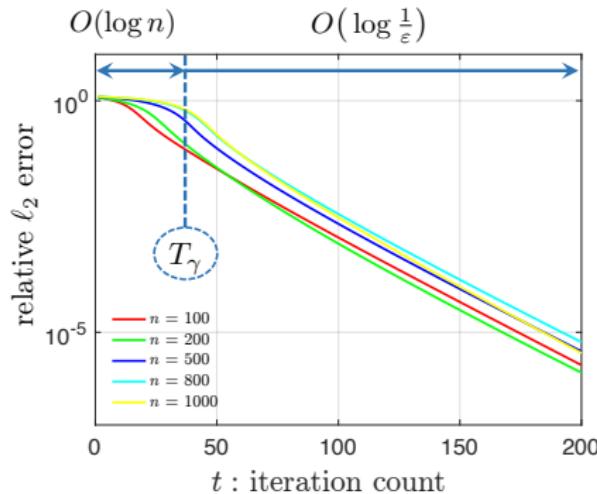
- Stage 1: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \underbrace{\gamma}_{\text{small}}$
- Stage 2: linear convergence

Our theory



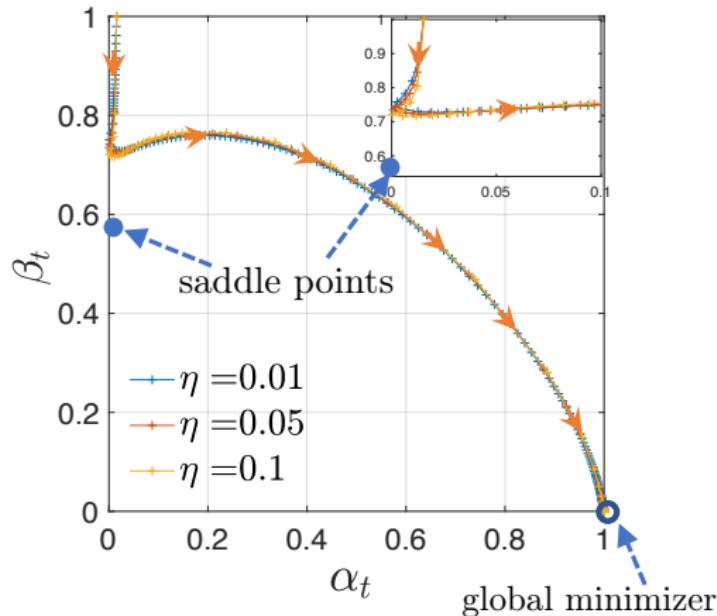
- *near-optimal computational cost:*
 - $O(\log n + \log \frac{1}{\varepsilon})$ iterations to yield ε accuracy

Our theory



- *near-optimal computational cost:*
 - $O(\log n + \log \frac{1}{\epsilon})$ iterations to yield ϵ accuracy
- *near-optimal sample size:* $m \gtrsim n \text{poly} \log m$

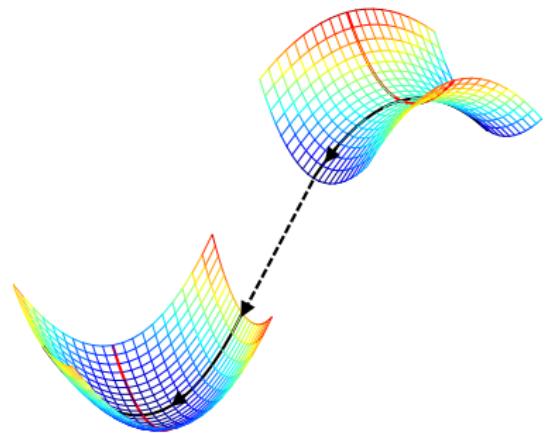
Saddle-escaping schemes?



Randomly initialized GD never hits saddle points!

Other saddle-escaping schemes based on generic landscape analysis

	iteration complexity
trust-region (Sun et al. '16)	$n^7 + \log \log \frac{1}{\varepsilon}$
perturbed GD (Jin et al. '17)	$n^3 + n \log \frac{1}{\varepsilon}$
perturbed accelerated GD (Jin et al. '17)	$n^{2.5} + \sqrt{n} \log \frac{1}{\varepsilon}$
GD (ours) (Chen et al. '18)	$\log n + \log \frac{1}{\varepsilon}$



Generic optimization theory yields highly suboptimal convergence guarantees

A little analysis

What if we have infinite samples?

Gaussian designs: $\mathbf{a}_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n), \quad 1 \leq k \leq m$

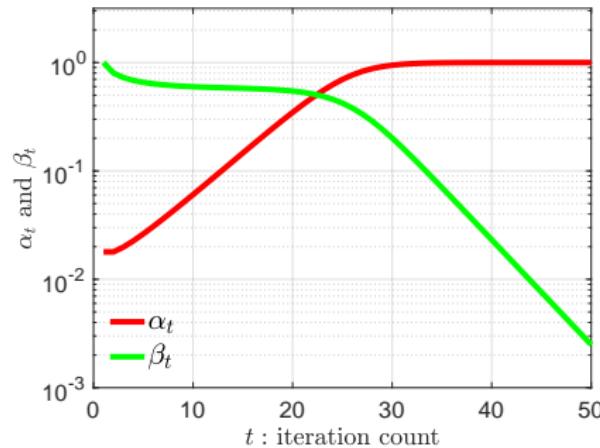
Population level (infinite samples)

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t),$$

where

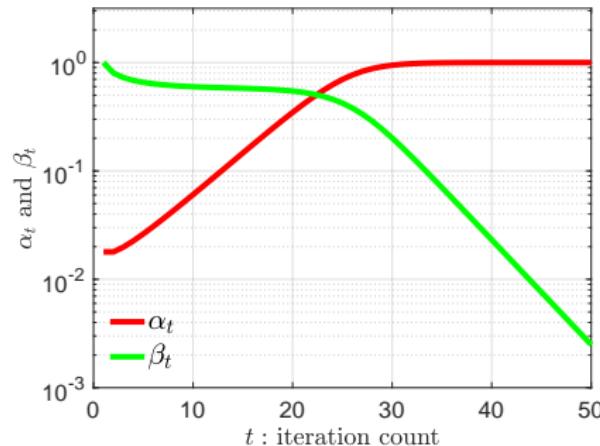
$$\nabla F(\mathbf{x}) := \mathbb{E}[\nabla f(\mathbf{x})] = (3\|\mathbf{x}\|_2^2 - 1)\mathbf{x} - 2(\mathbf{x}^\natural{}^\top \mathbf{x})\mathbf{x}^\natural$$

Population-level state evolution



Let $\alpha_t := \underbrace{|\langle \mathbf{x}^t, \mathbf{x}^\natural \rangle|}_{\text{signal strength}}$ and $\beta_t = \underbrace{\|\mathbf{x}^t - \langle \mathbf{x}^t, \mathbf{x}^\natural \rangle \mathbf{x}^\natural\|_2}_{\text{size of residual component}}$, then

Population-level state evolution



Let $\alpha_t := \underbrace{|\langle \mathbf{x}^t, \mathbf{x}^\natural \rangle|}_{\text{signal strength}}$ and $\beta_t = \underbrace{\|\mathbf{x}^t - \langle \mathbf{x}^t, \mathbf{x}^\natural \rangle \mathbf{x}^\natural\|_2}_{\text{size of residual component}}$, then

$$\alpha_{t+1} = \{1 + 3\eta[1 - (\alpha_t^2 + \beta_t^2)]\}\alpha_t$$

$$\beta_{t+1} = \{1 + \eta[1 - 3(\alpha_t^2 + \beta_t^2)]\}\beta_t$$

2-parameter dynamics

Back to finite-sample analysis

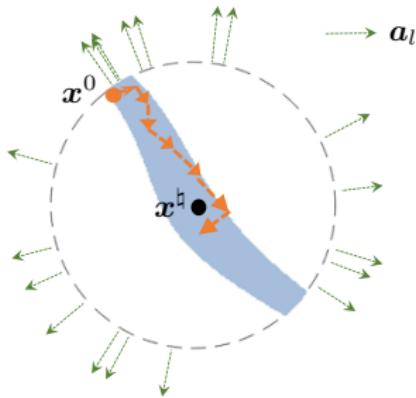
$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \eta \nabla f(\boldsymbol{x}^t)$$

Back to finite-sample analysis

$$\boldsymbol{x}^{t+1} = \boldsymbol{x}^t - \eta \nabla f(\boldsymbol{x}^t) = \boldsymbol{x}^t - \eta \nabla F(\boldsymbol{x}^t) - \underbrace{\eta (\nabla f(\boldsymbol{x}^t) - \nabla F(\boldsymbol{x}^t))}_{\text{residual}}$$

Back to finite-sample analysis

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) = \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t) - \underbrace{\eta (\nabla f(\mathbf{x}^t) - \nabla F(\mathbf{x}^t))}_{\text{residual}}$$

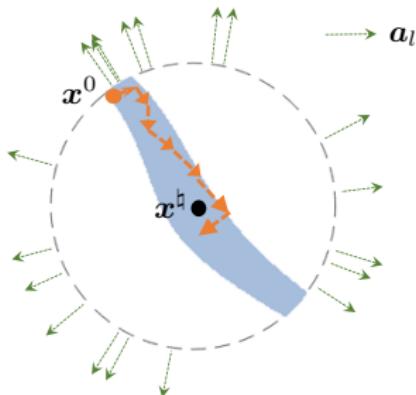


a region with
well-controlled residual

- population-level analysis holds
approximately if \mathbf{x}^t is independent of $\{a_l\}$

Back to finite-sample analysis

$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) = \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t) - \underbrace{\eta (\nabla f(\mathbf{x}^t) - \nabla F(\mathbf{x}^t))}_{\text{residual}}$$



a region with
well-controlled residual

- population-level analysis holds approximately if \mathbf{x}^t is independent of $\{a_l\}$
- **key analysis ingredient:** show \mathbf{x}^t is “nearly-independent” of each a_l

Key proof idea: leave-one-out analysis

Leave out a small amount of information from data and run GD

Key proof idea: leave-one-out analysis

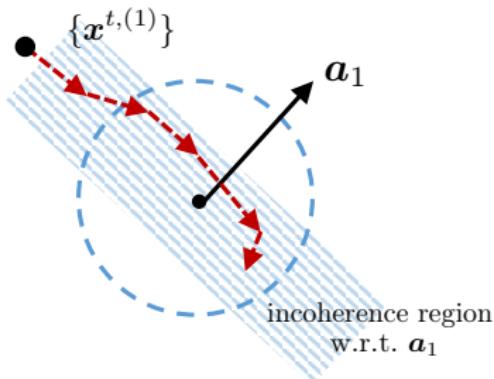
Leave out a small amount of information from data and run GD

The diagram shows the computation of a layer's output $y^{(l)}$ from its input x^h and weight matrix $A^{(l)}$. The input x^h is a vertical vector of size 4. The weight matrix $A^{(l)}$ is a 4x4 matrix with values ranging from -4 to 4. The product $A^{(l)}x^h$ is a vertical vector of size 4. The final output $y^{(l)}$ is the squared magnitude of this product, represented by a red horizontal bar above the result.

$A^{(l)}$	x^h	$A^{(l)}x^h$	$y^{(l)} = A^{(l)}x^h ^2$
$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \\ 3 & 4 & 1 & 2 \\ 4 & 1 & 2 & 3 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}$	$\begin{bmatrix} 1 \\ -3 \\ 2 \\ -1 \end{bmatrix}$	$\begin{bmatrix} 1 \\ 9 \\ 4 \\ 16 \end{bmatrix}$
a_l^\top	$=$	\rightarrow	
$\begin{bmatrix} 1 & 2 & 3 & 4 \\ 2 & 3 & 4 & 1 \\ 3 & 4 & 1 & 2 \\ 4 & 1 & 2 & 3 \end{bmatrix}$	$\begin{bmatrix} -2 \\ -1 \\ 3 \\ 4 \end{bmatrix}$	$\begin{bmatrix} 4 \\ 1 \\ 9 \\ 16 \end{bmatrix}$	

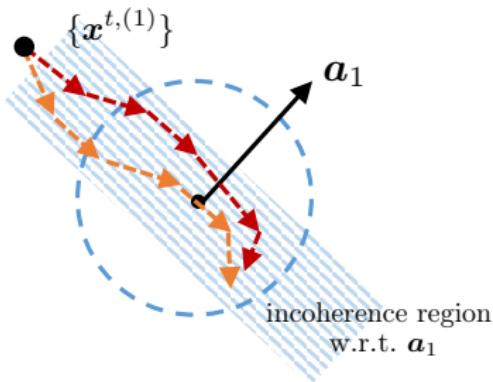
e.g. introduce leave-one-out iterates $x^{t,(l)}$ by running GD without l th sample

Key proof idea: leave-one-out analysis



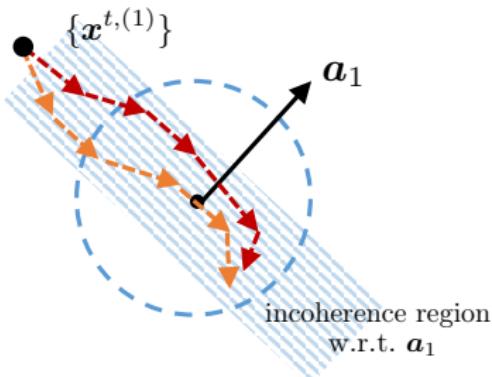
- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l

Key proof idea: leave-one-out analysis



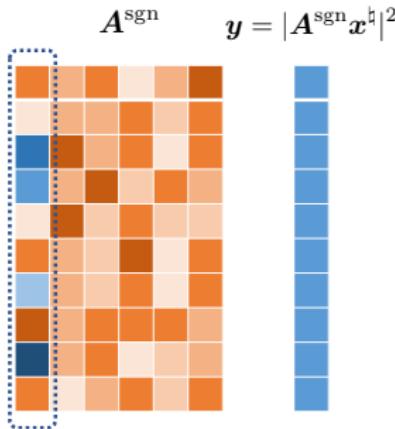
- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t

Key proof idea: leave-one-out analysis

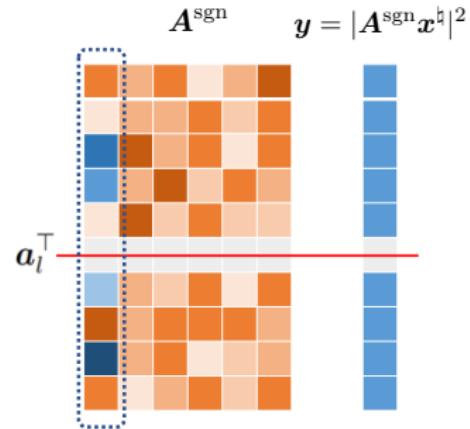


- Leave-one-out iterate $x^{t,(l)}$ is independent of a_l
- Leave-one-out iterate $x^{t,(l)} \approx$ true iterate x^t
 $\implies x^t$ is nearly independent of a_l
nearly orthogonal to

Other leave-one-out sequences



$\mathbf{x}^{t,\text{sgn}}$: indep. of sign info of $\{a_{i,1}\}$



$\mathbf{x}^{t,\text{sgn},(l)}$: indep. of both sign info of $\{a_{i,1}\}$ and a_l

Concluding remarks

Even **simplest** nonconvex methods
are remarkably **efficient** under suitable statistical models

smart initialization	sample splitting	saddle escaping
		

"Gradient Descent with Random Initialization: Fast Global Convergence for Nonconvex Phase Retrieval", Y. Chen, Y. Chi, J. Fan, C. Ma, arXiv:1803.07726