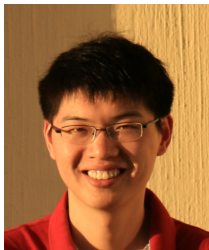


Random Initialization in Nonconvex Statistical Estimation

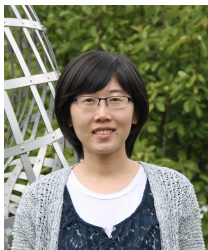


Yuxin Chen

Electrical Engineering, Princeton University



Cong Ma
Princeton ORFE



Yuejie Chi
CMU ECE

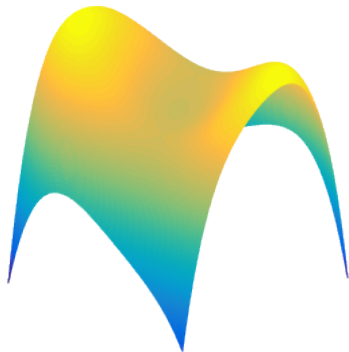


Jianqing Fan
Princeton ORFE

Nonconvex problems are everywhere

Empirical risk minimization is usually nonconvex

$$\text{minimize}_x \quad f(x; \text{data})$$

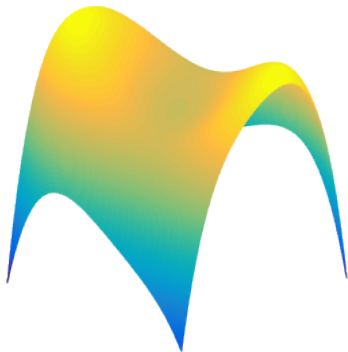


Nonconvex problems are everywhere

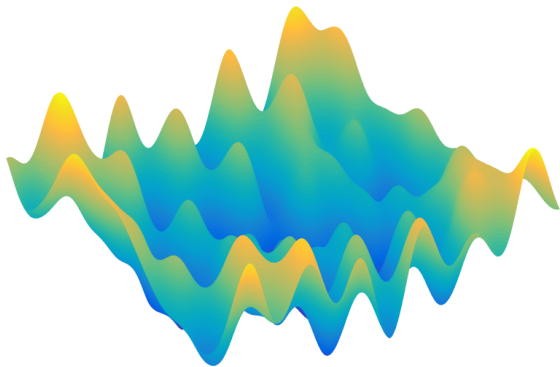
Empirical risk minimization is usually nonconvex

$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}; \text{data})$$

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



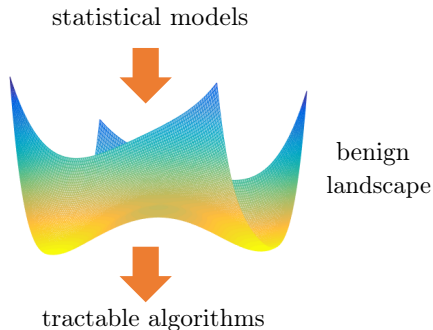
Nonconvex optimization may be super scary



There may be bumps everywhere and exponentially many local optima

e.g. 1-layer neural net (Auer, Herbster, Warmuth '96; Vu '98)

Statistical models come to rescue



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

Example: low-rank matrix recovery

$$\underset{U \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(U) := \sum_{i=1}^m (\langle \mathbf{A}_i, UU^\top \rangle - \langle \mathbf{A}_i, U^*U^{*\top} \rangle)^2$$

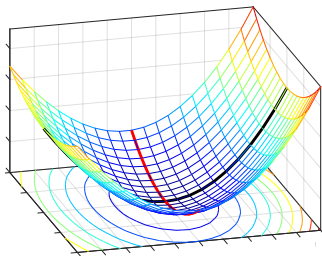
where entries of \mathbf{A}_i are i.i.d. Gaussian

Example: low-rank matrix recovery

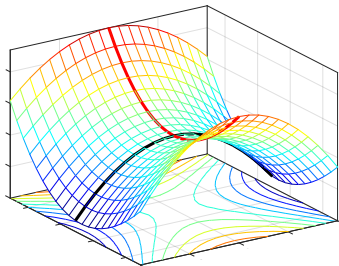
$$\underset{U \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(U) := \sum_{i=1}^m (\langle \mathbf{A}_i, UU^\top \rangle - \langle \mathbf{A}_i, U^*U^{*\top} \rangle)^2$$

where entries of \mathbf{A}_i are i.i.d. Gaussian

- *no spurious local minima* under large enough sample size (Bhojanapalli, Srebro '16)



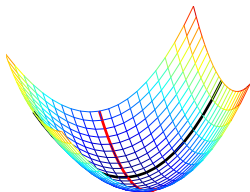
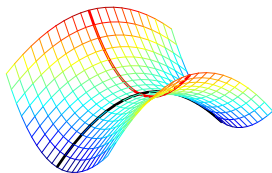
global minimum



saddle point

Separation of landscape analysis and generic algorithm design

landscape analysis
(statistics)

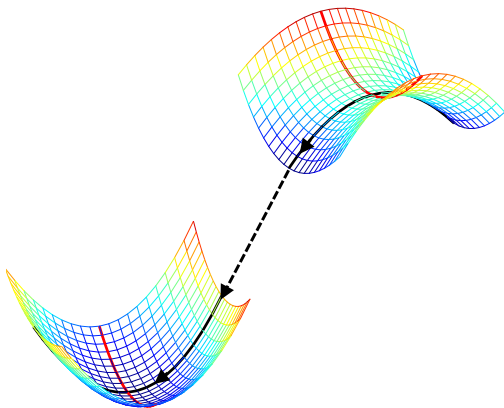


Separation of landscape analysis and generic algorithm design

landscape analysis
(statistics)



generic algorithms
(optimization)



Separation of landscape analysis and generic algorithm design

landscape analysis (statistics)



generic algorithms (optimization)

- 2-layer linear neural network (Baldi, Hornik '89)
 - dictionary learning (Sun et al. '15)
 - phase retrieval (Sun et al. '16, Davis et al. '17)
 - matrix completion (Ge et al. '16, Chen et al. '17)
 - matrix sensing (Bhojanapalli et al. '16, Li et al. '16)
 - empirical risk minimization (Mei et al. '16)
 - synchronization (Bandeira et al. '16)
 - robust PCA (Ge et al. '17)
 - inverting deep neural nets (Hand et al. '17)
 - 1-hidden-layer neural nets (Ge et al. '17)
 - blind deconvolution (Zhang et al. '18, Li et al. '18)
 - ...
- cubic regularization (Nesterov, Polyak '06)
 - gradient descent (Lee et al. '16)
 - trust region method (Sun et al. '16)
 - Carmon et al. '16
 - perturbed GD (Jin et al. '17)
 - perturbed accelerated GD (Jin et al. '17)
 - Agarwal et al. '17
 - Natasha (Allen-Zhu '17)
 - ...

Separation of landscape analysis and generic algorithm design

landscape analysis (statistics)



generic algorithms (optimization)

- 2-layer linear neural network (Baldi, Hornik '89)
- dictionary learning (Sun et al. '15)
- phase retrieval (Sun et al. '16, Davis et al. '17)
- matrix completion (Ge et al. '16, Chen et al. '17)
- matrix sensing (Bhojanapalli et al. '16, Li et al. '16)
- empirical risk minimization (Mei et al. '16)
- synchronization (Bandeira et al. '16)
- robust PCA (Ge et al. '17)
- inverting deep neural nets (Hand et al. '17)
- 1-hidden-layer neural nets (Ge et al. '17)
- blind deconvolution (Zhang et al. '18, Li et al. '18)
- ...
- cubic regularization (Nesterov, Polyak '06)
- gradient descent (Lee et al. '16)
- trust region method (Sun et al. '16)
- Carmon et al. '16
- perturbed GD (Jin et al. '17)
- perturbed accelerated GD (Jin et al. '17)
- Agarwal et al. '17
- Natasha (Allen-Zhu '17)
- ...

Issue: conservative computational guarantees for specific problems
(e.g. solving quadratic systems, matrix completion)

This talk: blending landscape and convergence analysis

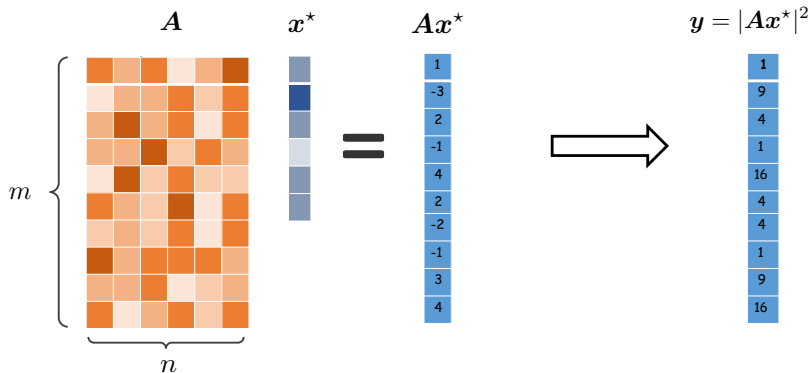
This talk: blending landscape and convergence analysis



Even **simplest** possible nonconvex methods
can be remarkably **efficient** under suitable statistical models

A case study: solving random quadratic systems of equations

Solving quadratic systems of equations



Estimate $\mathbf{x}^* \in \mathbb{R}^n$ from m random quadratic measurements

$$y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2 + \text{noise}, \quad k = 1, \dots, m$$

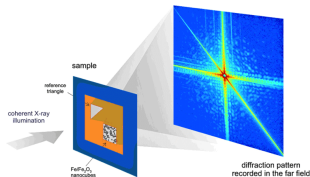
assume w.l.o.g. $\|\mathbf{x}^*\|_2 = 1$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

- electric field $x(t_1, t_2) \rightarrow$ Fourier transform $\hat{x}(f_1, f_2)$

Fig credit: Stanford SLAC



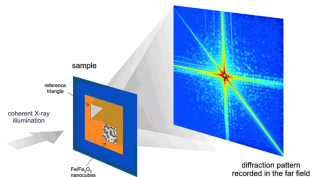
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

Motivation: phase retrieval

Detectors record **intensities** of diffracted rays

- electric field $x(t_1, t_2) \rightarrow$ Fourier transform $\hat{x}(f_1, f_2)$

Fig credit: Stanford SLAC

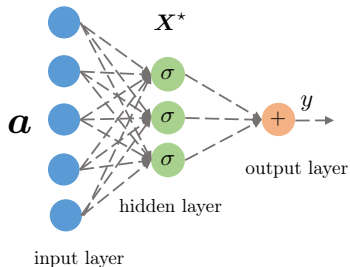


intensity of electrical field: $|\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$

Phase retrieval: recover signal $x(t_1, t_2)$ from intensity $|\hat{x}(f_1, f_2)|^2$

Motivation: learning neural nets with quadratic activation

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17

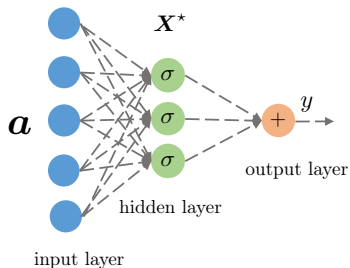


input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*)$$

Motivation: learning neural nets with quadratic activation

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17

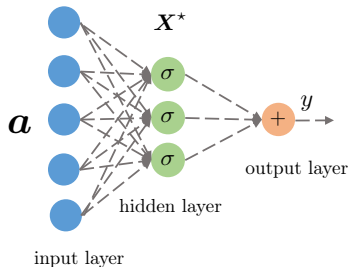


input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*) \stackrel{\sigma(z)=z^2}{:=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i^*)^2$$

Motivation: learning neural nets with quadratic activation

— Soltanolkotabi, Javanmard, Lee '17, Li, Ma, Zhang '17



input features: \mathbf{a} ; weights: $\mathbf{X}^* = [\mathbf{x}_1^*, \dots, \mathbf{x}_r^*]$

$$\text{output: } y = \sum_{i=1}^r \sigma(\mathbf{a}^\top \mathbf{x}_i^*) \stackrel{\sigma(z)=z^2}{:=} \sum_{i=1}^r (\mathbf{a}^\top \mathbf{x}_i^*)^2$$

We consider simplest model when $r = 1$

A natural least squares formulation

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

A natural least squares formulation

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

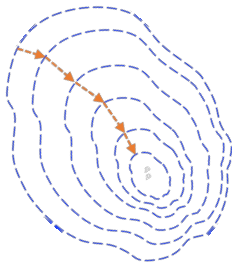
- **issue:** $f(\cdot)$ is highly nonconvex
→ *computationally challenging!*

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

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

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

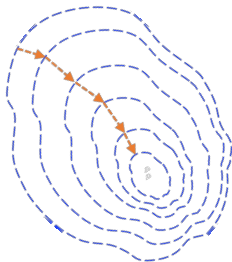
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



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

Wirtinger flow (Candès, Li, Soltanolkotabi '14)

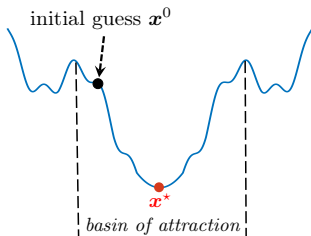
$$\text{minimize}_{\mathbf{x}} \quad f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m \left[(\mathbf{a}_k^\top \mathbf{x})^2 - y_k \right]^2$$



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

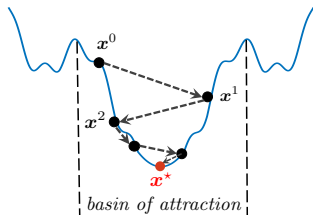
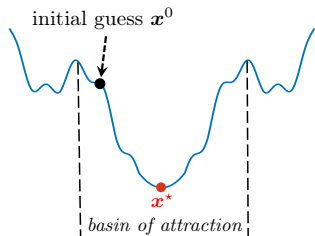
$$\mathbf{x}^{t+1} = \mathbf{x}^t - \eta_t \nabla f(\mathbf{x}^t), \quad t = 0, 1, \dots$$

Rationale of two-stage approach



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

Rationale of two-stage approach



1. initialize within local basin sufficiently close to x^*
(restricted) strongly convex; no saddles / spurious local mins
2. iterative refinement

A highly incomplete list of two-stage methods

phase retrieval:

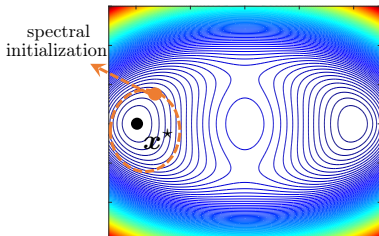
- Netrapalli, Jain, Sanghavi '13
- Candès, Li, Soltanolkotabi '14
- Chen, Candès '15
- Cai, Li, Ma '15
- Wang, Giannakis, Eldar '16
- Zhang, Zhou, Liang, Chi '16
- Kolte, Ozgur '16
- Zhang, Chi, Liang '16
- Soltanolkotabi '17
- Vaswani, Nayer, Eldar '16
- Chi, Lu '16
- Wang, Zhang, Giannakis, Akcakaya, 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, Simchowicz, 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
- ...

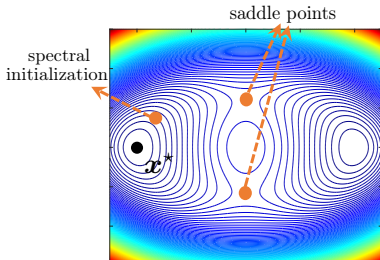
*Is carefully-designed initialization necessary
for fast convergence?*

Initialization



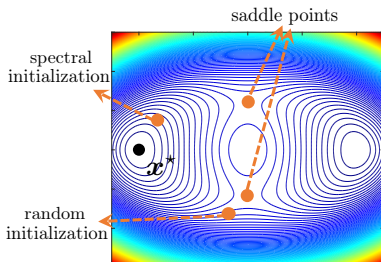
- spectral initialization gets us to (restricted) strongly cvx region

Initialization



- spectral initialization gets us to (restricted) strongly cvx region
- cannot initialize GD anywhere, e.g. might get stuck at saddles

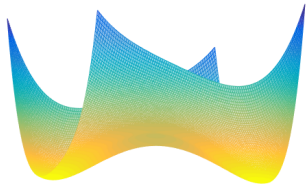
Initialization



- spectral initialization gets us to (restricted) strongly cvx region
- cannot initialize GD anywhere, e.g. might get stuck at saddles

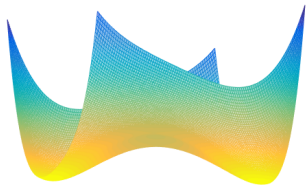
Can we initialize GD randomly, which is **simpler** and **model-agnostic**?

What does prior theory say?



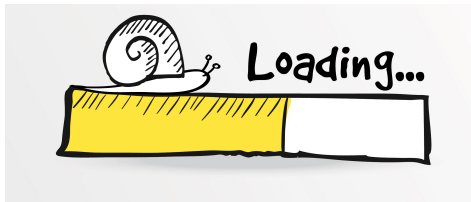
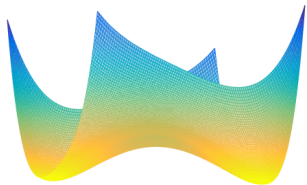
- **landscape:** no spurious local mins (Sun, Qu, Wright '16)

What does prior theory say?



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

What does prior theory say?

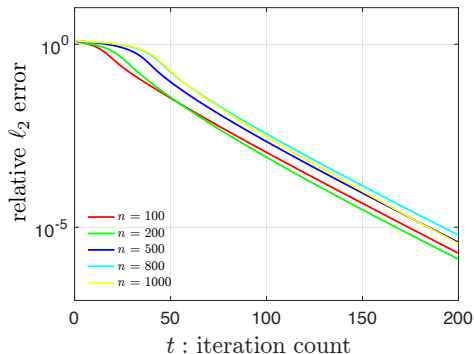


- **landscape:** no spurious local mins (Sun, Qu, Wright '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

“almost surely” might mean “take forever”

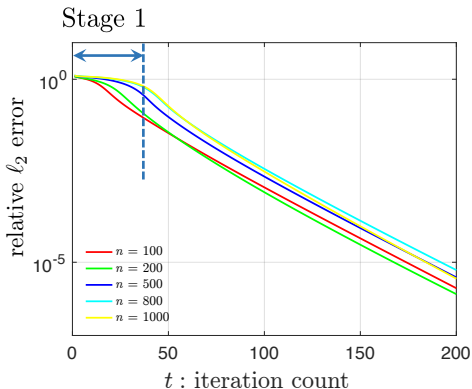
Numerical efficiency of randomly initialized GD

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



Numerical efficiency of randomly initialized GD

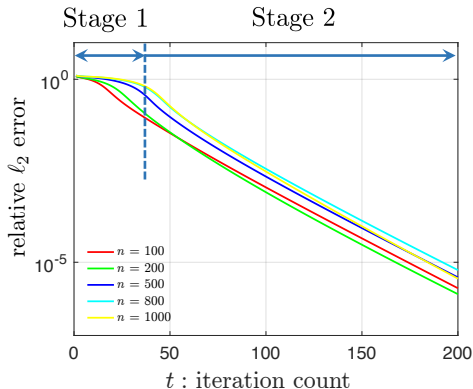
$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



Randomly initialized GD enters local basin within **tens of iterations**

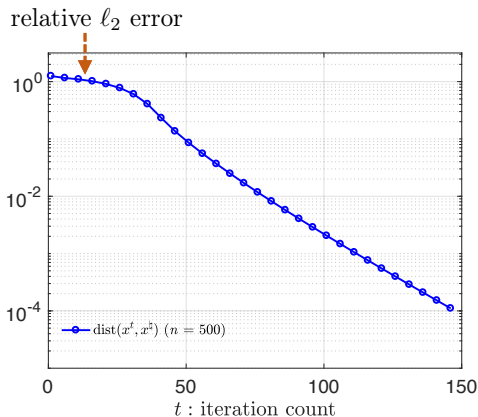
Numerical efficiency of randomly initialized GD

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

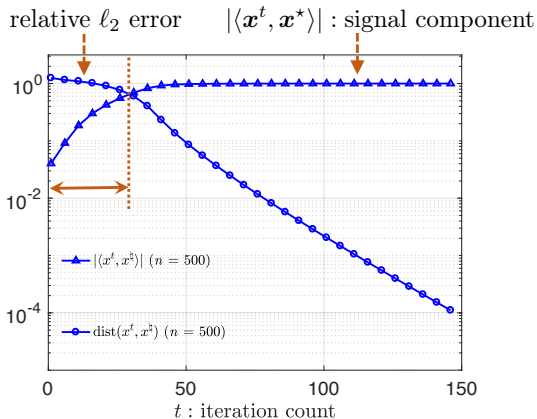


Randomly initialized GD enters local basin within **tens of iterations**

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: noiseless case

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: noiseless case

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

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

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

Our theory: noiseless case

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

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

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

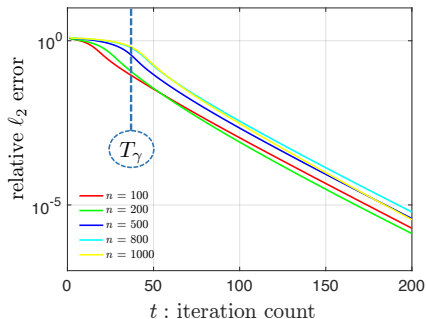
Under i.i.d. Gaussian design, GD with $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1}\mathbf{I}_n)$ achieves

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

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

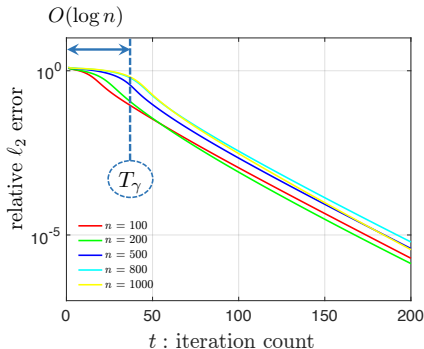
Our theory: noiseless case

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



Our theory: noiseless case

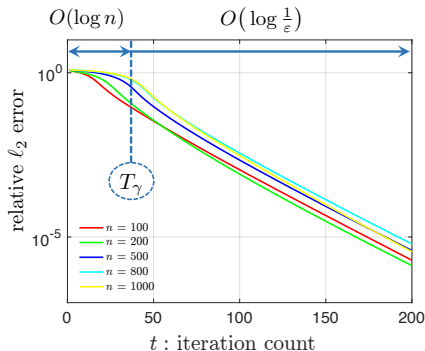
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *Stage 1*: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$ (e.g. $\gamma = 0.1$)

Our theory: noiseless case

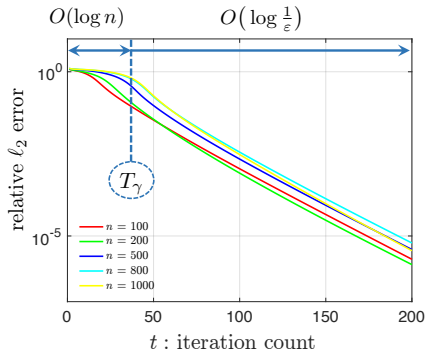
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



- *Stage 1*: takes $O(\log n)$ iterations to reach $\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma$ (e.g. $\gamma = 0.1$)
- *Stage 2*: linear (geometric) convergence

Our theory: noiseless case

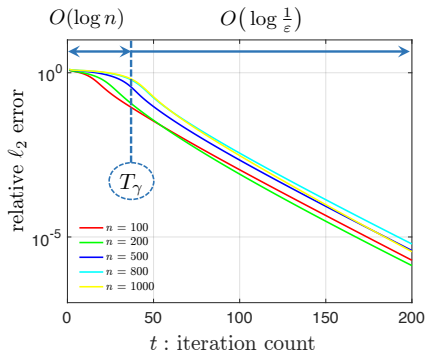
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^*) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^*\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

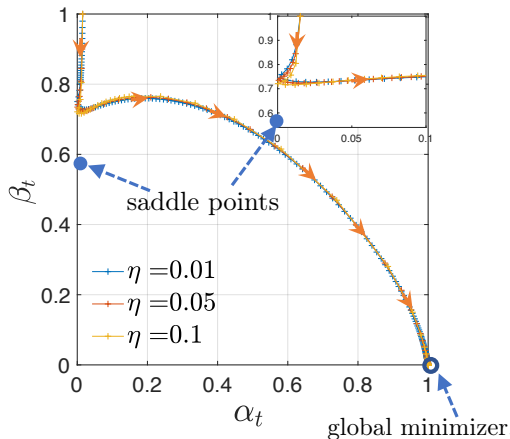
Our theory: noiseless case

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



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

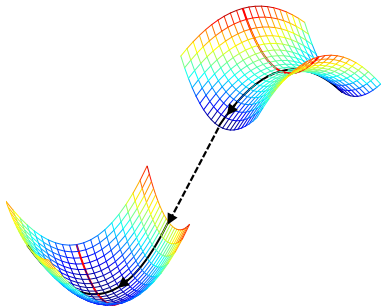
Automatic saddle avoidance



Randomly initialized GD never hits saddle points!

Other saddle-escaping schemes based on generic landscape analysis

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



Generic optimization theory yields highly suboptimal convergence guarantees

A bit of analysis

What if we have infinite samples?

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

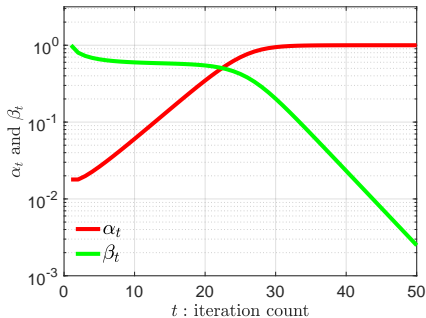
Population level (infinite samples)

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

where

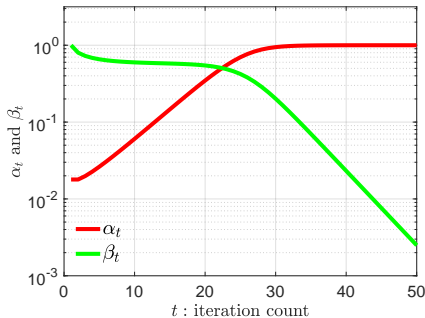
$$\nabla F(\mathbf{x}) := \mathbb{E}[\nabla f(\mathbf{x})] = (3\|\mathbf{x}\|_2^2 - 1)\mathbf{x} - 2(\mathbf{x}^{\star\top} \mathbf{x})\mathbf{x}^{\star}$$

Population-level state evolution



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

Population-level state evolution



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

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

Back to finite-sample analysis

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

Back to finite-sample analysis

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

— take one term in $\mathbf{x}^{\star\top} (\nabla f(\mathbf{x}^t) - \nabla F(\mathbf{x}^t))$ for example:

$$r_1 := \frac{1}{m} \sum_i (\mathbf{a}_{i,\perp}^\top \mathbf{x}_\perp^t)^3 a_{i,1}$$

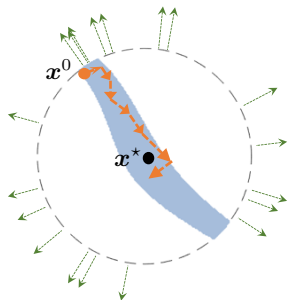
Back to finite-sample analysis

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

— take one term in $\mathbf{x}^{*\top} (\nabla f(\mathbf{x}^t) - \nabla F(\mathbf{x}^t))$ for example:

$$r_1 := \frac{1}{m} \sum_i (\mathbf{a}_{i,\perp}^\top \mathbf{x}_\perp^t)^2 a_{i,1}$$

- $r_1 \asymp \frac{1}{\sqrt{m}}$ if \mathbf{x}^t is independent of $\{\mathbf{a}_l\}$
desired level



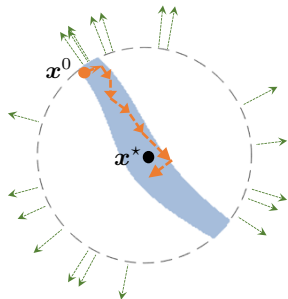
a region with
well-controlled residual

Back to finite-sample analysis

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

— take one term in $\mathbf{x}^{*\top}(\nabla f(\mathbf{x}^t) - \nabla F(\mathbf{x}^t))$ for example:

$$r_1 := \frac{1}{m} \sum_i (\mathbf{a}_{i,\perp}^\top \mathbf{x}_\perp^t)^2 a_{i,1}$$



a region with
well-controlled residual

- $r_1 \asymp \underbrace{\frac{1}{\sqrt{m}}}_{\text{desired level}}$ if \mathbf{x}^t is independent of $\{\mathbf{a}_l\}$
- **key analysis ingredient:** show \mathbf{x}^t is “nearly-independent” of (some part of) $\{\mathbf{a}_l\}$

Key proof idea: leave-one-out analysis

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

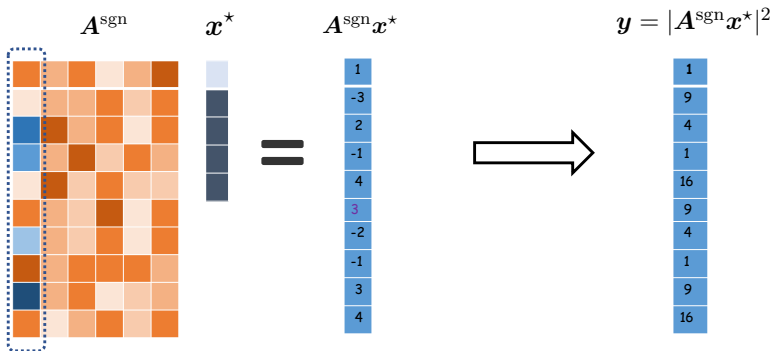
Key proof idea: leave-one-out analysis

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

- Stein '72
- El Karoui, Bean, Bickel, Lim, Yu '13
- El Karoui '15
- Javanmard, Montanari '15
- Zhong, Boumal '17
- Lei, Bickel, El Karoui '17
- Sur, Chen, Candès '17
- Abbe, Fan, Wang, Zhong '17
- Chen, Fan, Ma, Wang '17

Key proof idea: leave-one-out analysis

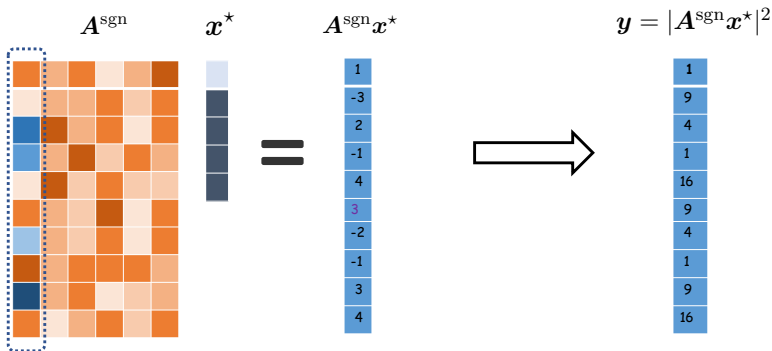
Leave out a small amount of information from data and re-run GD



- generate A^{sgn} by randomly flipping $\text{sgn}(a_{i,1})$, $\forall i$

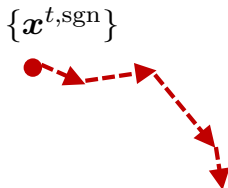
Key proof idea: leave-one-out analysis

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



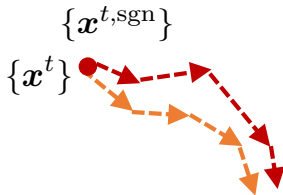
- generate A^{sgn} by randomly flipping $\text{sgn}(a_{i,1})$, $\forall i$
- generate auxiliary iterates $\{x^{t,\text{sgn}}\}$ by re-running GD w.r.t. A^{sgn}

Key proof idea: leave-one-out analysis



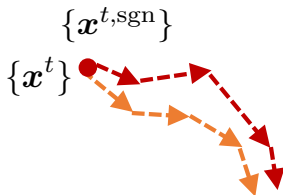
- Auxiliary iterate $\mathbf{x}^{t,\text{sgn}}$ is independent of $\{\text{sgn}(a_{i,1})\}$

Key proof idea: leave-one-out analysis



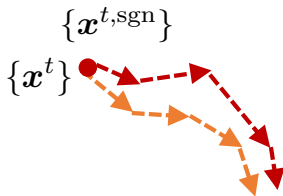
- Auxiliary iterate $x^{t,sgn}$ is independent of $\{sgn(a_{i,1})\}$
- Auxiliary iterate $x^{t,sgn} \approx$ true iterate x^t

Key proof idea: leave-one-out analysis



- Auxiliary iterate $x^{t,\text{sgn}}$ is independent of $\{\text{sgn}(a_{i,1})\}$
- Auxiliary iterate $x^{t,\text{sgn}} \approx$ true iterate x^t
 $\implies x^t$ is **nearly independent of** $\{\text{sgn}(a_{i,1})\}$

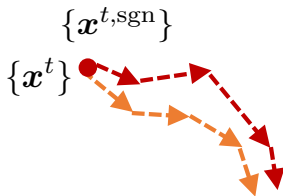
Key proof idea: leave-one-out analysis



- Auxiliary iterate $x^{t,\text{sgn}}$ is independent of $\{\text{sgn}(a_{i,1})\}$
- Auxiliary iterate $x^{t,\text{sgn}} \approx$ true iterate x^t
 $\implies x^t$ is **nearly independent of** $\{\text{sgn}(a_{i,1})\}$
- This makes it easy to control

$$r_1 = \frac{1}{m} \sum_i (\mathbf{a}_{i,\perp}^\top \mathbf{x}_\perp^t)^3 a_{i,1}$$

Key proof idea: leave-one-out analysis

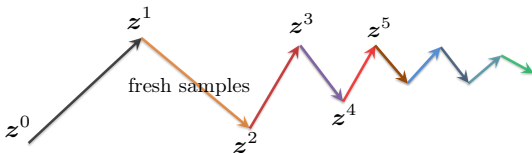


- Auxiliary iterate $x^{t,\text{sgn}}$ is independent of $\{\text{sgn}(a_{i,1})\}$
- Auxiliary iterate $x^{t,\text{sgn}} \approx$ true iterate x^t
 $\implies x^t$ is **nearly independent of** $\{\text{sgn}(a_{i,1})\}$
- This makes it easy to control

$$r_1 = \frac{1}{m} \sum_i (\mathbf{a}_{i,\perp}^\top \mathbf{x}_\perp^t)^3 |a_{i,1}| \text{sgn}(a_{i,1})$$

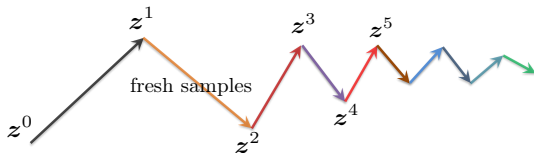
No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis

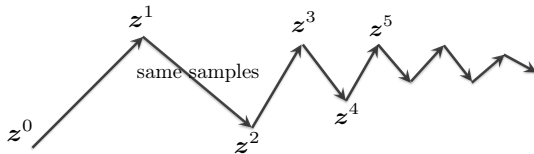


No need of sample splitting

- Several prior works use sample-splitting: require **fresh samples** at each iteration; not practical but helps analysis






- **This work:** reuses all samples in all iterations



Concluding remarks

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

smart initialization	sample splitting	saddle escaping
		

1. “Gradient descent with random initialization: ...”, Y. Chen, Y. Chi, J. Fan, C. Ma, accepted to Mathematical Programming
2. “Implicit regularization in nonconvex statistical estimation: ...”, C. Ma, K. Wang, Y. Chi, Y. Chen, arXiv:1711.10467
3. “Nonconvex optimization meets low-rank matrix factorization: An overview”, Y. Chi, Y. Lu, Y. Chen, arXiv:1809.09573