

# Random Initialization and Implicit Regularization in Nonconvex Statistical Estimation



Cong Ma

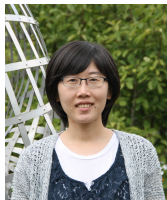
ORFE, Princeton University



Yuxin Chen  
Princeton EE



Kaizheng Wang  
Princeton ORFE



Yuejie Chi  
CMU ECE



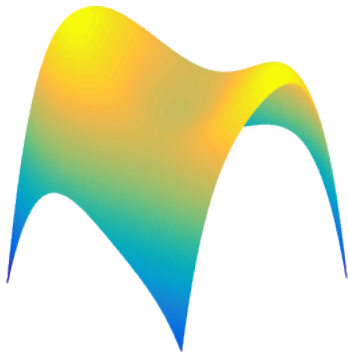
Jianqing Fan  
Princeton ORFE

# Nonconvex estimation problems are everywhere

---

Empirical risk minimization is usually nonconvex

$\text{minimize}_x \quad f(\mathbf{x}; \mathbf{y}) \quad \rightarrow \quad \text{loss function may be nonconvex}$



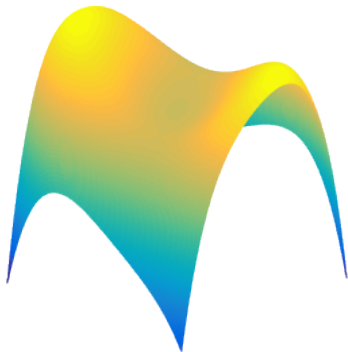
# Nonconvex estimation problems are everywhere

---

Empirical risk minimization is usually nonconvex

$\text{minimize}_x \quad f(\mathbf{x}; \mathbf{y}) \quad \rightarrow \quad \text{loss function may be nonconvex}$

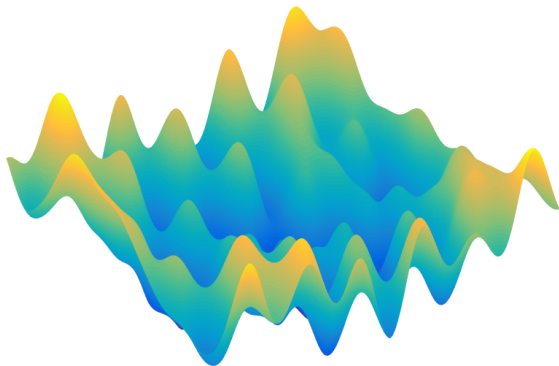
- low-rank matrix completion
- blind deconvolution
- dictionary learning
- mixture models
- deep learning
- ...





# 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)

# Nonconvex optimization may be super scary

---



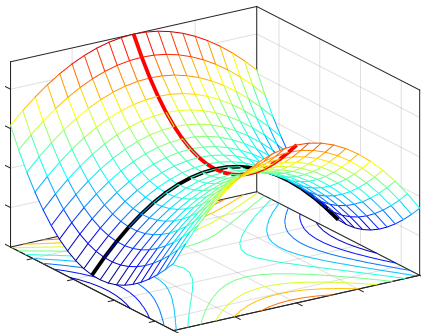
There may be bumps everywhere and exponentially many local optima

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

... but is sometimes much nicer than we think

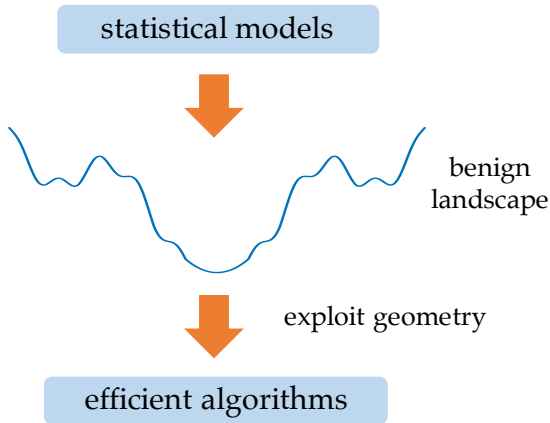
---

Under certain **statistical models**,  
we see benign global geometry: **no spurious local optima**



... but is sometimes much nicer than we think

---



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

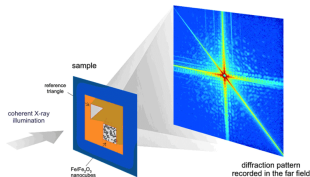
This talk: a case study — phase retrieval

# Missing phase problem

Detectors record **intensities** of diffracted rays

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

*Fig credit: Stanford SLAC*



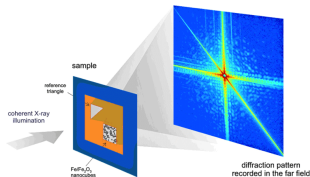
$$\text{intensity of electrical field: } |\hat{x}(f_1, f_2)|^2 = \left| \int x(t_1, t_2) e^{-i2\pi(f_1 t_1 + f_2 t_2)} dt_1 dt_2 \right|^2$$

# Missing phase problem

Detectors record **intensities** of diffracted rays

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

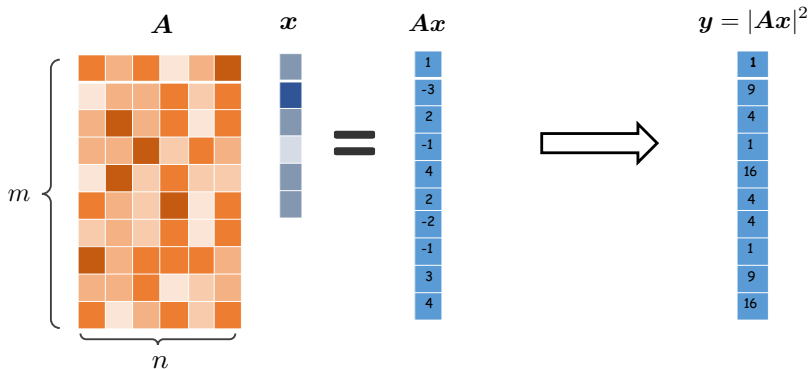
*Fig credit: Stanford SLAC*



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  $|\hat{x}(f_1, f_2)|^2$

# Solving quadratic systems of equations



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

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

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



# A natural least squares formulation

---

$$\text{given:} \quad y_k = |\mathbf{a}_k^\top \mathbf{x}|^2, \quad 1 \leq k \leq m$$

$\Downarrow$

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

# A natural least squares formulation

---

$$\text{given:} \quad y_k = |\mathbf{a}_k^\top \mathbf{x}|^2, \quad 1 \leq k \leq m$$

$\Downarrow$

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

- **pros:** often exact as long as sample size is sufficiently large

# A natural least squares formulation

---

$$\text{given:} \quad y_k = |\mathbf{a}_k^\top \mathbf{x}|^2, \quad 1 \leq k \leq m$$

$\Downarrow$

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

- **pros:** often exact as long as sample size is sufficiently large
- **cons:**  $f(\cdot)$  is highly nonconvex  
 $\longrightarrow$  *computationally challenging!*

## Wirtinger flow (Candès, Li, Soltanolkotabi '14)

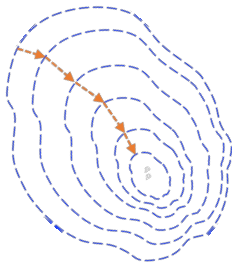
---

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

# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

---

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

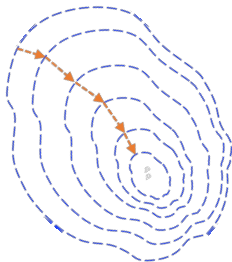


- **spectral initialization:**  $\mathbf{x}^0 \leftarrow$  leading eigenvector of certain data matrix

# Wirtinger flow (Candès, Li, Soltanolkotabi '14)

---

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

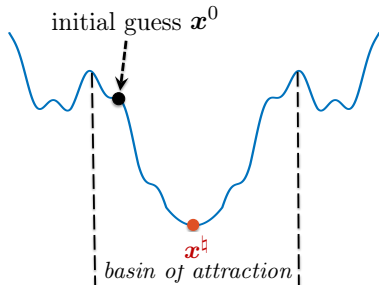


- **spectral initialization:**  $\mathbf{x}^0 \leftarrow$  leading eigenvector of certain data matrix
- **gradient descent:**

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

# Rationale of two-stage approach

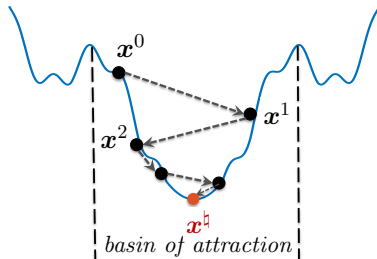
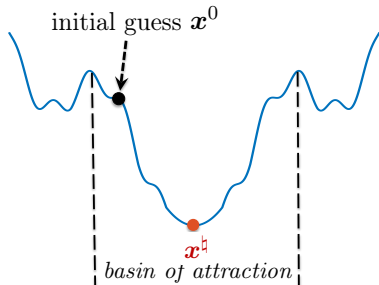
---



1. find an initial point within a local basin sufficiently close to  $x^*$

# Rationale of two-stage approach

---



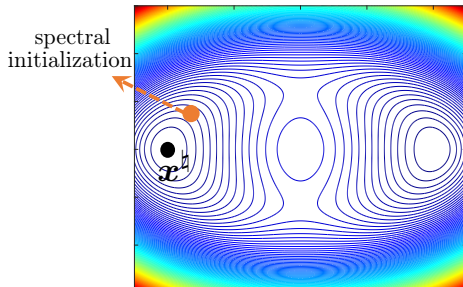
1. find an initial point within a local basin sufficiently close to  $x^*$
2. careful iterative refinement without leaving this local basin



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

# Initialization

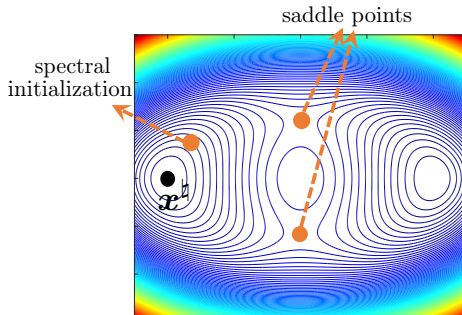
---



- spectral initialization gets us reasonably close to truth

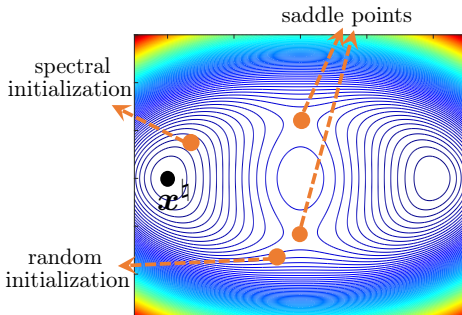
# Initialization

---



- spectral initialization gets us reasonably close to truth
- cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

# Initialization

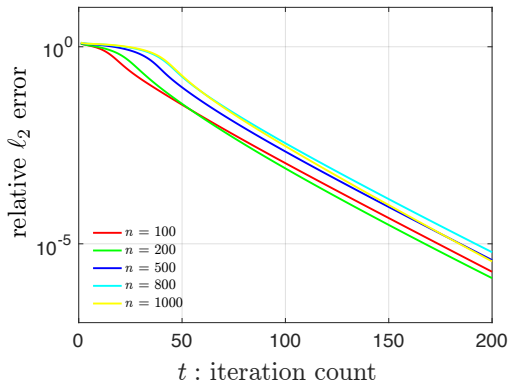


- spectral initialization gets us reasonably close to truth
- cannot initialize GD from anywhere, e.g. it might get stucked at local stationary points (e.g. saddle points)

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

# Numerical efficiency of randomly initialized GD

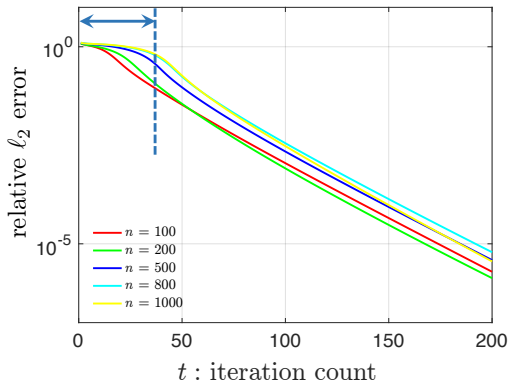
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



# Numerical efficiency of randomly initialized GD

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

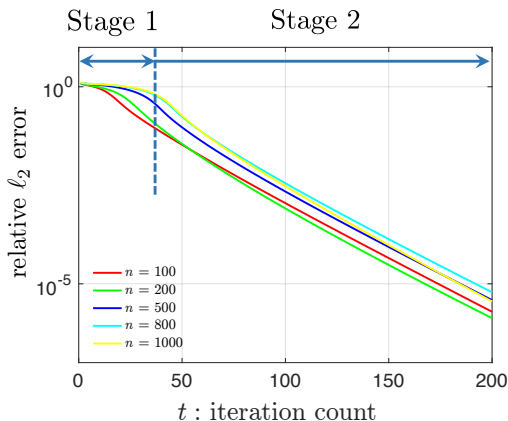
Stage 1



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

# Numerical efficiency of randomly initialized GD

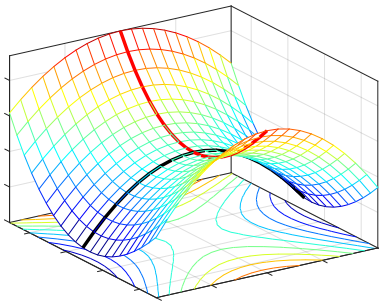
$$\eta_t = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



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

# What does prior theory say?

---

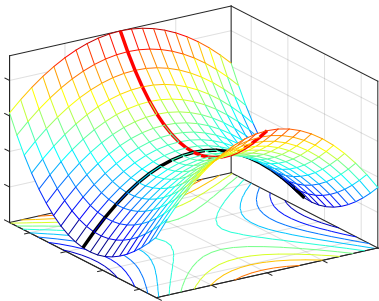


- no spurious local mins (Sun et al. '16)



# What does prior theory say?

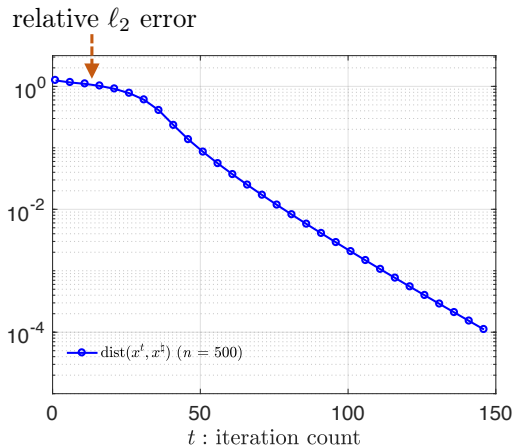
---



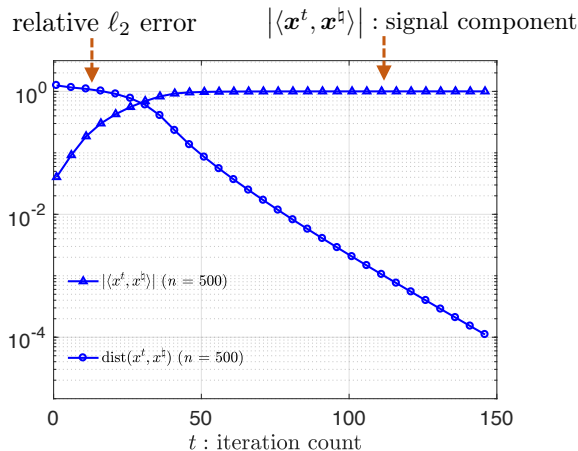
- no spurious local mins (Sun et al. '16)
- GD with random initialization converges to global min **almost surely** (Lee et al. '16)

No convergence rate guarantees for vanilla GD!

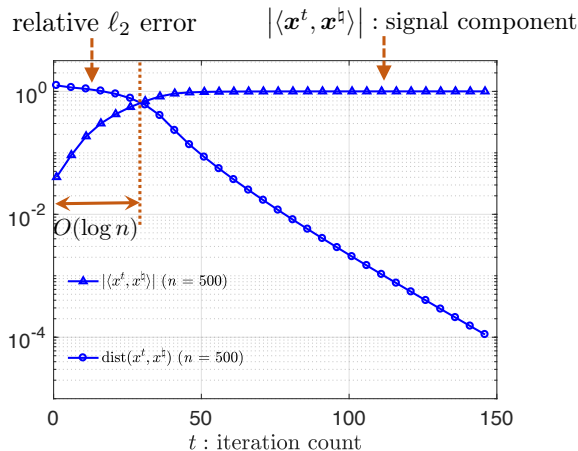
# Exponential growth of signal strength in Stage 1



# Exponential growth of signal strength in Stage 1

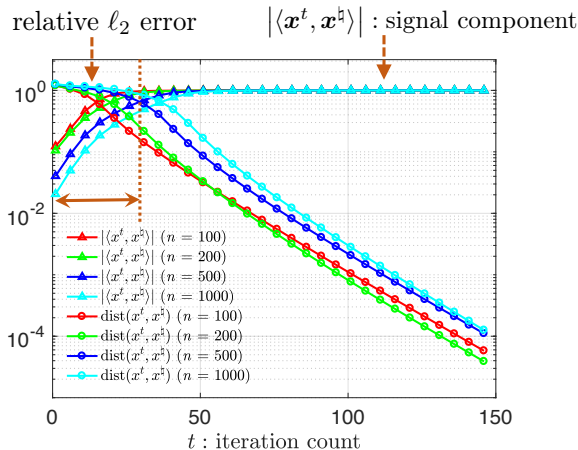


# Exponential growth of signal strength in Stage 1



Numerically,  $O(\log n)$  iterations are enough to enter local region

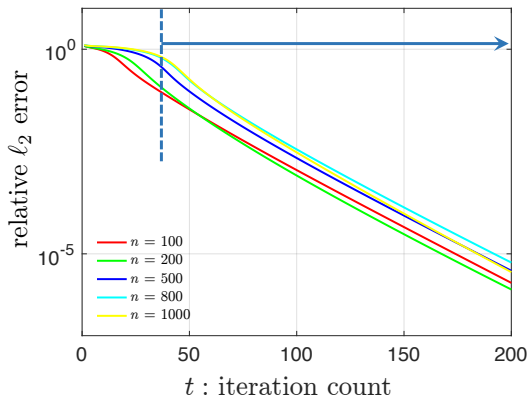
# Exponential growth of signal strength in Stage 1



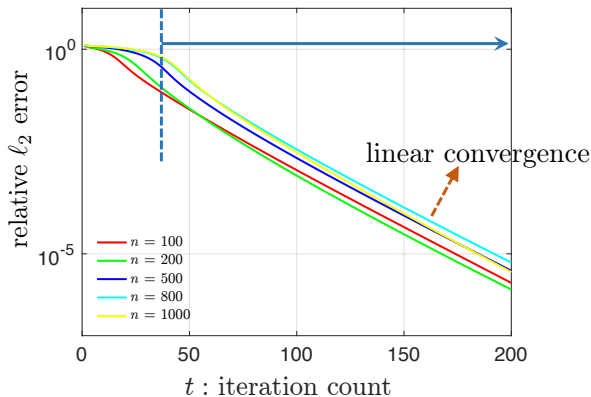
Numerically,  $O(\log n)$  iterations are enough to enter local region

## Linear / geometric convergence in Stage 2

---



## Linear / geometric convergence in Stage 2



Numerically, GD converges linearly within local region

# Experiments on images

---



- coded diffraction patterns
- $\mathbf{x}^b \in \mathbb{R}^{256 \times 256}$
- $m/n = 12$



# GD with random initialization

---

$\mathbf{x}^t$   
GD iterate

$\langle \mathbf{x}^t, \mathbf{x}^{\natural} \rangle \mathbf{x}^{\natural}$   
signal component

$\mathbf{x}^t - \langle \mathbf{x}^t, \mathbf{x}^{\natural} \rangle \mathbf{x}^{\natural}$   
perpendicular component

*use Adobe to view the animation*

## Exponential growth of “signal-to-noise” ratio

---

$$\frac{|\langle \mathbf{x}^t, \mathbf{x}^\natural \rangle|}{\|\mathbf{x}^t\|} \rightarrow \text{signal component}$$

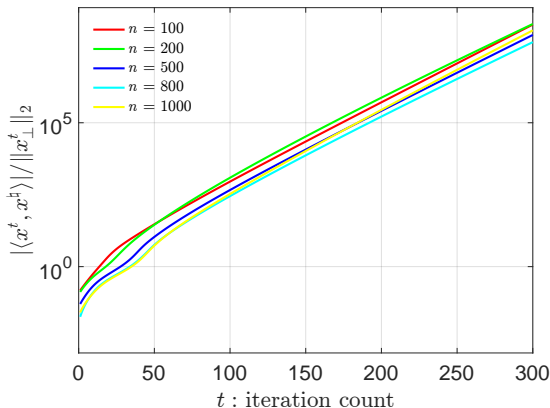
# Exponential growth of “signal-to-noise” ratio

---

$$\frac{|\langle \mathbf{x}^t, \mathbf{x}^q \rangle|}{\underbrace{\|\mathbf{x}^t - \langle \mathbf{x}^t, \mathbf{x}^q \rangle \mathbf{x}^q\|_2}_{:= \mathbf{x}_\perp^t}} \quad \begin{array}{l} \rightarrow \text{ signal component} \\ \rightarrow \text{ residual component} \end{array}$$

# Exponential growth of “signal-to-noise” ratio

$$\frac{|\langle \mathbf{x}^t, \mathbf{x}^\natural \rangle|}{\underbrace{\|\mathbf{x}^t - \langle \mathbf{x}^t, \mathbf{x}^\natural \rangle \mathbf{x}^\natural\|_2}_{:= \mathbf{x}_\perp^t}} \quad \begin{array}{l} \rightarrow \text{signal component} \\ \rightarrow \text{residual component} \end{array}$$



# Theoretical guarantees

---

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

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

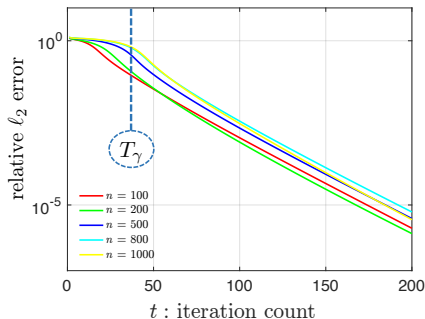
*Under i.i.d. Gaussian design, GD with  $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1}\mathbf{I}_n)$  achieves*

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

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

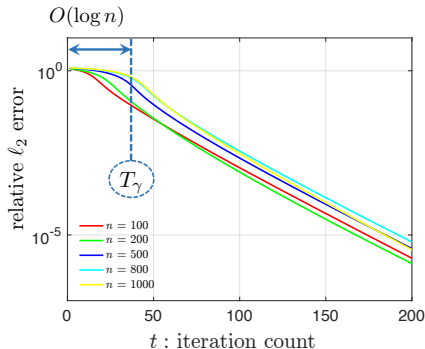
# Theoretical guarantees

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



# Theoretical guarantees

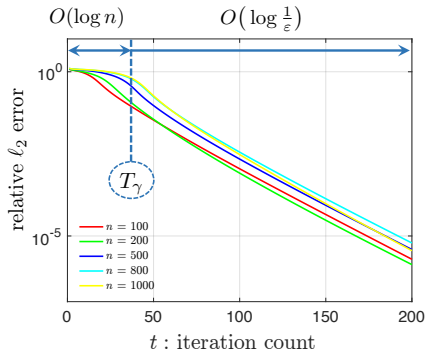
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\natural\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Theoretical guarantees

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

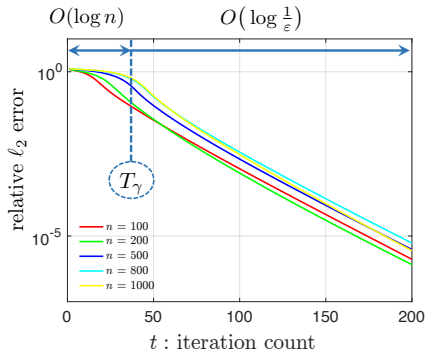


- *Stage 1*: takes  $O(\log n)$  iterations to reach  $\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \gamma$
- *Stage 2*: linear convergence



# Theoretical guarantees

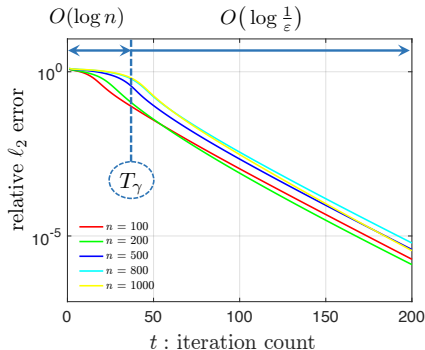
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\natural\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Theoretical guarantees

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



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

# Comparison with prior theory

---

Iteration complexity:

	prior theory	our theory
<b>Stage 1:</b> random init $\rightarrow$ local region	almost surely (Lee et al. '16)	$O(\log n)$
<b>Stage 2:</b> local refinement		

# Comparison with prior theory

---

## Iteration complexity:

	<b>prior theory</b>	<b>our theory</b>
<b>Stage 1:</b> random init $\rightarrow$ local region	<b>almost surely</b> (Lee et al. '16)	$O(\log n)$
<b>Stage 2:</b> local refinement	$O(n \log \frac{1}{\varepsilon})$ (Candes et al. '14)	$O(\log \frac{1}{\varepsilon})$

**Stage 1: random initialization → local region**

# What if we have infinite samples?

---

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

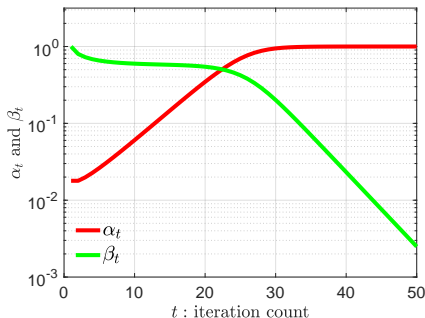
**Population level (infinite samples)**

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

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

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

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

2-parameter dynamics

## Back to finite-sample analysis

---

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



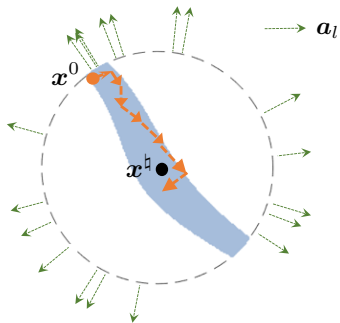
## Back to finite-sample analysis

---

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

# Back to finite-sample analysis

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



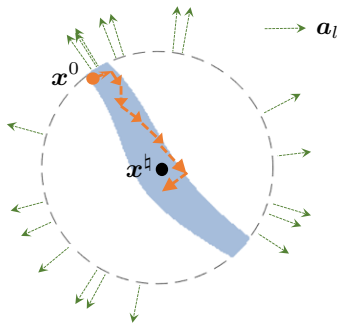
a region with well-controlled

$\mathbf{r}(\mathbf{x})$

- population-level analysis holds *approximately* if  $\mathbf{r}(\mathbf{x}^t) \ll \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t)$

# Back to finite-sample analysis

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

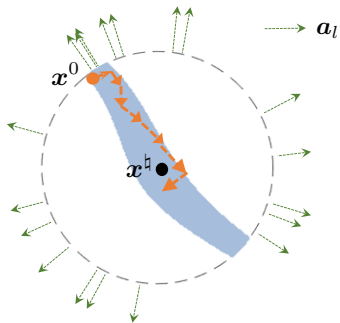


a region with well-controlled  
 $\mathbf{r}(\mathbf{x})$

- population-level analysis holds *approximately* if  $\mathbf{r}(\mathbf{x}^t) \ll \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t)$
- $\mathbf{r}(\mathbf{x}^t)$  is well-controlled if  $\mathbf{x}^t$  is independent of  $\{\mathbf{a}_k\}$

# Back to finite-sample analysis

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



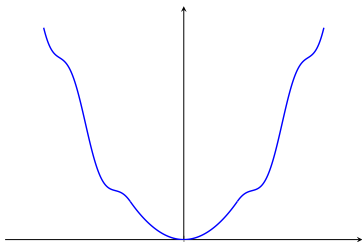
a region with well-controlled  
 $\mathbf{r}(\mathbf{x})$

- population-level analysis holds *approximately* if  $\mathbf{r}(\mathbf{x}^t) \ll \mathbf{x}^t - \eta \nabla F(\mathbf{x}^t)$
- $\mathbf{r}(\mathbf{x}^t)$  is well-controlled if  $\mathbf{x}^t$  is independent of  $\{\mathbf{a}_k\}$
- **key analysis ingredient:** show  $\mathbf{x}^t$  is “nearly-independent” of each  $\mathbf{a}_k$

## **Stage 2: local refinement**

# Gradient descent theory revisited

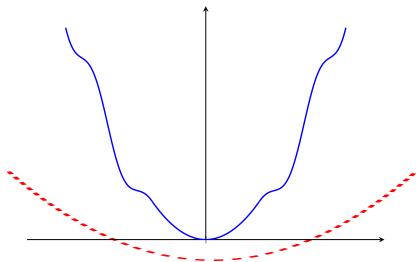
---



Two standard conditions that enable geometric convergence of GD

# Gradient descent theory revisited

---

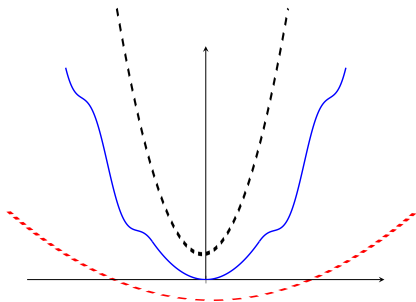


Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)

# Gradient descent theory revisited

---



Two standard conditions that enable geometric convergence of GD

- (local) restricted strong convexity (or regularity condition)
- (local) smoothness

$$\nabla^2 f(\mathbf{x}) \succcurlyeq \mathbf{0} \quad \text{and} \quad \text{is well-conditioned}$$



# Gradient descent theory revisited

---

$f$  is said to be  $\alpha$ -strongly convex and  $\beta$ -smooth if

$$\mathbf{0} \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

$\ell_2$  **error contraction:** GD with  $\eta = 1/\beta$  obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

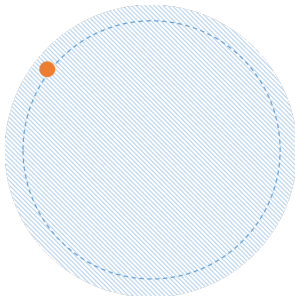
# Gradient descent theory revisited

---

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$



region of local strong convexity + smoothness



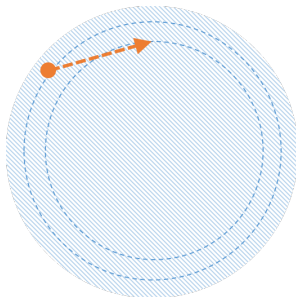
# Gradient descent theory revisited

---

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$



region of local strong convexity + smoothness



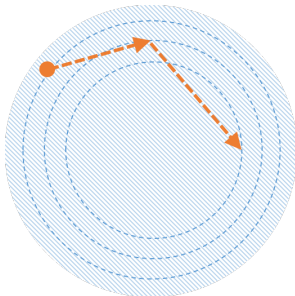
# Gradient descent theory revisited

---

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{\natural}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^{\natural}\|_2$$



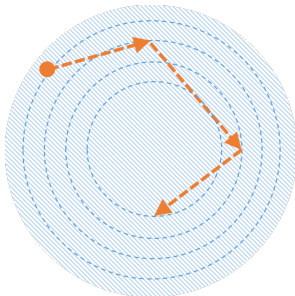
region of local strong convexity + smoothness



# Gradient descent theory revisited

$$\|\mathbf{x}^{t+1} - \mathbf{x}^{\natural}\|_2 \leq (1 - \alpha/\beta) \|\mathbf{x}^t - \mathbf{x}^{\natural}\|_2$$

- region of local strong convexity + smoothness



# Gradient descent theory revisited

---

$$0 \preceq \alpha I \preceq \nabla^2 f(\mathbf{x}) \preceq \beta I, \quad \forall \mathbf{x}$$

$\ell_2$  error contraction: GD with  $\eta = 1/\beta$  obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

- Condition number  $\beta/\alpha$  determines rate of convergence

# Gradient descent theory revisited

---

$$0 \preceq \alpha \mathbf{I} \preceq \nabla^2 f(\mathbf{x}) \preceq \beta \mathbf{I}, \quad \forall \mathbf{x}$$

$\ell_2$  **error contraction:** GD with  $\eta = 1/\beta$  obeys

$$\|\mathbf{x}^{t+1} - \mathbf{x}^\natural\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right) \|\mathbf{x}^t - \mathbf{x}^\natural\|_2$$

- Condition number  $\beta/\alpha$  determines rate of convergence
- Attains  $\varepsilon$ -accuracy within  $O(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon})$  iterations

## What does this optimization theory say about WF?

---

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



# What does this optimization theory say about WF?

---

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)

$$\mathbb{E}[\nabla^2 f(\mathbf{x})] = \underbrace{3 \left( \|\mathbf{x}\|_2^2 \mathbf{I} + 2\mathbf{x}\mathbf{x}^\top \right) - \left( \|\mathbf{x}^\natural\|_2^2 \mathbf{I} + 2\mathbf{x}^\natural \mathbf{x}^{\natural\top} \right)}_{\text{locally positive definite and well-conditioned}}$$

**Consequence:** Given good initialization, WF converges within  $O(\log \frac{1}{\varepsilon})$  iterations if  $m \rightarrow \infty$

# What does this optimization theory say about WF?

---

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

**Finite-sample level** ( $m \asymp n \log n$ )

$$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$$

# What does this optimization theory say about WF?

---

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

**Finite-sample level** ( $m \asymp n \log n$ )

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$  but ill-conditioned (even locally)  
condition number  $\asymp n$

# What does this optimization theory say about WF?

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

Finite-sample level ( $m \asymp n \log n$ )

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$  but ill-conditioned (even locally)  
condition number  $\asymp n$

**Consequence (Candès et al '14):** WF attains  $\varepsilon$ -accuracy within  $O(n \log \frac{1}{\varepsilon})$  iterations if  $m \asymp n \log n$

# What does this optimization theory say about WF?

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

Finite-sample level ( $m \asymp n \log n$ )

$\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$  but ill-conditioned (even locally)  
condition number  $\asymp n$

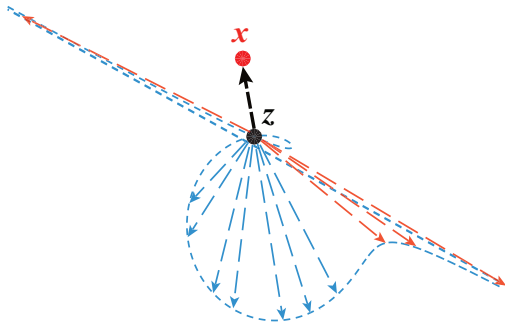
**Consequence (Candès et al '14):** WF attains  $\varepsilon$ -accuracy within  $O(n \log \frac{1}{\varepsilon})$  iterations if  $m \asymp n \log n$

*Too slow ... can we accelerate it?*

## Improvement: truncated WF (Chen, Candès '15)

---

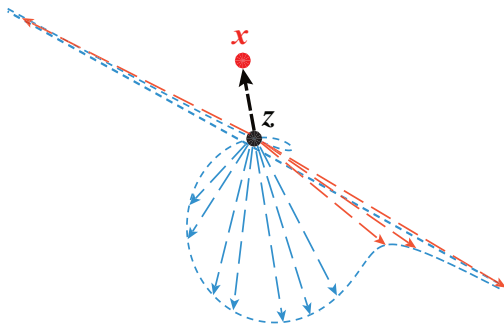
Regularize / trim gradient components to accelerate convergence



## Improvement: truncated WF (Chen, Candès '15)

---

Regularize / trim gradient components to accelerate convergence



But it still needs certain spectral initialization ...

# Recall

---

WF converges in  $O(n)$  iterations



# Recall

---

WF converges in  $O(n)$  iterations



Step size taken to be  $\eta_t = O(1/n)$

# Recall

---

WF converges in  $O(n)$  iterations



Step size taken to be  $\eta_t = O(1/n)$



This choice is suggested by **generic** optimization theory

# Recall

---

WF converges in  $O(n)$  iterations



Step size taken to be  $\eta_t = O(1/n)$



This choice is suggested by **worst-case** optimization theory

# Recall

---

WF converges in  $O(n)$  iterations



Step size taken to be  $\eta_t = O(1/n)$



This choice is suggested by **worst-case** optimization theory



Does it capture what really happens?

## A second look at gradient descent theory

---

Which region enjoys both strong convexity and smoothness?

$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[ 3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x})^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

## A second look at gradient descent theory

---

Which region enjoys both strong convexity and smoothness?

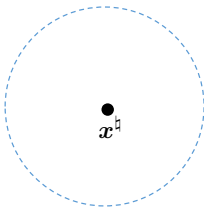
$$\nabla^2 f(\mathbf{x}) = \frac{1}{m} \sum_{k=1}^m \left[ 3(\mathbf{a}_k^\top \mathbf{x})^2 - (\mathbf{a}_k^\top \mathbf{x})^2 \right] \mathbf{a}_k \mathbf{a}_k^\top$$

- Not smooth if  $\mathbf{x}$  and  $\mathbf{a}_k$  are too close (coherent)

## A second look at gradient descent theory

---

Which region enjoys both strong convexity and smoothness?

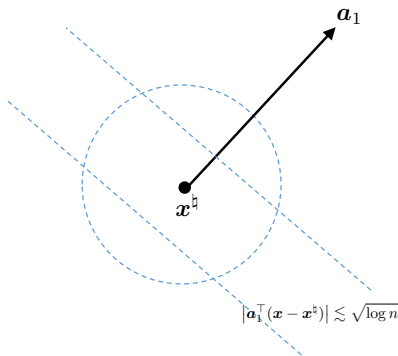


- $x$  is not far away from  $x^h$

## A second look at gradient descent theory

---

Which region enjoys both strong convexity and smoothness?

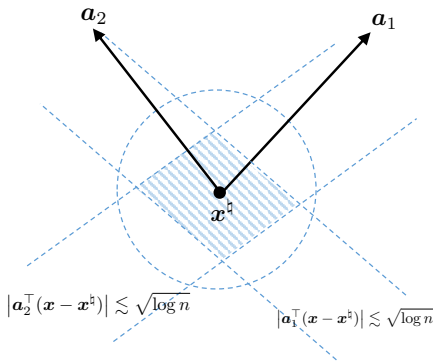


- $x$  is not far away from  $x^{\natural}$
- $x$  is incoherent w.r.t. sampling vectors (incoherence region)



## A second look at gradient descent theory

Which region enjoys both strong convexity and smoothness?



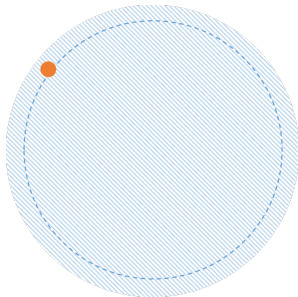
- $x$  is not far away from  $x^\dagger$
- $x$  is incoherent w.r.t. sampling vectors (incoherence region)

## A second look at gradient descent theory

---



region of local strong convexity + smoothness



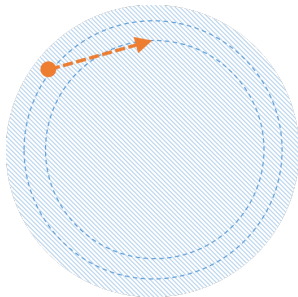
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

## A second look at gradient descent theory

---



region of local strong convexity + smoothness



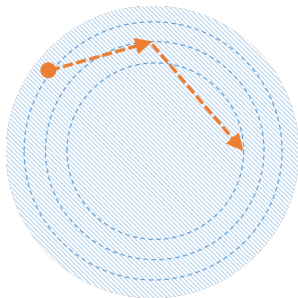
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

## A second look at gradient descent theory

---



region of local strong convexity + smoothness



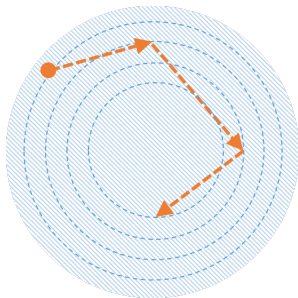
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

## A second look at gradient descent theory

---



region of local strong convexity + smoothness



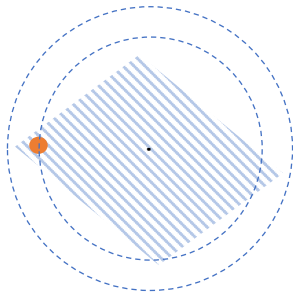
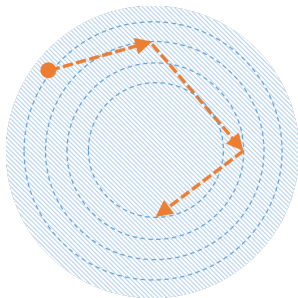
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# A second look at gradient descent theory

---



region of local strong convexity + smoothness



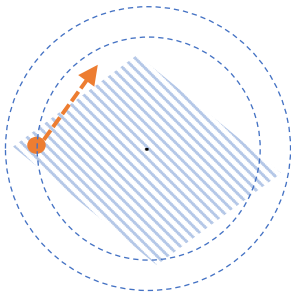
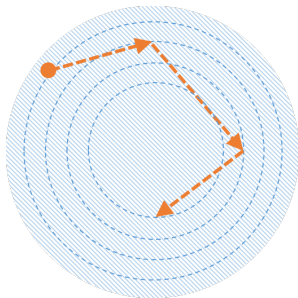
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# A second look at gradient descent theory

---



region of local strong convexity + smoothness

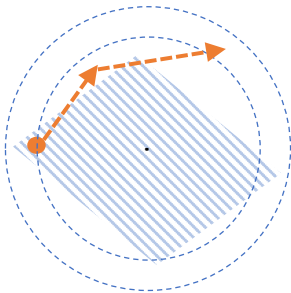
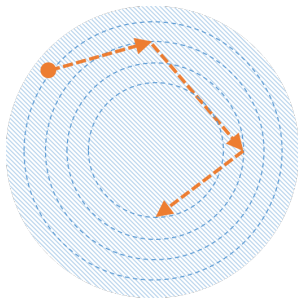


- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# A second look at gradient descent theory



region of local strong convexity + smoothness



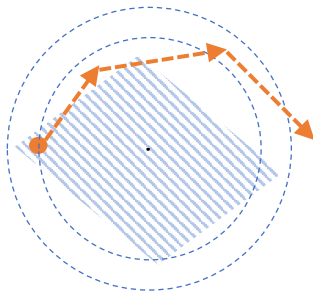
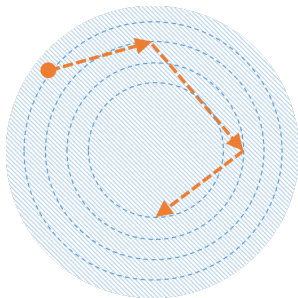
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region



# A second look at gradient descent theory



region of local strong convexity + smoothness

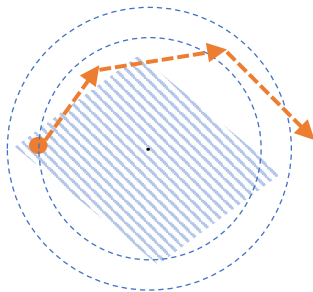
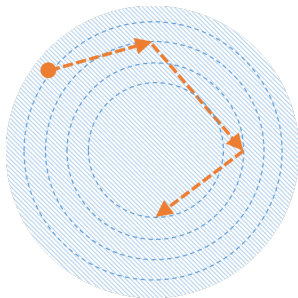


- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region

# A second look at gradient descent theory



region of local strong convexity + smoothness



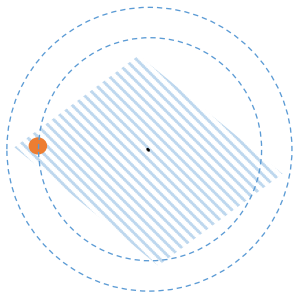
- Prior theory only ensures that iterates remain in  $\ell_2$  ball but not incoherence region
- *Prior theory enforces regularization to promote incoherence*

# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

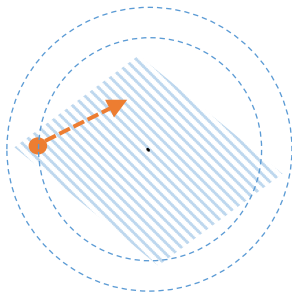


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

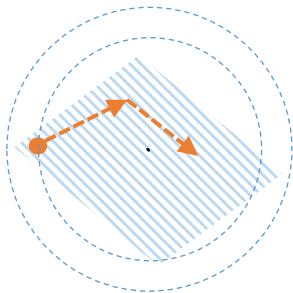


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

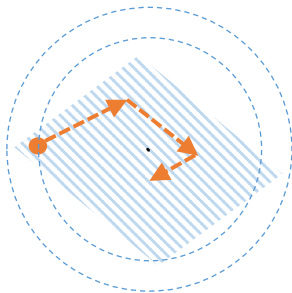


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness

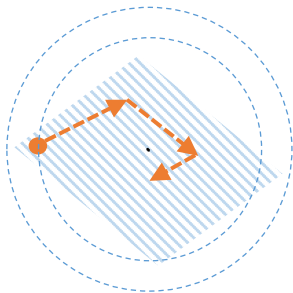


# Our findings: GD is implicitly regularized

---



region of local strong convexity + smoothness



GD implicitly forces iterates to remain **incoherent**

## Theoretical guarantees for Stage 2

---

### Theorem 2 (Phase retrieval)

*Under i.i.d. Gaussian design, GD with random initialization achieves for  $t \geq T_\gamma + 1$*

- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$  (incoherence)



## Theoretical guarantees for Stage 2

---

### Theorem 2 (Phase retrieval)

*Under i.i.d. Gaussian design, GD with random initialization achieves for  $t \geq T_\gamma + 1$*

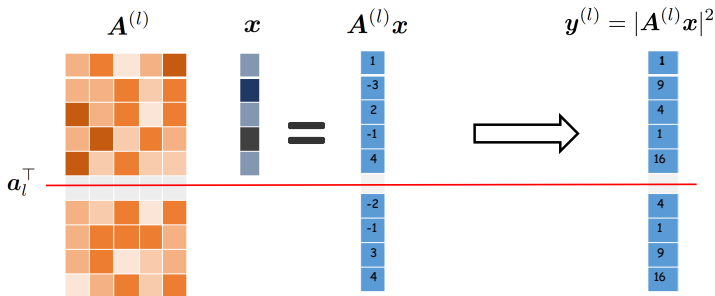
- $\max_k |\mathbf{a}_k^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \lesssim \sqrt{\log n} \|\mathbf{x}^\natural\|_2$  (incoherence)
- $\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \lesssim (1 - \frac{\eta}{2})^{t-T_\gamma} \cdot \gamma \|\mathbf{x}^\natural\|_2$  (linear convergence)

*provided that step size  $\eta \asymp c$  and sample size  $m \gtrsim n \text{ poly log } m$ .*

# Key ingredient: leave-one-out analysis

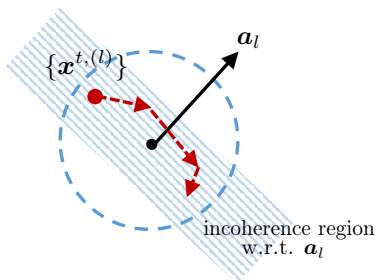
---

For each  $1 \leq l \leq m$ , introduce leave-one-out iterates  $x^{t,(l)}$  by dropping  $l$ th measurement



# Key ingredient: leave-one-out analysis

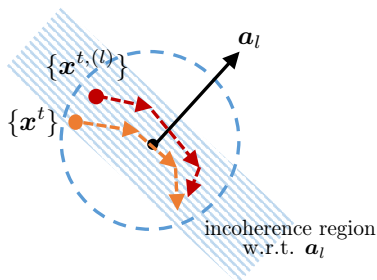
---



- Leave-one-out iterates  $\{x^{t,(l)}\}$  are independent of  $a_l$ , and are hence **incoherent** w.r.t.  $a_l$  with high prob.

# Key ingredient: leave-one-out analysis

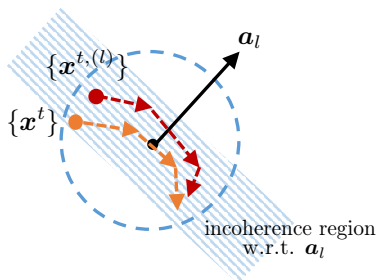
---



- Leave-one-out iterates  $\{x^{t,(l)}\}$  are independent of  $a_l$ , and are hence **incoherent** w.r.t.  $a_l$  with high prob.
- Leave-one-out iterates  $x^{t,(l)} \approx$  true iterates  $x^t$

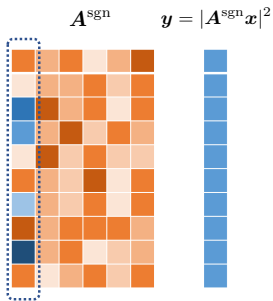
# Key ingredient: leave-one-out analysis

---

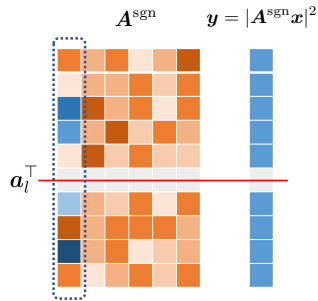


- Leave-one-out iterates  $\{\mathbf{x}^{t,(l)}\}$  are independent of  $\mathbf{a}_l$ , and are hence **incoherent** w.r.t.  $\mathbf{a}_l$  with high prob.
- Leave-one-out iterates  $\mathbf{x}^{t,(l)} \approx$  true iterates  $\mathbf{x}^t$
- $|\mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^\natural)| \leq |\mathbf{a}_l^\top (\mathbf{x}^{t,(l)} - \mathbf{x}^\natural)| + |\mathbf{a}_l^\top (\mathbf{x}^t - \mathbf{x}^{t,(l)})|$

## Other leave-one-out sequences



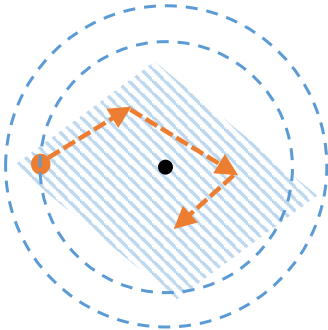
$x^{t,\text{sgn}}$ : indep. of sign info of  $\{a_{i,1}\}$



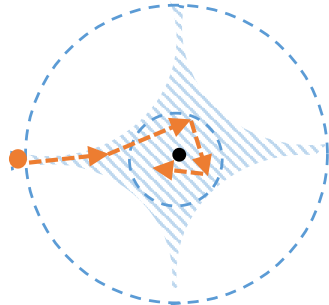
$x^{t,\text{sgn},(l)}$ : indep. of both sign info of  $\{a_{i,1}\}$  and  $a_l$

# Incoherence region in high dimensions

---



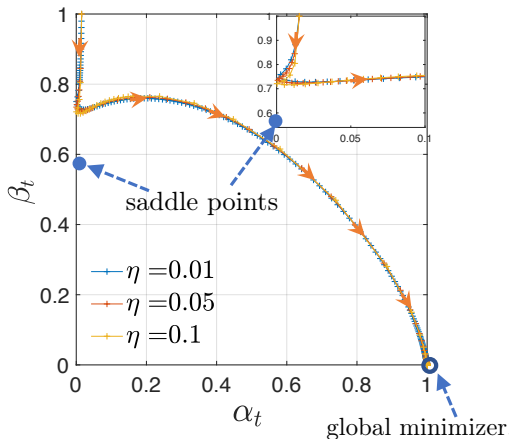
2-dimensional



high-dimensional (mental representation)

incoherence region is vanishingly small

# Saddle-escaping schemes?



Randomly initialized GD never hits saddle points in phase retrieval!



## Other saddle-escaping schemes

---

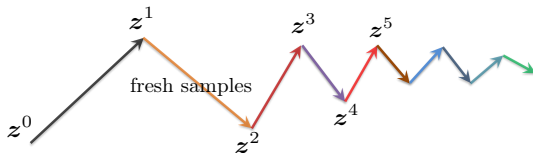
	iteration complexity	num of iterations needed to escape saddles	local iteration complexity
<b>Trust-region</b> (Sun et al. '16)	$n^7 + \log \log \frac{1}{\epsilon}$	$n^7$	$\log \log \frac{1}{\epsilon}$
<b>Perturbed GD</b> (Jin et al. '17)	$n^3 + n \log \frac{1}{\epsilon}$	$n^3$	$n \log \frac{1}{\epsilon}$
<b>Perturbed accelerated GD</b> (Jin et al. '17)	$n^{2.5} + \sqrt{n} \log \frac{1}{\epsilon}$	$n^{2.5}$	$\sqrt{n} \log \frac{1}{\epsilon}$
<b>GD (ours)</b> (Chen et al. '18)	$\log n + \log \frac{1}{\epsilon}$	$\log n$	$\log \frac{1}{\epsilon}$

Generic optimization theory yields highly suboptimal convergence guarantees

# No need of sample splitting

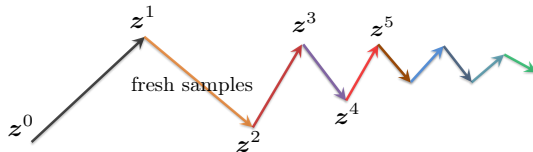
---

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

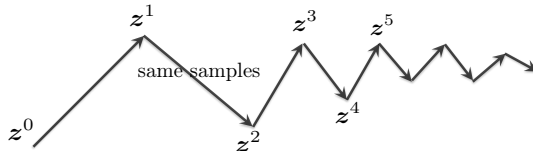


# No need of sample splitting

- 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



# Summary

---

- **Blessings of statistical models:** GD with random initialization converges fast
- **Implicit regularization:** vanilla gradient descent automatically focuses iterates to stay *incoherent*

## Paper:

“Implicit regularization in nonconvex statistical estimation: Gradient descent converges linearly for phase retrieval, matrix completion, and blind deconvolution”, Cong Ma, Kaizheng Wang, Yuejie Chi, Yuxin Chen, arXiv:1711.10467

“Gradient Descent with Random Initialization: Fast Global Convergence for Nonconvex Phase Retrieval”, Yuxin Chen, Yuejie Chi, Jianqing Fan, Cong Ma  
arXiv:XXXXXX