

# **Analysis of global convergence: random initialization**



Cong Ma

University of Chicago, Winter 2024

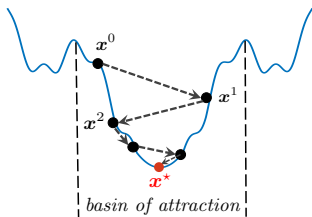
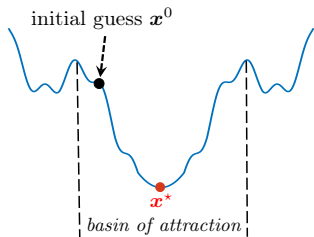
# Outline

---

- Strict saddle property
- Global landscape analysis: matrix sensing
- Gradient descent with random initialization: phase retrieval
- Generic saddle-escaping algorithms

# Rationale of two-stage approach

---

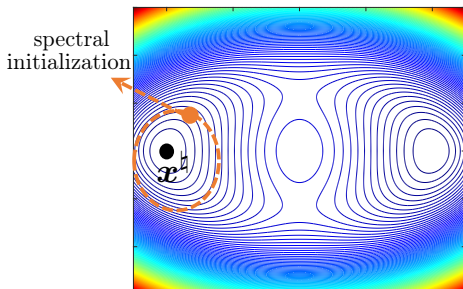


1. initialize within local basin sufficiently close to  $x^*$   
(restricted) strongly convex
2. iterative refinement

Is careful initialization necessary for fast convergence?

# Initialization

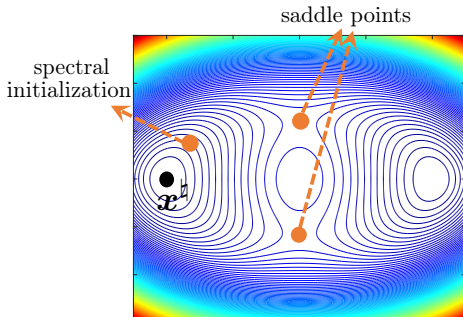
---



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

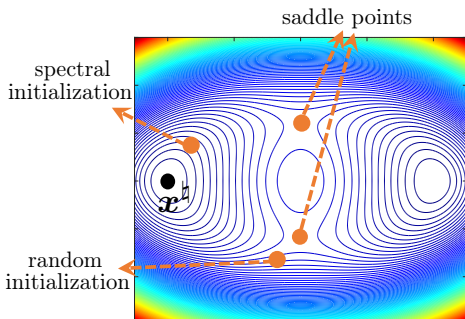
# Initialization

---



- 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**?

## **Generic saddle-escaping algorithms**



# Strict saddle property: qualitative version

---

All critical points can be classified into two categories

- local minimizers
- strict saddle points: Hessian has a strictly negative eigenvalue

Let  $\mathbf{x}$  be a critical point. Taylor expansion yields

$$f(\mathbf{x} + \Delta) \approx f(\mathbf{x}) + \frac{1}{2} \Delta^\top \nabla^2 f(\mathbf{x}) \Delta$$

# GD converges to local minimizers

---

## Theorem 10.1

*Consider any twice continuously differentiable function  $f$  that satisfies the strict saddle property. If  $\eta < 1/\beta$  with  $\beta$  the smoothness parameter, then GD with a random initialization converges to a local minimizer or  $-\infty$  almost surely.*

- This also holds for other optimization algorithms
- Exponential time for GD to converge in the worst case

**An example: low-rank matrix sensing**

# Low-rank matrix sensing

---

- Groundtruth: rank- $r$  psd matrix  $\mathbf{M}^* = \mathbf{X}^* \mathbf{X}^{*\top} \in \mathbb{R}^{n \times n}$
- Observations:

$$y_i = \langle \mathbf{A}_i, \mathbf{M}^* \rangle, \quad \text{for } 1 \leq i \leq m$$

- Goal: recover  $\mathbf{M}^*$  based on linear measurements  $\{\mathbf{A}_i, y_i\}_{1 \leq i \leq m}$

# Restricted isometry property (RIP)

---

Define linear operator  $\mathcal{A} : \mathbb{R}^{n_1 \times n_2} \mapsto \mathbb{R}^m$  to be

$$\mathcal{A}(\mathbf{M}) = [m^{-1/2} \langle \mathbf{A}_i, \mathbf{M} \rangle]_{1 \leq i \leq m}$$

## Definition 10.2

The operator  $\mathcal{A}$  is said to satisfy  $r$ -RIP with RIP constant  $\delta_r < 1$  if

$$(1 - \delta_r) \|\mathbf{M}\|_F^2 \leq \|\mathcal{A}(\mathbf{M})\|_2^2 \leq (1 + \delta_r) \|\mathbf{M}\|_F^2$$

holds simultaneously for all  $\mathbf{M}$  of rank at most  $r$ .

# An optimization-based method

---

Then least-squares estimation yields

$$\underset{\mathbf{X} \in \mathbb{R}^{n \times r}}{\text{minimize}} \quad f(\mathbf{X}) = \frac{1}{4m} \sum_{i=1}^m \left( \langle \mathbf{A}_i, \mathbf{X} \mathbf{X}^\top \rangle - y_i \right)^2$$

# Global landscape

---

## Theorem 10.3

Assume that the measurement operator  $\mathcal{A}$  satisfies  $2r$ -RIP with RIP constant  $\delta_{2r} \leq 1/10$ . Then for the matrix sensing objective, one has

- For any critical point  $\mathbf{U}$  that is not a local minimum, one has  $\lambda_{\min}(\nabla^2 f(\mathbf{U})) \leq -2/5\sigma_r(\mathbf{M}^*)$ ;
- All local minimizers are global.

- Matrix sensing obeys strict saddle property
- In addition, all local minimizers are global — GD converges to global minimizer

# Strict saddle property: quantitative version

---

## Definition 10.4

A function  $f(\cdot)$  is said to satisfy the  $(\varepsilon, \gamma, \xi)$ -strict saddle property for some positive  $\varepsilon, \gamma, \xi$ , if for each  $\mathbf{x}$ , at least one of the following is true

- **(strong gradient)**  $\|\nabla f(\mathbf{x})\|_2 \geq \varepsilon$ ;
- **(negative curvature)**  $\lambda_{\min}(\nabla^2 f(\mathbf{x})) \leq -\gamma$ ;
- **(local minimum)** there exists a local minimum  $\mathbf{x}_\star$  such that  $\|\mathbf{x} - \mathbf{x}_\star\|_2 \leq \xi$ .



# Saddle-escaping algorithms

---

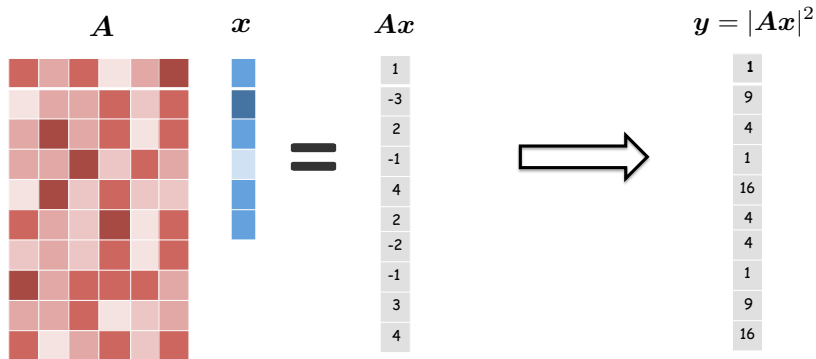
$$f(\mathbf{x} + \Delta) \approx f(\mathbf{x}) + \frac{1}{2} \Delta^\top \nabla^2 f(\mathbf{x}) \Delta$$

A rough categorization:

- Hessian-based algorithms
- Gradient-based algorithms

**Another example: phase retrieval**

# Solving quadratic systems of equations



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

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

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

# A natural least-squares formulation

---

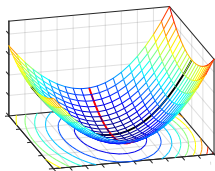
$$\text{given: } y_k = (\mathbf{a}_k^\top \mathbf{x}^*)^2, \quad 1 \leq k \leq m$$

↓

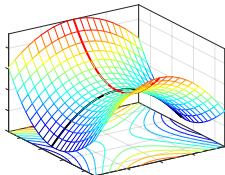
$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4m} \sum_{k=1}^m [(\mathbf{a}_k^\top \mathbf{x})^2 - y_k]^2$$

# What does prior theory say?

---



global minimum

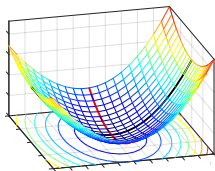


saddle point

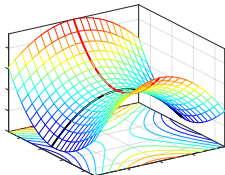
- **landscape:** no spurious local mins (Sun et al. '16)

# What does prior theory say?

---



global minimum

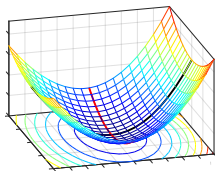


saddle point

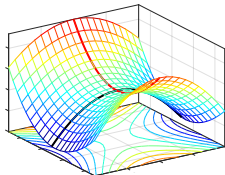
- **landscape:** no spurious local mins (Sun et al. '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

# What does prior theory say?

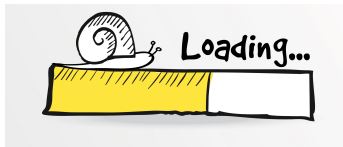
---



global minimum



saddle point

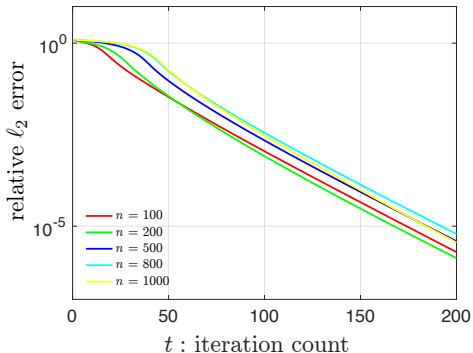


- **landscape:** no spurious local mins (Sun et al. '16)
- randomly initialized GD converges **almost surely** (Lee et al. '16)

“almost surely” might mean “takes forever”

# Numerical efficiency of randomly initialized GD

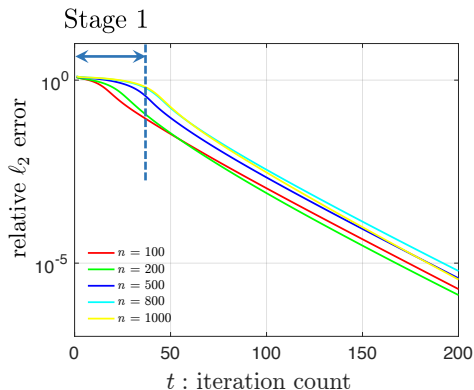
$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$





# Numerical efficiency of randomly initialized GD

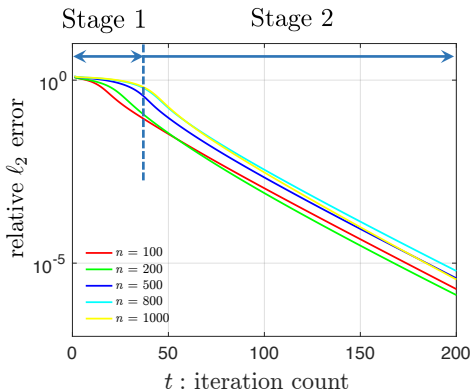
$$\eta = 0.1, \mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_n), m = 10n, \mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1} \mathbf{I}_n)$$



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**

## 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  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ :

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

### Theorem 10.5 (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  $\mathbf{a}_i \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(\mathbf{0}, \mathbf{I}_n)$ :

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

### Theorem 10.5 (Chen, Chi, Fan, Ma '18)

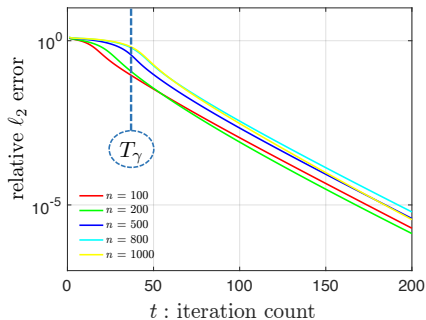
*Under i.i.d. Gaussian design, GD with  $\mathbf{x}^0 \sim \mathcal{N}(\mathbf{0}, n^{-1}\mathbf{I}_n)$  achieves*

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

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

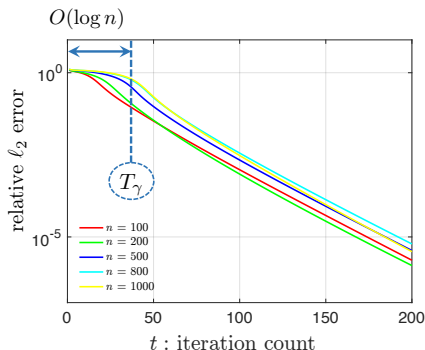
# Our theory

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



# Our theory

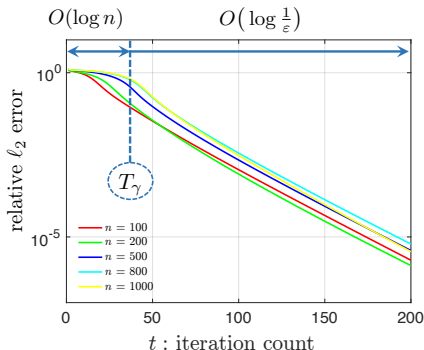
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\natural) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\natural\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Our theory

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

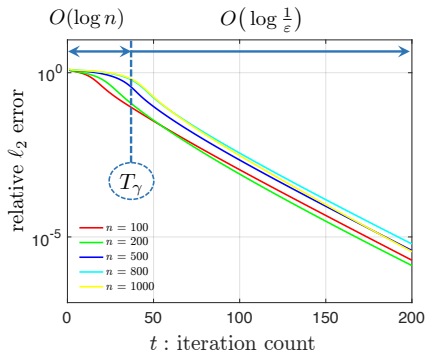


- *Stage 1*: takes  $O(\log n)$  iterations to reach  $\text{dist}(\mathbf{x}^t, \mathbf{x}^\dagger) \leq \gamma$
- *Stage 2*: linear (geometric) convergence



# Our theory

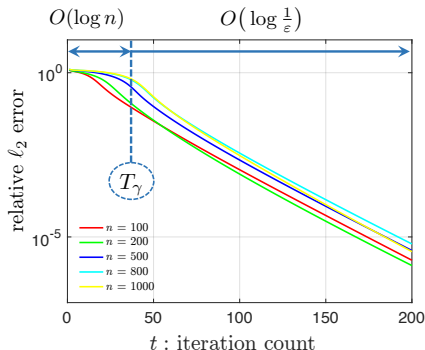
$$\text{dist}(\mathbf{x}^t, \mathbf{x}^\dagger) \leq \gamma(1 - \rho)^{t - T_\gamma} \|\mathbf{x}^\dagger\|_2, \quad t \geq T_\gamma \asymp \log n$$



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

# Our theory

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

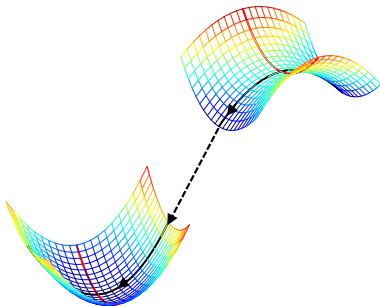


- *near-optimal computational cost:*
  - $O(\log n + \log \frac{1}{\epsilon})$  iterations to yield  $\epsilon$  accuracy
- *near-optimal sample size:*  $m \gtrsim n \text{poly} \log m$

# Generic algorithm design and analysis

---

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



Generic optimization theory yields highly suboptimal convergence guarantees

# What we have not discussed so far

---

- A lot of interesting problems that nonconvex optimization could work, e.g., robust PCA, tensor estimation, mixture models, etc.
- A lot of algorithms, e.g., expectation maximization, alternating minimization, scaledGD, etc.
- Inference for nonconvex estimators
- Connections between nonconvex and convex estimators