

Introduction to Nonconvex Optimization



Cong Ma

University of Chicago, Winter 2026

Roadmap

- Motivation: where nonconvexity appears and what guarantees we seek
- Part I: a concrete running example (rank-1 matrix factorization)
- Part II: generic local convergence tools (strong convexity, smoothness, RC)
- Part III: return to the example for local and global geometry

Where Nonconvex Optimization Appears

- Low-rank estimation: matrix factorization, PCA, matrix sensing
- Inverse problems: phase retrieval, blind deconvolution
- Representation learning: dictionary learning, tensor decomposition
- Deep learning: training neural networks

Different applications, but recurring questions: geometry, initialization, and algorithmic guarantees.

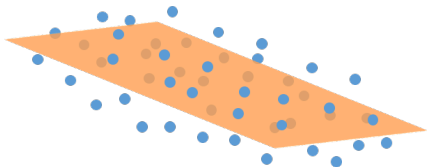
What Guarantees Do We Usually Seek?

- First-order guarantee: find x with $\|\nabla f(x)\|$ small
- Second-order guarantee: avoid strict saddles and reach local minima
- Global guarantee: reach a global minimizer (often under benign geometry)

- Typical intro hierarchy: stationary point \rightarrow local minimum \rightarrow global optimum
- Assumptions become stronger as guarantees get stronger
- In this lecture: understand when simple methods are still effective

A running example: rank-1 matrix factorization

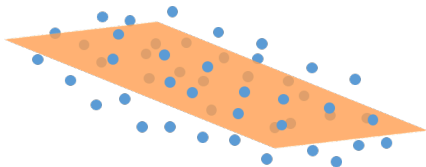
Principal Component Analysis



Given $\mathbf{M} \succeq \mathbf{0} \in \mathbb{R}^{n \times n}$ (not necessarily low-rank), find its best rank- r approximation:

$$\underbrace{\widehat{\mathbf{M}} = \arg \min_{\mathbf{Z}} \|\mathbf{Z} - \mathbf{M}\|_{\text{F}}^2 \quad \text{s.t.} \quad \text{rank}(\mathbf{Z}) \leq r}_{\text{nonconvex optimization!}}$$

Principal Component Analysis



This problem admits a closed-form solution.

- Let $\mathbf{M} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$ be eigen-decomposition of \mathbf{M} ($\lambda_1 \geq \dots \geq \lambda_r > \lambda_{r+1} \geq \lambda_n$), then

$$\widehat{\mathbf{M}} = \sum_{i=1}^r \lambda_i \mathbf{u}_i \mathbf{u}_i^\top$$

— *nonconvex, but tractable*

Optimization Viewpoint

If we factorize $\mathbf{Z} = \mathbf{X}\mathbf{X}^\top$ with $\mathbf{X} \in \mathbb{R}^{n \times r}$, then it leads to a nonconvex problem:

$$\text{minimize}_{\mathbf{X} \in \mathbb{R}^{n \times r}} f(\mathbf{X}) = \frac{1}{4} \|\mathbf{X}\mathbf{X}^\top - \mathbf{M}\|_{\text{F}}^2$$

To simplify exposition, set $r = 1$:

$$\text{minimize}_{\mathbf{x}} f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2$$

Why Is the Global Minimizer the Leading Eigenvector? (I)

For the rank-1 objective

$$f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2,$$

expand the Frobenius norm:

$$\begin{aligned} 4f(\mathbf{x}) &= \|\mathbf{x}\mathbf{x}^\top\|_{\text{F}}^2 - 2\langle \mathbf{x}\mathbf{x}^\top, \mathbf{M} \rangle + \|\mathbf{M}\|_{\text{F}}^2 \\ &= \|\mathbf{x}\|_2^4 - 2\mathbf{x}^\top \mathbf{M} \mathbf{x} + \|\mathbf{M}\|_{\text{F}}^2. \end{aligned}$$

Let $\mathbf{x} = r\mathbf{v}$ with $\|\mathbf{v}\|_2 = 1$. Then

$$4f(r\mathbf{v}) = r^4 - 2r^2\mathbf{v}^\top \mathbf{M} \mathbf{v} + \|\mathbf{M}\|_{\text{F}}^2.$$

For fixed r , minimizing f means maximizing $\mathbf{v}^\top \mathbf{M} \mathbf{v}$.

Why Is the Global Minimizer the Leading Eigenvector? (II)

By the Rayleigh quotient bound for $M \succeq 0$,

$$\mathbf{v}^\top M \mathbf{v} \leq \lambda_1,$$

with equality iff \mathbf{v} is a leading eigenvector \mathbf{u}_1 .

Thus the best direction is $\mathbf{v} = \pm \mathbf{u}_1$. Substituting gives

$$4f(r\mathbf{u}_1) = r^4 - 2\lambda_1 r^2 + \|M\|_F^2,$$

whose minimum is attained at $r^2 = \lambda_1$.

Therefore the global minimizers are

$$\mathbf{x}_{\text{opt}} = \pm \sqrt{\lambda_1} \mathbf{u}_1$$

(unique up to sign when $\lambda_1 > \lambda_2$).

Part I Takeaway

- Nonconvex objectives can still have transparent global structure
- In this example, the global minimizer is explicit: $\mathbf{x}_{\text{opt}} = \pm\sqrt{\lambda_1}\mathbf{u}_1$

- Remaining question: when do iterative algorithms provably reach such solutions?
- Next: build a generic toolkit to answer this

Part II: Generic Nonconvex Toolkit

Local geometry + algorithmic guarantees

Unconstrained Optimization

Consider an unconstrained optimization problem

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

- For simplicity, we assume $f(\mathbf{x})$ is twice differentiable
- We assume the minimizer \mathbf{x}_{opt} exists, i.e.,

$$\mathbf{x}_{\text{opt}} := \arg \min_{\mathbf{x}} f(\mathbf{x})$$

Critical/Stationary Points

Definition 1

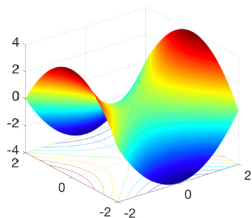
A first-order critical point of f satisfies

$$\nabla f(\mathbf{x}) = \mathbf{0}$$

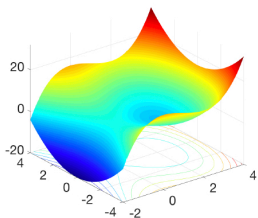
- If f is convex, any 1st-order critical point is a global minimizer
- Finding 1st-order stationary point is sufficient for convex optimization
- Example: gradient descent (GD)

What Changes in Nonconvex Optimization?

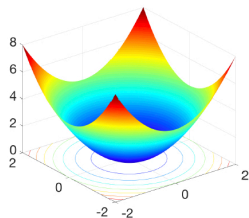
First-order critical points could be global min, local min, local max, saddle points...



(a) strict saddle



(b) local minimum



(c) global minimum

figure credit: Li et al. '16

Simple algorithms like GD can get stuck at undesired stationary points.

Types of Critical Points

Definition 2

A second-order critical point \mathbf{x} satisfies

$$\nabla f(\mathbf{x}) = \mathbf{0} \quad \text{and} \quad \nabla^2 f(\mathbf{x}) \succeq \mathbf{0}$$

For any first-order critical point \mathbf{x} :

- $\nabla^2 f(\mathbf{x}) \prec \mathbf{0}$ \rightarrow local maximum
- $\nabla^2 f(\mathbf{x}) \succ \mathbf{0}$ \rightarrow local minimum
- $\lambda_{\min}(\nabla^2 f(\mathbf{x})) < 0$ \rightarrow *strict saddle point*

What Structure Makes Nonconvex Problems Solvable?

- Good local curvature near the target (strong convexity + smoothness)
- Benign global geometry (no spurious minima / strict saddles)

Local Strong Convexity and Smoothness

Definition 3

A twice differentiable function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is said to be α -strongly convex in a set \mathcal{B} if for all $\mathbf{x} \in \mathcal{B}$

$$\nabla^2 f(\mathbf{x}) \succeq \alpha \mathbf{I}_n.$$

Definition 4

A twice differentiable function $f : \mathbb{R}^n \mapsto \mathbb{R}$ is said to be β -smooth in a set \mathcal{B} if for all $\mathbf{x} \in \mathcal{B}$

$$\|\nabla^2 f(\mathbf{x})\| \leq \beta.$$

Gradient Descent Theory Revisited

Gradient descent method with step size $\eta > 0$

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

Lemma 5

Suppose f is α -strongly convex and β -smooth in the local ball $\mathcal{B}_\delta(\mathbf{x}_{\text{opt}}) := \{\mathbf{x} \mid \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2 \leq \delta\}$. Running gradient descent from $\mathbf{x}^0 \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$ with $\eta = 1/\beta$ achieves linear convergence

$$\|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2 \leq \left(1 - \frac{\alpha}{\beta}\right)^t \|\mathbf{x}^0 - \mathbf{x}_{\text{opt}}\|_2, \quad t = 0, 1, 2, \dots$$

Implications

- Condition number β/α determines rate of convergence
- Attains ε -accuracy (i.e., $\|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2 \leq \varepsilon \|\mathbf{x}_{\text{opt}}\|_2$) within

$$O\left(\frac{\beta}{\alpha} \log \frac{1}{\varepsilon}\right)$$

iterations

- Needs initialization $\mathbf{x}^0 \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$: basin of attraction

Proof of Lemma 5

Since $\nabla f(\mathbf{x}_{\text{opt}}) = \mathbf{0}$, we can rewrite GD as

$$\begin{aligned}\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}} &= \mathbf{x}^t - \eta \nabla f(\mathbf{x}^t) - [\mathbf{x}_{\text{opt}} - \eta \nabla f(\mathbf{x}_{\text{opt}})] \\ &= \left[\mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) d\tau \right] (\mathbf{x}^t - \mathbf{x}_{\text{opt}}),\end{aligned}$$

where $\mathbf{x}(\tau) := \mathbf{x}_{\text{opt}} + \tau(\mathbf{x}^t - \mathbf{x}_{\text{opt}})$. By local strong convexity and smoothness, one has

$$\alpha \mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}(\tau)) \preceq \beta \mathbf{I}_n, \quad \text{for all } 0 \leq \tau \leq 1$$

Therefore $\eta = 1/\beta$ yields

$$\mathbf{0} \preceq \mathbf{I}_n - \eta \int_0^1 \nabla^2 f(\mathbf{x}(\tau)) d\tau \preceq \left(1 - \frac{\alpha}{\beta}\right) \mathbf{I}_n,$$

which further implies

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

Regularity Condition

More generally, for update rule

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

where $\mathbf{g}(\cdot) : \mathbb{R}^n \mapsto \mathbb{R}^n$

Definition 6

$\mathbf{g}(\cdot)$ is said to obey $\text{RC}(\mu, \lambda, \delta)$ for some $\mu, \lambda, \delta > 0$ if

$$2\langle \mathbf{g}(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \mu \|\mathbf{g}(\mathbf{x})\|_2^2 + \lambda \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad \forall \mathbf{x} \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$$

- Negative search direction \mathbf{g} is positively correlated with error $\mathbf{x} - \mathbf{x}_{\text{opt}} \implies$ one-step improvement
- $\mu\lambda \leq 1$ by Cauchy-Schwarz

RC = One-Point Strong Convexity + Smoothness

- One-point α -strong convexity:

$$f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) \geq \langle \nabla f(\mathbf{x}), \mathbf{x}_{\text{opt}} - \mathbf{x} \rangle + \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 \quad (1)$$

- β -smoothness:

$$\begin{aligned} f(\mathbf{x}_{\text{opt}}) - f(\mathbf{x}) &\leq f\left(\mathbf{x} - \frac{1}{\beta} \nabla f(\mathbf{x})\right) - f(\mathbf{x}) \\ &\leq \left\langle \nabla f(\mathbf{x}), -\frac{1}{\beta} \nabla f(\mathbf{x}) \right\rangle + \frac{\beta}{2} \left\| \frac{1}{\beta} \nabla f(\mathbf{x}) \right\|_2^2 \\ &= -\frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2 \end{aligned} \quad (2)$$

RC = One-Point Strong Convexity + Smoothness

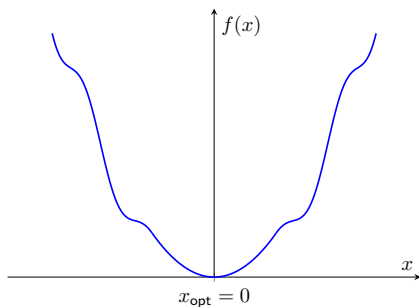
Combining relations (1) and (2) yields

$$\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{x}_{\text{opt}} \rangle \geq \frac{\alpha}{2} \|\mathbf{x} - \mathbf{x}_{\text{opt}}\|_2^2 + \frac{1}{2\beta} \|\nabla f(\mathbf{x})\|_2^2$$

— *RC holds with $\mu = 1/\beta$ and $\lambda = \alpha$*

Example of Nonconvex Functions

When $\mathbf{g}(\mathbf{x}) = \nabla f(\mathbf{x})$, f is not necessarily convex



$$f(x) = \begin{cases} x^2, & |x| \leq 6, \\ x^2 + 1.5|x|(\cos(|x| - 6) - 1), & |x| > 6 \end{cases}$$

Convergence Under RC

Lemma 7

Suppose $g(\cdot)$ obeys $\text{RC}(\mu, \lambda, \delta)$. The update rule ($\mathbf{x}^{t+1} = \mathbf{x}^t - \eta g(\mathbf{x}^t)$) with $\eta = \mu$ and $\mathbf{x}^0 \in \mathcal{B}_\delta(\mathbf{x}_{\text{opt}})$ obeys

$$\|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 \leq (1 - \mu\lambda)^t \|\mathbf{x}^0 - \mathbf{x}_{\text{opt}}\|_2^2$$

- $g(\cdot)$: more general search directions
 - example: in vanilla GD, $g(\mathbf{x}) = \nabla f(\mathbf{x})$
- The product $\mu\lambda$ determines the rate of convergence
- Attains ε -accuracy within $O\left(\frac{1}{\mu\lambda} \log \frac{1}{\varepsilon}\right)$ iterations

Proof of Lemma 7

By definition, one has

$$\begin{aligned}\|\mathbf{x}^{t+1} - \mathbf{x}_{\text{opt}}\|_2^2 &= \|\mathbf{x}^t - \eta \mathbf{g}(\mathbf{x}^t) - \mathbf{x}_{\text{opt}}\|_2^2 \\ &= \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta^2 \|\mathbf{g}(\mathbf{x}^t)\|_2^2 - 2\eta \langle \mathbf{g}(\mathbf{x}^t), \mathbf{x}^t - \mathbf{x}_{\text{opt}} \rangle \\ &\leq \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta^2 \|\mathbf{g}(\mathbf{x}^t)\|_2^2 - \eta \left(\lambda \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \mu \|\mathbf{g}(\mathbf{x}^t)\|_2^2 \right) \\ &= (1 - \eta\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2 + \eta(\eta - \mu) \|\mathbf{g}(\mathbf{x}^t)\|_2^2 \\ &\leq (1 - \mu\lambda) \|\mathbf{x}^t - \mathbf{x}_{\text{opt}}\|_2^2\end{aligned}$$

Initialization Matters in Nonconvex Problems

- Local convergence results require starting in a basin of attraction
- Good initializations reduce time to reach the fast local regime
- Common strategies:
 - spectral initialization
 - random initialization + perturbation
 - warm start from a convex/relaxed estimator

Limitations and Failure Modes

- Spurious local minima can appear in some models
- Flat/degenerate saddles can slow practical convergence
- Ill-conditioning (large condition number) hurts rates
- Optimization success does not automatically imply statistical generalization

Takeaway: always match algorithms and guarantees to structure in the specific problem.

Part III: Return to the Running Example

Apply the toolkit: local rates and global critical-point geometry

Interesting Questions

$$\text{minimize}_{\mathbf{x} \in \mathbb{R}^n} f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2$$

- How does the curvature behave, at least locally around the global minimizer?
- Where / what are the critical points? (Global geometry)

Gradient and Hessian via $f(\mathbf{x} + \mathbf{h})$ Expansion

For

$$f(\mathbf{x}) = \frac{1}{4} \|\mathbf{x}\mathbf{x}^\top - \mathbf{M}\|_{\text{F}}^2, \quad \mathbf{A} := \mathbf{x}\mathbf{x}^\top - \mathbf{M},$$

expand

$$f(\mathbf{x} + \mathbf{h}) = \frac{1}{4} \|\mathbf{A} + \mathbf{x}\mathbf{h}^\top + \mathbf{h}\mathbf{x}^\top + \mathbf{h}\mathbf{h}^\top\|_{\text{F}}^2.$$

Collecting terms up to second order in \mathbf{h} :

$$f(\mathbf{x} + \mathbf{h}) = f(\mathbf{x}) + \underbrace{\langle \mathbf{A}\mathbf{x}, \mathbf{h} \rangle}_{\text{1st order}} + \frac{1}{2} \underbrace{\mathbf{h}^\top (2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}) \mathbf{h}}_{\text{2nd order}} + O(\|\mathbf{h}\|_2^3).$$

Hence

$$\nabla f(\mathbf{x}) = (\mathbf{x}\mathbf{x}^\top - \mathbf{M})\mathbf{x}, \quad \nabla^2 f(\mathbf{x}) = 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I} - \mathbf{M}.$$

Local Linear Convergence of GD

Theorem 8

Suppose that $\|\mathbf{x}_0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$ and set $\eta = \frac{1}{4.5\lambda_1}$, GD obeys

$$\|\mathbf{x}^t - \sqrt{\lambda_1}\mathbf{u}_1\|_2 \leq \left(1 - \frac{\lambda_1 - \lambda_2}{18\lambda_1}\right)^t \|\mathbf{x}^0 - \sqrt{\lambda_1}\mathbf{u}_1\|_2, \quad t \geq 0,$$

- condition number/eigengap determines rate of convergence
- Requires initialization: use spectral method?

Proof of Theorem 8

It suffices to show that for all \mathbf{x} obeying $\underbrace{\|\mathbf{x} - \sqrt{\lambda_1}\mathbf{u}_1\|_2}_{\text{basin of attraction}} \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$,

$$0.25(\lambda_1 - \lambda_2)\mathbf{I}_n \preceq \nabla^2 f(\mathbf{x}) \preceq 4.5\lambda_1\mathbf{I}_n$$

Express gradient and Hessian as

$$\begin{aligned}\nabla f(\mathbf{x}) &= (\mathbf{x}\mathbf{x}^\top - \mathbf{M})\mathbf{x} \\ \nabla^2 f(\mathbf{x}) &= 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2\mathbf{I}_n - \mathbf{M}\end{aligned}$$

Preliminary Facts

Let $\Delta := \mathbf{x} - \sqrt{\lambda_1} \mathbf{u}_1$. It is seen that when $\|\Delta\|_2 \leq \frac{\lambda_1 - \lambda_2}{15\sqrt{\lambda_1}}$, one has

$$\lambda_1 - 0.25(\lambda_1 - \lambda_2) \leq \|\mathbf{x}\|_2^2 \leq 1.15\lambda_1;$$

$$\|\Delta\|_2 \leq \|\mathbf{x}\|_2;$$

$$\|\Delta\|_2 \|\mathbf{x}\|_2 \leq (\lambda_1 - \lambda_2)/12$$

Triangle inequality gives

$$\begin{aligned}\|\nabla^2 f(\mathbf{x})\| &\leq \|2\mathbf{x}\mathbf{x}^\top\| + \|\mathbf{x}\|_2^2 + \|\mathbf{M}\| \\ &\leq 3\|\mathbf{x}\|_2^2 + \lambda_1 < 4.5\lambda_1,\end{aligned}$$

where the last relation follows from $\|\mathbf{x}\|_2^2 \leq 1.15\lambda_1$

Local Strong Convexity

Recall that $\Delta = x - \sqrt{\lambda_1}u_1$

$$\begin{aligned}xx^\top &= \lambda_1 u_1 u_1^\top + \Delta x^\top + x \Delta^\top - \Delta \Delta^\top \\ &\succeq \lambda_1 u_1 u_1^\top - 3\|\Delta\|_2 \|x\|_2 \mathbf{I}_n \quad (\|\Delta\|_2 \leq \|x\|_2) \\ &\succeq \lambda_1 u_1 u_1^\top - 0.25(\lambda_1 - \lambda_2) \mathbf{I}_n,\end{aligned}$$

where last line relies on $\|\Delta\|_2 \|x\|_2 \leq (\lambda_1 - \lambda_2)/12$. Consequently,

$$\begin{aligned}\nabla^2 f(x) &= 2xx^\top + \|x\|_2^2 \mathbf{I}_n - \lambda_1 u_1 u_1^\top - \sum_{i=2}^n \lambda_i u_i u_i^\top \\ &\succeq 2\lambda_1 u_1 u_1^\top + (\|x\|_2^2 - 0.5(\lambda_1 - \lambda_2)) \mathbf{I}_n - \lambda_1 u_1 u_1^\top - \sum_{i=2}^n \lambda_i u_i u_i^\top \\ &\succeq (\|x\|_2^2 - 0.5(\lambda_1 - \lambda_2) + \lambda_1) u_1 u_1^\top \\ &\quad + \sum_{i=2}^n (\|x\|_2^2 - 0.5(\lambda_1 - \lambda_2) - \lambda_i) u_i u_i^\top \\ &\succeq (\|x\|_2^2 - 0.5(\lambda_1 - \lambda_2) - \lambda_2) \mathbf{I}_n \\ &\succeq 0.25(\lambda_1 - \lambda_2) \mathbf{I}_n \quad (\lambda_1 - 0.25(\lambda_1 - \lambda_2) \leq \|x\|_2^2)\end{aligned}$$

Critical Points Characterization

\mathbf{x} is a critical point, i.e., $\nabla f(\mathbf{x}) = (\mathbf{x}\mathbf{x}^\top - \mathbf{M})\mathbf{x} = \mathbf{0}$

\Leftrightarrow

$$\mathbf{M}\mathbf{x} = \|\mathbf{x}\|_2^2 \mathbf{x}$$

\Leftrightarrow

\mathbf{x} aligns with an eigenvector of \mathbf{M} or $\mathbf{x} = \mathbf{0}$

Since $\mathbf{M}\mathbf{u}_i = \lambda_i \mathbf{u}_i$, the set of critical points is given by

$$\{\mathbf{0}\} \cup \{\pm\sqrt{\lambda_i} \mathbf{u}_i, \quad i = 1, \dots, n\}$$

Categorization of Critical Points

The critical points can be further categorized based on the **Hessian**:

$$\nabla^2 f(\mathbf{x}) = 2\mathbf{x}\mathbf{x}^\top + \|\mathbf{x}\|_2^2 \mathbf{I}_n - \mathbf{M}$$

- For any non-zero critical point $\mathbf{x}_k = \pm\sqrt{\lambda_k}\mathbf{u}_k$:

$$\begin{aligned}\nabla^2 f(\mathbf{x}_k) &= 2\lambda_k \mathbf{u}_k \mathbf{u}_k^\top + \lambda_k \mathbf{I} - \mathbf{M} \\ &= 2\lambda_k \mathbf{u}_k \mathbf{u}_k^\top + \lambda_k \left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^\top \right) - \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top \\ &= \sum_{i:i \neq k} (\lambda_k - \lambda_i) \mathbf{u}_i \mathbf{u}_i^\top + 2\lambda_k \mathbf{u}_k \mathbf{u}_k^\top\end{aligned}$$

Categorization of Critical Points (Cont.)

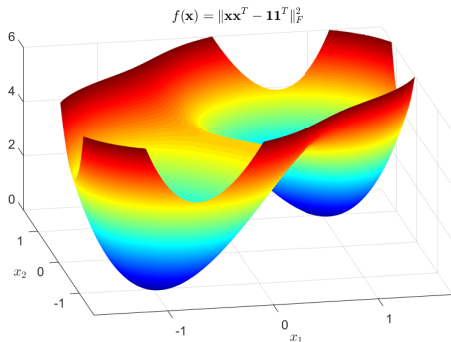
If $\lambda_1 > \lambda_2 \geq \dots \geq \lambda_n \geq 0$, then

- $\nabla^2 f(\mathbf{x}_1) \succ \mathbf{0}$ \rightarrow local minima
- $1 < k \leq n$: $\lambda_{\min}(\nabla^2 f(\mathbf{x}_k)) < 0$, $\lambda_{\max}(\nabla^2 f(\mathbf{x}_k)) > 0$
 \rightarrow strict saddle
- $\mathbf{x} = \mathbf{0}$: $\nabla^2 f(\mathbf{0}) = -\mathbf{M} \preceq \mathbf{0}$ \rightarrow local maxima, strict saddle

all local minima are global; all saddles are strict

A Pictorial Example

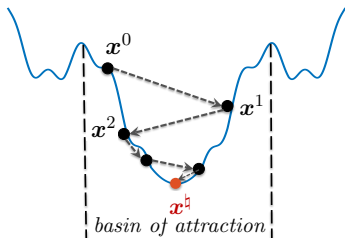
For example, for 2-dimensional case $f(\mathbf{x}) = \left\| \mathbf{x}\mathbf{x}^\top - \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix} \right\|_F^2$



global minima: $\mathbf{x} = \pm \begin{bmatrix} 1 \\ 1 \end{bmatrix}$; strict saddles: $\mathbf{x} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$, and $\pm \begin{bmatrix} 1 \\ -1 \end{bmatrix}$
— No “spurious” local minima!

Two Vignettes

Two-stage approach:



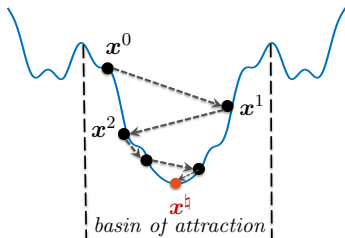
smart initialization

+

local refinement

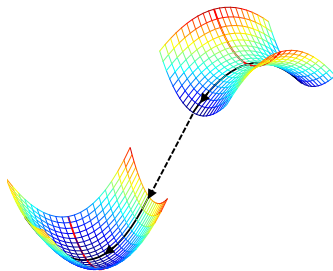
Two Vignettes

Two-stage approach:



smart initialization
+
local refinement

Global landscape:



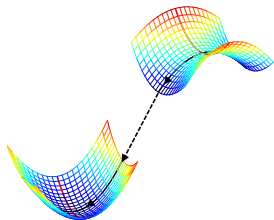
benign landscape
+
saddle-point escaping

Benign landscape:

- all local minima = global minima
- other critical points = strict saddle points

Saddle-point escaping algorithms:

- trust-region methods
- perturbed gradient descent
- perturbed SGD
- ...



Next Steps

- Generic local analysis of (regularized) gradient descent
- Refined local analysis for gradient descent
- Global landscape analysis
- Gradient descent with random initialization
- (Maybe) Gradient descent with arbitrary initialization