

535520: Optimization Algorithms

Lecture 4 — Gradient Descent

Ping-Chun Hsieh (謝秉均)

September 23, 2024

This Lecture

1. Gradient Descent for Convex Problems

2. Gradient Descent for Non-Convex Problems

- Reading Material:
 - Chapter 5 of Amir Beck's textbook "First-Order Methods in Optimization"
 - Chapters 2 & 5 of Dimitri Bertsekas's textbook "Nonlinear Programming"
 - Chapter 3 of Jorge Nocedal and Stephen Wright's textbook "Numerical Optimization"
 - Yuxin Chen's lecture note: https://yuxinchen2020.github.io/ele522_optimization/lectures/grad_descent_unconstrained.pdf

Convergence Rates of GD

	Convergence Rate (under constant step sizes)
Quadratic problem	$\ x_t - x^*\ _2 = \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \right)^t \cdot \ x_0 - x^*\ _2$
Strongly convex and L-smooth	$\ x_t - x^*\ \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^t \cdot \ x_0 - x^*\ $
PL condition and L-smooth	$f(x_t) - f(x^*) \leq \left(1 - \frac{\mu}{L} \right)^t \cdot (f(x_0) - f(x^*))$
Convex and L-smooth	$f(x_t) - f(x^*) \leq \frac{L}{2t} \cdot \ x_0 - x^*\ ^2 = O\left(\frac{1}{t}\right)$
Non-convex and L-smooth	$\min_{0 \leq k \leq T} \ \nabla f(x_k)\ \leq \sqrt{\frac{2L(f(x_0) - f(x^*))}{T}}, \quad \lim_{t \rightarrow \infty} \ \nabla f(x_t)\ = 0$

Intuition: Why can GD converge under convexity?

Ideally, if GD converges (to some point), then we expect $\|\nabla f(x_t)\| \rightarrow 0$ as $t \rightarrow \infty$

Recall from the definition of convexity:

$$f(x) - f(x^*) \leq \nabla f(x)^\top (x - x^*) \leq \|\nabla f(x)\| \cdot \|x - x^*\|$$

Suppose the domain is of finite radius, then $\|\nabla f(x)\| \rightarrow 0$ implies that $(f(x) - f(x^*)) \rightarrow 0$

Quick Recap: GD for Quadratic Problems

Let's motivate the convergence of GD using a quadratic objective function

$$\min_{x \in \mathbb{R}^n} f(x) := \frac{1}{2}(x - x^*)^\top Q(x - x^*), \quad Q \text{ is pd } (\lambda_1(Q) \geq \dots \geq \lambda_n(Q) > 0)$$

GD update: $x_{k+1} = x_k - \eta_k \nabla f(x_k) =$

Theorem (Convergence Rate of GD for Quadratic Problems):

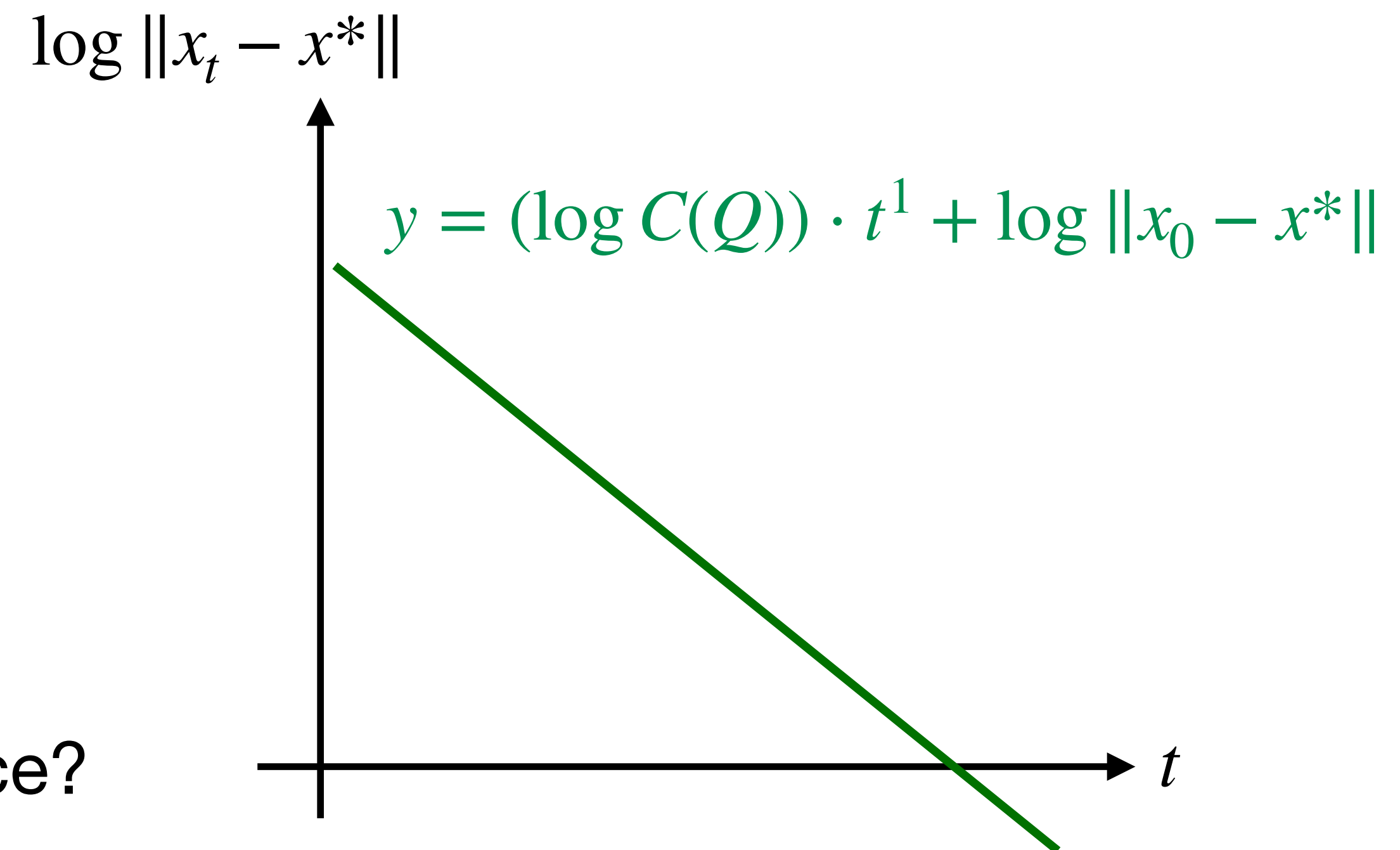
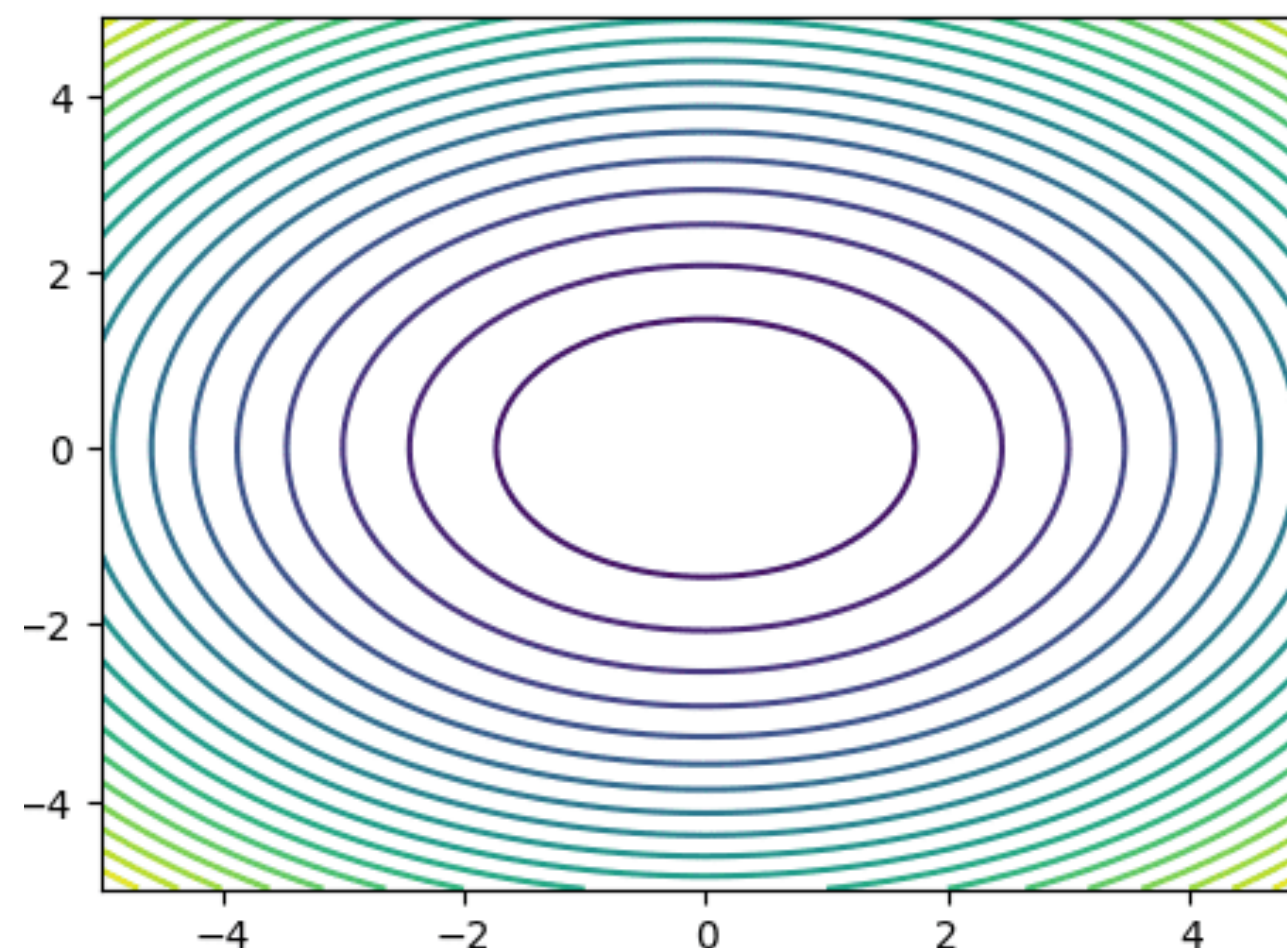
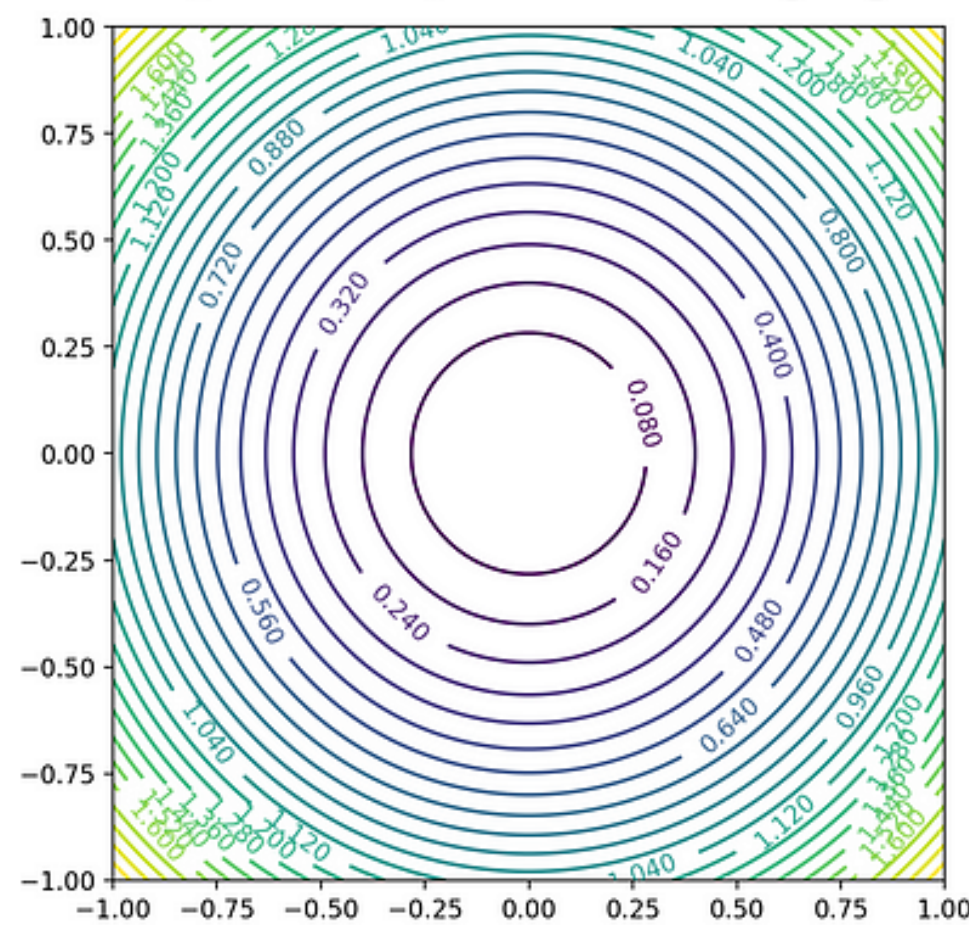
Under the step size $\eta_t \equiv \eta = 2/(\lambda_1(Q) + \lambda_n(Q))$, then we have

$$\|x_t - x^*\|_2 = \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \right)^t \cdot \|x_0 - x^*\|_2, \quad \forall t \in \mathbb{N}$$

Quick Recap: GD for Quadratic Problems

$$\begin{aligned}\|x_t - x^*\|_2 &= \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \right)^t \cdot \|x_0 - x^*\|_2, \quad \forall t \in \mathbb{N} \\ &= \left(\frac{1 - C(Q)}{1 + C(Q)} \right)^t \cdot \|x_0 - x^*\|_2 \quad \text{where } C(Q) := \frac{\lambda_n(Q)}{\lambda_1(Q)} \text{ is “condition number”}\end{aligned}$$

- **Question:** This is often called “geometric convergence” or “linear convergence”



How about “sub-linear” or “super-linear” convergence?

Terminology of Convergence Rates

- Let $e(x)$ denote the distance from optimality
 - Example: $e(x) = \|x - x^*\|$
 - Example: $e(x) = |f(x) - f(x^*)|$
- **Rate of convergence**: The limit of the ratio of successive errors

$$\lim_{k \rightarrow \infty} \frac{e(x_{k+1})}{e(x_k)} = \beta$$

- If $\beta = 1$: We call it a **sub-linear rate** of convergence
- If $\beta \in (0,1)$: We call it a **linear rate** of convergence
- If $\beta = 0$: We call it a **super-linear rate** of convergence

Proof: Convergence of GD for Quadratic Problems

Step 1: Consider the GD update

$$x_{t+1} - x^* = (x_t - \eta \cdot \nabla f(x^t)) - x^* = (I_n - \eta Q) \cdot (x_t - x^*)$$

This implies that

$$\|x_{t+1} - x^*\|_2 \leq \underbrace{\|I_n - \eta \cdot Q\|}_{\text{What norm?}} \cdot \|x_t - x^*\|_2 \quad \dots\dots (\quad)$$

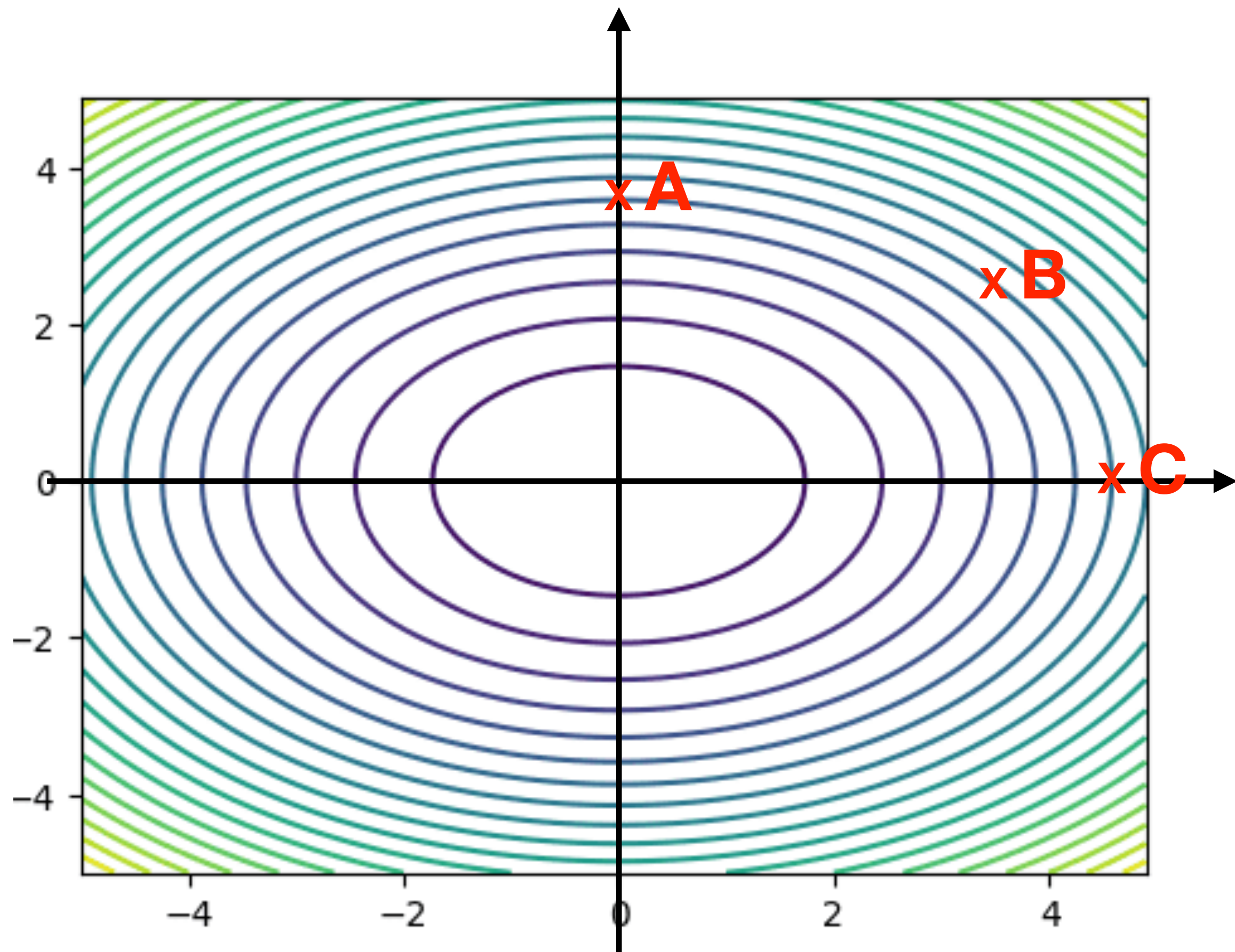
Step 2: Moreover, we have

$$\begin{aligned} \|I_n - \eta Q\| &= \max \left\{ |1 - \eta \lambda_1(Q)|, |1 - \eta \lambda_n(Q)| \right\} \\ &= 1 - \frac{2\lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \end{aligned}$$

Step 3: By repeating Step 1 recursively, we complete the proof.

Observations: GD with Constant Step Sizes

Consider $f(x) := \frac{1}{2}x^\top Qx$ with $Q = \begin{bmatrix} 1 & 0 \\ 0 & 10 \end{bmatrix}$



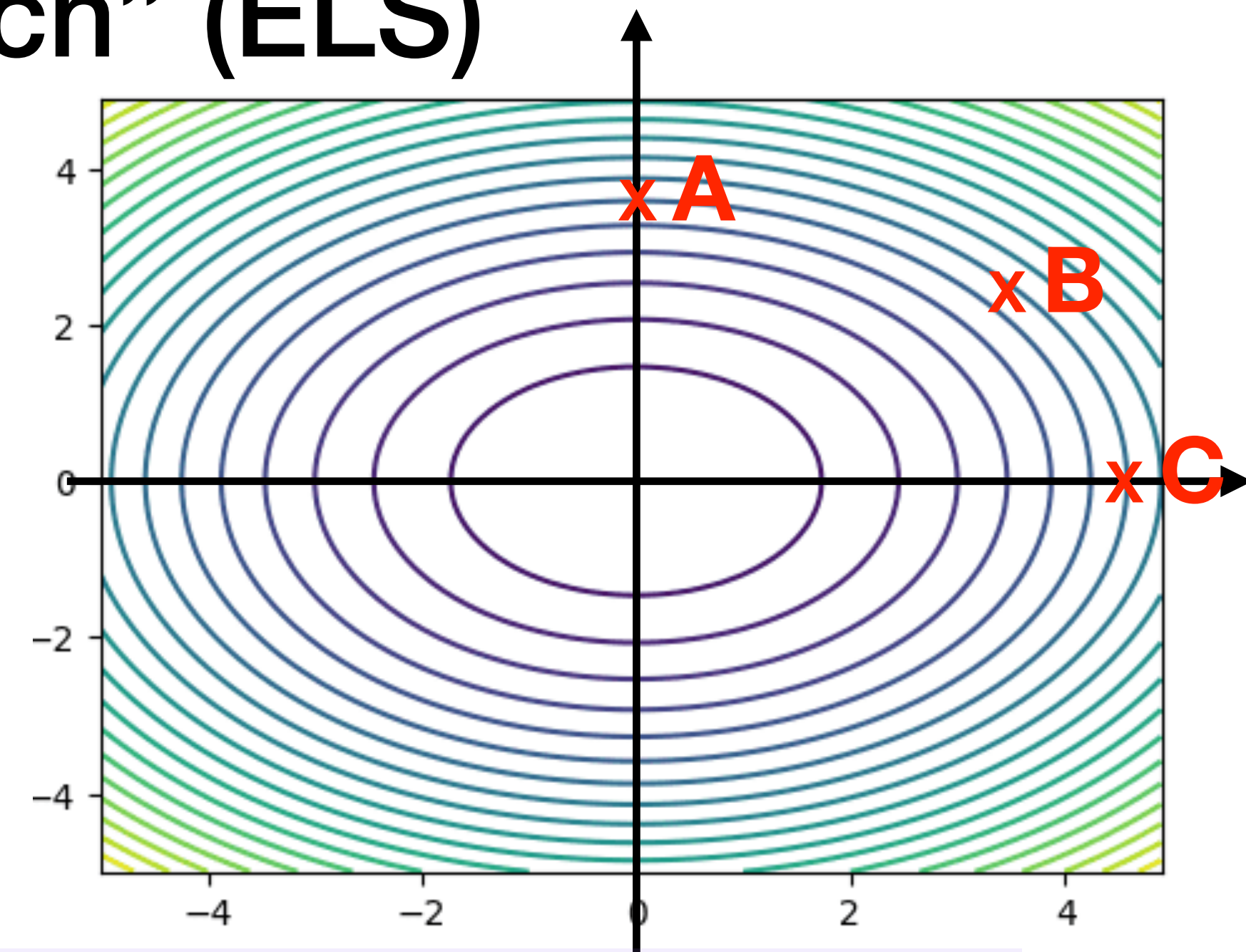
$$\eta = \frac{2}{\lambda_1(Q) + \lambda_n(Q)} = \frac{2}{11}$$

Can you find any special behavior of GD with constant step sizes?

Another Variant: GD with “Exact Line Search” (ELS)

To accelerate GD, we can choose the step sizes by

$$\eta_t = \arg \min_{\eta \geq 0} f\left(x_t - \eta \nabla f(x_t)\right)$$



Theorem (Convergence Rate of GD with ELS):

By applying GD with ELS to the quadratic problems, we have

$$f(x_t) - f(x^*) \leq \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \right)^{2t} \cdot f(x_0) - f(x^*), \quad \forall t \in \mathbb{N}$$

Remark: The convergence rate is actually the same of GD with constant η

Proof: Convergence of GD for Quadratic Problems (1/2)

Step 1: Under ELS, we have

Notation: $g_t \equiv \nabla f(x_t) = Q(x_t - x^*)$

$$\eta_t = \arg \min_{\eta \geq 0} f\left(x_t - \eta \nabla f(x_t)\right) = \frac{g_t^\top g_t}{g_t^\top Q g_t} \quad (\text{Why?})$$

Step 2: Then, we can find $f(x_{t+1})$ as

$$\begin{aligned} f(x_{t+1}) &= \frac{1}{2} \left(x_t - \eta_t Q (x_t - x^*) - x^* \right)^\top Q \left(x_t - \eta_t Q (x_t - x^*) - x^* \right) \\ &= \frac{1}{2} (x_t - x^*)^\top Q (x_t - x^*) - \frac{\|g_t\|^4}{2g_t^\top Q g_t} \\ &= \underbrace{\frac{1}{2} (x_t - x^*)^\top Q (x_t - x^*)}_{f(x_t)} \cdot \left(1 - \frac{\|g_t\|^4}{(g_t^\top Q g_t) \cdot (g_t^\top Q^{-1} \cdot Q \cdot Q^{-1} g)} \right) \end{aligned}$$

Proof: Convergence of GD for Quadratic Problems (2/2)

Step 3: To bound (A), we can use the “Kantorovich’s inequality”

Lemma (Kantorovich’s inequality):

Let Q be a symmetric and pd matrix. Then, for any $y \in \mathbb{R} \setminus \{0\}$,

$$\frac{\|y\|^4}{(y^\top Q y) \cdot (y^\top Q^{-1} y)} \geq \frac{4\lambda_1(Q) \cdot \lambda_n(Q)}{(\lambda_1(Q) + \lambda_n(Q))^2}$$

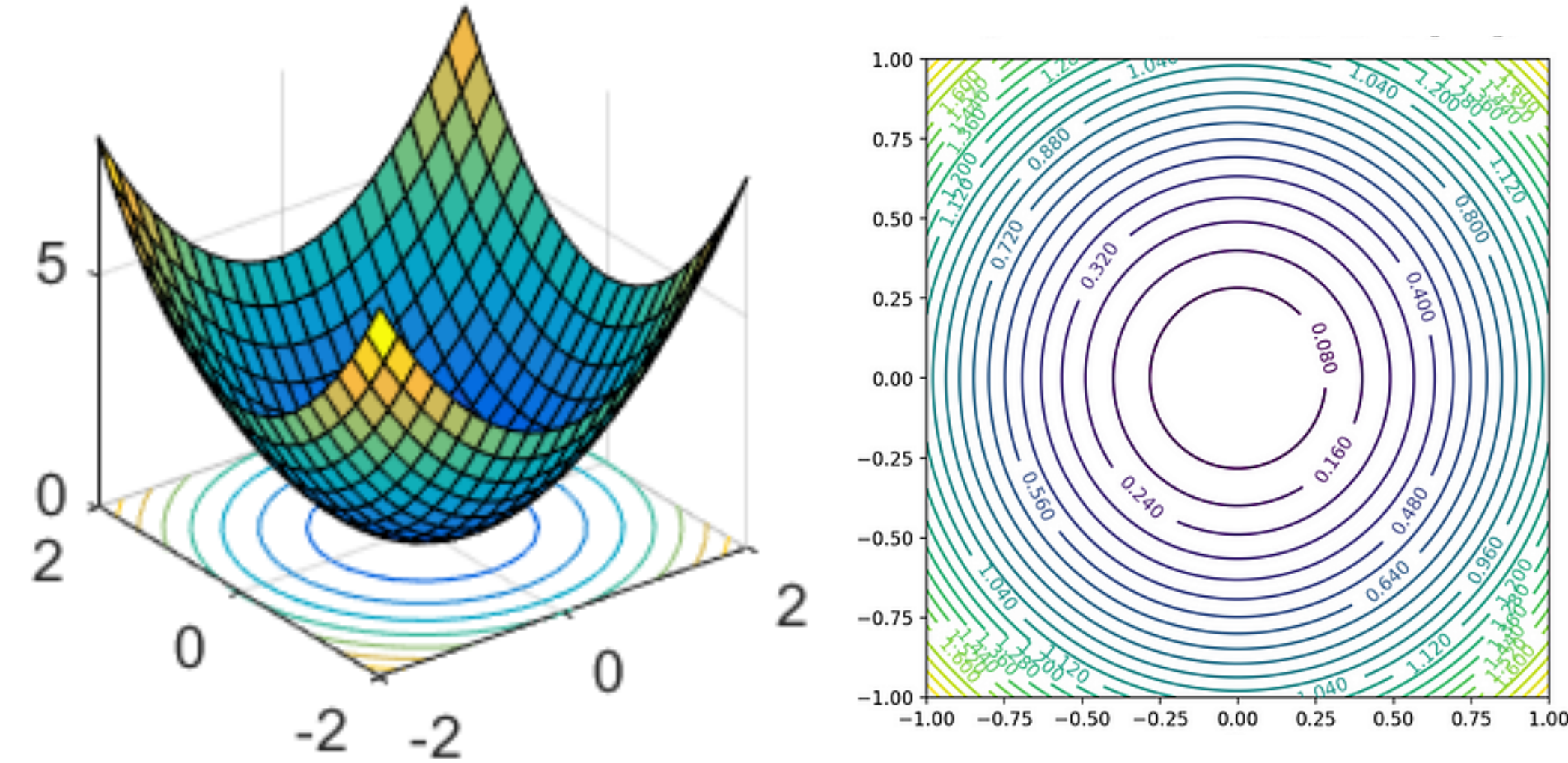
As a result, we have

$$f(x_{t+1}) \leq \left(1 - \frac{4 \cdot \lambda_1(Q) \cdot \lambda_n(Q)}{(\lambda_1(Q) + \lambda_n(Q))^2} \right) \cdot f(x_t) = \left(\frac{\lambda_1(Q) - \lambda_n(Q)}{\lambda_1(Q) + \lambda_n(Q)} \right)^2 f(x_t)$$

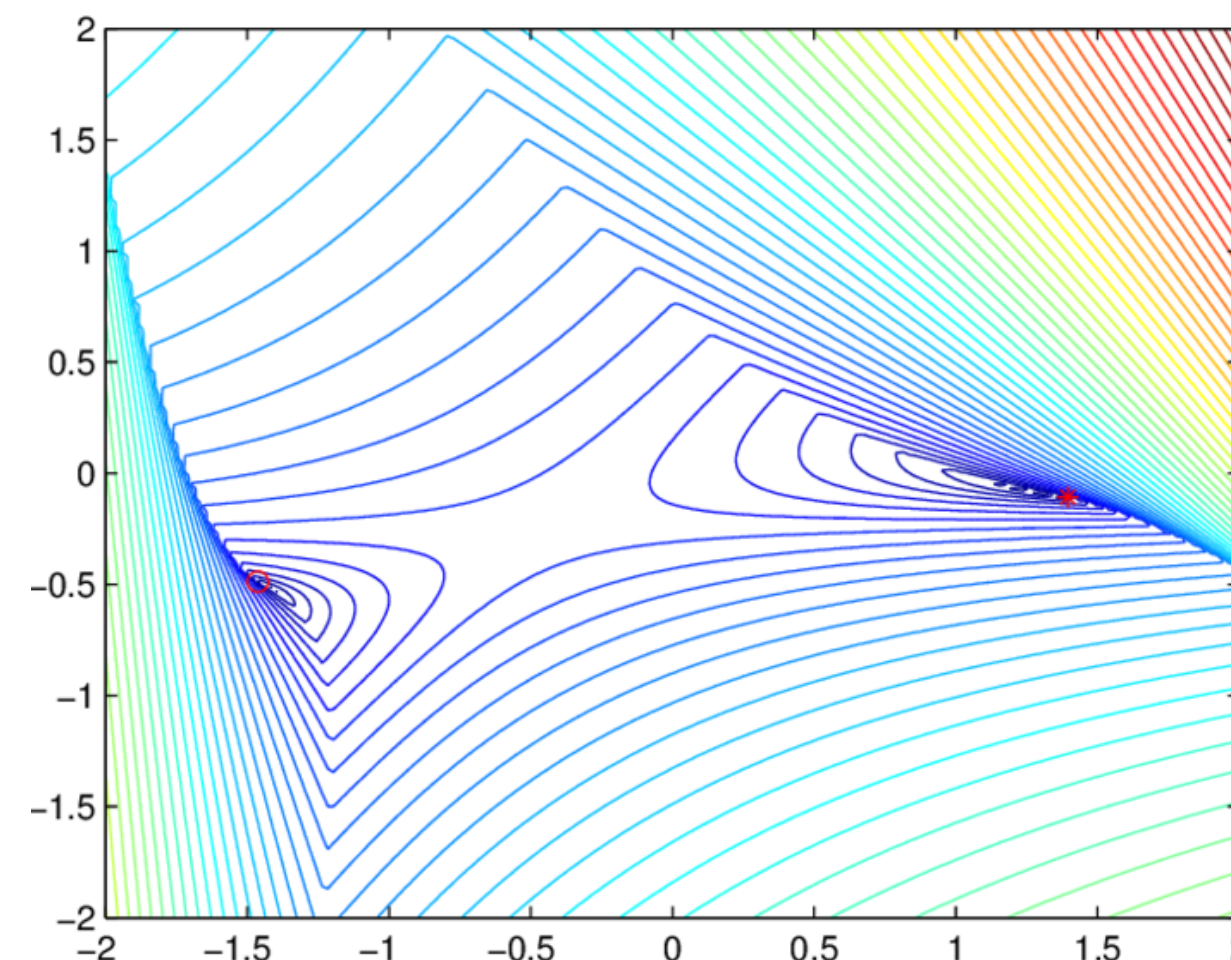
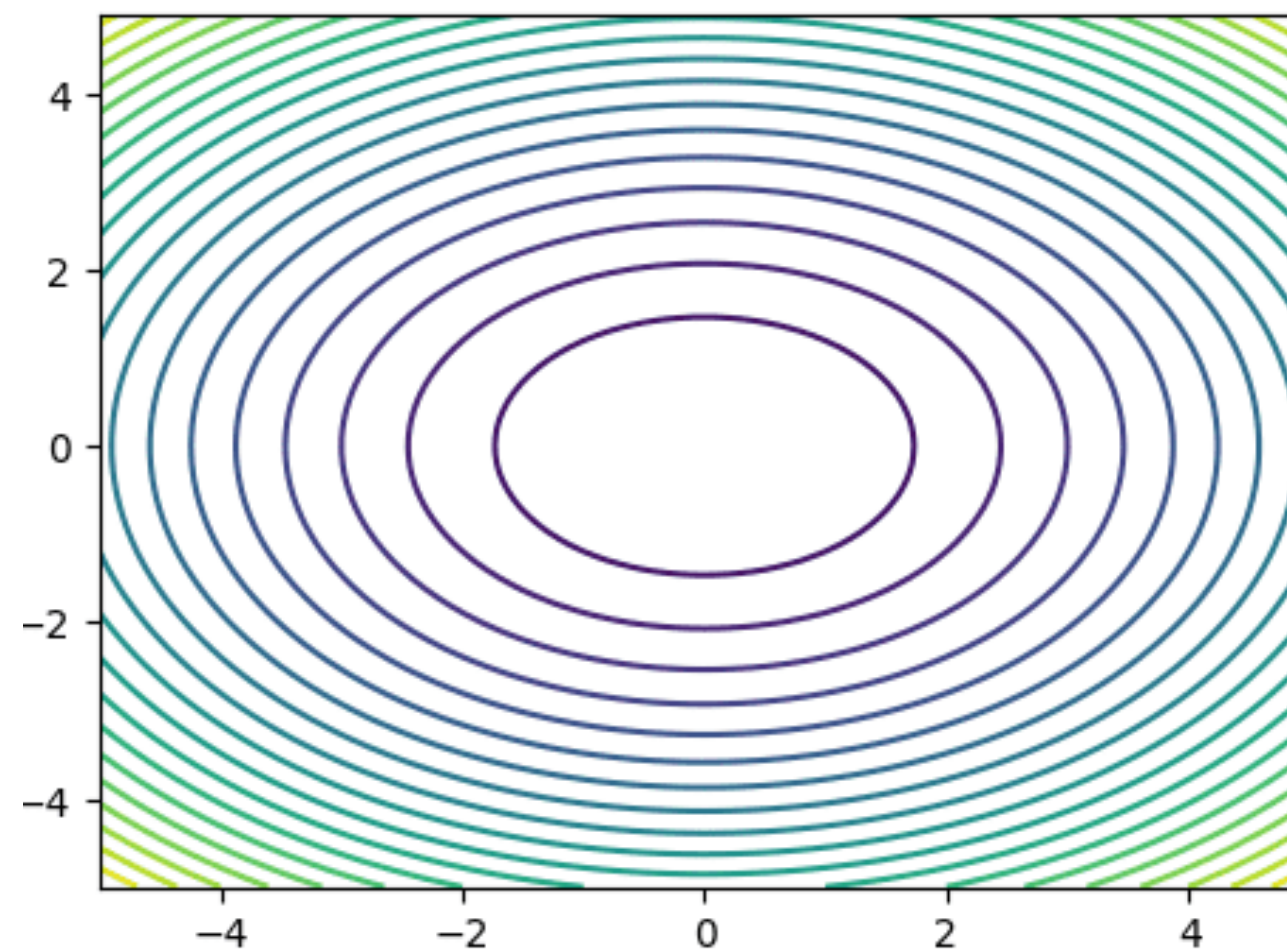
Let's go beyond quadratic problems:
Strongly-convex and *smooth* problems

Why Strong Convexity and Smoothness?

Strong convexity: GD can always attain sufficient per-step improvement



Smoothness: Gradient serves as a useful direction for improvement



Strict Convexity vs Strong Convexity

Definition: A function $f: X \rightarrow \mathbb{R}$ is called **strictly convex** if its domain X is a convex set and for any $x, y \in X$ with $x \neq y$ and any $\alpha \in [0, 1]$, we have

$$f(\alpha x + (1 - \alpha)y) < \alpha f(x) + (1 - \alpha)f(y)$$

Intuition: “The segment lies strictly above the function”

Definition: A function $f: X \rightarrow \mathbb{R}$ is called **μ -strongly convex** if its domain X is a convex set and there exists some $\mu > 0$ such that for any $x, y \in X$

$$f(\alpha x + (1 - \alpha)y) \leq \alpha f(x) + (1 - \alpha)f(y) - \frac{\mu}{2}\alpha(1 - \alpha)\|x - y\|^2$$

Intuition: 1-dimensional case

An Alternative Definition of Strong Convexity

Theorem 1: Let $f : X \rightarrow \mathbb{R}$ be a *continuously differentiable* function. Then, the following are equivalent characterization of **strong convexity**.

(1) There exists some $\mu > 0$ such that for any $x, y \in X$,

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{\mu}{2} \|y - x\|^2$$

(2) There exists some $\mu > 0$ such that for any $x, y \in X$,

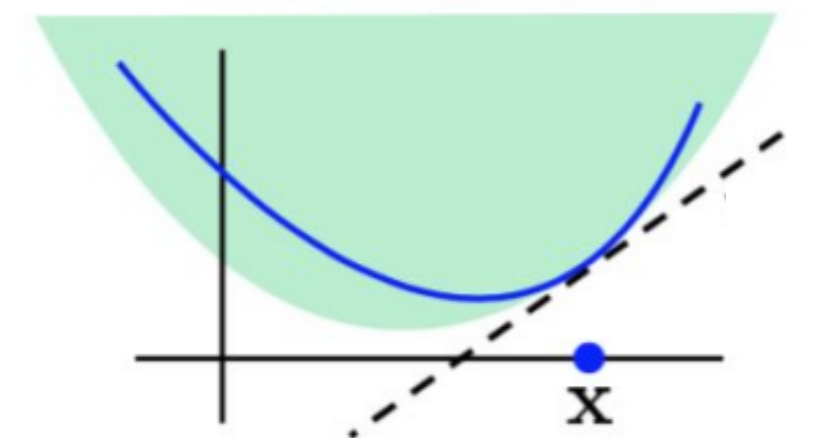
$$(\nabla f(x) - \nabla f(y))^\top (x - y) \geq \mu \|x - y\|^2$$

(3) Moreover, if f is twice continuously differentiable, then there exists some $\mu > 0$ such that for any $x \in X$,

$$\nabla^2 f(x) - \mu I \succ 0$$

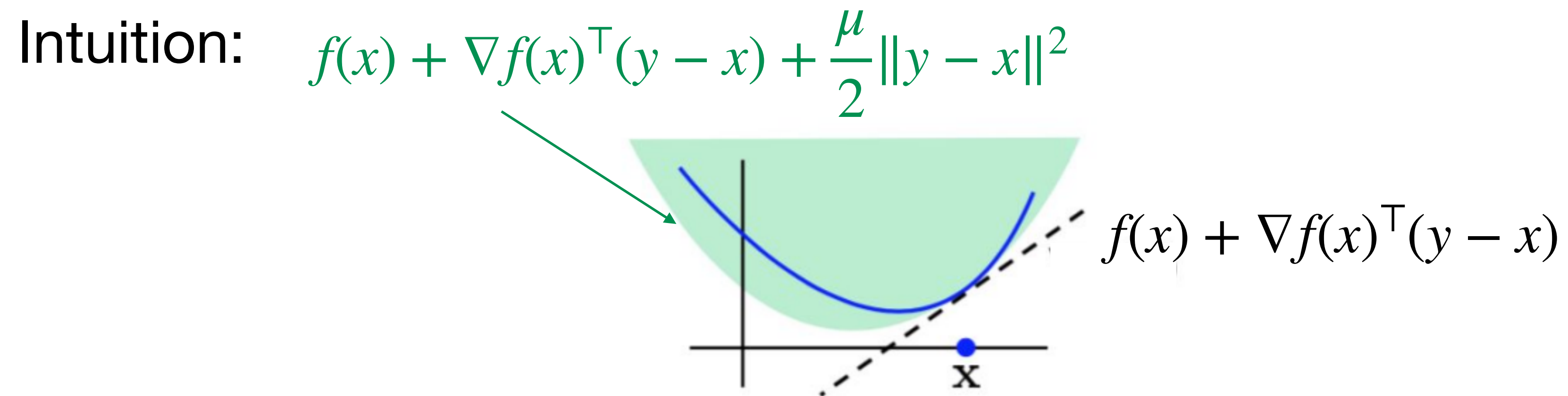
Intuition: Taylor expansion

(Proof: HW1 Problem)



Connecting “Strict Convexity” and “Strong Convexity”

Theorem 2: Let $f : X \rightarrow \mathbb{R}$ be a *continuously differentiable* function with an open convex domain X . If f is **strongly convex**, then f is also **strictly convex**.



Connecting “Strict Convexity” and “Strong Convexity”

Theorem 2: Let $f : X \rightarrow \mathbb{R}$ be a *continuously differentiable* function with an open convex domain X . If f is **strongly convex**, then f is also **strictly convex**.

Proof: Define $h(t) := f(x + t(y - x))$, $t \in \mathbb{R}$

Step 1: Consider $t, t' \in [0, 1]$ such that $t < t'$

$$\begin{aligned} & \underbrace{\left(\nabla f(x + t'(y - x)) - \nabla f(x + t(y - x)) \right)^\top}_{= \left(\frac{dh(t')}{dt} - \frac{dh(t)}{dt} \right) (t' - t)} \left((t' - t)(y - x) \right) \geq \alpha(t' - t)^2 \|y - x\|^2 > 0 \\ & = \left(\frac{dh(t')}{dt} - \frac{dh(t)}{dt} \right) (t' - t) \end{aligned}$$

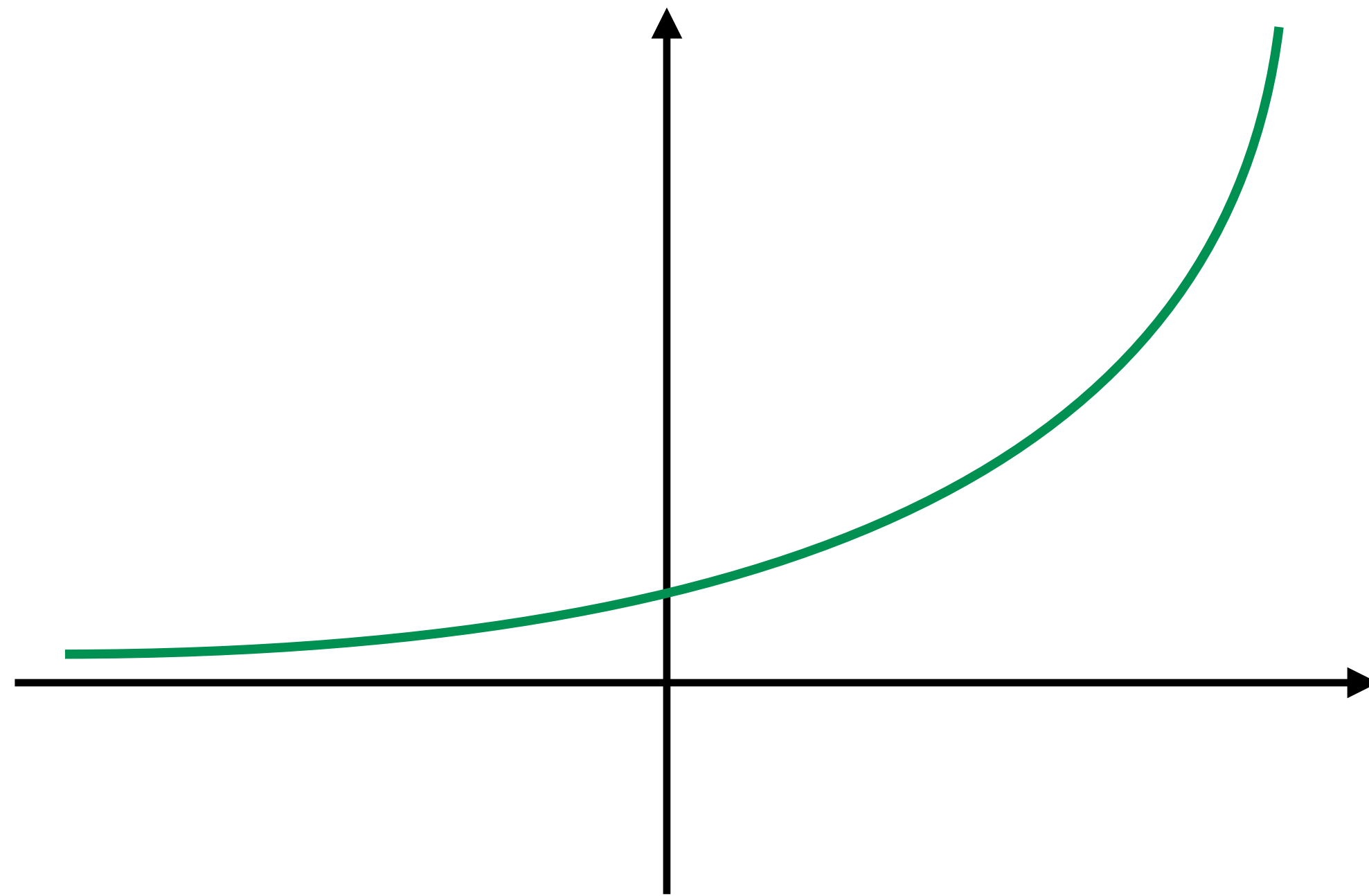
Step 2: By Step 1, we know $\frac{dh}{dt}$ is strictly increasing. As a result,

$$\frac{h(t) - h(0)}{t} = \frac{1}{t} \int_0^t \frac{dh(s)}{ds} ds < \frac{1}{1 - t} \int_t^1 \frac{dh(s)}{ds} ds = \frac{h(1) - h(t)}{1 - t} \quad (\text{Why?})$$

Step 3: Hence, we have $t \cdot h(1) + (1 - t)h(0) > h(t)$

“Strict Convexity” does NOT imply “Strong Convexity”

- A strictly convex function is NOT necessarily strongly convex
- Example: $f(x) = \exp(x)$ is not strongly convex on \mathbb{R}



Check $f''(x)$:

Lipschitz Smoothness and L -Smoothness

Definition: $f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ is called **Lipschitz continuous** if there exists $L < \infty$ such that for all $x, y \in \mathbb{R}^n$

$$\|f(x) - f(y)\| \leq L\|x - y\|$$

Definition: $f: \mathbb{R}^n \rightarrow \mathbb{R}$ is called **L -smooth** if it has *Lipschitz continuous gradients*, i.e., there exists $L < \infty$ such that for all $x, y \in \mathbb{R}^n$

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$$

Theorem 3: Let $f: X \rightarrow \mathbb{R}$ be *twice differentiable*. Then, f is **L -smooth** if and only if

$$\nabla^2 f(x) \preceq LI$$

Equivalent Characterization of L -Smoothness for Convex Functions

Theorem 4: Let $f : X \rightarrow \mathbb{R}$ be a convex and differentiable function. Then, the following are equivalent characterization of L -smoothness:

$$(1) f(y) \leq f(x) + \nabla f(x)^\top (y - x) + \frac{L}{2} \|y - x\|^2, \text{ for all } x, y \in X$$

$$(2) f(y) \geq f(x) + \nabla f(x)^\top (y - x) + \frac{1}{2L} \|\nabla f(y) - \nabla f(x)\|^2, \text{ for all } x, y \in X$$

$$(3) (\nabla f(x) - \nabla f(y))^\top (y - x) \geq \frac{1}{L} \|\nabla f(y) - \nabla f(x)\|^2, \text{ for all } x, y \in X$$

(For the details, please see Chapter 5.1.2 of Amir Beck's textbook)

In the next few slides, we focus on GD for
 μ -strongly convex and *L -smooth* objective functions

Convergence of GD for μ -Strongly Convex and Smooth Functions

Theorem (Convergence of GD under strong convexity and smoothness):

Let f be μ -strongly convex and L -smooth. Under GD with constant step sizes $\eta = 2/(\mu + L)$, we have

$$\|x_t - x^*\| \leq \left(\frac{\kappa - 1}{\kappa + 1} \right)^t \cdot \|x_0 - x^*\|$$

($\kappa := L/\mu$ is called the condition number)

Comparison:

• Step size:	$\frac{2}{\mu + L}$	vs	$\frac{2}{\lambda_1(Q) + \lambda_n(Q)}$
• Contraction:	$\frac{\kappa - 1}{\kappa + 1}$	vs	$\frac{1 - C(Q)}{1 + C(Q)}$

Proof: Convergence of GD for μ -Strongly Convex and Smooth Functions

Step 1: Let's rewrite

$$\nabla f(x_t) = \nabla f(x_t) - \nabla f(\underline{x_t + 1 \cdot (x^* - x_t)}) = \left(\int_0^1 \nabla^2 f(x_t + s \cdot (x^* - x_t)) \cdot ds \right) (x_t - x^*)$$

Step 2:

$$\|x_{t+1} - x^*\| = \|x_t - x^* - \eta \cdot \nabla f(x_t)\|$$

$$\begin{aligned} &= \left\| \left(I - \eta \cdot \int_0^1 \nabla^2 f(x_t + s(x^* - x_t)) ds \right) (x_t - x^*) \right\| \\ &= \underbrace{\sup_{0 \leq s \leq 1} \left\| I - \eta \cdot \int_0^1 \nabla^2 f(x_t + s \cdot (x^* - x_t)) ds \right\|}_{\leq |1 - \eta L|} \cdot \|x_t - x^*\| \end{aligned}$$

Next Question: Do we still get “linear convergence” while relaxing the strong convexity condition?

Polyak-Łojasiewicz (PL) Condition in Non-Convex Optimization

Question: When can GD succeed under non-convex objective functions?

Polyak-Łojasiewicz Condition

Gradient norm

Sub-optimality gap

$$\|\nabla f(\theta)\|^2 \geq 2\mu \cdot (f(\theta^*) - f(\theta)) \quad \text{for some } \mu > 0$$

(aka “gradient dominance”)



Boris
Polyak



Stanisław
Łojasiewicz

Interpretation:

- PL ensures that gradient grows fast as it moves away from the optimum
- PL ensures that every stationary point is a global optimum

Convergence of GD Under PL Condition

Theorem (Convergence of GD under PL and smoothness):

Let f satisfies PL condition and is L -smooth. Under GD with constant step sizes $\eta = 1/L$, we have

$$f(x_t) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^t \cdot (f(x_0) - f(x^*)), \quad \forall t \in \mathbb{N}$$

Proof: $f(x_{t+1}) - f(x^*)$

$$\leq f(x_t) - f(x^*) - \frac{1}{2L} \|\nabla f(x_t)\|^2 \quad \dots\dots (\quad)$$

$$\leq f(x_t) - f(x^*) - \frac{\mu}{L} \cdot (f(x_t) - f(x^*)) \quad \dots\dots (\quad)$$

$$= \left(1 - \frac{\mu}{L}\right) \cdot (f(x_t) - f(x^*)) \quad \dots\dots (\quad)$$

Question: Any known problem that satisfies a PL-like condition?

Example 1: Overparametrized Linear Regression

Linear regression: Given N data samples $\{a_i \in \mathbb{R}^m, y_i \in \mathbb{R}\}$, find a linear model by minimizing

$$f(x) = \frac{1}{2} \sum_{i=1}^n (a_i^\top x - y_i)^2$$

Overparameterization: Model dimension $m >$ sample size n

(This regime occurs frequently in deep learning)

Remark: This is a convex but not strongly convex problem (why?)

$$\nabla^2 f(x) = \sum_{i=1}^n a_i a_i^\top = \mathbf{A} \mathbf{A}^\top$$

Notation:

$$\mathbf{A} = [a_1, \dots, a_n]^\top \in \mathbb{R}^{n \times m}$$

Question: Does $f(x)$ satisfy the PL condition?

Example 1: Overparametrized Linear Regression

Let's show that $\|\nabla f(x)\|_2^2 \geq 2\lambda_{\min}(AA^\top)f(x)$

$$\nabla f(x) = A^\top(Ax - y)$$

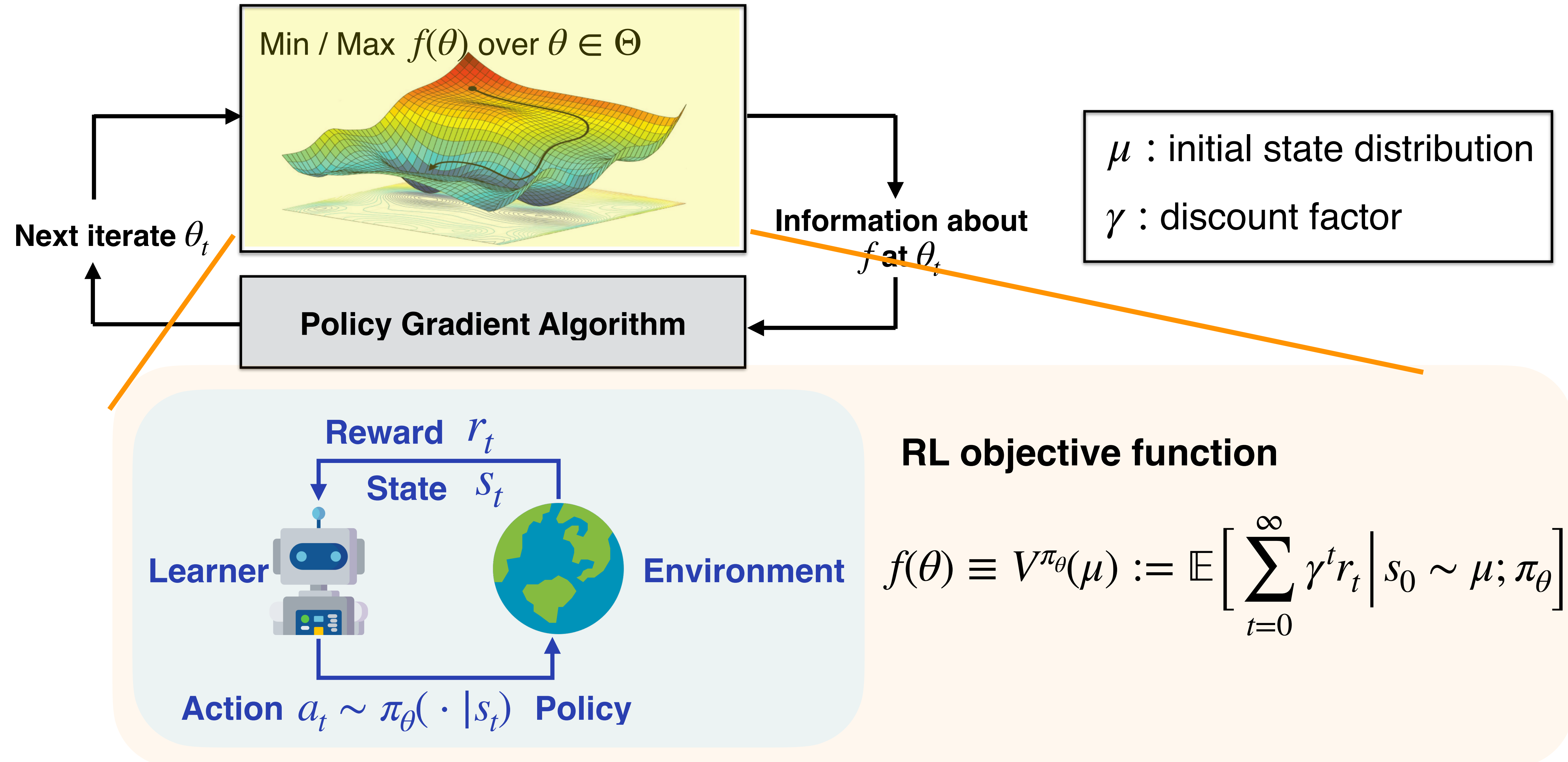
$$\|\nabla f(x)\|_2^2 = (Ax - y)^\top AA^\top(Ax - y) \dots\dots (\hspace{10em})$$

$$\geq \lambda_{\min}(AA^\top)\|Ax - y\|^2 \dots\dots (\hspace{10em})$$

$$= 2\lambda_{\min}(AA^\top)f(x) \dots\dots (\hspace{10em})$$

Notation:
 $A = [a_1, \dots, a_n]^\top \in \mathbb{R}^{n \times m}$
 $y = [y_1, \dots, y_n]^\top \in \mathbb{R}^{n \times 1}$

Example 2: Policy Gradient in Reinforcement Learning



Example 2: Non-Uniform PL in Reinforcement Learning

Non-Uniform PL-like Condition (Mei et al., ICML 2020)

Gradient norm

Non-uniformity

Sub-optimality gap

$$\left\| \frac{\partial V^{\pi_\theta}(\mu)}{\partial \theta} \right\|_2 \geq \frac{\min_s \pi_\theta(a^*(s)|s)}{\sqrt{S} \cdot \|d_\rho^{\pi^*} / d_\mu^{\pi_\theta}\|_\infty} \cdot [V^*(\rho) - V^{\pi_\theta}(\rho)] .$$

Nuance:

- Gradient could be extremely small if π_θ is far from an optimal one
- This “non-uniformity” results in complicated convergence analysis

Recent Breakthrough on Policy Gradient Theory in RL

(Agarwal et al., 2019)

(Mei et al., 2020)

(Xiao, 2022)

(Chen et al., 2024)

On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift

Alekh Agarwal* Sham M. Kakade† Jason D. Lee‡ Gaurav Mahajan§

Abstract

Policy gradient methods are among the most effective methods in challenging reinforcement learning problems with large state and/or action spaces. However, little is known about even their most basic theoretical convergence properties, including: if and how fast they converge to a globally optimal solution or how they cope with approximation error due to using a restricted class of parametric policies. This work provides provable characterizations of the computational, approximation, and sample size properties of policy gradient methods in the context of discounted Markov Decision Processes (MDPs). We focus on both: “tabular” policy parameterizations, where the optimal policy is contained in the class and where we show global convergence to the optimal policy; and parametric policy classes (considering both log-linear and neural policy classes), which may not contain the optimal policy and where we provide agnostic learning results. One central contribution of this work is in providing approximation guarantees that are average case — which avoid explicit worst-case dependencies on the size of state space — by making a formal connection to supervised learning under *distribution shift*. This characterization shows an important interplay between estimation error, approximation error, and exploration (as characterized through a precisely defined condition number).

On the Global Convergence Rates of Softmax Policy Gradient Methods

Jincheng Mei✱✱✱ Chenjun Xiao✱ Csaba Szepesvári♥✱ Dale Schuurmans✱✱

✱University of Alberta ♥DeepMind ✱Google Research, Brain Team

Abstract

We make three contributions toward better understanding policy gradient methods in the tabular setting. First, we show that with the true gradient, policy gradient with a softmax parametrization converges at a $O(1/t)$ rate, with constants depending on the problem and initialization. This result significantly expands the recent asymptotic convergence results. The analysis relies on two findings: that the softmax policy gradient satisfies a Łojasiewicz inequality, and the minimum probability of a policy gradient method to achieve effective empirical performance (e.g., Schulman et al., 2015; 2017).

On the Convergence Rates of Policy Gradient Methods

Lin Xiao

Meta AI Research

Seattle, WA 98109, USA

Editor: Alekh Agarwal

Abstract

We consider infinite-horizon discounted Markov decision problems with finite state and action spaces and study the convergence rates of the projected policy gradient method and a general class of policy mirror descent methods, all with direct parametrization in the policy space. First, we develop a theory of weak gradient-mapping dominance and use it to prove sharp sublinear convergence rate of the projected policy gradient method. Then we show that with geometrically increasing step sizes, a general class of policy mirror descent methods, including the natural policy gradient method and a projected Q-descent method, all enjoy a linear rate of convergence without relying on entropy or other strongly convex

Accelerated Policy Gradient: On the Convergence Rates of the Nesterov Momentum for Reinforcement Learning

Yen-Ju Chen*†1 Nai-Chieh Huang*†1 Ching-pei Lee2 Ping-Chun Hsieh1

Abstract

Various acceleration approaches for Policy Gradient (PG) have been analyzed within the realm of Reinforcement Learning (RL). However, the theoretical understanding of the widely used momentum-based acceleration method on PG remains largely open. In response to this gap, we adapt the celebrated Nesterov’s accelerated gradient (NAG) method to policy optimization in RL, termed *Accelerated Policy Gradient* (APG). To demonstrate the potential of APG in achieving fast convergence, we formally prove that with the true gradient and under the softmax policy parametrization, APG converges to an optimal policy at rates: (i) $\tilde{O}(1/t^2)$ with constant step sizes; (ii) $O(e^{-ct})$ with exponentially-growing step sizes. To the best of our knowledge, this is the first characterization of the convergence rates

1. Introduction

Policy gradient (PG) is a fundamental technique utilized in the field of reinforcement learning (RL) for policy optimization. It operates by directly optimizing the RL objectives to determine the optimal policy, employing first-order derivatives similar to the gradient descent algorithm in conventional optimization problems. Notably, PG has demonstrated empirical success (Mnih et al., 2016; Wang et al., 2016; Silver et al., 2014; Lillicrap et al., 2016; Schulman et al., 2017; Espeholt et al., 2018) and is supported by strong theoretical guarantees (Agarwal et al., 2021; Fazeli et al., 2018; Liu et al., 2020; Bhandari & Russo, 2019; Mei et al., 2020; Wang et al., 2021; Mei et al., 2021a, 2022; Xiao, 2022). In a recent study, Mei et al. (2020) characterize the $O(1/t)$ convergence rate of PG in the non-regularized tabular softmax setting. This convergence behavior aligns with that of the gradient descent algorithm for convex minimization problems, despite that the RL objectives lack concave characteristics for maximization. Additionally, advance-

COLT 2019

ICML 2020

JMLR 2022

ICML 2024

Asymptotic convergence to optimum under PG

The first convergence rate of $O(1/t)$ under PG

Similar rates for a large class of PG

Our recent result: $\tilde{O}(1/t^2)$ under Accelerated PG

GD for *convex* and *L-smooth* objective functions

Convergence of GD for Convex and Smooth Functions

Theorem (Convergence of GD under convexity and smoothness):

Let f be convex and L -smooth. Under GD with constant step sizes $\eta = 1/L$, we have

$$f(x_T) - f(x^*) \leq \frac{L}{2T} \cdot \|x_0 - x^*\|^2, \quad \forall T \in \mathbb{N}$$

How is this convergence rate compared to the strongly convex case?

A Useful Tool: “Descent Lemma”

Descent Lemma:

Let f be an L -smooth function (and not necessarily convex). Then, under GD with step size $\eta \leq 1/L$, we have

$$f(x_{t+1}) \leq f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2$$

Proof:

$$f(x_{t+1}) \leq f(x_t) + \nabla f(x_t)^\top (x_{t+1} - x_t) + \frac{L}{2} \|x_{t+1} - x_t\|^2 \quad \dots\dots (\quad)$$

$$= f(x_t) + \nabla f(x_t)^\top (- \eta \nabla f(x_t)) + \frac{L}{2} \left\| \eta \nabla f(x_t) \right\|^2 \quad \dots\dots (\quad)$$

$$= f(x_t) - \frac{\eta}{2} \|\nabla f(x_t)\|^2 \quad \dots\dots (\quad)$$

Proof: Convergence of GD for Convex and Smooth Functions

Step 1: Let's quantify the distance from x^*

$$\begin{aligned}\|x_t - x^*\|^2 &= \left\| (x_{t-1} - \eta \nabla f(x_{t-1})) - x^* \right\|^2 \\ &= \|x_{t-1} - x^*\|^2 - 2\eta \nabla f(x_{t-1})^\top (x_{t-1} - x^*) + \eta^2 \|\nabla f(x_{t-1})\|^2\end{aligned}$$

By reorganizing the terms, we have

$$\nabla f(x_{t-1})^\top (x_{t-1} - x^*) = \frac{1}{2\eta} (\|x_{t-1} - x^*\|_2^2 - \|x_t - x^*\|_2^2) + \frac{\eta}{2} \|\nabla f(x_{t-1})\|_2^2$$

Step 2:

$$\begin{aligned}\underline{f(x_{t-1}) - f(x^*)} &\leq \nabla f(x_{t-1})^\top (x_{t-1} - x^*) = \frac{1}{2\eta} (\|x_{t-1} - x^*\|_2^2 - \|x_t - x^*\|_2^2) + \frac{\eta}{2} \|\nabla f(x_{t-1})\|_2^2 \\ \text{(Why?)} &\leq \frac{1}{2\eta} (\|x_{t-1} - x^*\|_2^2 - \|x_t - x^*\|_2^2) + \underline{(f(x_{t-1}) - f(x_t))} \\ &\hspace{15em} \text{(Why?)}\end{aligned}$$

This implies

$$(f(x_t) - f(x^*)) \leq \frac{1}{2\eta} (\|x_{t-1} - x^*\|_2^2 - \|x_t - x^*\|_2^2)$$

Proof: Convergence of GD for Convex and Smooth Functions (Cont.)

Step 3: By taking the summation over $t = 1, \dots, T$

$$\sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{1}{2\eta} \|x_0 - x^*\|_2^2$$

Since GD is a descent algorithm, we have $f(x_0) \geq f(x_1) \geq \dots \geq f(x_T)$

$$f(x_T) - f(x^*) \leq \frac{1}{T} \sum_{t=1}^T (f(x_t) - f(x^*)) \leq \frac{1}{2T\eta} \|x_0 - x^*\|_2^2 = \frac{L}{2T} \|x_0 - x^*\|_2^2$$

Remarks

- Why selecting the step size $\eta = 1/L$?

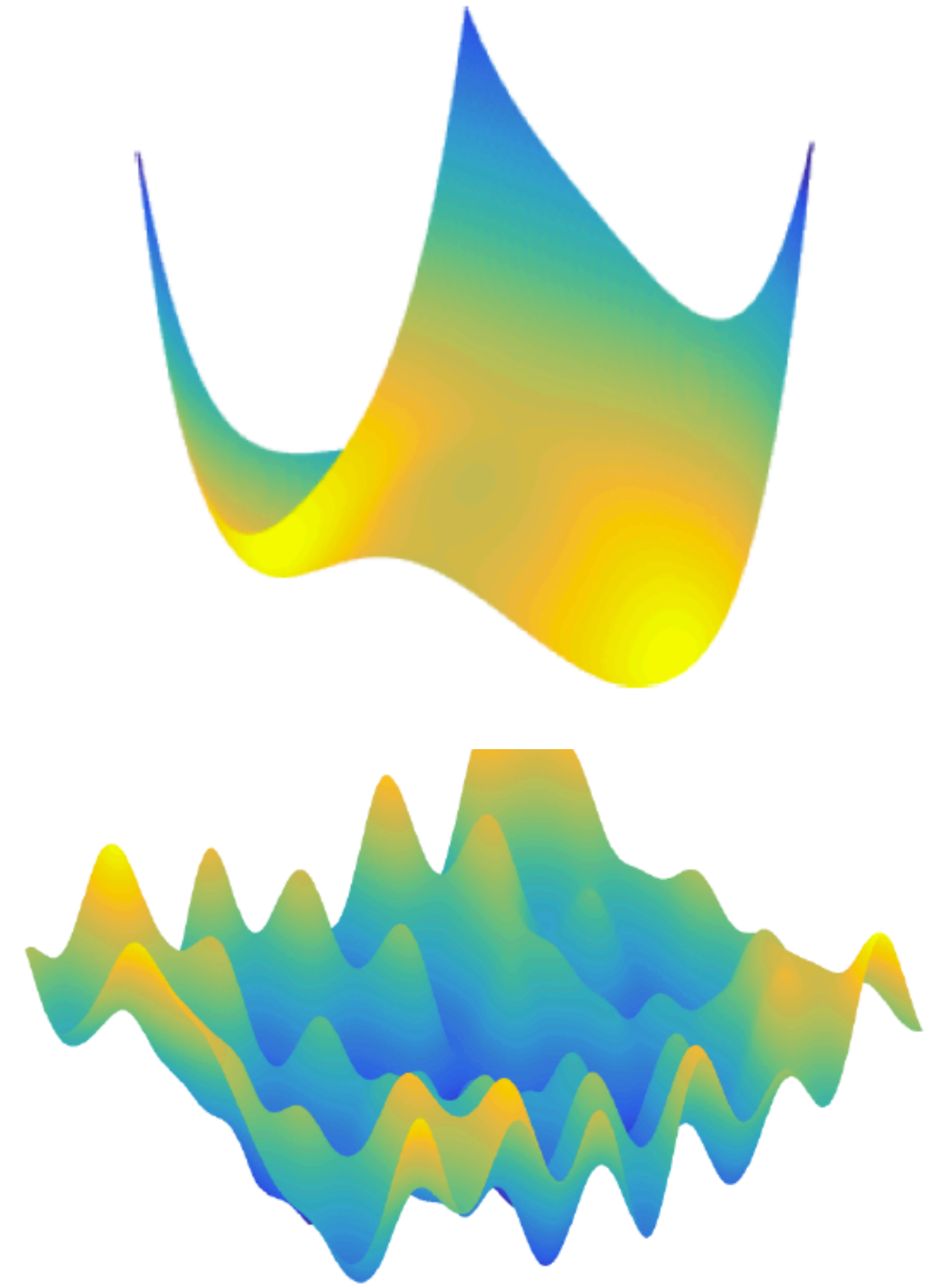
- What's the hidden assumption about x^* when we state the result

$$f(x_T) - f(x^*) \leq \frac{L}{2T} \cdot \|x_0 - x^*\|^2 ?$$

GD for non-convex and smooth objective functions

GD for General Non-Convex Problems?

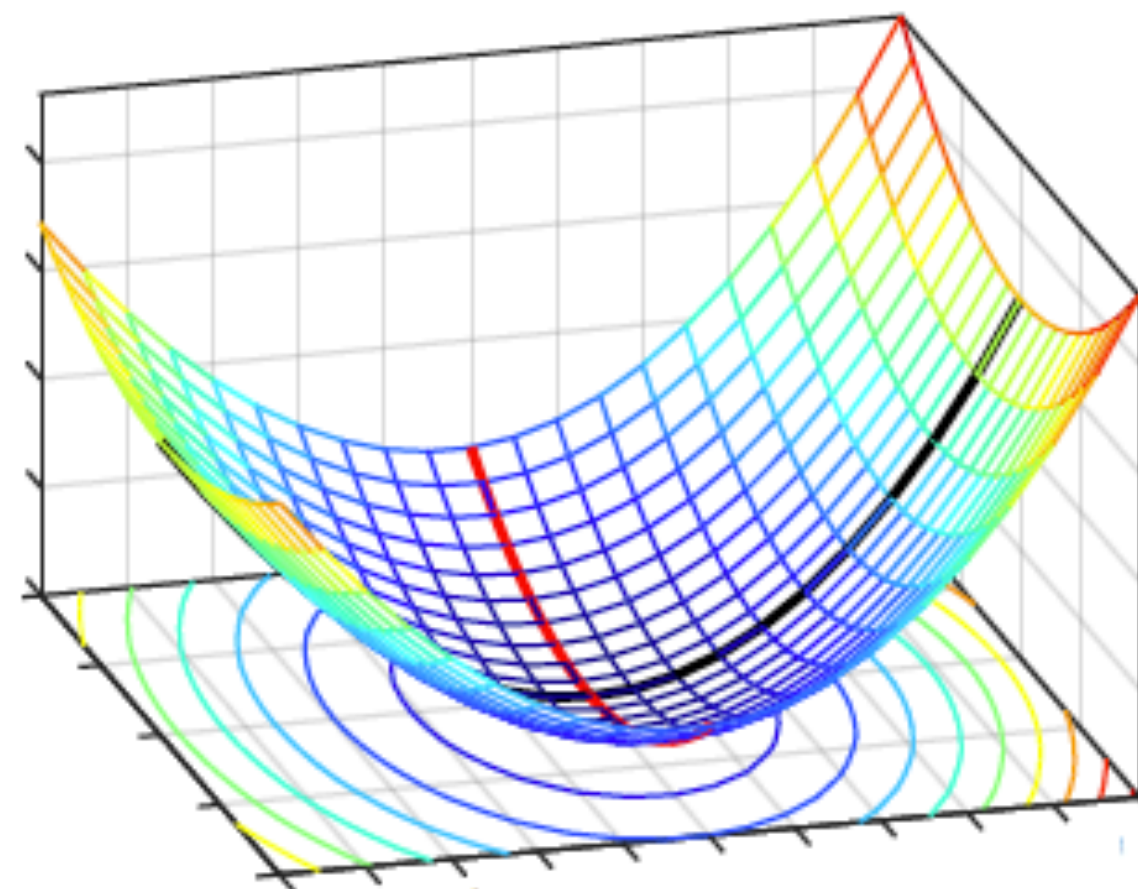
- ▶ Many machine learning problems are non-convex
 - ▶ Mixture models
 - ▶ Learning deep neural networks
 - ▶ Meta learning (e.g., MAML)
- ▶ Challenge
 - ▶ Bumps and local minima everywhere
 - ▶ No algorithm can solve non-convex problems efficiently in all cases



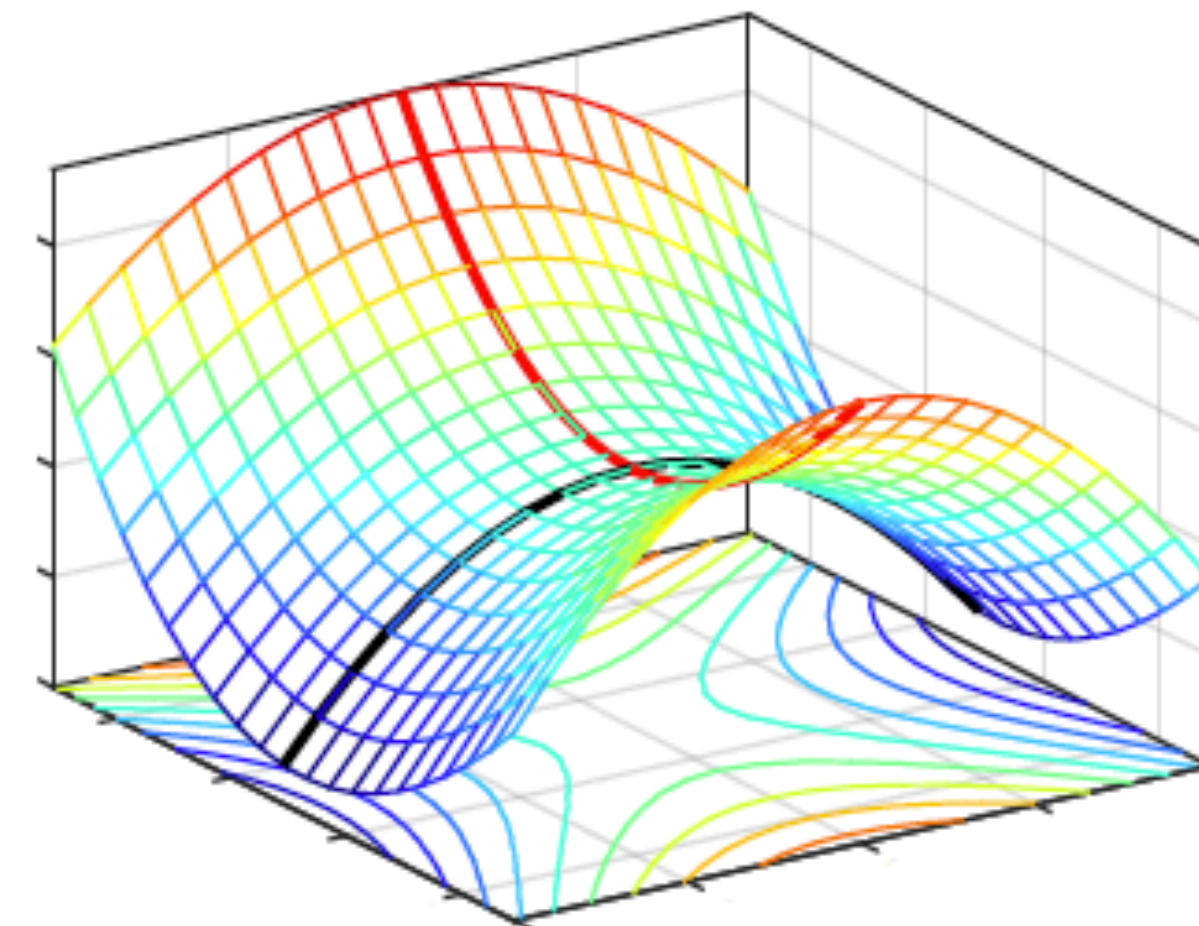
Typical Convergence Guarantees for Non-Convex Problems

- ▶ No efficient global convergence to global minima in general
- ▶ But, we may still hope for convergence to (nearly-)stationary points (i.e. $\|\nabla f(x)\| \leq \epsilon$)

There are at least 2 types of stationary points



global and local minimum



saddle point

Escaping Saddle Points under GD for Non-Convex Problems

- ▶ GD cannot always escape saddle points
- ▶ **Example**: if x^0 happens to be a saddle point, then GD is trapped (as $\nabla f(x^0) = 0$)
- ▶ **Existing results**: Under mild conditions (strict saddle property), **randomly initialized** GD converges to local minimum with probability 1

Lee et al., Gradient Descent Only Converges to Minimizers (COLT 2016)

Pascanu et al., On the saddle point problem for non-convex optimization (NIPS 2014)

- ▶ As a result, we are happy with finding a **(nearly-)stationary points** with

$$\|\nabla f(x)\| \leq \epsilon$$

Convergence of GD for Non-Convex and Smooth Functions

Theorem (Convergence of GD under only smoothness):

Let f be L -smooth. Under GD with constant step sizes $\eta = 1/L$, we have

$$(1) \quad \|\nabla f(x_t)\| \rightarrow 0, \quad \text{as } t \rightarrow \infty$$

$$(2) \quad \min_{0 \leq k \leq T} \|\nabla f(x_k)\| \leq \sqrt{\frac{2L(f(x_0) - f(x^*))}{T}}$$

Implication: GD reaches a nearly-stationary point after sufficiently many iterations

Question: Does (1) imply (2)? And how about vice versa?

Proof: Convergence of GD for Non-Convex and Smooth Functions

Step 1: By the descent lemma under GD with $\eta = 1/L$, we have

$$f(x_{t+1}) \leq f(x_t) - \frac{1}{2L} \|\nabla f(x_t)\|^2$$

Step 2: By taking the telescoping sum of the above,

$$\underbrace{\frac{1}{2L} \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2}_{LHS} \leq f(x_0) - f(x_T)$$

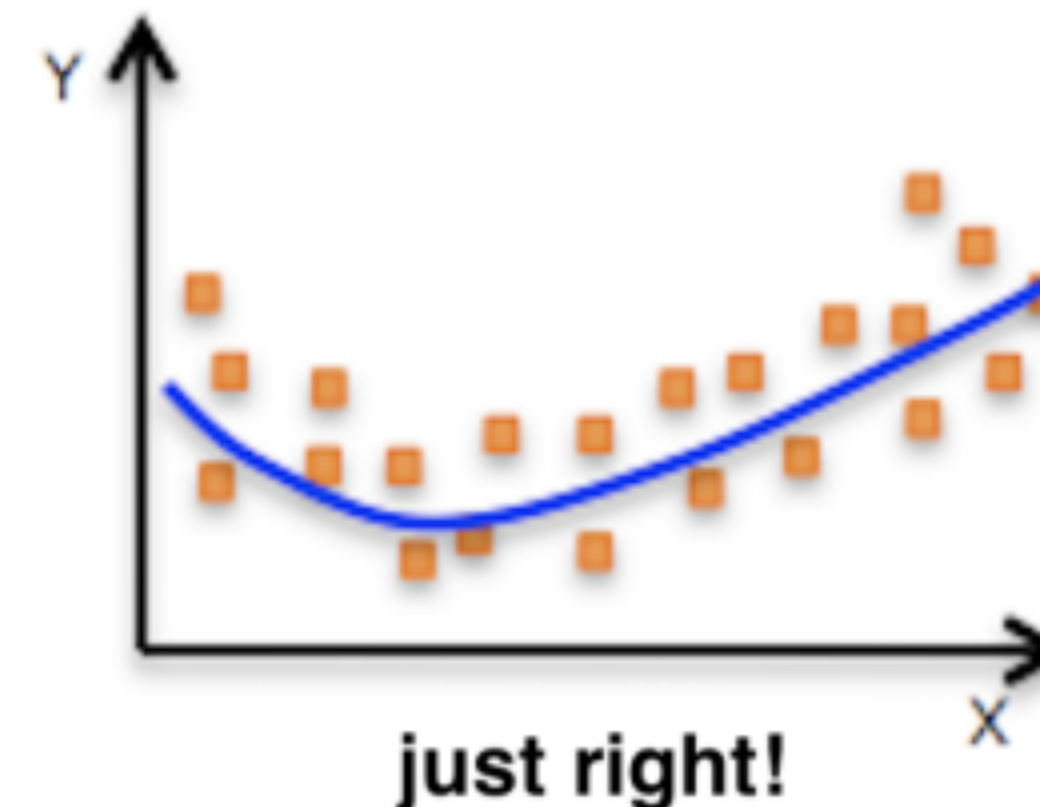
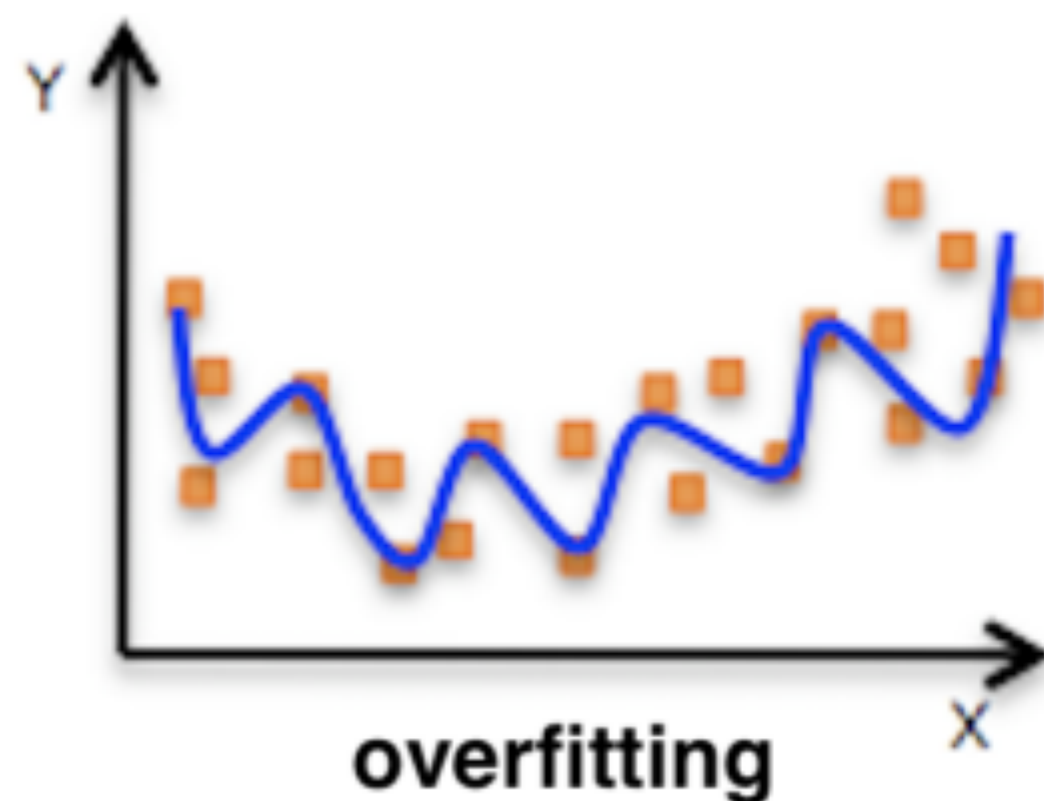
Step3: As LHS is non-negative and upper bounded, we know $\lim_{t \rightarrow \infty} \|\nabla f(x_t)\| = 0$

Moreover,
$$\min_{0 \leq k \leq T} \|\nabla f(x_k)\|^2 \leq \frac{1}{T} \sum_{k=0}^{T-1} \|\nabla f(x_k)\|^2 \leq \frac{2L(f(x_0) - f(x^*))}{T}$$

GD for Fully-Connected Neural Nets

Gradient Descent vs Neural Networks

- ▶ Back in 2010-2016, deep neural nets has shown significant empirical success in supervised learning
- ▶ Since 2016, a lot of research interests are focused on the question: “Why can GD converge for neural nets?”
- ▶ Answer: “overparameterization” (model size \gg training samples)



Example: Overparametrized Non-Linear Least-Squares

Non-linear least-squares regression: Given N data samples $\{x_i \in \mathbb{R}^n, y_i \in \mathbb{R}\}$, find a nonlinear model $f(\theta)$ by minimizing

$$L(\theta) = \frac{1}{2} \sum_{i=1}^n (f(x_i; \theta) - y_i)^2 = \frac{1}{2} \|f(\theta) - \mathbf{y}\|^2$$

where $\mathbf{y} := [y_1, \dots, y_n]^\top$, $f(\theta) := [f(x_1, \theta), \dots, f(x_n, \theta)]^\top$

Overparameterization: Model dimension $p >$ sample size n

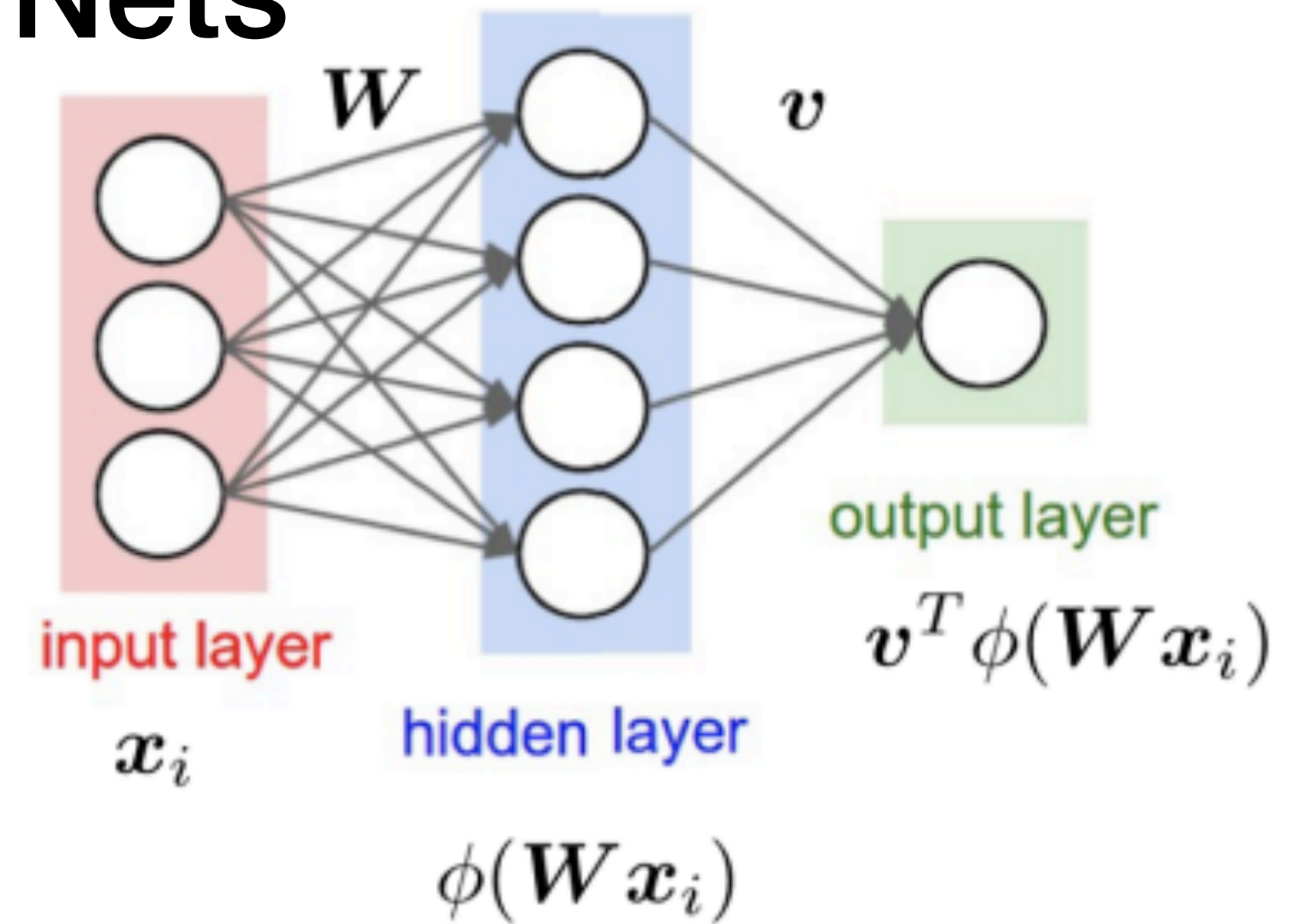
Run GD on this nonlinear least-squares problem: $\theta_{t+1} = \theta_t - \eta_t \nabla L(\theta_t)$

Gradient and Jacobian: $\nabla L(\theta) = J(\theta)^\top (f(\theta) - \mathbf{y})$ (Compared to linear case?)

where $J(\theta) = \frac{\partial f(\theta)}{\partial \theta} \in \mathbb{R}^{n \times p}$ is called the Jacobian

Example: One-Hidden Layer Neural Nets

- Training data:
 $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$
- Loss:
 $\mathcal{L}(\mathbf{v}, \mathbf{W}) := \sum_{i=1}^n (\mathbf{v}^T \phi(\mathbf{W} \mathbf{x}_i) - y_i)^2$
- Algorithm: gradient descent
with random Gaussian initialization



Theorem (Oymak and Soltanolkotabi 2019)

As long as

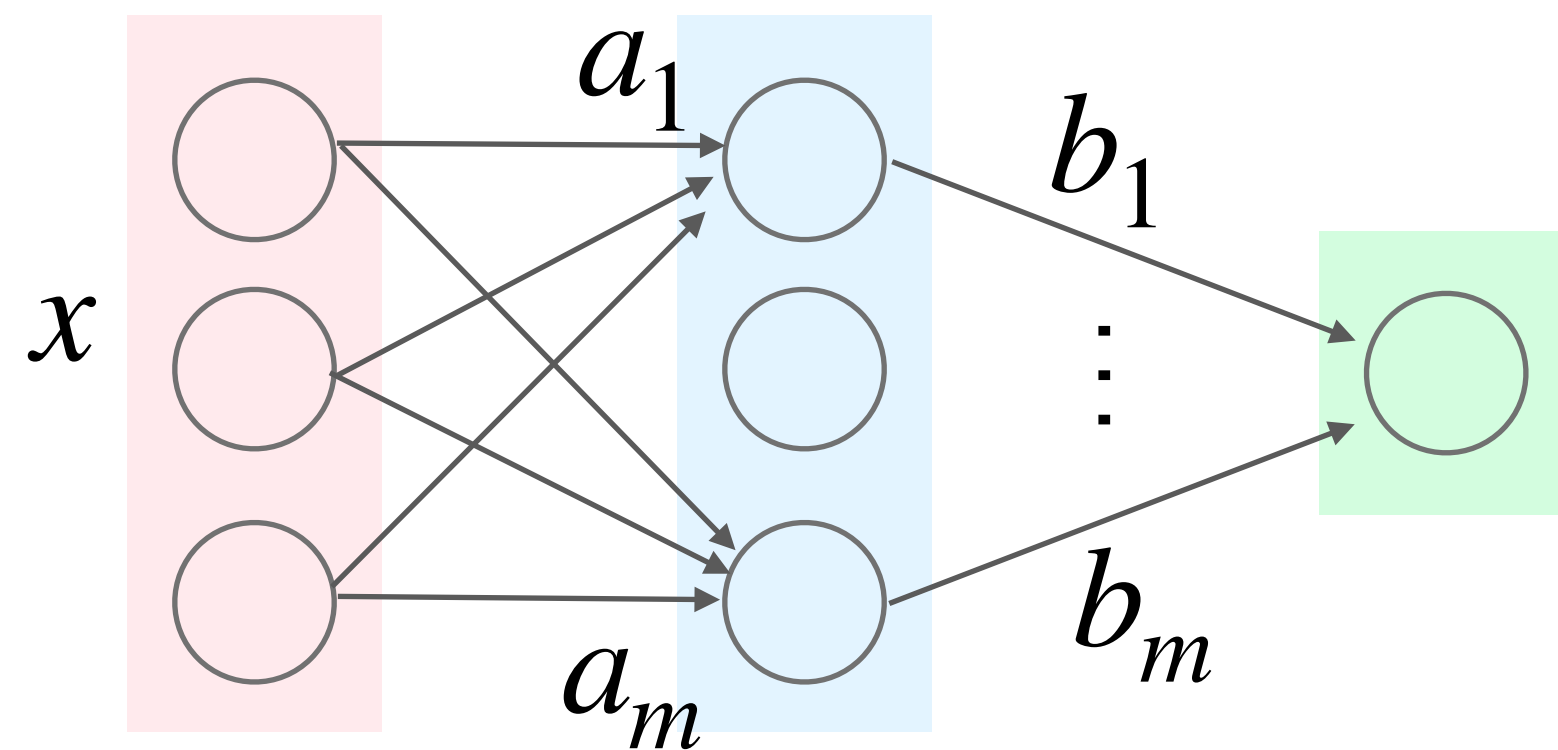
$$\#parameters \gtrsim (\#of\ training\ data)^2$$

Then, with high probability

- Zero training error: $\mathcal{L}(\mathbf{v}_\tau, \mathbf{W}_\tau) \leq (1 - \rho)^\tau \mathcal{L}(\mathbf{v}_0, \mathbf{W}_0)$
- Iterates remain close to initialization

An Alternative Explanation: Neural Tangent Kernel

An empirical observation about NNs:



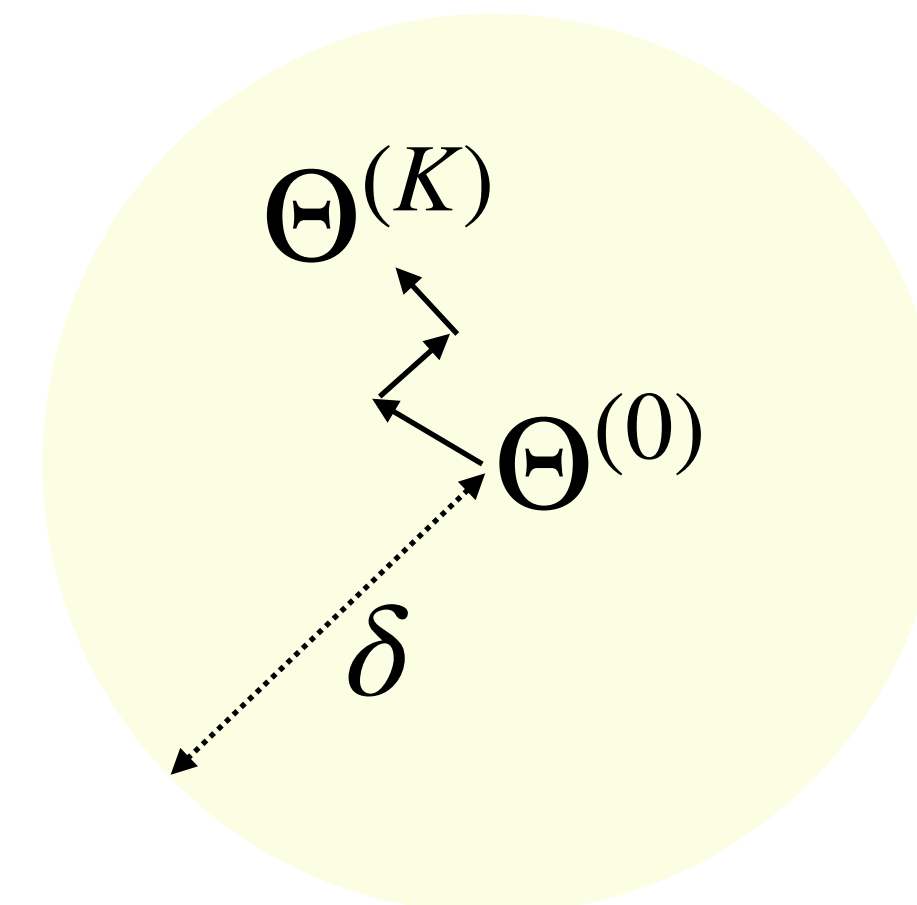
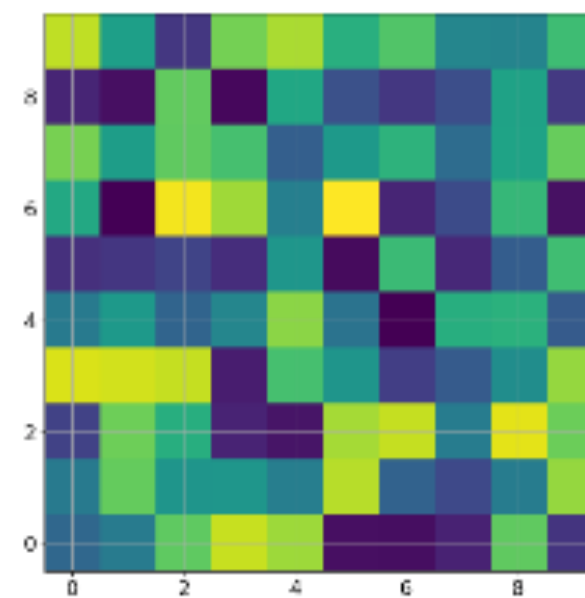
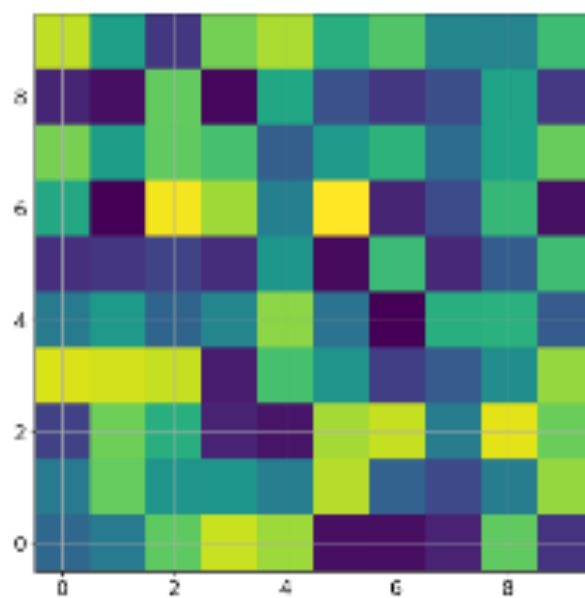
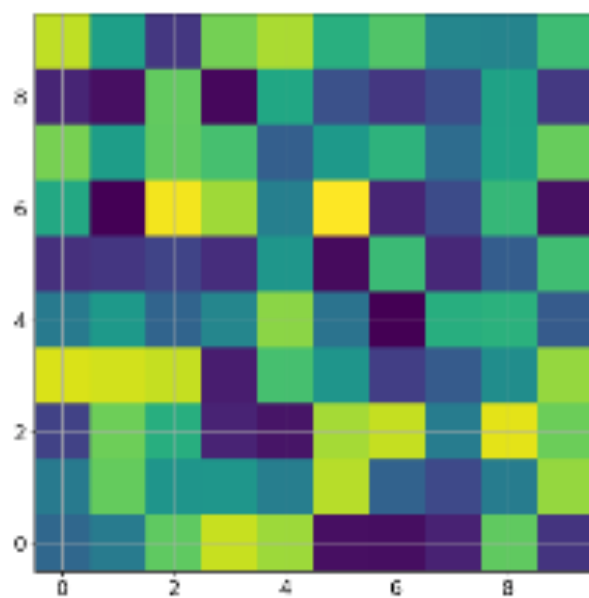
$$f_m(x; \Theta) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \cdot \sigma(a_i^\top x)$$

(σ : ReLU activation)

$$(\Theta \equiv \{[a_i]_{i=1}^m, [b_i]_{i=1}^m\})$$

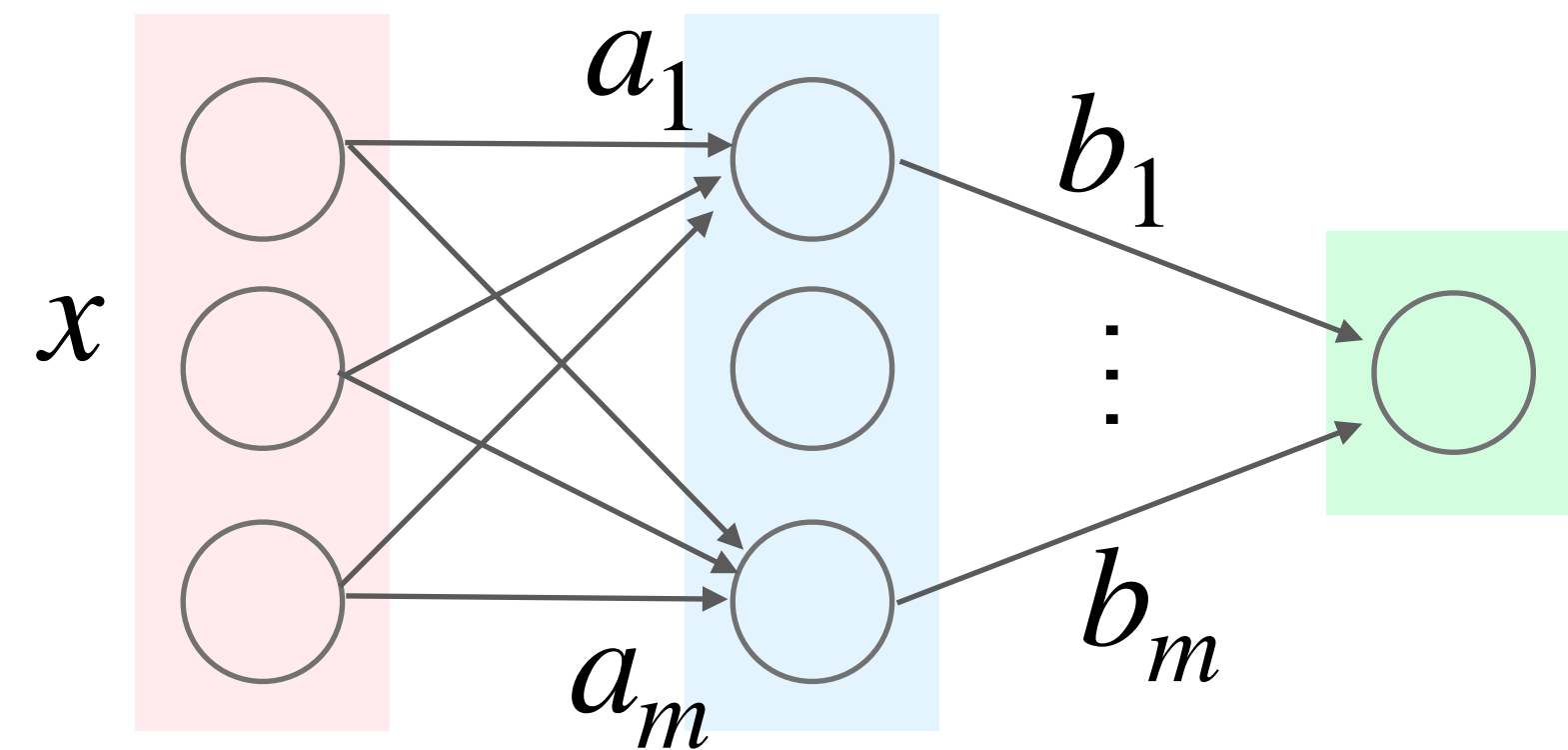
(Random Initialization)

$$\Theta^{(0)} \xrightarrow{\text{GD}} \Theta^{(1)} \xrightarrow{\text{GD}} \Theta^{(2)}$$



NN parameters “almost” static (under large m)

A Primer on Neural Tangent Kernel (NTK)



$$f_m(x; \Theta) = \frac{1}{\sqrt{m}} \sum_{i=1}^m b_i \cdot \sigma(a_i^\top x)$$

(σ : ReLU activation)

($\Theta \equiv \{[a_i]_{i=1}^m, [b_i]_{i=1}^m\}$)

First-order Taylor expansion:

$$f_m(x; \Theta) = f_m(x; \Theta^{(0)}) + \underbrace{\nabla_{\Theta} f_m(x; \Theta^{(0)})^\top}_{\text{Viewed as a feature map in kernel methods}} (\Theta - \Theta^{(0)}) + O(\|\Theta - \Theta^{(0)}\|^2)$$

Viewed as a feature map in kernel methods

Neural tangent kernel function [Jacot et al., 2018]:

$$\mathbf{H}_m(x, x') := \langle \nabla_{\Theta} f_m(x; \Theta^{(0)}), \nabla_{\Theta} f_m(x'; \Theta^{(0)}) \rangle$$

$\downarrow m \rightarrow \infty$

$$\mathbf{H}(x, x') = \mathbb{E}_{a,b}[b^2 \sigma'(a^\top x) \sigma'(a^\top x') \langle x, x' \rangle] + \mathbb{E}_a[\sigma(a^\top x) \sigma(a^\top x')]$$

References: GD in Overparameterization Regime

(ICLR 2019)

GRADIENT DESCENT PROVABLY OPTIMIZES OVER-PARAMETERIZED NEURAL NETWORKS

Simon S. Du*

Machine Learning Department
Carnegie Mellon University
ssdu@cs.cmu.edu

Xiyu Zhai*

Department of EECS
Massachusetts Institute of Technology
xiyuzhai@mit.edu

Barnabás Póczos

Machine Learning Department
Carnegie Mellon University
bapozos@cs.cmu.edu

Aarti Singh

Machine Learning Department
Carnegie Mellon University
aartisinh@cmu.edu

(NeurIPS 2018)

Neural Tangent Kernel: Convergence and Generalization in Neural Networks

Arthur Jacot

École Polytechnique Fédérale de Lausanne
arthur.jacot@netopera.net

Franck Gabriel

Imperial College London and École Polytechnique Fédérale de Lausanne
franckrgabriel@gmail.com

Clément Hongler

École Polytechnique Fédérale de Lausanne
clement.hongler@gmail.com

(ICML 2019)

Overparameterized Nonlinear Learning: Gradient Descent Takes the Shortest Path?

Samet Oymak¹ Mahdi Soltanolkotabi²