

# Oxford Machine Learning Summer School 2022

## Optimisation for Machine Learning

Yali Du

Assistant professor

Department of Informatics, King's College London

# People



Instructor Yali Du  
[yali.du@kcl.ac.uk](mailto:yali.du@kcl.ac.uk)



KCL Cristina Gava  
[cristina.gava@kcl.ac.uk](mailto:cristina.gava@kcl.ac.uk)



KCL Ziyan Wang:  
[wangziyan1998@gmail.com](mailto:wangziyan1998@gmail.com)



CAS Xue Yan  
[yanxuesdu@gmail.com](mailto:yanxuesdu@gmail.com)



CAS Runji Lin:  
[linprophet@icloud.com](mailto:linprophet@icloud.com)

# Content

Part 0: Introduction to introduction ...

Part 1: Theory of Convexity

Part 2: Gradient Descent and convergence

Part 3: Stochastic Gradient Descent, Non-Convex Optimization, Adam, and Machine Learning Applications

Practical work

<https://colab.research.google.com/drive/1r-FTu17Utca7JjblyubSnFEC0n5ph3V6?usp=sharing>

# Optimization

General optimization problem (unconstrained minimization)

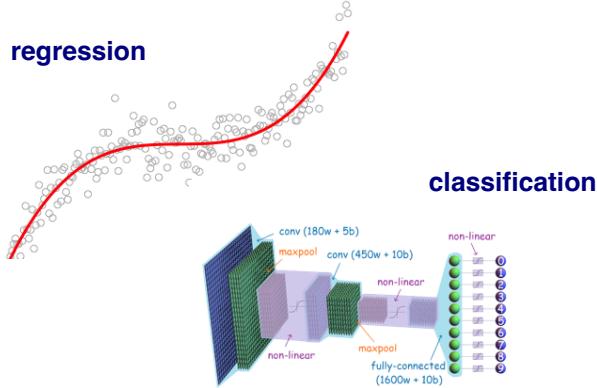
$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{with} & \mathbf{x} \in \mathbb{R}^d\end{array}$$

candidate solutions, variables, parameters  $\mathbf{x} \in \mathbb{R}^d$

objective function  $f : \mathbb{R}^d \mapsto \mathbb{R}$

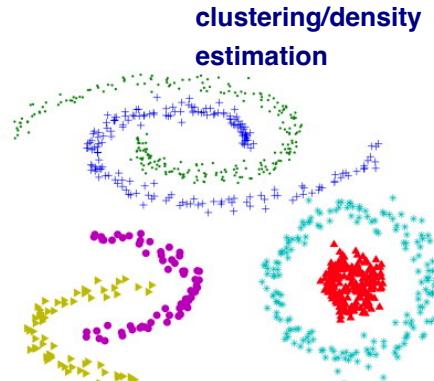
typically: technical assumption:  $f$  is continuous and differentiable

# why optimization?



**Supervised  
Learning**

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}, \mathbf{a}_i, b_i)$$



**Unsupervised  
Learning**

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\mathbf{x}, \mathbf{a}_i)$$

**computer games**



**Reinforcement  
Learning**

$$\max_{\mathbf{x} \in \mathbb{R}^d} \mathbb{E}_{\tau \sim p_{\mathbf{x}(\tau)}}[R(\tau)]$$

**... all these involve a minimisation of some function ...**

$$\min_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x})$$

# Optimization for machine learning

Mathematical Modeling:

- defining & and measuring the machine learning model

Computational Optimization:

- learning the model parameters

Theory vs. practice:

- libraries are available, algorithms treated as “black box” by most practitioners
- **Not here:** we look inside the algorithms and try to understand why and how fast they work!

# Solving optimization problems

Optimization at large scale:

Main approaches:

- Gradient Descent
- Stochastic Gradient Descent (SGD)
- ...

History:

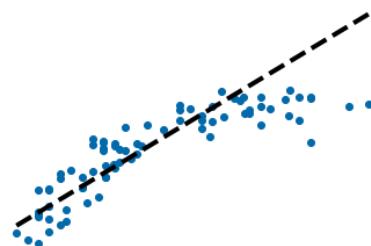
- 1847: Cauchy proposes gradient descent
- 1950s: Linear Programs, soon followed by non-linear, SGD
- 1980s: General optimization, convergence theory
- 2005-2015: Large scale optimization (mostly convex), convergence of SGD
- 2015-today: Improved understanding of SGD for deep learning

# Function types, and what one can hope for ...

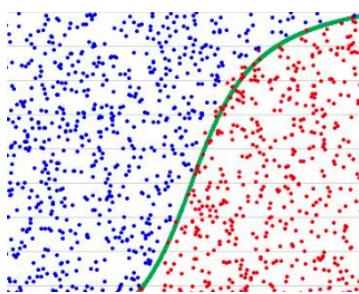
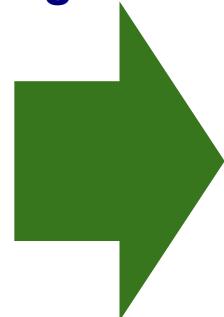
... optimising for unknown parameters depends on the type of function under study ...

**Convex**

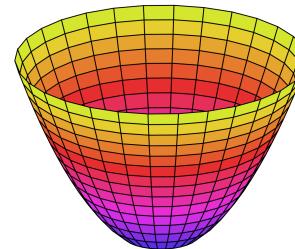
$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}), \lambda \in [0,1], \forall \mathbf{x}, \mathbf{y} \in \mathbb{R}^d$$



**Linear  
Regression**



**Classification with  
Hinge Loss**

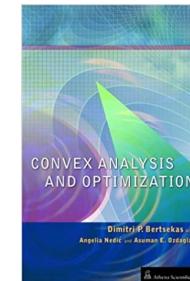


**Unique global  
minimum**

$$\begin{matrix} n \\ m \end{matrix} \begin{matrix} A \\ \approx \\ W \end{matrix} \begin{matrix} k \\ m \end{matrix} \quad \begin{matrix} n \\ H \end{matrix} \begin{matrix} k \\ \end{matrix}$$

$A \geq 0, W \geq 0, H \geq 0$

**Non-Negative Matrix  
Factorisation**

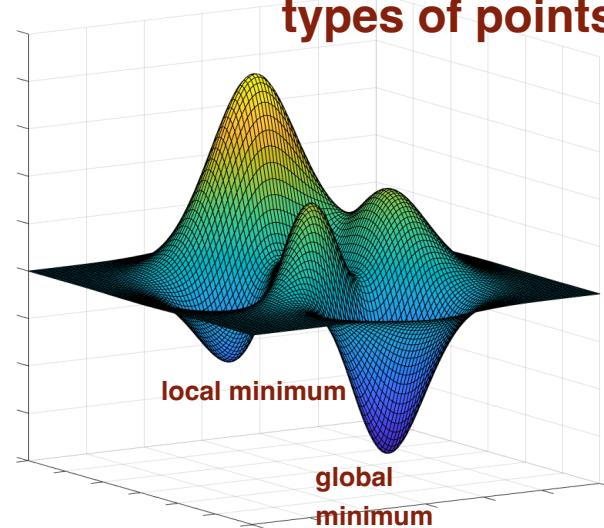


**... admits polynomial time  
algorithms**

# Function types, and what one can hope for ...

... optimising for unknown parameters depends on the type of function under study ...

Non-Convex



... global and local minima (checking) are NP-Hard, we look for other types of points ...

$$\nabla_{\mathbf{x}} f(\mathbf{x}_{\text{stationary}}) = 0$$

... so instead, the community is fetching for stationary points ...

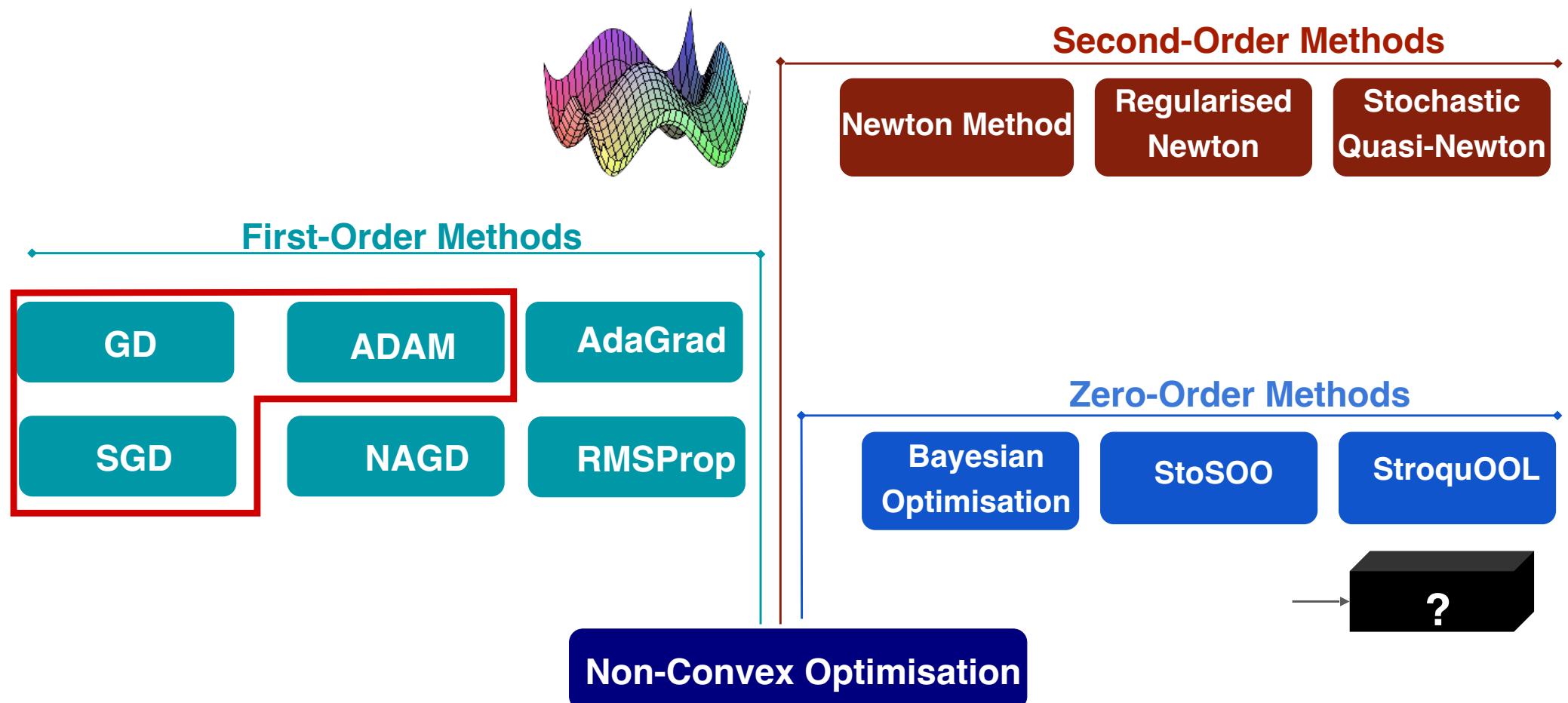
1.  $\epsilon$ -First-Order-Stationary Point (FOSP):  $\|\nabla_{\mathbf{x}} f(\mathbf{x}_{\text{FOSP}})\|_2 \leq \epsilon$   
[e.g., all global and local minima, saddle points, plateau points]

2.  $\epsilon$ -Second-Order-Stationary Point (SOSP):

$$\|\nabla_{\mathbf{x}} f(\mathbf{x}_{\text{SOSP}})\|_2 \leq \epsilon \quad \text{and} \quad \lambda_{\min}(\nabla_{\mathbf{x}}^2 f(\mathbf{x}_{\text{SOSP}})) \geq -\sqrt{\epsilon}$$

[e.g., all global and local minima, plateau points]

# Algorithms vary in type of information used ...



# References

This tutorial is based on

Books:

- Convex Optimization by Stephen Boyd & Lieven Vandenberghe, 2003
- Convex Optimization: Algorithms and Complexity, by Sébastien Bubeck 2015

Course:

- Lecture notes: Optimization for Machine Learning, by Martin Jaggi & Nicolas Flammarion
- UTDallas Course - Convex Optimization- CS 7301, by Nicholas Ruozzi
- ADAM's Story and Proof - [Machine Learning and AI Academy](#), by Haitham Bou-Ammar

# Notes on Notations

$\mathbb{R}$  : reals

$\mathbb{R}^d$  :  $d$ -dimensional Euclidean space

$\mathbf{x} \in \mathbb{R}^d$  : the variable or unknown parameters

$f : \mathbb{R}^d \mapsto \mathbb{R}$  the objective or cost function, and  $\mathbb{R}^d$  is called the domain of  $f$ , or

$\text{dom}(f) = \mathbb{R}^d$

$\nabla f(\mathbf{x}) = \left( \frac{\partial f(\mathbf{x}_1)}{\partial x_1}, \dots, \frac{\partial f(\mathbf{x}_d)}{\partial x_d} \right)^T$  : the vector of partial derivatives, also the gradient of  $f$  at  $\mathbf{x}$

# Content

Part 0: Introduction to introduction ...

Part 1: Theory of Convexity

Part 2: Gradient Descent and convergence

Part 3: Stochastic Gradient Descent, Non-Convex Optimization, Adam, and Machine Learning Applications

Practical work

<https://colab.research.google.com/drive/1r-FTu17Utca7JjblyubSnFEC0n5ph3V6?usp=sharing>

# Theory of convexity

# Convex Sets

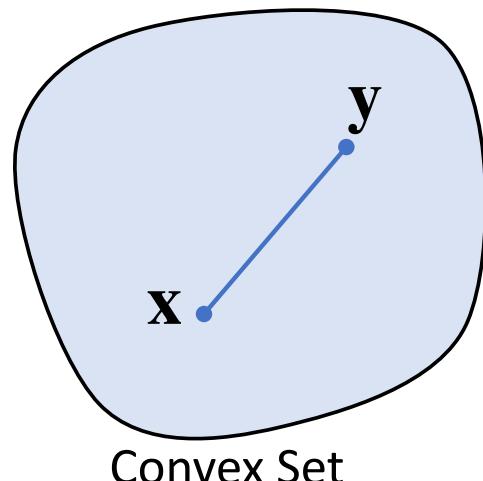
**line segment** between  $\mathbf{x}_1$  and  $\mathbf{x}_2$  : all points

$$\mathbf{x} = t\mathbf{x}_1 + (1 - t)\mathbf{x}_2$$

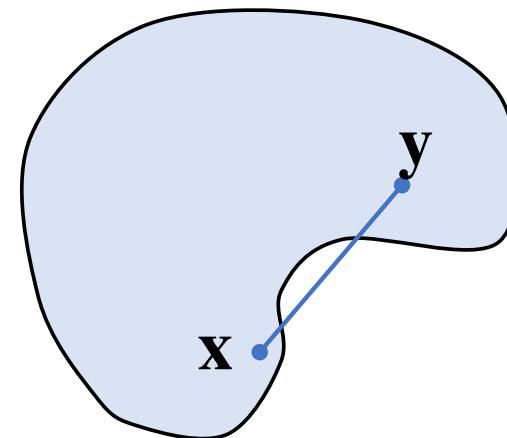
with  $0 \leq t \leq 1$ .

**convex set**: contains line segment between any two points in the set

That is, if  $\mathbf{x}, \mathbf{y} \in C$ , then  $t\mathbf{x} + (1 - t)\mathbf{y} \in C$  for all  $t \in [0,1]$



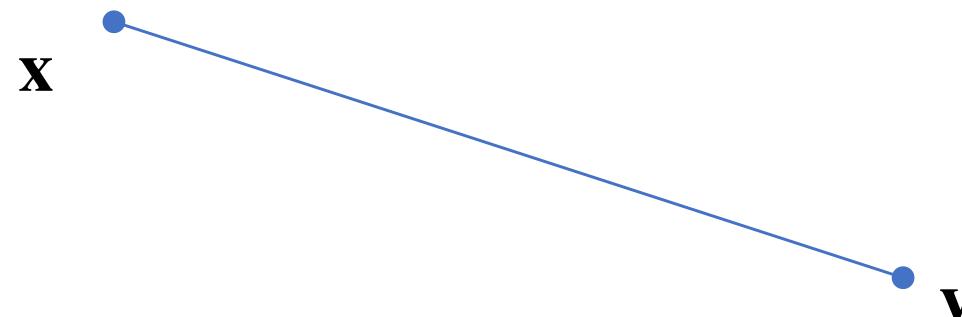
Convex Set



Not a Convex Set

# Examples of Convex Sets

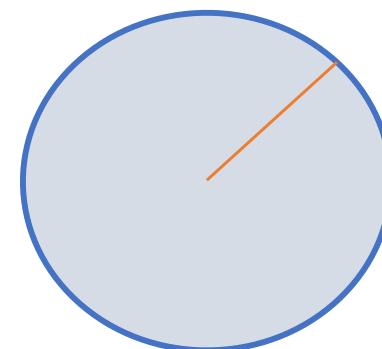
**Line Segments:**  $C = \{t\mathbf{x} + (1 - t)\mathbf{y} : t \in [0,1]\}$  for some  $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$



Lines, planes, hyperplanes, etc. also define convex sets

**Balls of Radius  $\epsilon$ :**  $C = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\|_2 \leq \epsilon\}$  for some  $\epsilon \geq 0 \in \mathbb{R}$

$$\|\mathbf{x}\|_2 = \sqrt{\mathbf{x}^T \mathbf{x}} = \sqrt{\sum_{i=1}^n x_i^2}$$



# Properties of Convex Sets

Intersections of convex sets are convex

**Observation.** Let  $C_i, i \in I$  be convex sets, where  $I$  is a (possibly infinite) index set. Then  $C = \cap_{i \in I} C_i$  is a convex set.

Projections onto convex sets are unique, and often efficient to compute

$$P_C(\mathbf{x}) := \operatorname{argmin}_{\mathbf{y} \in C} \|\mathbf{y} - \mathbf{x}\|$$

# Convex functions

# Convex Functions

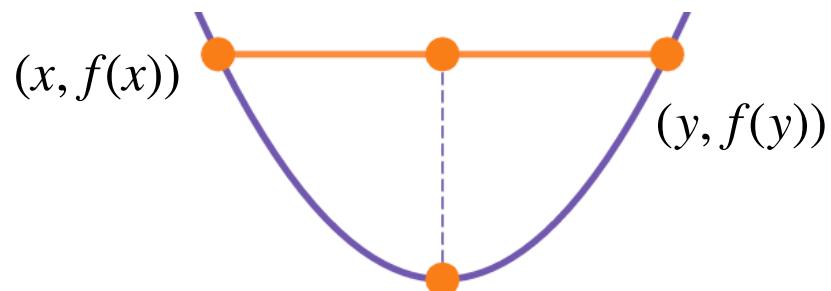
A **function**  $f: C \rightarrow \mathbb{R}^n$  is **convex** if  $C$  is a convex set and

$$f(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y})$$

for all  $\mathbf{x}, \mathbf{y} \in C$  and  $\lambda \in [0,1]$

**Geometrically:** The line segment between  $(\mathbf{x}, f(\mathbf{x}))$  and  $(\mathbf{y}, f(\mathbf{y}))$  lies above the graph of  $f$ .

$f$  is called **concave** if  $-f$  is convex



Smiley face



# Examples

convex:

- affine function:  $\mathbf{a}^T \mathbf{x} + b$
- exponential:  $e^{ax}$ , for any  $a \in \mathbb{R}$
- powers:  $x^\alpha$  on  $\mathbb{R}_+$ , for  $\alpha \geq 1$  or  $\alpha \leq 0$
- negative entropy:  $x \log x$  on  $\mathbb{R}_+$

concave:

- affine:  $\mathbf{a}^T \mathbf{x} + b$
- powers:  $x^\alpha$  on  $\mathbb{R}_+$ , for  $0 \leq \alpha \leq 1$
- logarithm:  $\log x$  on  $\mathbb{R}_+$

Note:

- $\mathbb{R}_+$  indicate positive reals.
- affine functions are convex and concave.

# First-order condition

$f$  is differentiable if  $C$  is open and the gradient

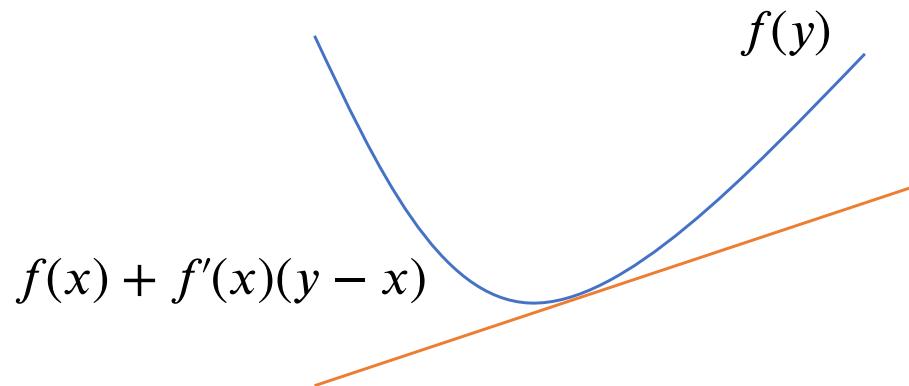
$$\nabla f(\mathbf{x}) := \left( \frac{\partial f}{\partial x_1}(\mathbf{x}), \dots, \frac{\partial f}{\partial x_d}(\mathbf{x}) \right)^T$$

exist at each  $\mathbf{x} \in C$

**1st-order condition:** differentiable  $f$  with convex domain is convex iff

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \mathbf{x}, \mathbf{y} \in \text{dom}(f)$$

i.e. first-order approximation of  $f$  is global underestimator



# First-order condition

i.e. first-order approximation of  $f$  is global underestimator

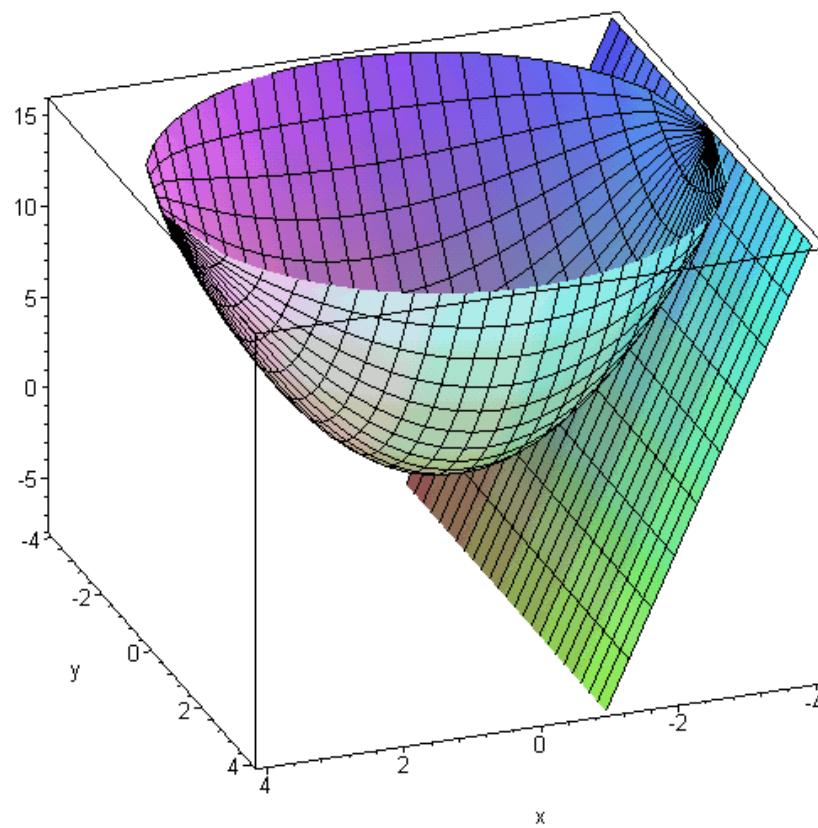


Image: Lane Vosbury, Seminole State College

# Second-order condition

$f$  is twice differentiable if  $C$  is open and the Hessian  $\nabla^2 f(\mathbf{x})$ ,

$$\nabla^2 f(\mathbf{x})_{ij} = \frac{\partial^2 f}{\partial x_i \partial x_j}(\mathbf{x}), i, j = 1, \dots, d$$

exist at each  $\mathbf{x} \in C$

**2nd-order conditions:** for twice differentiable  $f$  with convex domain  
 $f$  is convex if and only if

$$\nabla^2 f(\mathbf{x}) \geq 0 \text{ for all } \mathbf{x} \in C$$

i.e. the hessian is positive semi-definite (all of the eigenvalues are nonnegative).

# Second-order conditions

Special case when  $C \subseteq \mathbb{R}$

A twice continuously differentiable **function**  $f: C \rightarrow \mathbb{R}$  with  $C \subseteq \mathbb{R}$  is **convex** if  $C$  is a convex set and

$$\frac{d^2 f}{dx^2}(x) \geq 0$$

for all  $x \in C$

# Operations Preserving Convexity

**Nonnegative weighted sums** of convex functions are convex, i.e., if  $f_1: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $f_2: \mathbb{R}^n \rightarrow \mathbb{R}$  are convex functions and  $c_1, c_2 \geq 0$ , then

$$g(\mathbf{x}) = c_1 f_1(\mathbf{x}) + c_2 f_2(\mathbf{x})$$

is a convex function.

**Proof:**

Let  $\mathbf{z} = t\mathbf{x} + (1 - t)\mathbf{y}$ ,

$$\begin{aligned} g(\mathbf{z}) &= c_1 f_1(\mathbf{z}) + c_2 f_2(\mathbf{z}) \\ &\leq c_1(t f_1(\mathbf{x}) + (1 - t) f_1(\mathbf{y})) + c_2(t f_2(\mathbf{x}) + (1 - t) f_2(\mathbf{y})) \\ &= t(c_1 f_1(\mathbf{x}) + c_2 f_2(\mathbf{x})) + (1 - t)c_1(f_1(\mathbf{y})) + c_2(f_2(\mathbf{y})) \\ &= t g(\mathbf{x}) + (1 - t) g(\mathbf{y}) \end{aligned}$$

The statement follows.

# Operations Preserving Convexity

**Composition with an affine function**, i.e., if  $f: \mathbb{R}^n \rightarrow \mathbb{R}$  is a convex function and  $M \in \mathbb{R}^{n \times m}$ ,  $b \in \mathbb{R}^n$ , then  $g: \mathbb{R}^m \rightarrow \mathbb{R}$  given by

$$g(\mathbf{x}) = f(M\mathbf{x} + b)$$

is a convex function.

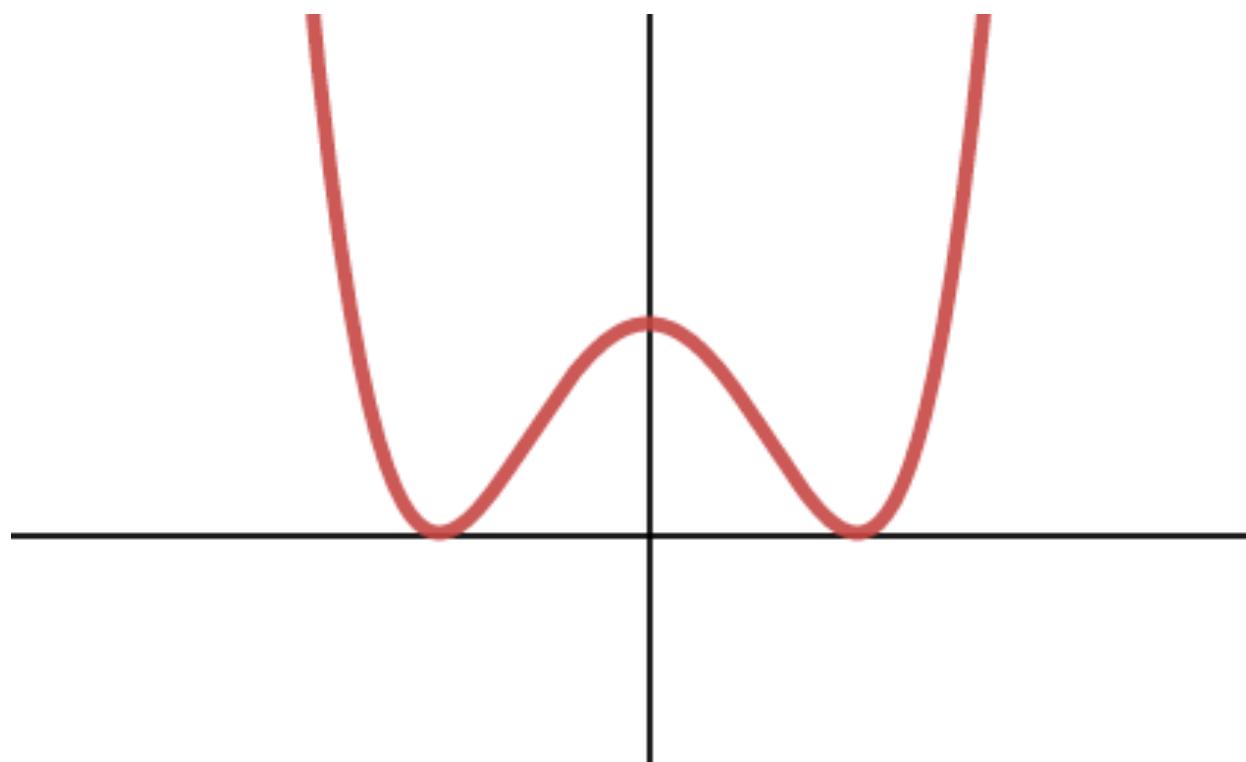
Questions:

Would any composition of two convex functions be convex?

# Counterexample: composite function

**Counterexample:** Consider two convex functions  $f(x) = x^2$  and  $g(x) = x^2 - 1$ . Then, the composition

$$(f \circ g)(x) = x^4 - 2x^2 + 1.$$

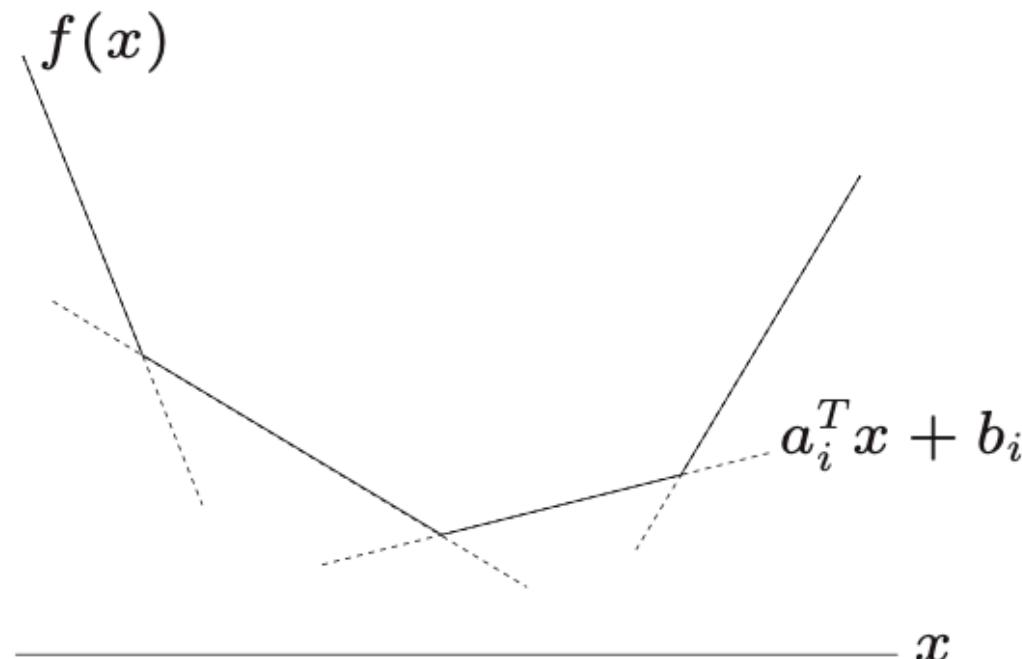


# Operations Preserving Convexity

**Pointwise maximum** of convex functions are convex, i.e., if  $f_1: \mathbb{R}^n \rightarrow \mathbb{R}$  and  $f_2: \mathbb{R}^n \rightarrow \mathbb{R}$  are convex functions, then

$$g(\mathbf{x}) = \max(f_1(\mathbf{x}), f_2(\mathbf{x}))$$

is a convex function.



# Convex optimisation

# Optimization

Convex Optimization Problem is of the form

$$\begin{array}{ll}\text{minimize} & f(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in X\end{array}$$

Where both

- $f$  is a convex function
- $X \subseteq \text{dom}(f)$  is a convex set

Crucial Property of Convex Optimization Problems

- Every local minimum is a **global minimum**, see later...

# Why do we favor convex optimization ?

For convex optimization problems, all algorithms

- Coordinate Descent, Gradient Descent, Stochastic Gradient Descent, Projected and Proximal Gradient Descent

do **converge** to the global optimum! (assuming  $f$  differentiable)

**Example Theorem:** The convergence rate is proportional to  $\frac{1}{t}$ ,  
i.e.

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{c}{t}$$

(where  $\mathbf{x}^*$  is some optimal solution to the problem.)

Meaning: **Approximation error** converges to 0 over time.

# Local Minima are Global Minima

## Definition

A **local minimum** of  $f : \text{dom}(f) \mapsto \mathbb{R}$  is a point  $\mathbf{x}$  such that there exists  $\epsilon > 0$  with

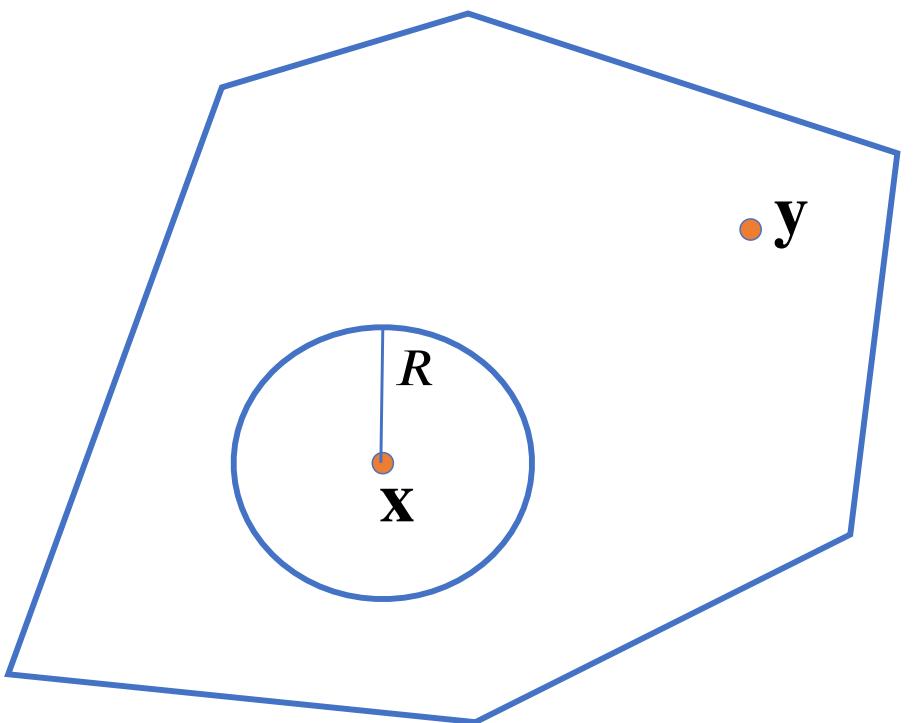
$$f(\mathbf{x}) \leq f(\mathbf{y}) \quad \forall \mathbf{y} \in \text{dom}(f) \text{ satisfying } \|\mathbf{y} - \mathbf{x}\| < \epsilon$$

**Lemma:** Let  $\mathbf{x}^*$  be a local minimum of a convex function  $f : \text{dom}(f) \mapsto \mathbb{R}$ . Then  $\mathbf{x}^*$  is a global minimum, meaning that  $f(\mathbf{x}^*) \leq f(\mathbf{y}) \quad \forall \mathbf{y} \in \text{dom}(f)$ .

# Local Minima are Global Minima

Proof.

Let  $\mathbf{x}$  be a local minimum, assume there exists  $\mathbf{y} \in \text{dom}(f)$  such that  $f(\mathbf{x}) > f(\mathbf{y})$



$\mathbf{x}$  is a local minimum implies that there exists a radius  $R$  such that for all points  $\mathbf{x}'$  within the ball of radius  $R$  of  $\mathbf{x}$ ,

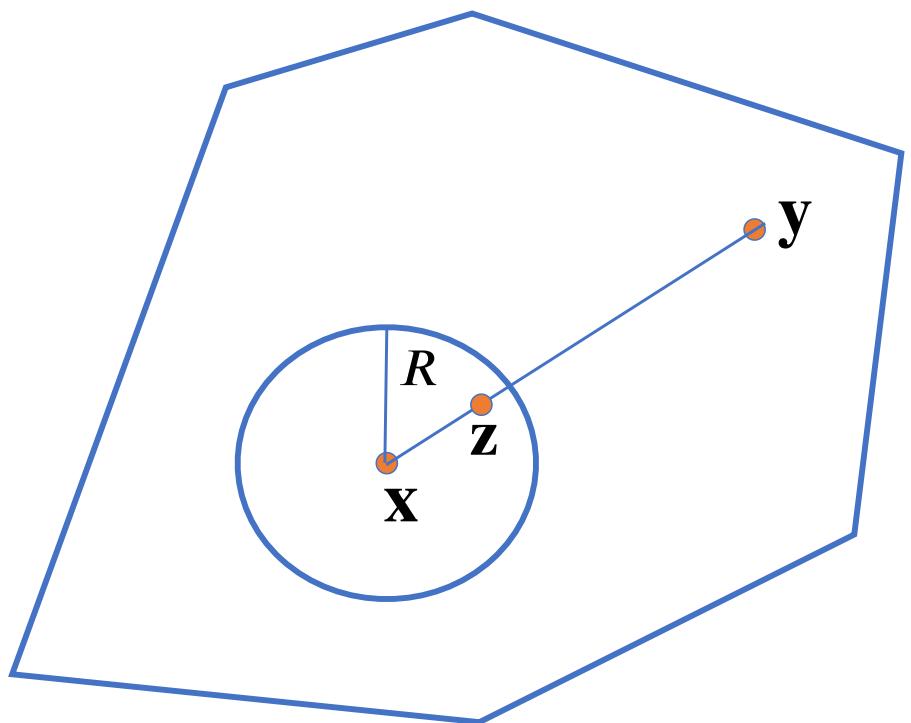
$$f(\mathbf{x}) \leq f(\mathbf{x}')$$

# Local Minima are Global Minima

Let  $\mathbf{z} = \theta\mathbf{y} + (1 - \theta)\mathbf{x}$ ,

(1) By convexity,

$$f(\mathbf{z}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) < f(\mathbf{x})$$



(2) Let  $\theta = \frac{R}{2\|\mathbf{x} - \mathbf{y}\|_2}$ ,

Then

$$\|\mathbf{x} - \mathbf{z}\|_2 = \theta\|\mathbf{x} - \mathbf{y}\|_2 = \frac{R}{2} < R$$

This implies that  $f(\mathbf{x}) \leq f(\mathbf{z})$ .

(1) and (2) cannot hold at the same time, which is a contradiction to  $\mathbf{x}$  being a local minimum.

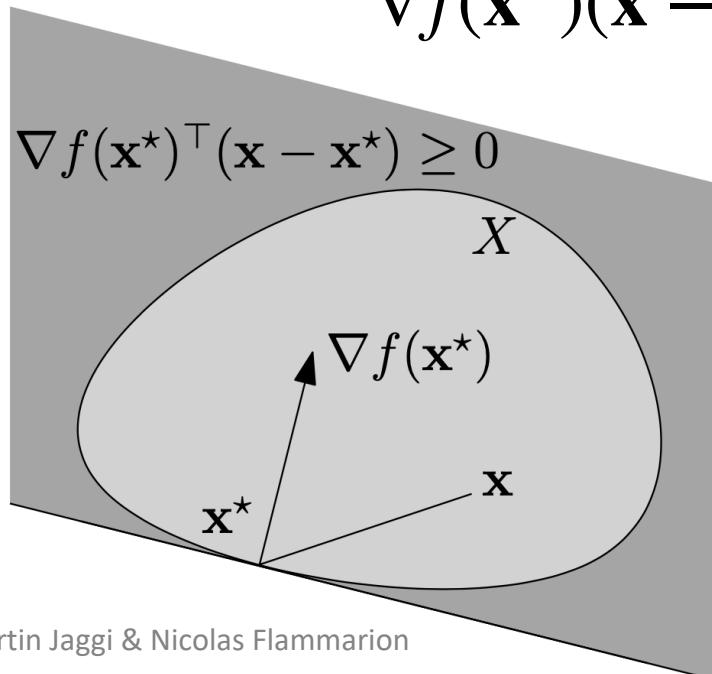
# Optimality criterion for differentiable $f$

Let  $f : \text{dom}(f) \mapsto \mathbb{R}$  be convex and let  $X \in \text{dom}(f)$  be a convex set. A point  $\mathbf{x}$  is a minimizer of  $f$  over  $X$  if

$$f(\mathbf{x}) \leq f(\mathbf{y}) \quad \forall \mathbf{y} \in X$$

- $\mathbf{x}^*$  is optimal if and only if it is feasible and

$$\nabla f(\mathbf{x}^*)(\mathbf{x} - \mathbf{x}^*) \geq 0 \text{ for all feasible } \mathbf{x}$$



In words: all feasible directions from  $\mathbf{x}$  are aligned with gradient  $\nabla f(\mathbf{x})$

# Critical Points are Global Minima

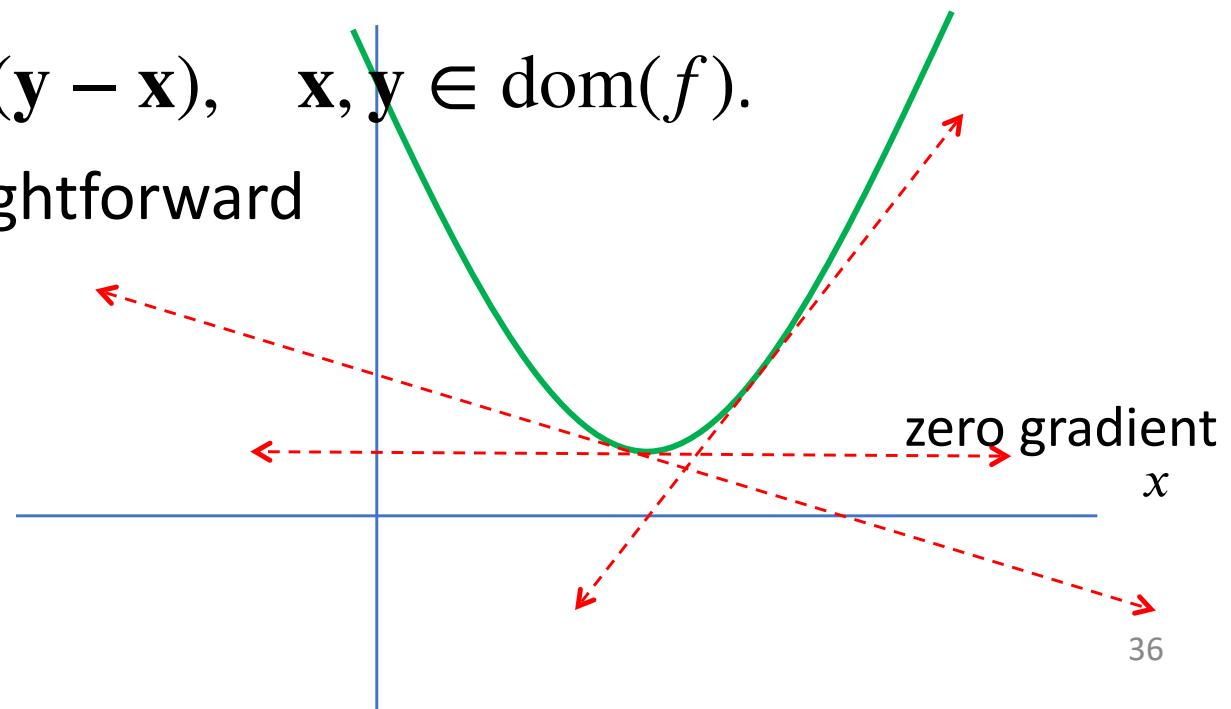
Suppose that  $f$  is convex and differentiable over an open domain  $\text{dom}(f)$ . Let  $\mathbf{x} \in \text{dom}(f)$ . If  $\nabla f(\mathbf{x}) = 0$  (critical point), then  $\mathbf{x}$  is a global minimum.

## Proof

Recall the first-order condition of convexity,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \quad \mathbf{x}, \mathbf{y} \in \text{dom}(f).$$

The conclusion is straightforward



# Questions?

# Content

Part 0: Introduction to introduction ...

Part 1: Theory of Convexity

Part 2: Gradient Descent and convergence

Part 3: Stochastic Gradient Descent, Non-Convex Optimization, Adam, and Machine Learning Applications

Practical work

<https://colab.research.google.com/drive/1r-FTu17Utca7JjblyubSnFEC0n5ph3V6?usp=sharing>

# Gradient method

# Gradient Descent

Consider unconstrained convex optimization

$$\min_{\mathbf{x}} f(\mathbf{x})$$

where  $f$  is convex, differentiable with  $\text{dom}(f) = \mathbb{R}^d$  and has a global minimum  $\mathbf{x}^\star$ .

**Goal:** Find  $\mathbf{x} \in \mathbb{R}^d$  such that

$$f(\mathbf{x}) - f(\mathbf{x}^\star) \leq \epsilon.$$

Note that there can be several global minima  $\mathbf{x}_1^\star \neq \mathbf{x}_2^\star$  with  $f(\mathbf{x}_1^\star) = f(\mathbf{x}_2^\star)$ .

## Gradient Descent Algorithm:

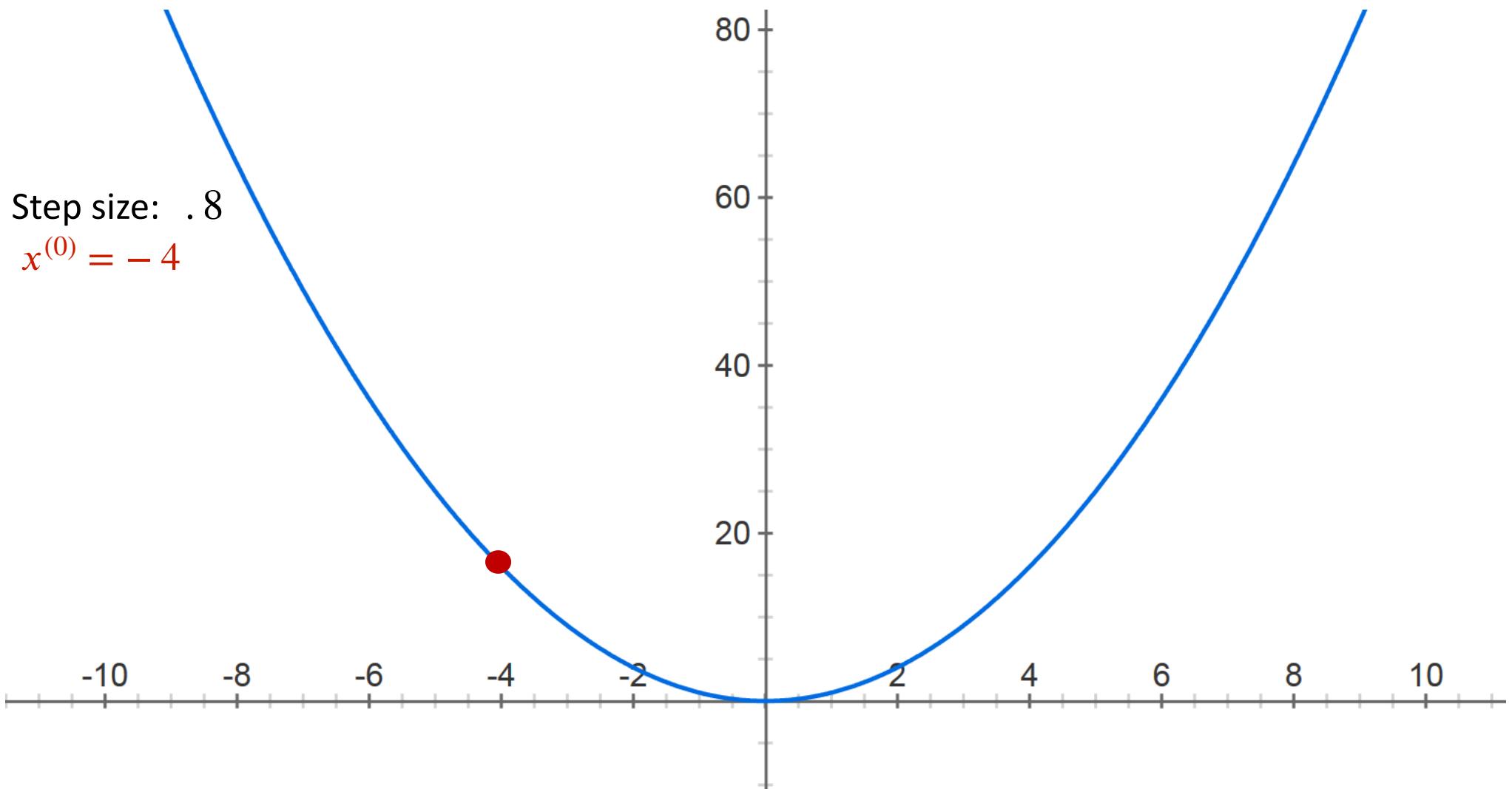
Choose initial point  $\mathbf{x}_0 \in \mathbb{R}^d$ , repeat

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f(\mathbf{x}_t), t = 0, 1, 2, \dots$$

for **step size**  $\gamma$  ( or learning rate)

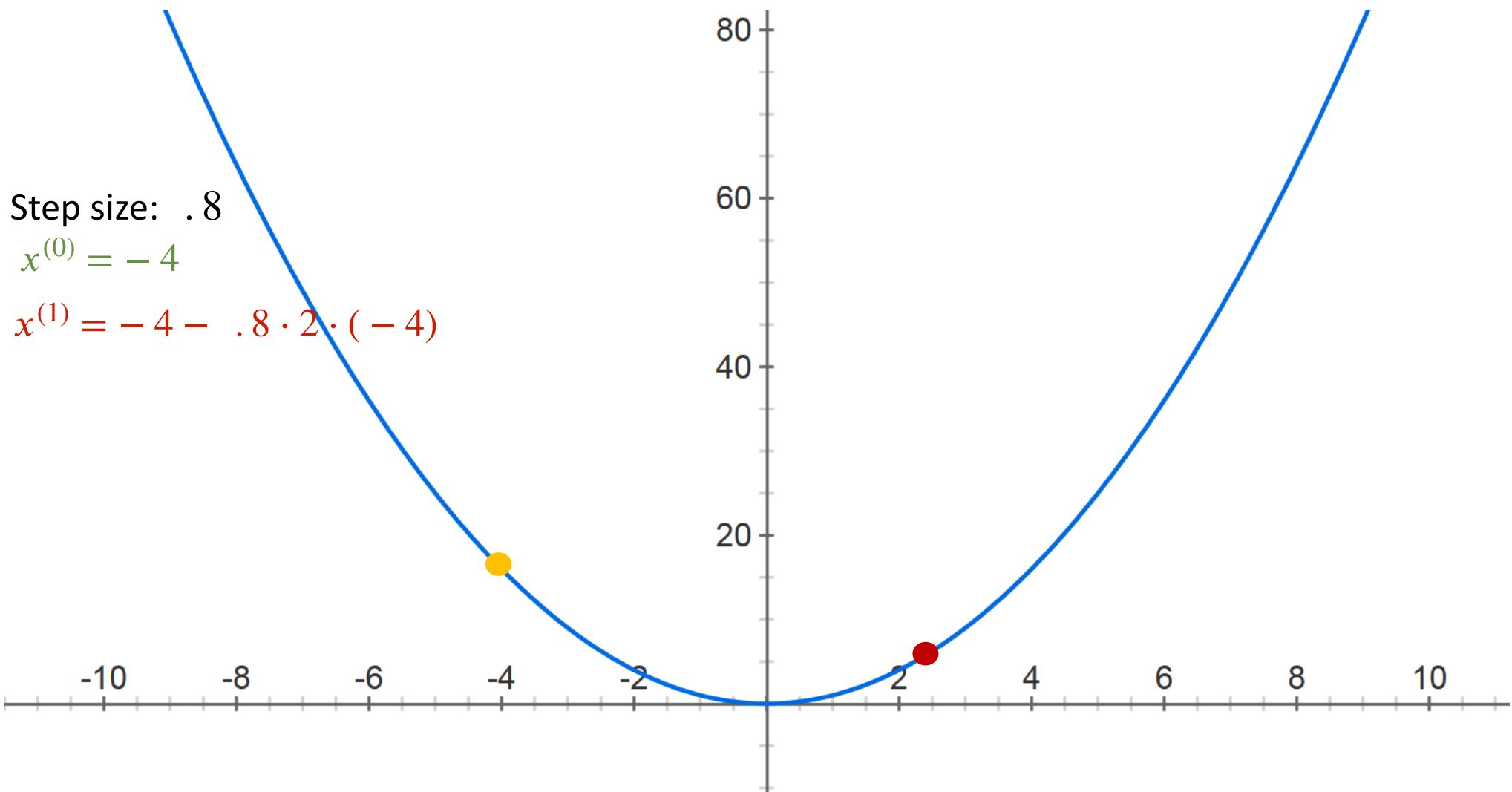
# Gradient Descent

$$f(x) = x^2$$



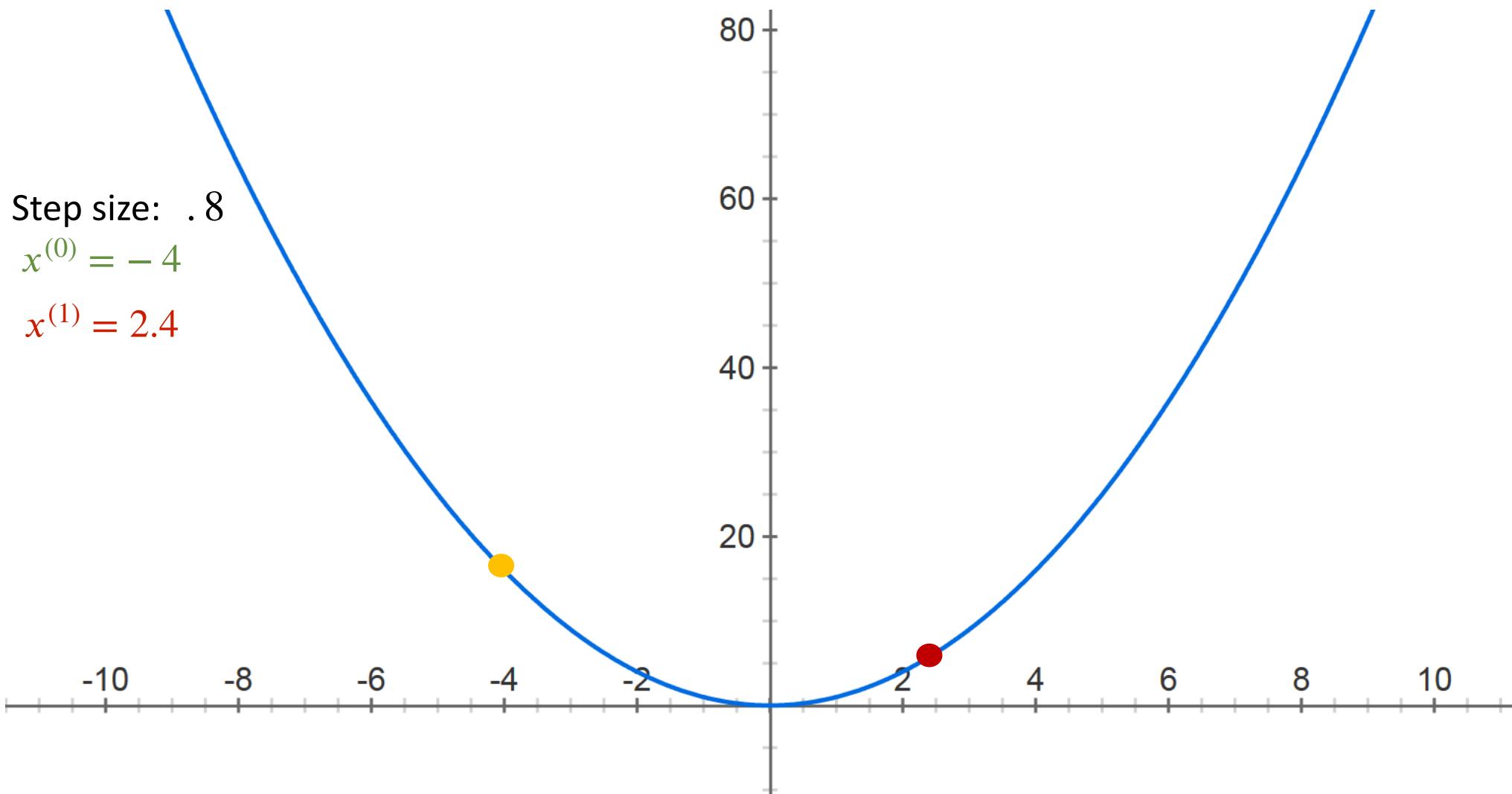
# Gradient Descent

$$f(x) = x^2$$



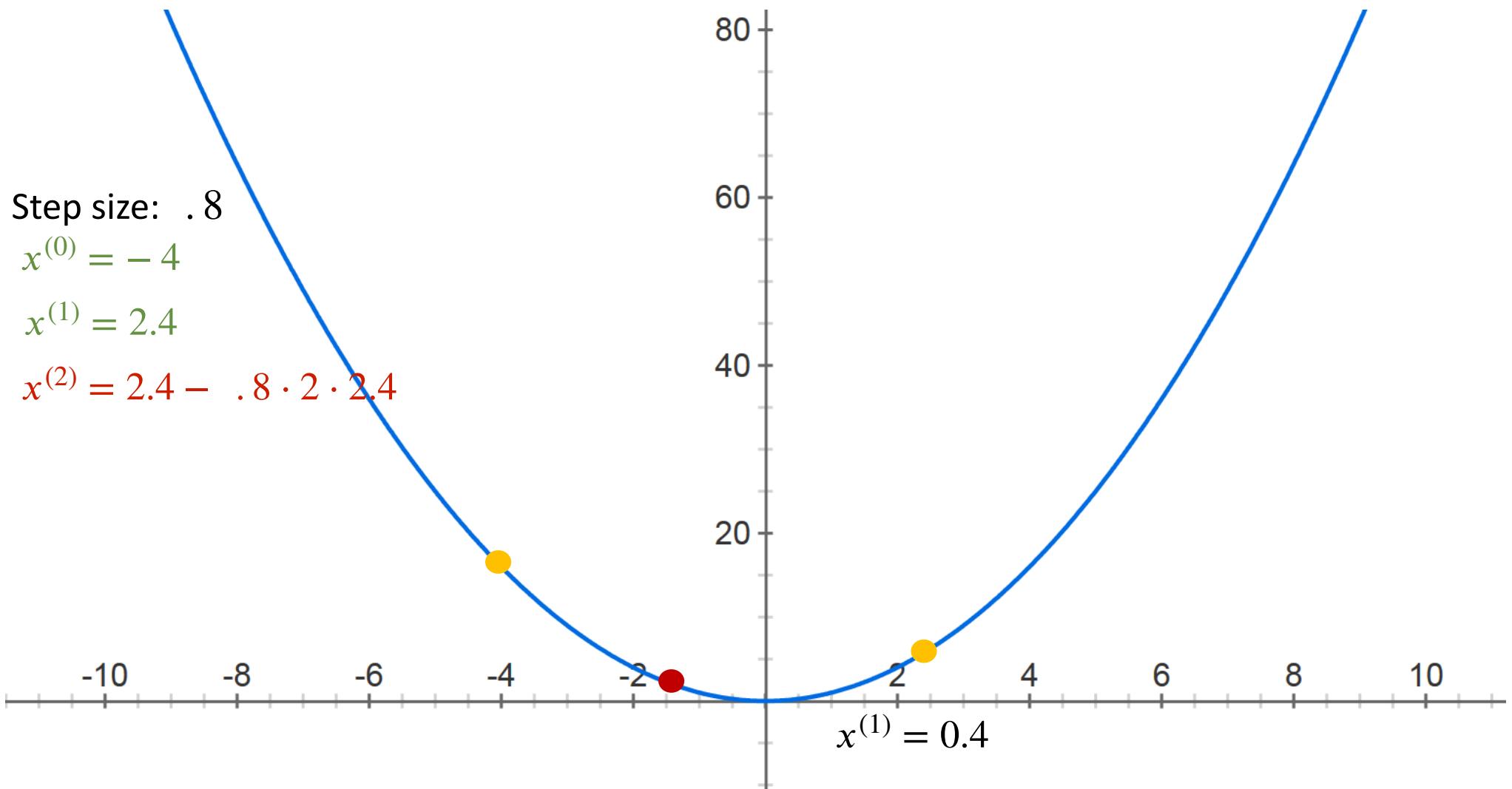
# Gradient Descent

$$f(x) = x^2$$



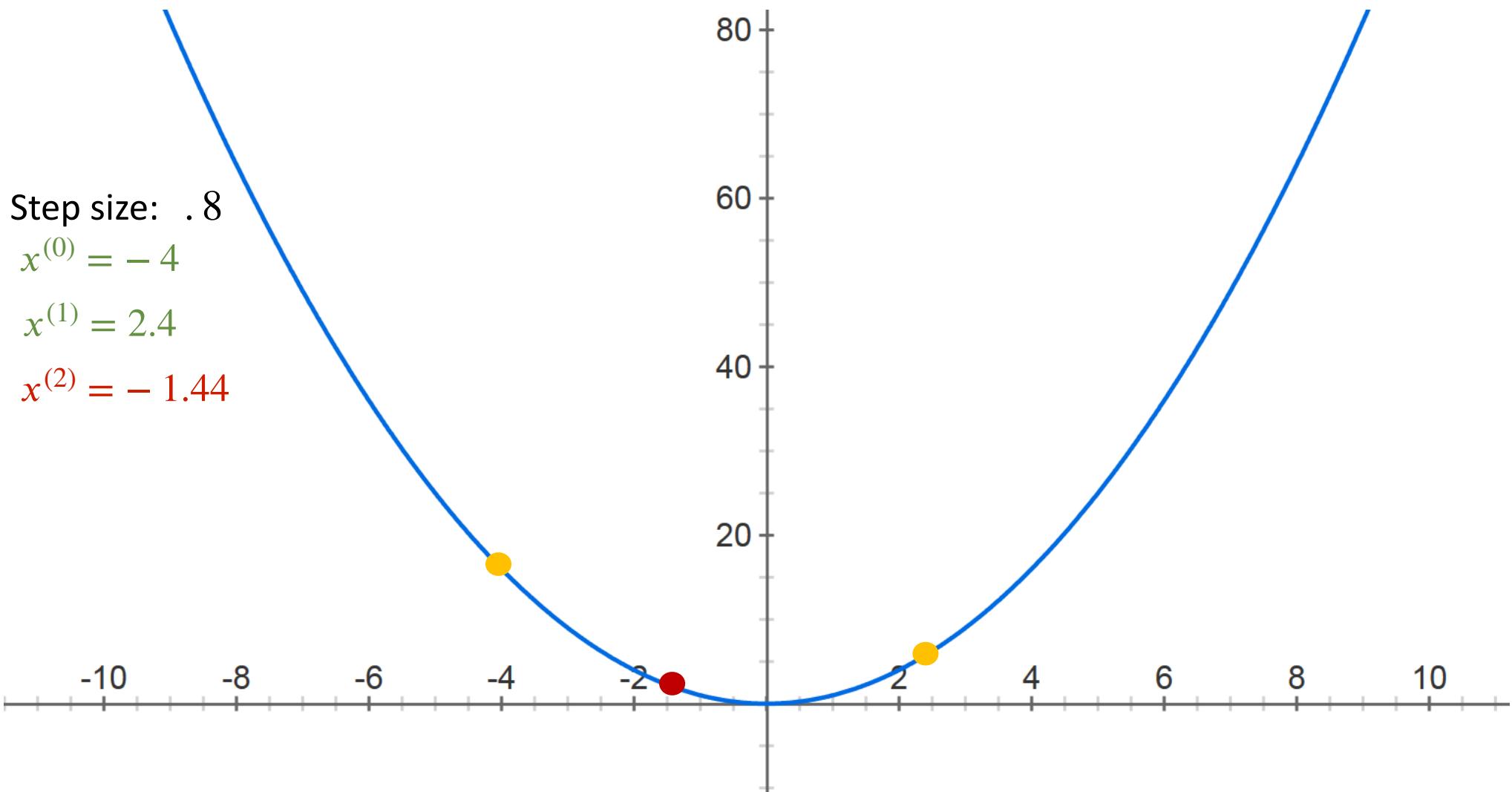
# Gradient Descent

$$f(x) = x^2$$



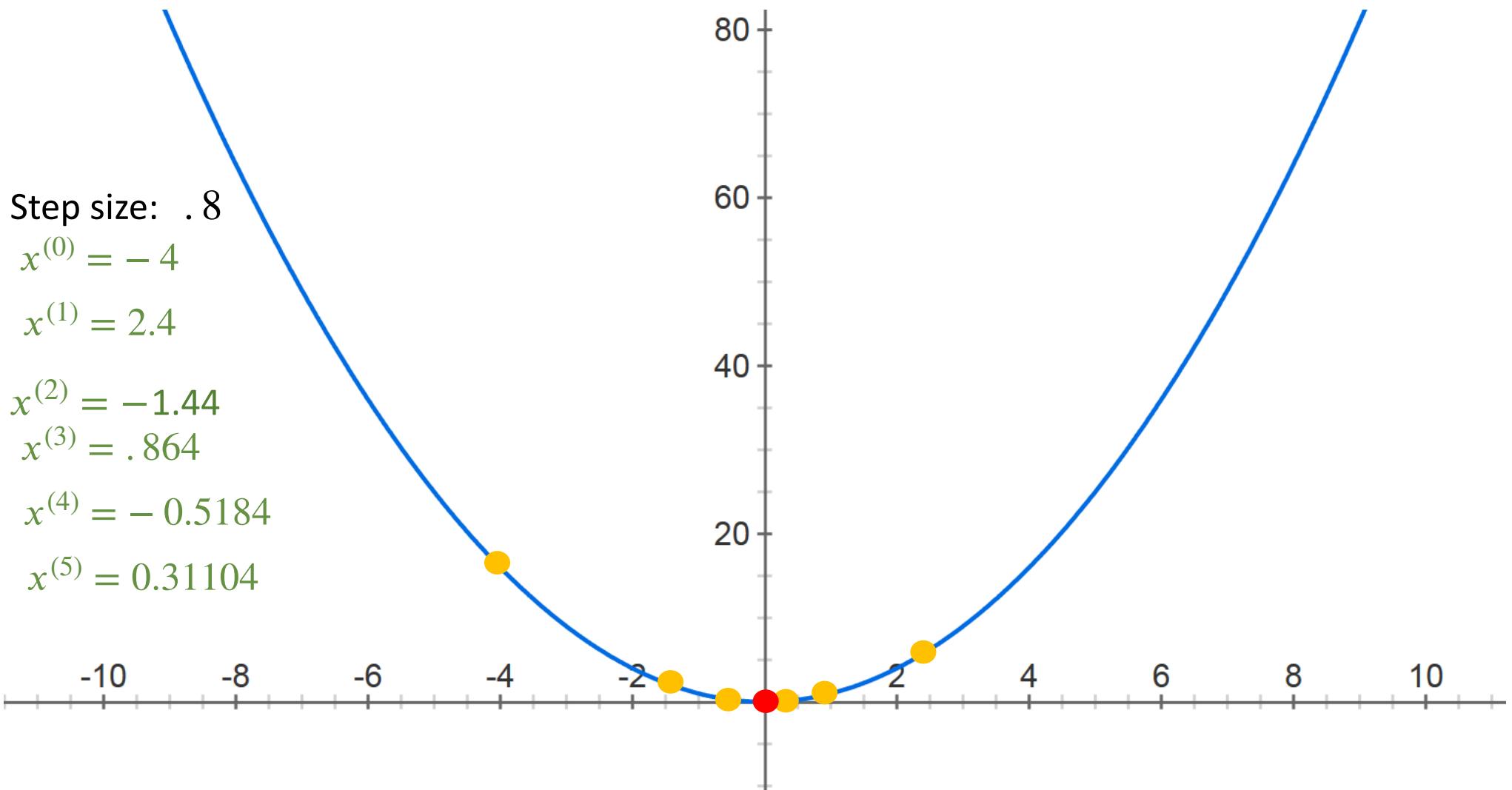
# Gradient Descent

$$f(x) = x^2$$



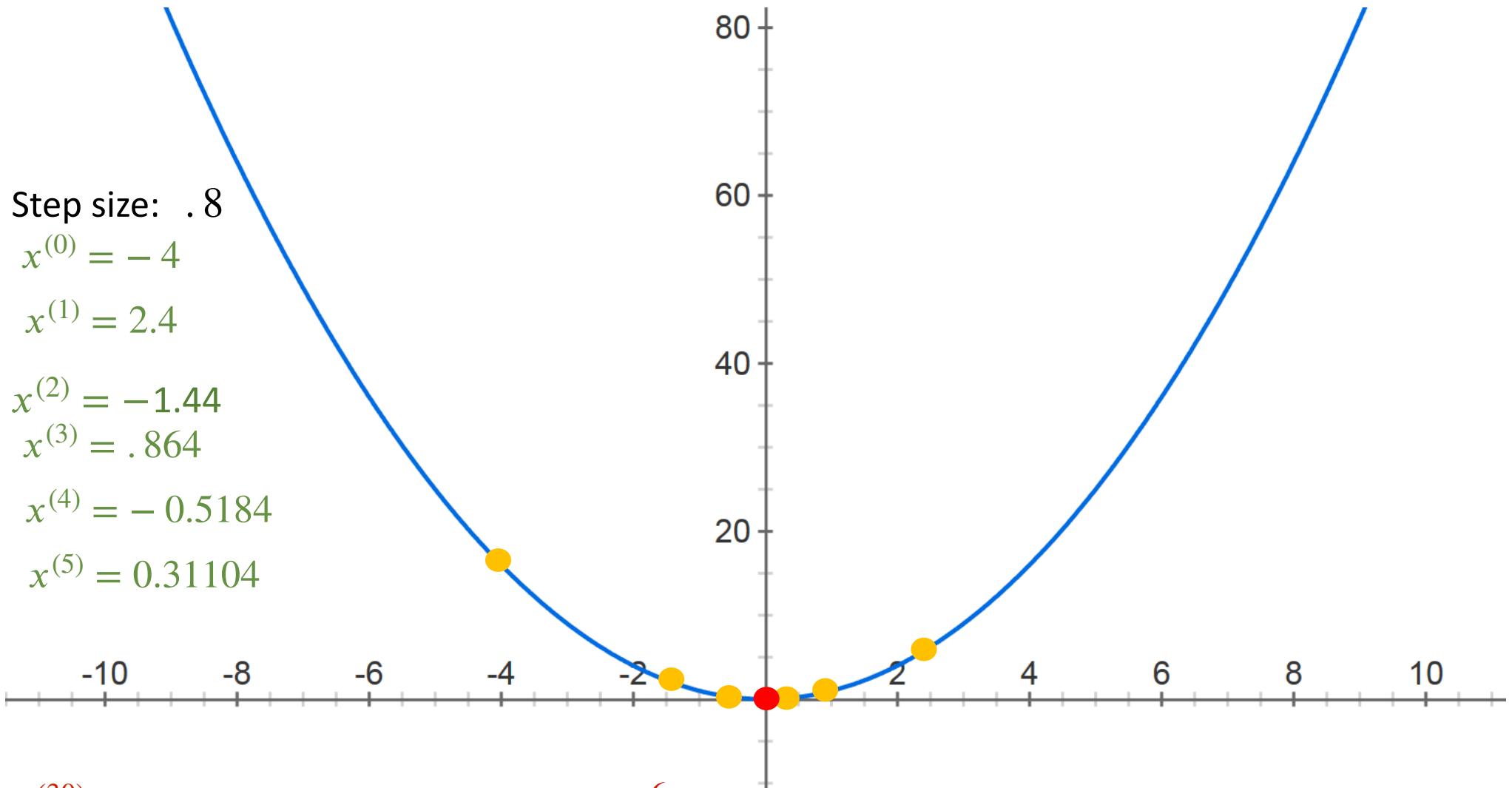
# Gradient Descent

$$f(x) = x^2$$

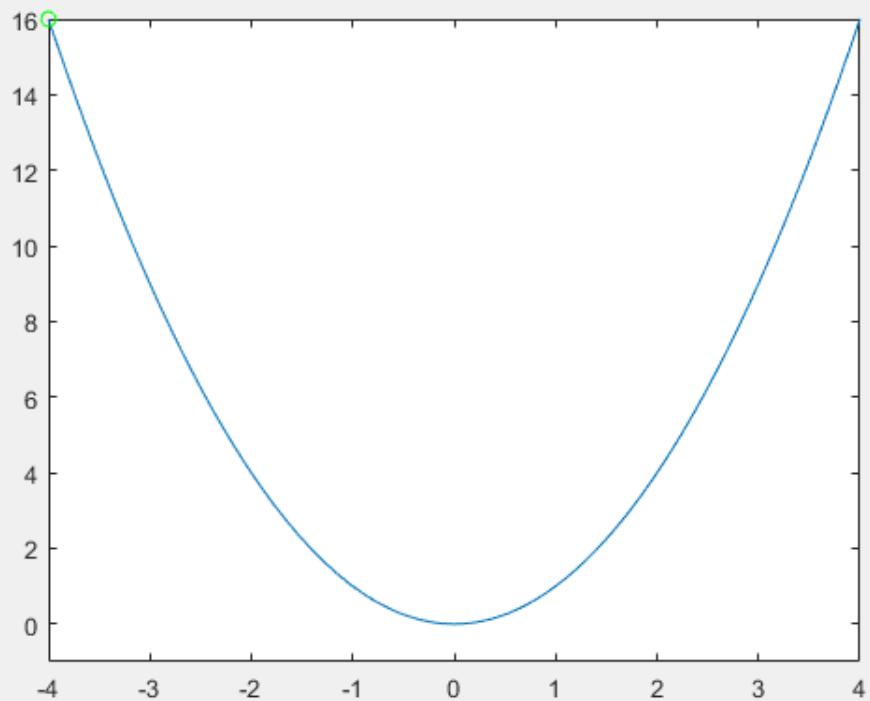


# Gradient Descent

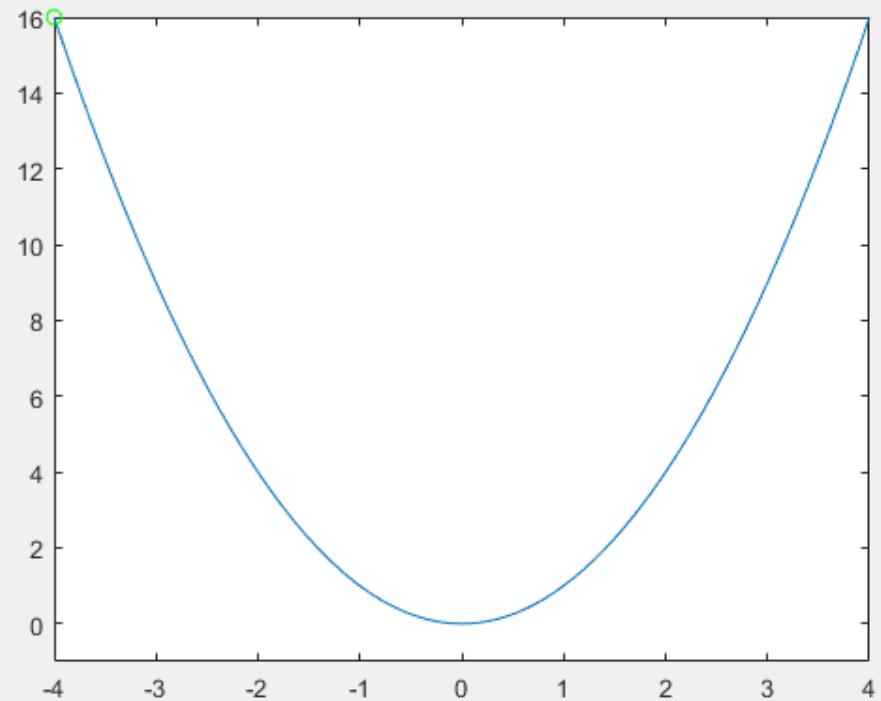
$$f(x) = x^2$$



# Gradient Descent



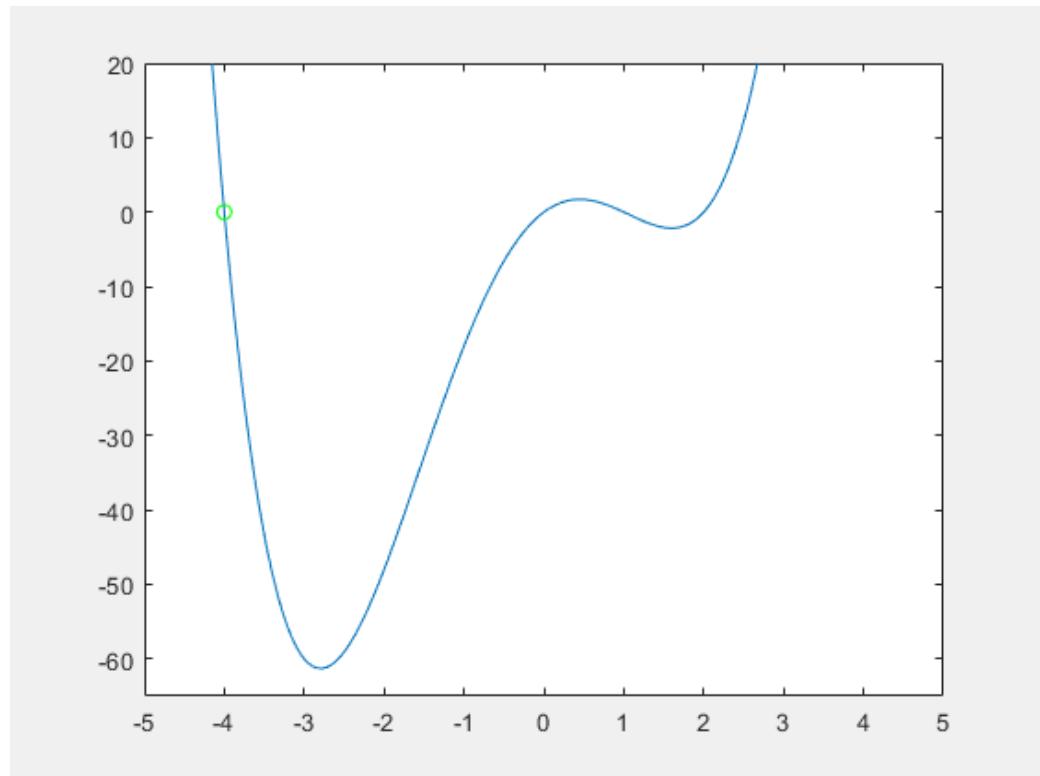
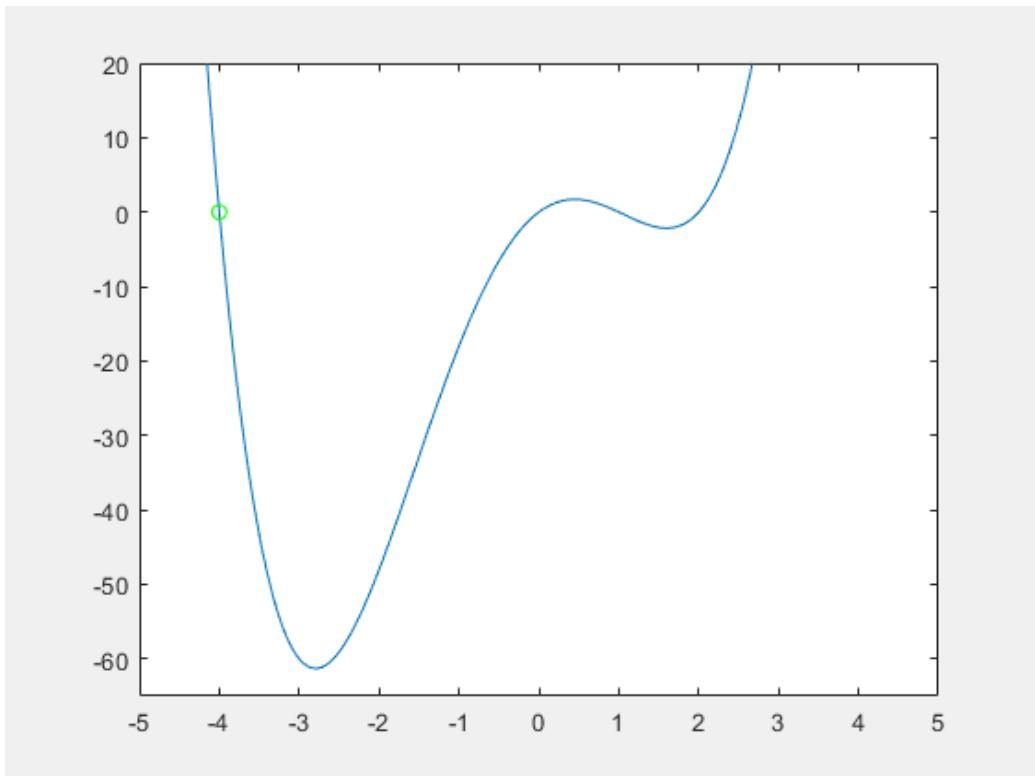
Step size: .2



Step size: .9

# Gradient Descent

Step size matters!



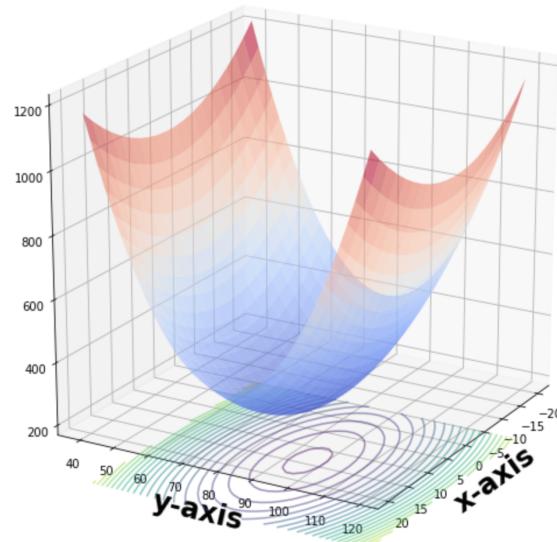
# Example: linear regression(MSE loss)

**Problem:** Given data points  $\mathbf{x}^1, \dots, \mathbf{x}^n \in \mathbb{R}^d$  and  $y^1, \dots, y^n \in \mathbb{R}$ , find the best fit line. Consider a simple linear regression model with  $\hat{y} = \mathbf{w}^T \mathbf{x} + b$ . The formulation can be simplified to  $\hat{y} = \boldsymbol{\theta}^T \tilde{\mathbf{x}}$ , in which  $\boldsymbol{\theta} = [\mathbf{w}^T, b]^T$ ,  $\tilde{\mathbf{x}} = [\mathbf{x}^T, 1]^T$ .

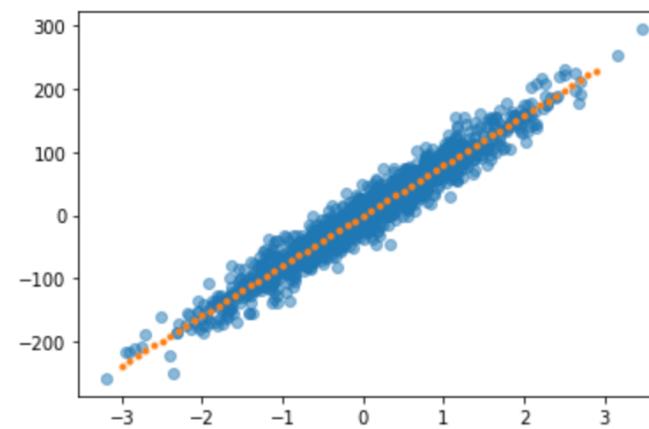
**Objective:**

$$\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (\boldsymbol{\theta}^T \tilde{\mathbf{x}}_i - y_i)^2$$

**Loss Surface:**



**Data Distribution:**



# GD result with different step size!

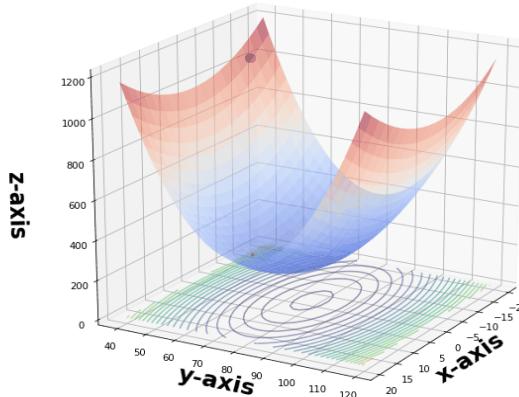
$$\theta_{t+1} := \theta_t - \gamma \sum_{i=1}^n \tilde{\mathbf{x}}_i^T (\theta^T \tilde{\mathbf{x}}_i - y_i), t = 0, 1, 2, \dots$$

```
pred_y = np.dot(x, theta)
error = pred_y - y
gradient = x.T.dot(error)/sample_number
theta = theta - step_size * gradient
```

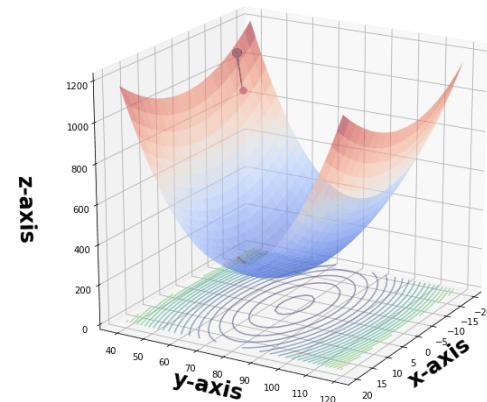
step\_size = 0.5

step\_size = 0.1

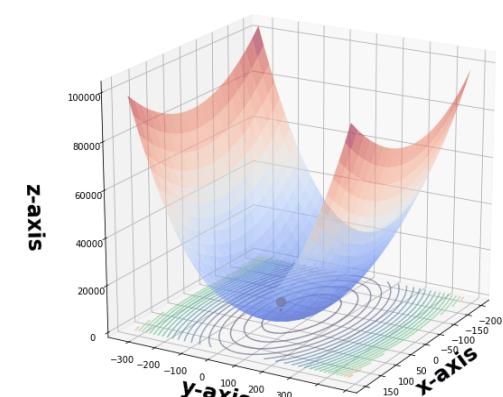
step\_size = 2.3



Converge



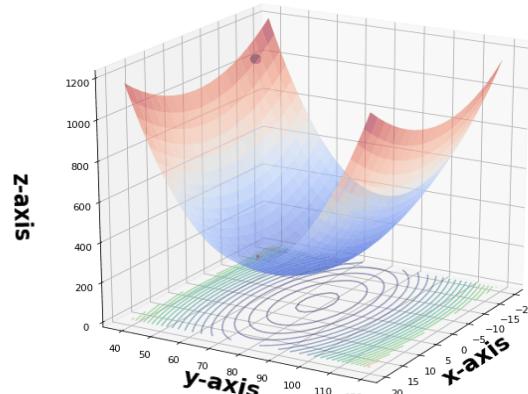
Slow Converge



Diverge

# GD result with different step size!

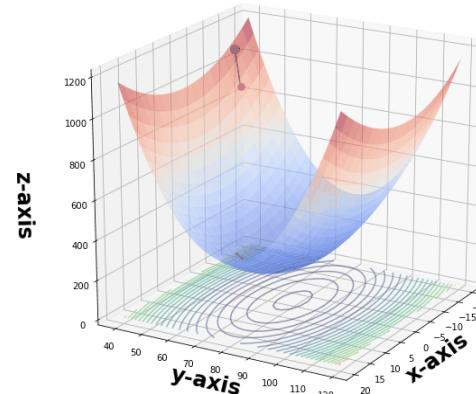
step\_size = 0.5



```
## Gradient Descent ##  
converge to the:  
[-0.34,79.17]  
iterations :  
14 steps
```

**Converge**

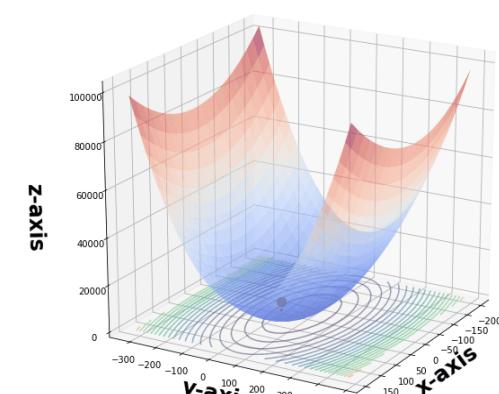
step\_size = 0.1



```
## Gradient Descent ##  
converge to the:  
[-0.35,79.15]  
iterations :  
69 steps
```

**Slow Converge**

step\_size = 2.3



```
## Gradient Descent ##  
Diverge to the:  
[7.55e64,7.31e64]  
iterations :  
200+ steps
```

**Diverge** 52

# Convergence of gradient descent

# Gradient Descent

**Assumption:**  $f : \mathbb{R}^d \mapsto \mathbb{R}$  is convex, differentiable and has a global minimum  $\mathbf{x}^\star$

**Goal:** Find  $\mathbf{x} \in \mathbb{R}^d$  such that

$$f(\mathbf{x}) - f(\mathbf{x}^\star) \leq \epsilon.$$

Note that there can be several global minima  $\mathbf{x}_1^\star \neq \mathbf{x}_2^\star$  with  $f(\mathbf{x}_1^\star) = f(\mathbf{x}_2^\star)$ .

**Average error:**

$$\text{Averaged error} = \frac{1}{T} \sum_{t=0}^{T-1} \left( f(\mathbf{x}_t) - f(\mathbf{x}^\star) \right)$$

How can we bound the average error?

# Warmup

- First let's look at  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$ , abbreviating  $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ .
- With  $\mathbf{x} = \mathbf{x}_t, \mathbf{y} = \mathbf{x}^*$

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$$

- giving

$$\sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$$

Then?

first-order characterization of convexity  
 $f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}), \forall \mathbf{x}, \mathbf{y}$

# Warmup

As  $\mathbf{g}_t := \nabla f(\mathbf{x}_t)$ , from gradient descent,  $\mathbf{g}_t = (\mathbf{x}_t - \mathbf{x}_{t+1})/\gamma$ .

$$\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) = \frac{1}{\gamma} (\mathbf{x}_t - \mathbf{x}_{t+1})^\top (\mathbf{x}_t - \mathbf{x}^*).$$

- Rewrite according to  $2\mathbf{v}^\top \mathbf{w} = \|\mathbf{v}\|^2 + \|\mathbf{w}\|^2 - \|\mathbf{v} - \mathbf{w}\|^2$

$$\begin{aligned} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) &= \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 + \|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \\ &= \frac{\gamma}{2} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2) \end{aligned}$$

- Sum this up over the first  $T$  iterations:

$$\begin{aligned} \sum_{t=0}^{T-1} \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*) &= \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \underbrace{\sum_{t=0}^{T-1} \frac{1}{2\gamma} (\|\mathbf{x}_t - \mathbf{x}^*\|^2 - \|\mathbf{x}_{t+1} - \mathbf{x}^*\|^2)}_{\text{Telescoping sum}} \\ &= \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2) \end{aligned}$$

# Observation 1

This gives

$$\begin{aligned}\sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2 - \|\mathbf{x}_T - \mathbf{x}^*\|^2) \\ &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2)\end{aligned}$$

an upper bound for the **total error**  $f(\mathbf{x}_t) - f(\mathbf{x}^*)$  over the steps

- last iterate is not necessarily the best one.
- step size is crucial

# Lipschitz convex functions: $O(1/\epsilon^2)$ steps

**Theorem** Let  $f : \mathbb{R}^d \mapsto \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^\star$ ; furthermore, suppose that  $\|\mathbf{x}_0 - \mathbf{x}^\star\| \leq R$  and  $\|\nabla f(\mathbf{x})\| \leq B$  for all  $\mathbf{x}$ . Choosing the stepsize

$$\gamma := \frac{R}{B\sqrt{T}}$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^\star) \leq \frac{RB}{\sqrt{T}}$$

Assume that all gradients of  $f$  are bounded in norm.

$f$  is  $B$ -Lipschitz with

$$\|f(\mathbf{x}) - f(\mathbf{y})\| \leq B\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

$\Leftrightarrow$

$$\|\nabla f(\mathbf{x})\| \leq B \quad \forall \mathbf{x} \in X$$

# Lipschitz convex functions: $O(1/\epsilon^2)$ steps

Proof

- plug  $\|\mathbf{x}_0 - \mathbf{x}^\star\| \leq R$  and  $\|\nabla f(\mathbf{x})\| \leq B$  into observation 1:

$$\begin{aligned}\sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^\star) &\leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\mathbf{g}_t\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^\star\|^2) \\ &\leq \frac{\gamma}{2} B^2 T + \frac{1}{2\gamma} R^2\end{aligned}$$

- choose  $\gamma$  such that

$$q(\gamma) = \frac{\gamma}{2} B^2 T + \frac{1}{2\gamma} R^2$$

is minimized.

- giving  $\gamma = \frac{R}{B\sqrt{T}}$ , and  $q(R/B\sqrt{T}) = RB\sqrt{T}$ .

- the result follows by dividing  $q(\gamma)$  by  $T$ .

# Lipschitz convex functions: $O(1/\epsilon^2)$ steps

$$T \geq \frac{R^2 B^2}{\epsilon^2} \quad \Rightarrow \quad \text{average error} \leq \frac{RB}{\sqrt{T}} \leq \epsilon$$

Advantages:

- dimension-independent (no  $d$  in the bound)!
- holds for both average, or best iterate

# Smooth functions

# Smooth functions

## Definition

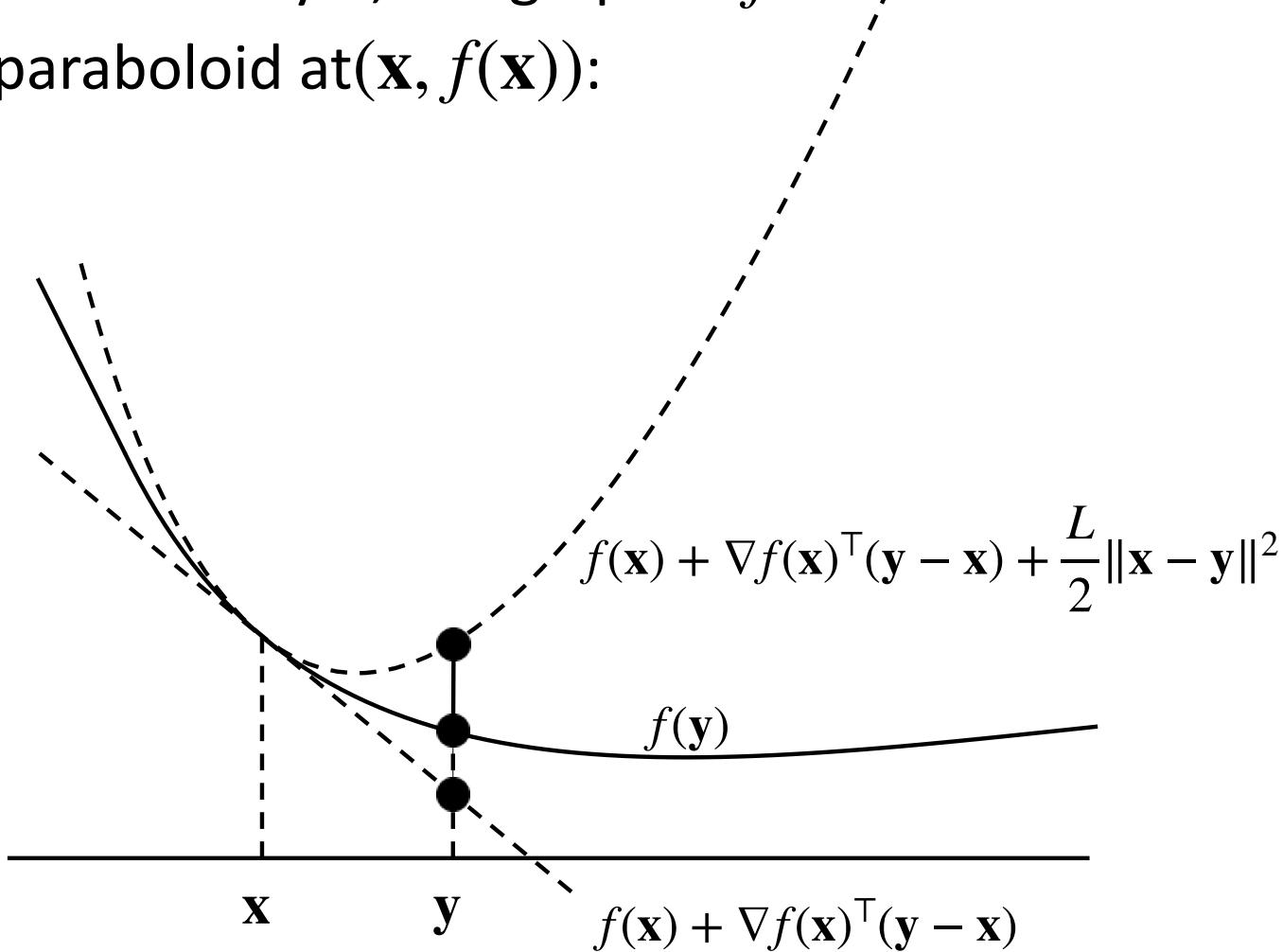
Let  $f : \text{dom}(f) \mapsto \mathbb{R}$  be differentiable,  $X \subseteq \text{dom}(f)$ ,  $L > 0$ .  $f$  is called smooth (with parameter  $L$ ) over  $X$  if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

- Definition does not require convexity (useful in non-convex world)

# Smooth functions

**Smoothness:** For any  $\mathbf{x}$ , the graph of  $f$  is below a not too steep tangent paraboloid at  $(\mathbf{x}, f(\mathbf{x}))$ :



# Observation 2: Sufficient decrease

**Lemma:** Let  $f : \mathbb{R}^d \mapsto \mathbb{R}$  be differentiable and smooth with parameter  $L$ . With stepsize

$$\gamma := \frac{1}{L}$$

gradient descent satisfies

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2, \quad t \geq 0$$

**Proof:**

Use smoothness and definition of gradient descent ( $\mathbf{x}_{t+1} - \mathbf{x}_t = -\nabla f(\mathbf{x})/L$ ):

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{L} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \\ &= f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \end{aligned}$$

# Smooth convex functions: $O(1/\epsilon)$ steps

**Theorem** Let  $f : \mathbb{R}^d \mapsto \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^\star$ ; furthermore, suppose that  $f$  is smooth with parameter  $L$ . Choosing the stepsize

$$\gamma := \frac{1}{L}$$

gradient descent yields

$$f(\mathbf{x}_T) - f(\mathbf{x}^\star) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^\star\|^2.$$

# Smooth convex functions: $O(1/\epsilon)$ steps

Proof.

Recall: observation 1

$$(1) \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{\gamma}{2} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{1}{2\gamma} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2)$$

The, based on observation 2 (sufficient decrease),

$$(2) \sum_{t=0}^{T-1} \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2 \leq \underbrace{\sum_{t=0}^{T-1} (f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}))}_{\text{Telescoping sum}} = f(\mathbf{x}_0) - f(\mathbf{x}_T)$$

Putting (1) and (2) together with  $\gamma = 1/L$ ,

$$\begin{aligned} \sum_{t=0}^{T-1} f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{1}{2L} \sum_{t=0}^{T-1} \|\nabla f(\mathbf{x}_t)\|^2 + \frac{L}{2} (\|\mathbf{x}_0 - \mathbf{x}^*\|^2) \\ &\leq f(\mathbf{x}_0) - f(\mathbf{x}_T) + \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \\ &\quad \Downarrow \\ \sum_{t=1}^T f(\mathbf{x}_t) - f(\mathbf{x}^*) &\leq \frac{L}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2 \end{aligned}$$

As last iterate is the best (sufficient decrease!)

$$f(\mathbf{x}_T) - f(\mathbf{x}^*) \leq \frac{1}{T} \sum_{t=0}^T f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \frac{L}{2T} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

The result follows!

Observation 2:

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2$$

# Smooth convex functions: $O(1/\epsilon)$ steps

Let  $R^2 := \| \mathbf{x}_0 - \mathbf{x}^\star \|^2$ , we have

$$T \geq \frac{R^2 L}{2\epsilon} \Rightarrow \text{error} \leq \frac{L}{2T} R^2 \leq \epsilon$$

Recall that for Lipschitz-convex

$$T \geq \frac{R^2 B^2}{\epsilon^2} \Rightarrow \text{average error} \leq \frac{RB}{\sqrt{T}} \leq \epsilon$$

To achieve an approximation error  $\epsilon \leq 0.01$ , we need

- $50 \cdot R^2 L$  iterations for smooth convex function
- $10,000 \cdot R^2 B^2$  iterations for Lipschitz convex function

# Summary

- Lipschitz convex functions:  $O(1/\epsilon^2)$  steps
- Smooth convex functions:  $O(1/\epsilon)$  steps

Next

- Stochastic gradient descent
- Non-convex optimization

# Questions?

# Content

Part 0: Introduction to introduction ...

Part 1: Theory of Convexity

Part 2: Gradient Descent and convergence

Part 3: Stochastic Gradient Descent, Non-Convex Optimization, Adam, and Machine Learning Applications

Practical work

<https://colab.research.google.com/drive/1r-FTu17Utca7JjblyubSnFEC0n5ph3V6?usp=sharing>

# Stochastic gradient descent

# Stochastic gradient descent

Consider minimizing an average of functions

$$\min f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x})$$

## Stochastic GD:

Choose initial point  $\mathbf{x}_0 \in \mathbb{R}^d$ , repeat

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \nabla f_i(\mathbf{x}_t), t = 0, 1, 2, \dots$$

for **step size**  $\gamma \geq 0$ , uniformly sampled index  $i \in [n]$ .

Gradient descent algorithm:

choose initial point  $\mathbf{x}_0 \in \mathbb{R}^d$ , repeat

$$\mathbf{x}_{t+1} := \mathbf{x}_t - \gamma \cdot \frac{1}{n} \sum_{i=1}^n \nabla f_i(\mathbf{x}_t), t = 0, 1, 2, \dots$$

for **step size**  $\gamma$  ( or learning rate)

# Unbiasedness

The vector  $\mathbf{g}_t := \nabla f_i(\mathbf{x}_t)$  is called a **stochastic gradient**

Can't use convexity

$$f(\mathbf{x}_t) - f(\mathbf{x}^*) \leq \mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)$$

on top of the **Observation 1**, as this may hold or not hold, depending on how the stochastic gradient  $\mathbf{g}_t$  turns out.

- We will show (and exploit): the inequality holds in expectation.

For this, we use that by definition,  $\mathbf{g}_t$  is an unbiased estimate of  $\nabla f(\mathbf{x}_t)$ :

$$\mathbb{E}_{i \in [n]} [\mathbf{g}_t \mid \mathbf{x}_t = \mathbf{x}] = \frac{1}{n} \sum^n \nabla f_i(\mathbf{x}) = \nabla f(\mathbf{x}), \quad \mathbf{x} \in \mathbb{R}^d$$

- Main appeal of SGD:

- Iteration cost is independent of  $n$  (number of functions)
- Can also be a big savings in terms of memory usage

- Further, the below equation holds

$$\mathbb{E}_{i \in [n], \mathbf{x}_t} [\mathbf{g}_t^\top (\mathbf{x}_t - \mathbf{x}^*)] = \mathbb{E}_{\mathbf{x}_t} [\nabla f(\mathbf{x}_t)^\top (\mathbf{x}_t - \mathbf{x}^*)] \geq \mathbb{E}_{\mathbf{x}_t} [f(\mathbf{x}_t) - f(\mathbf{x}^*)]$$

# Bounded stochastic gradients: $O(1/\epsilon^2)$ steps

**Theorem** Let  $f : \mathbb{R}^d \mapsto \mathbb{R}$  be convex and differentiable with a global minimum  $\mathbf{x}^\star$ ; furthermore, suppose that  $\|\mathbf{x}_0 - \mathbf{x}^\star\| \leq R$  and  $\mathbb{E}[\|\mathbf{g}_t\|^2] \leq B^2$  for all  $\mathbf{x}$ . Choosing the stepsize

$$\gamma := \frac{R}{B\sqrt{T}}$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[f(\mathbf{x}_t)] - f(\mathbf{x}^\star) \leq \frac{RB}{\sqrt{T}}.$$

Same procedure as the proof of gradient descent, except

- we assume bounded stochastic gradients **in expectation**;
- error bound holds **in expectation**.

# Mini-batch SGD

Instead of using a single element  $f_i$ , use an average of several of them:

$$\tilde{\mathbf{g}}_t := \frac{1}{m} \sum_{j=1}^m \mathbf{g}_t^j$$

Extreme cases:

$m = 1 \Leftrightarrow$  SGD as originally defined

$m = n \Leftrightarrow$  full gradient descent

Benefit: Gradient computation can be naively parallelized

# Mini-batch SGD

Variance Intuition: Taking an average of many independent random variables reduces the variance. So for larger size of the mini-batch  $m$ ,  $\tilde{\mathbf{g}}_t$  will be closer to the true gradient, in expectation:

$$\begin{aligned}\mathbb{E}[\|\tilde{\mathbf{g}}_t - \nabla f(\mathbf{x}_t)\|^2] &= \mathbb{E}\left[\left\|\frac{1}{m} \sum_{j=1}^m \mathbf{g}_t^j - \nabla f(\mathbf{x}_t)\right\|^2\right] \\ &= \frac{1}{m} \mathbb{E}[\|\mathbf{g}_t^1 - \nabla f(\mathbf{x}_t)\|^2] \\ \text{var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2 \quad \rightarrow \quad &= \frac{1}{m} \mathbb{E}[\|\mathbf{g}_t^1\|^2] - \frac{1}{m} \|\nabla f(\mathbf{x}_t)\|^2 \leq \frac{B^2}{m}\end{aligned}$$

Using a modification of the SGD analysis, can use this quantity to relate convergence rate to the rate of full gradient descent.

# Convergence rate

Recall:

- Lipschitz convex functions:  $O(1/\epsilon^2)$  steps
- Smooth convex functions:  $O(1/\epsilon)$  steps

What about SGD?

- Lipschitz convex functions:  $O(1/\epsilon^2)$  steps
- Unfortunately this **does not improve** when we further assume  $f$  has Lipschitz gradient

# Non-convex optimization

# Nonconvex problem

The problems solved in practice, especially in machine learning/statistics, are mostly convex.

- Linear regression, logistic regression;
- Kernel methods;
- Linear programming, semi-definite programming, SOS (Sum Of Squares programming);

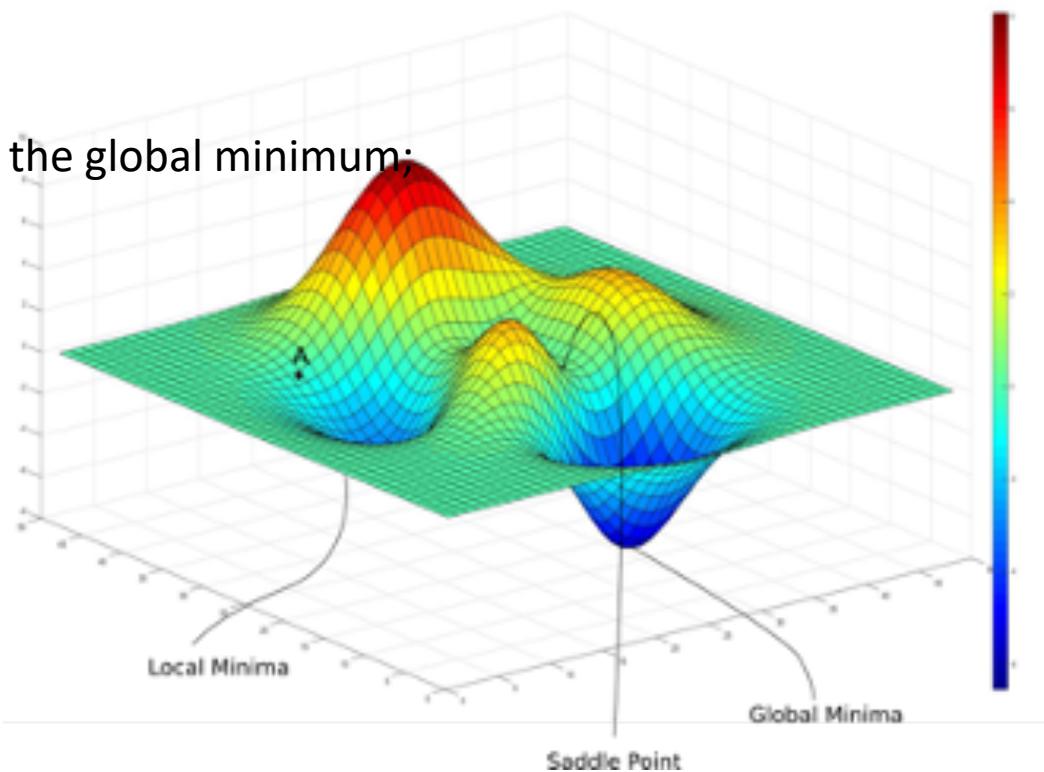
But now, they are mostly non-convex, mainly for one reason:

- Deep learning / Neural networks.

GD may get stuck in a local minimum and miss the global minimum;

- but GD still works well in practice. Why?

Let's look into theoretical explanations!



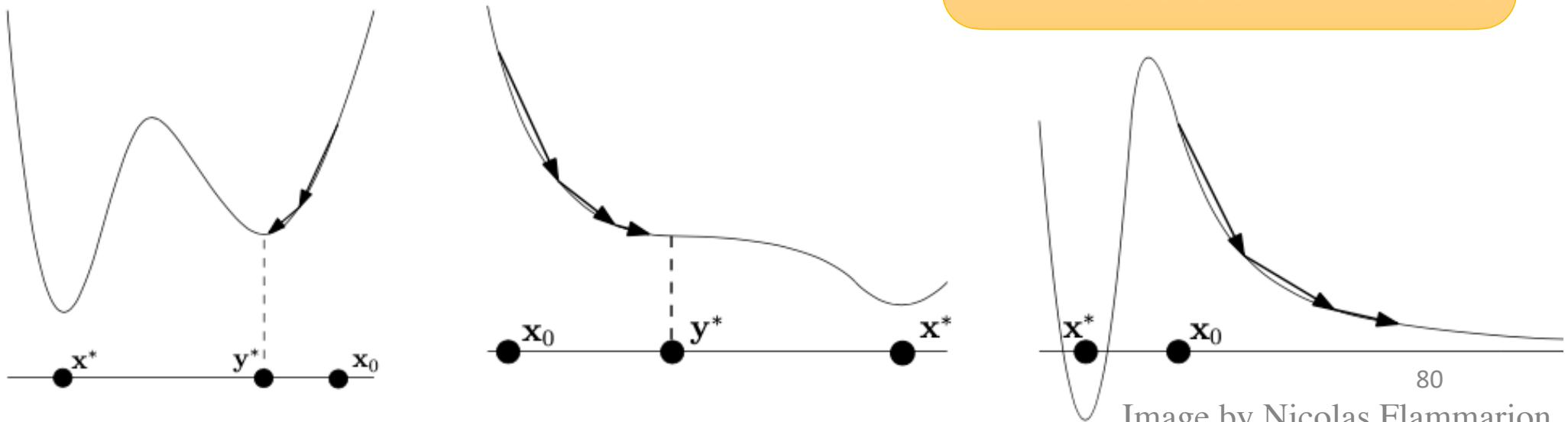
# Gradient descent on non-convex problems

Will prove  $\| \nabla f(\mathbf{x}_t) \|_2 \rightarrow 0$  for  $t \rightarrow \infty$  at the same rate as  $f(\mathbf{x}_t) - f(\mathbf{x}^*) \rightarrow 0$  in the convex case.

$f(\mathbf{x}_t) - f(\mathbf{x}^*)$  itself may not converge to 0 in the nonconvex case.

It may or may not mean that we converge to a critical point ( $\nabla f(\mathbf{y}^*) = 0$ ).

No convexity needed!!

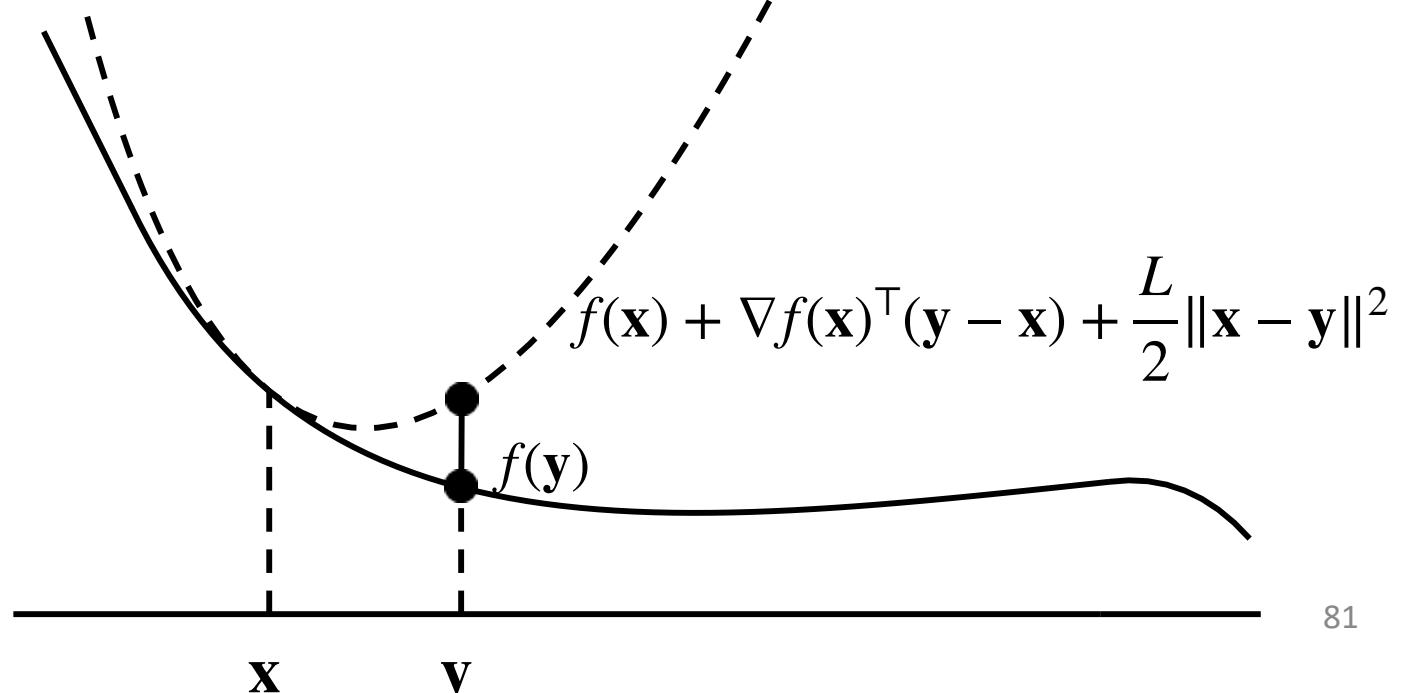


## Recall: Smooth (but not necessarily convex) functions

Let  $f : \text{dom}(f) \mapsto \mathbb{R}$  be differentiable,  $X \subseteq \text{dom}(f)$ ,  $L > 0$ .  $f$  is called smooth (with parameter  $L$ ) over  $X$  if

$$f(\mathbf{y}) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^\top (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in X$$

Definition does not require convexity.



# Smooth gradient descent

**Theorem:** Let  $f: \mathbb{R}^d \mapsto \mathbb{R}$  be differentiable with a global minimum  $x^\star$ ; furthermore, suppose that  $f$  is smooth with parameter  $L$ . Choosing stepsize

$$\gamma := \frac{1}{L}$$

gradient descent yields

$$\frac{1}{T} \sum_{t=0}^{T-1} \| \nabla f(\mathbf{x}_t) \|^2 \leq \frac{2L}{T} (f(\mathbf{x}_0) - f(\mathbf{x}^\star)), \quad T > 0.$$

# Smooth gradient descent

## Proof

Based on observation 2 (sufficient decrease)

$$f(\mathbf{x}_{t+1}) \leq f(\mathbf{x}_t) - \frac{1}{2L} \|\nabla f(\mathbf{x}_t)\|^2, \quad t \geq 0$$

Rewriting:

$$\|\nabla f(\mathbf{x}_t)\|^2 \leq 2L(f(\mathbf{x}_t) - f(\mathbf{x}_{t+1}))$$

Again, we have telescoping sum

$$\begin{aligned} \sum_{t=0}^T \|\nabla f(\mathbf{x}_t)\|^2 &\leq \sum_{t=0}^T 2L(f(\mathbf{x}_t) - f(\mathbf{x}_{t+1})) \\ &= 2L(f(\mathbf{x}_0) - f(\mathbf{x}_T)) \\ &\leq 2L(f(\mathbf{x}_0) - f(\mathbf{x}^\star)) \end{aligned}$$

Dividing by  $T$ , the proof is completed.

# Adaptive methods for SGD

Several methods exist

- AdaGrad
- AdaDelta
- RMSProp
- Adam
- ...

## Adam: A Method for Stochastic Optimization

Authors Diederik P. Kingma, Jimmy Ba

Publication date 2014/12/22

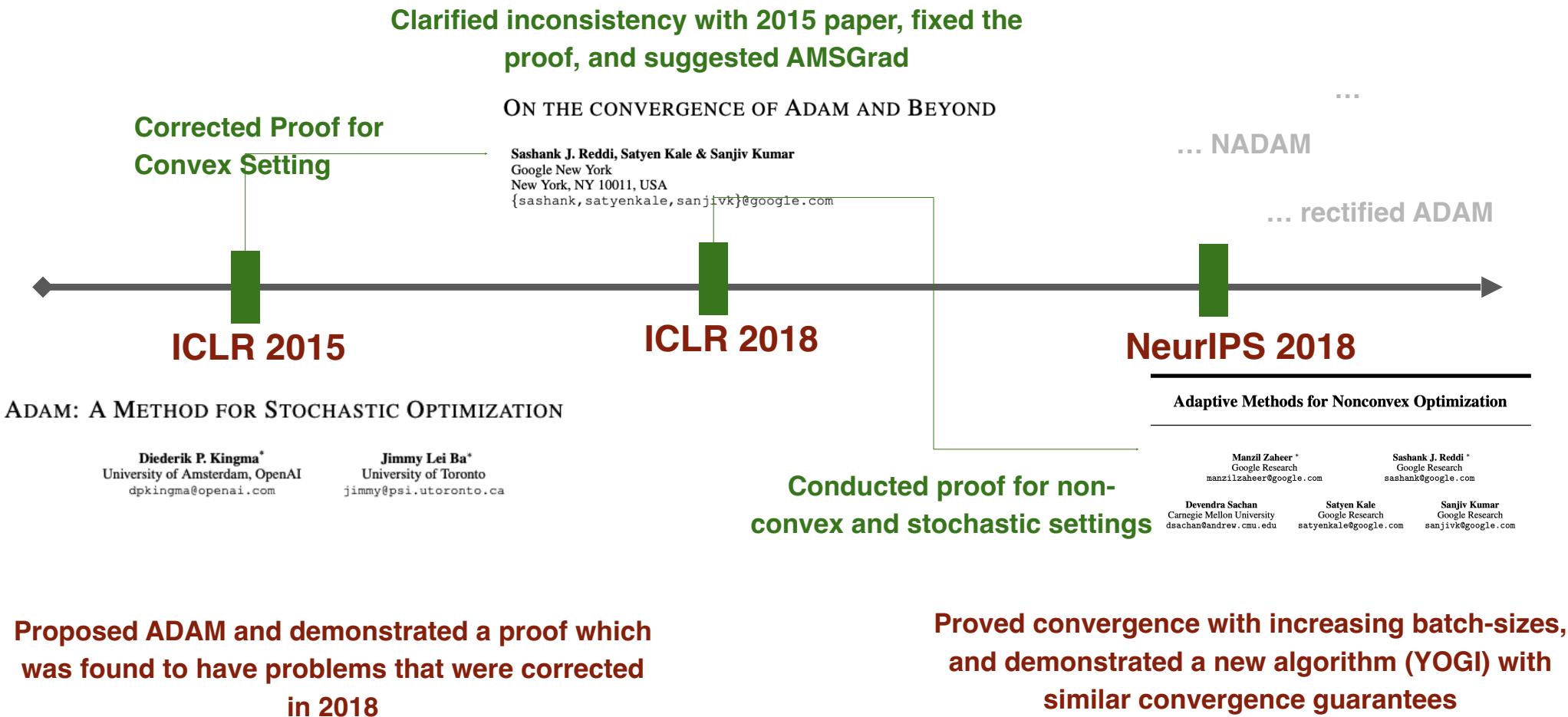
Journal Proceedings of the 3rd International Conference on Learning Representations (ICLR)

Description We introduce Adam, an algorithm for first-order gradient-based optimization of stochastic objective functions, based on adaptive estimates of lower-order moments. The method is straightforward to implement, is computationally efficient, has little memory requirements, is invariant to diagonal rescaling of the gradients, and is well suited for problems that are large in terms of data and/or parameters. The method is also appropriate for non-stationary objectives and problems with very noisy and/or sparse gradients. The hyperparameters have intuitive interpretations and typically require little tuning. Some connections to related algorithms, on which Adam was inspired, are discussed. We also analyze the theoretical convergence properties of the algorithm and provide a regret bound on the convergence rate that is comparable to the best known results under the online convex optimization framework. Empirical results demonstrate that Adam works well in practice and compares favorably to other stochastic optimization methods. Finally, we discuss AdaMax, a variant of Adam based on the infinity norm.

Total citations Cited by 113113



# History of Adam Optimizer



# Adam

**Theorem 1.** Let  $\gamma_t = \gamma$  for all  $t \in [T]$ . Furthermore, assume that  $\varepsilon, \beta_2$  and  $\gamma$  are chosen such that the following conditions satisfied:  $\gamma \leq \frac{\varepsilon}{2L}$  and  $1 - \beta_2 \leq \frac{\varepsilon^2}{16B^2}$ . Then for  $\mathbf{x}_t$  generated using ADAM, we have the following bound

$$\mathbb{E} \left\| \nabla f(\mathbf{x}_i) \right\|^2 \leq O \left( \frac{f(\mathbf{x}_1) - f(\mathbf{x}^*)}{\gamma T} + \sigma^2 \right)$$

where  $\mathbf{x}^*$  is an optimal solution to the problem in (1) and  $\mathbf{x}_i$  is an iterate uniformly randomly chosen from  $\{\mathbf{x}_1, \dots, \mathbf{x}_T\}$ .  $\sigma^2$  is the bound of the variance in stochastic gradients.  $B$  is upper bound of stochastic gradient norm.

## Adam

```

Set  $\mathbf{m}_0 = 0, \mathbf{v}_0 = 0$ 
for  $t = 1$  to  $T$  do
    Draw a sample  $\mathbf{s}_t$ .
    Compute  $\mathbf{g}_t = \nabla \ell(\mathbf{x}_t, \mathbf{s}_t)$ .
     $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$ 
     $\mathbf{v}_t = \mathbf{v}_{t-1} - (1 - \beta_2) (\mathbf{v}_{t-1} - \mathbf{g}_t^2)$ 
     $\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{m}_t / (\sqrt{\mathbf{v}_t} + \epsilon)$ 
end for

```



$$f(\mathbf{x}) = \frac{1}{n} \sum_{i=1}^n \ell(\mathbf{x}_i, \mathbf{a}_i, b_i), \text{ with } \ell(\mathbf{x}_i, \mathbf{a}_i, b_i) = (b_i - f_{\mathbf{x}}(\mathbf{a}_i))^2$$

$$\text{Sample } \xi_t = i_t \in \{1, \dots, n\}$$

$$\Rightarrow \ell(\mathbf{x}, i_t) = (b_{i_t} - f_{\mathbf{x}}(\mathbf{a}_{i_t}))^2$$

$$\nabla_{\mathbf{x}} \ell(\mathbf{x}, i_t) = \nabla_{\mathbf{a}} \left( b_{i_t} - f_{\mathbf{x}}(\mathbf{a}_{i_t}) \right)^2$$

$$= -2 \left( b_{i_t} - f_{\mathbf{x}}(\mathbf{a}_{i_t}) \right) \nabla f_{\mathbf{x}}(\mathbf{a}_{i_t})$$

# Adam proof

**Key steps in proof ( $b$  is minibatch size):**

1- assume that  $f$  is  $L$ -smooth, we have:

$$\begin{aligned} f(\mathbf{x}_{t+1}) &\leq f(\mathbf{x}_t) + \nabla f(\mathbf{x}_t)^\top (\mathbf{x}_{t+1} - \mathbf{x}_t) + \frac{L}{2} \|\mathbf{x}_t - \mathbf{x}_{t+1}\|^2 \\ &= f(\mathbf{x}_t) - \gamma_t \sum_{i=1}^d ([\nabla f(\mathbf{x}_t)]_i \times \frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon}) + \frac{L\gamma_t^2}{2} \sum_{i=1}^d \frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \end{aligned}$$

2- taking expectation

$$\begin{aligned} \mathbb{E}[f(\mathbf{x}_{t+1}) | \mathbf{x}_t] &\leq f(\mathbf{x}_t) - \gamma_t \sum_{i=1}^d ([\nabla f(\mathbf{x}_t)]_i \times \mathbb{E}_t \left[ \frac{g_{t,i}}{\sqrt{v_{t,i}} + \epsilon} \right]) + \frac{L\gamma_t^2}{2} \sum_{i=1}^d \mathbb{E}_t \left[ \frac{g_{t,i}^2}{(\sqrt{v_{t,i}} + \epsilon)^2} \right] \\ &\leq f(\mathbf{x}_t) - \frac{\gamma_t}{2(\sqrt{\beta_2}B + \epsilon)} \|\nabla f(\mathbf{x}_t)\|^2 + \left( \frac{\gamma_t B \sqrt{1-\beta_2}}{\epsilon^2} + \frac{L\gamma_t^2}{2\epsilon^2} \right) \frac{\sigma^2}{b} \end{aligned}$$

3-Using telescoping sum and rearranging the inequality, we obtain

$$\frac{\gamma}{2(\sqrt{\beta_2}B + \epsilon)} \sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 \leq f(\mathbf{x}_1) - \mathbb{E}[f(\mathbf{x}_{T+1})] + \left( \frac{\gamma B \sqrt{1-\beta_2}}{\epsilon^2} + \frac{L\gamma^2}{2\epsilon^2} \right) \frac{T\sigma^2}{b}$$

4- using the fact that  $f(\mathbf{x}^*) \leq f(\mathbf{x}_{t+1})$

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\nabla f(\mathbf{x}_t)\|^2 \leq 2(\sqrt{\beta_2}B + \epsilon) \times \left[ \frac{f(\mathbf{x}_1) - f(\mathbf{x}^*)}{\gamma T} + \left( \frac{B \sqrt{1-\beta_2}}{\epsilon^2} + \frac{L\gamma}{2\epsilon^2} \right) \frac{\sigma^2}{b} \right]$$

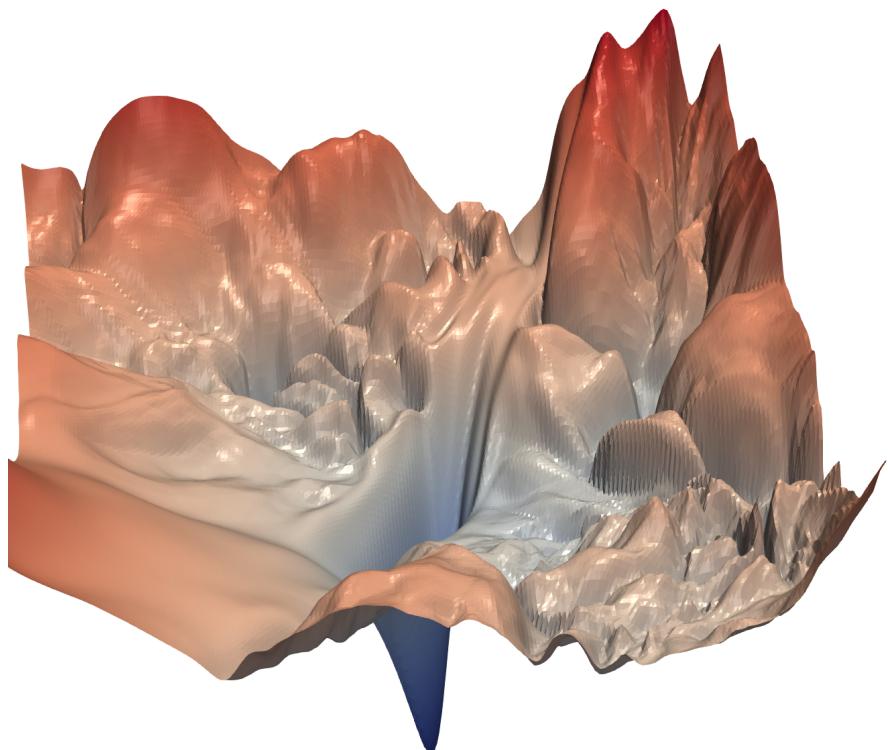
Notes:

- $\mathbb{E}[\|\nabla f(\mathbf{x}_i)\|^2] \leq O(1/b)$  as  $T \rightarrow \infty$ . Not necessarily imply that the  $\mathbf{x}_i$  is close to a stationary point, but sufficient for many ML applications
- ADAM with  $b = \Theta(T) \Rightarrow \mathbb{E}[\|\nabla f(\mathbf{x}_i)\|^2] \leq O(1/T)$ , implying converging to the stationary point

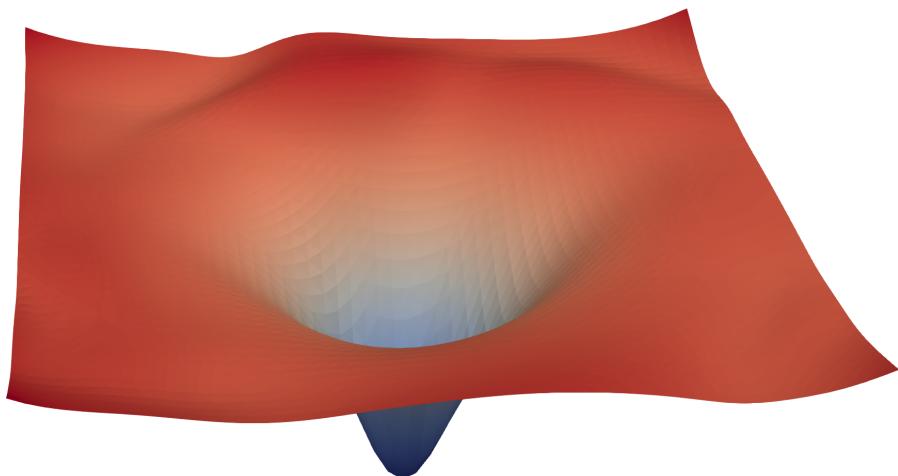
<sup>1</sup>Zaheer, Manzil, et al. "Adaptive methods for nonconvex optimization." *Advances in neural information processing systems* 31 (2018).

<sup>2</sup>[https://www.youtube.com/watch?v=n65pElMp6x8&ab\\_channel=MachineLearningandAIAcademy](https://www.youtube.com/watch?v=n65pElMp6x8&ab_channel=MachineLearningandAIAcademy)

# what one can hope for DNNs?



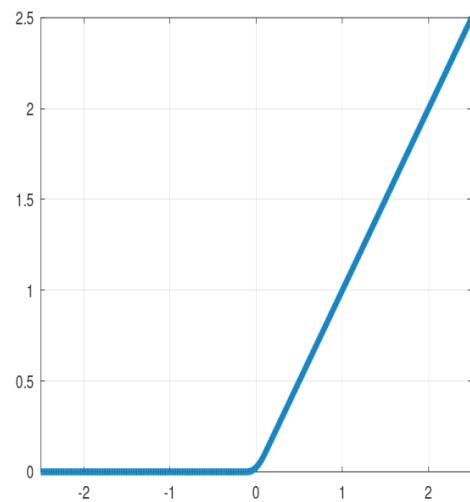
(a) without skip connections



(b) with skip connections

The loss surfaces of ResNet-56 with/without skip connections.

# what one can hope for DNNs?



Smoothed ReLU



Loss surfaces (as functions of a 2D input) for two sample loss functions (middle and right) as the activation function's transition region widens, going from from ReLU to an increasingly smoother SmeLU (left).

# Machine learning applications

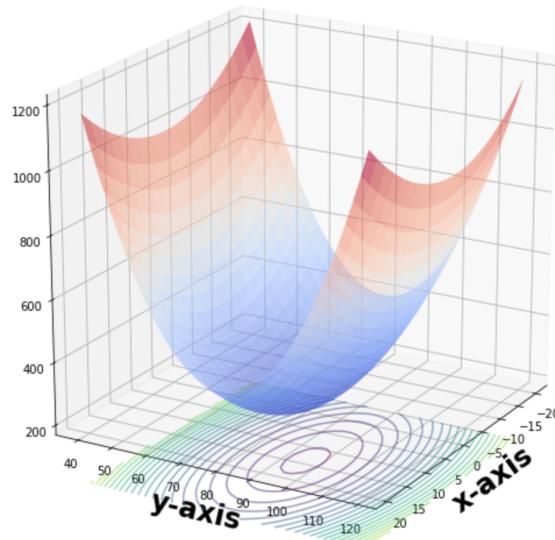
# Example: linear regression(MSE loss)

**Problem:** Given data points  $\mathbf{x}^1, \dots, \mathbf{x}^n \in \mathbb{R}^d$  and  $y^1, \dots, y^n \in \mathbb{R}$ , find the best fit line. Consider a simple linear regression model with  $\hat{y} = \mathbf{w}^T \mathbf{x} + b$ . The formulation can be simplified to  $\hat{y} = \boldsymbol{\theta}^T \tilde{\mathbf{x}}$ , in which  $\boldsymbol{\theta} = [\mathbf{w}^T, b]^T$ ,  $\tilde{\mathbf{x}} = [\mathbf{x}^T, 1]^T$ .

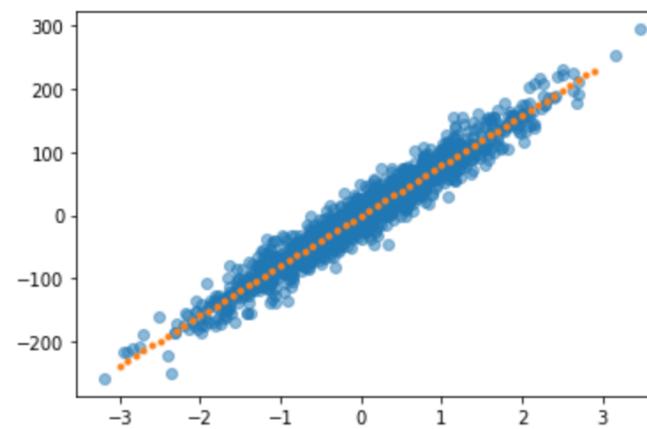
**Objective:**

$$\min_{\boldsymbol{\theta}} f(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n (\hat{y}_i - y_i)^2 = \frac{1}{2n} \sum_{i=1}^n (\boldsymbol{\theta}^T \mathbf{x}'_i - y_i)^2$$

**Loss Surface:**



**Data Distribution:**



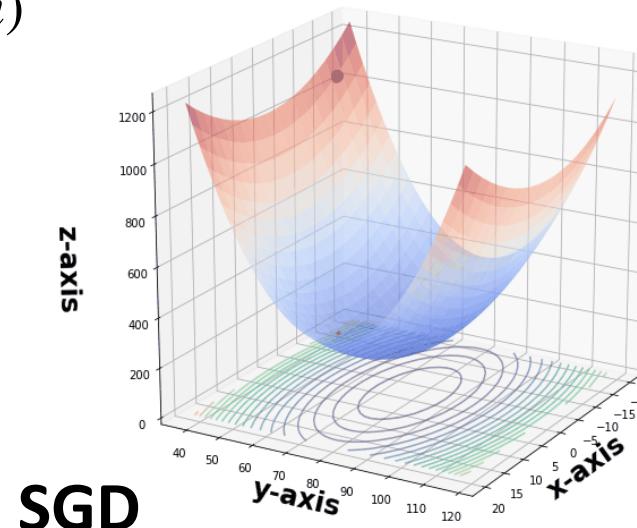
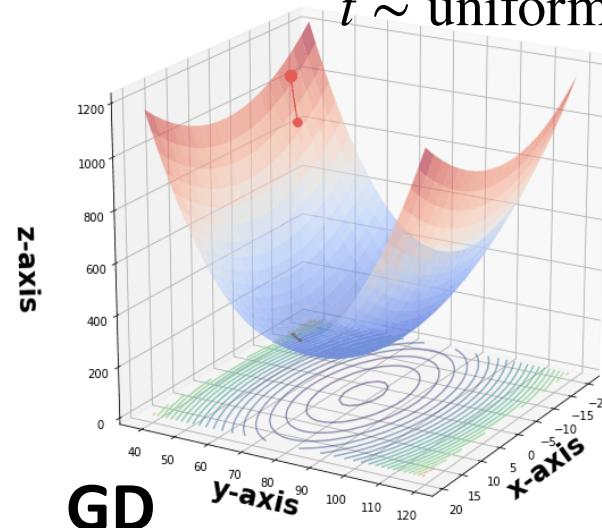
# Example: linear regression(MSE loss)

**GD:**  $\theta_{t+1} := \theta_t - \gamma \sum_{i=1}^n \tilde{\mathbf{x}}_i^T (\theta_t^T \tilde{\mathbf{x}}_i - y_i),$   
 $t = 0, 1, 2, \dots$

**SGD:**  $\theta_{t+1} := \theta_t - \gamma \tilde{\mathbf{x}}_i^T (\theta_t^T \tilde{\mathbf{x}}_i - y_i),$   
 $t \sim \text{uniform}(0, n)$

**Mini-SGD:**

$\theta_{t+1} := \theta_t - \gamma \sum_{i=1}^b \tilde{\mathbf{x}}_i^T (\theta_t^T \tilde{\mathbf{x}}_i - y_i),$   
 $t \sim \text{uniform}(0, n)$

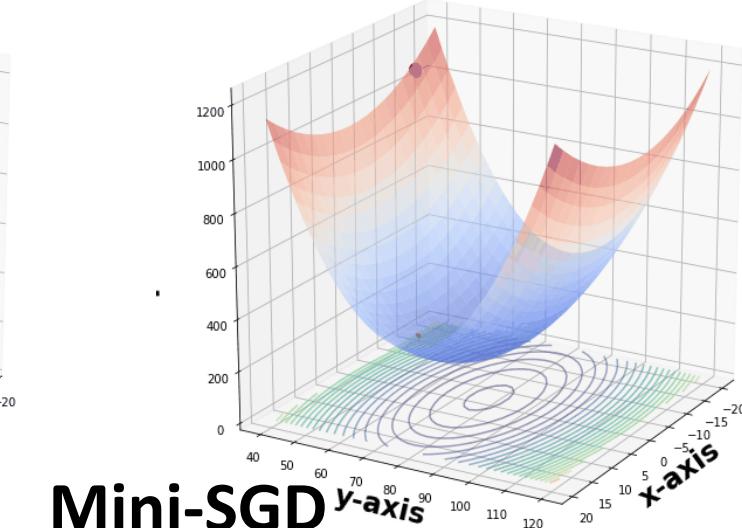


SGD

```
pred_y = np.dot(x, theta)
error = pred_y - y
gradient = x.T.dot(error)/sample_number
theta = theta - step * gradient
```

```
pred_y = np.dot(x, theta)
error = pred_y - y
idx = random.sample(range(0, m), 1)
gradient = x[idx].T.dot(error[idx])
theta = theta - step * gradient
```

```
pred_y = np.dot(x, theta)
error = pred_y - y
idx = random.sample(range(0, m), batch_size)
gradient = x[idx].T.dot(error[idx])/batch_size
theta = theta - step * gradient
```



Mini-SGD

# Example: linear regression(MSE loss)

Adam:

Draw a batch of sample .

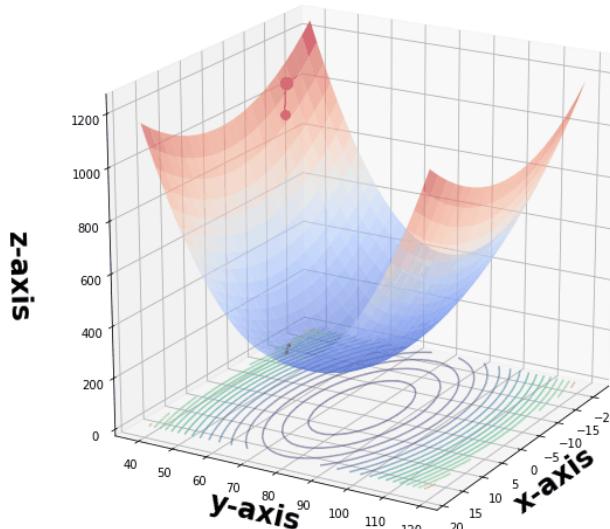
$$\text{Compute } \mathbf{g}_t = \sum_{i=1}^b \tilde{\mathbf{x}}_i^T (\theta_t^T \tilde{\mathbf{x}}_i - y_i)$$

$$\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$$

$$\mathbf{v}_t = \mathbf{v}_{t-1} - (1 - \beta_2) (\mathbf{v}_{t-1} - \mathbf{g}_t^2)$$

(Bias correction)

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \eta_t \mathbf{m}_t / \left( \sqrt{\mathbf{v}_t} + \epsilon \right)$$

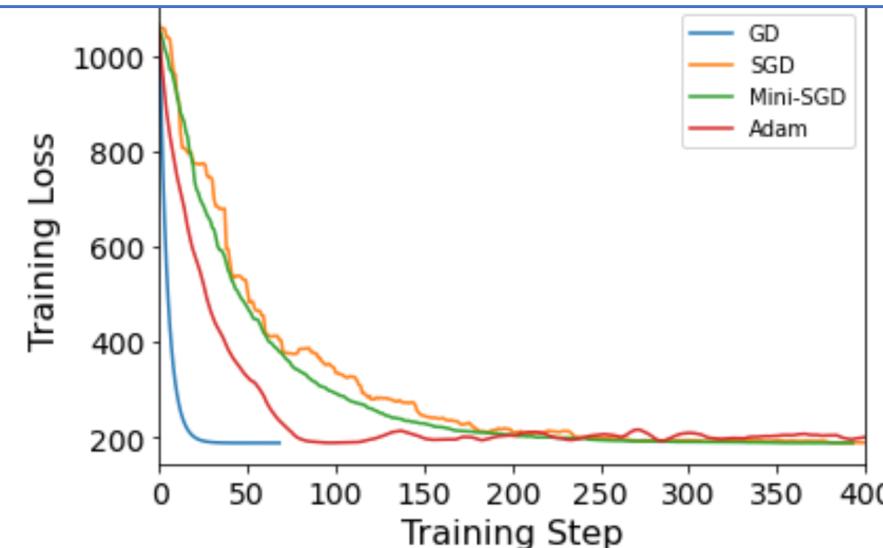


```
pred_y = np.dot(x, theta)
error = pred_y - y
idx = random.sample(range(0, m), batch_size)
batch_gradient = x[idx].T.dot(error[idx])/ batch_size

gred_m = beta1 * gred_m + (1-beta1) * batch_gradient
gred_n = beta2 * gred_n + (1-beta2) * batch_gradient**2

gred_m_adaptive = gred_m/(1-beta1**counter)
gred_n_adaptive = gred_n/(1-beta2**counter)
gradient = gred_m_adaptive/ (np.sqrt(gred_n_adaptive) + epsilon)

theta = theta - step * gradient
```



# Example: Logistic regression(MSE loss)

**Problem:** Given data points  $\mathbf{x}^1, \dots, \mathbf{x}^n \in \mathbb{R}^d$  and labels  $y^1, \dots, y^n \in \{0,1\}$ , find the best classifier.

Denote  $\mathbf{x} = [x_1, x_2]^T$ . Consider a sigmoid function that predict the soft label:

$$\sigma(\boldsymbol{\theta}^T \tilde{\mathbf{x}}) = \frac{1}{1 + \exp(-\boldsymbol{\theta}^T \tilde{\mathbf{x}})}, \text{ where } \boldsymbol{\theta} = [1, w, b]^T, \tilde{\mathbf{x}} = [x_1, x_2, 1]^T.$$

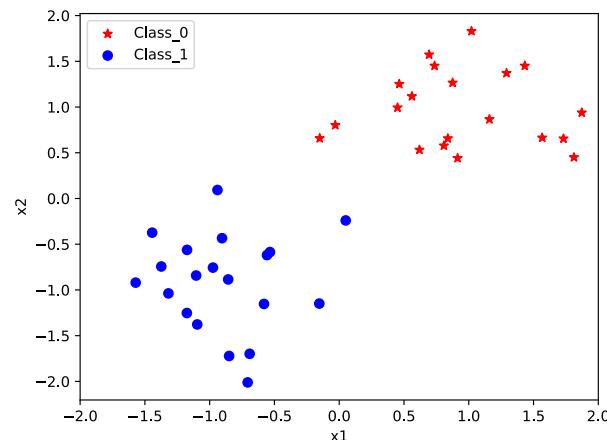
**Objective:**

$$\min_{w,b} f(w, b) = \frac{1}{2n} \sum_{i=1}^n (\sigma(\boldsymbol{\theta}^T \tilde{\mathbf{x}}_i) - y_i)^2$$

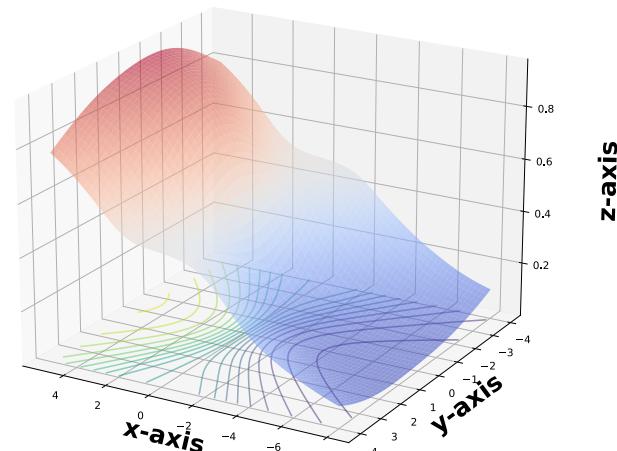
**Gradient**  $\nabla \ell_i(\boldsymbol{\theta}) = \tilde{\mathbf{x}}_i^T (\sigma(\boldsymbol{\theta}_t^T \tilde{\mathbf{x}}) - y_i) \sigma(\boldsymbol{\theta}_t^T \tilde{\mathbf{x}}_i) (1 - \sigma(\boldsymbol{\theta}_t^T \tilde{\mathbf{x}}_i))$

Derivative of the Sigmoid Function  
 $\sigma'(x) = \sigma(x)(1 - \sigma(x))$

**Data Distribution:**

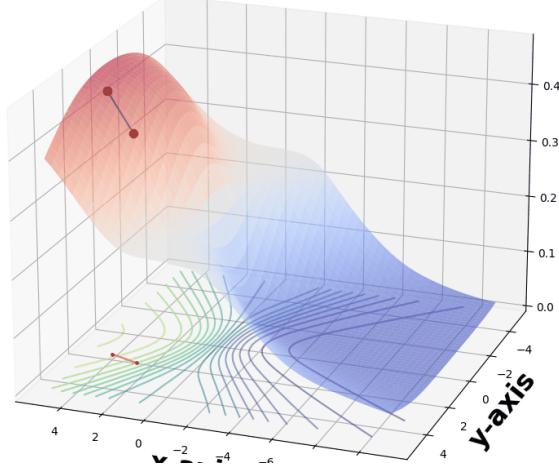


**Loss Surface:**

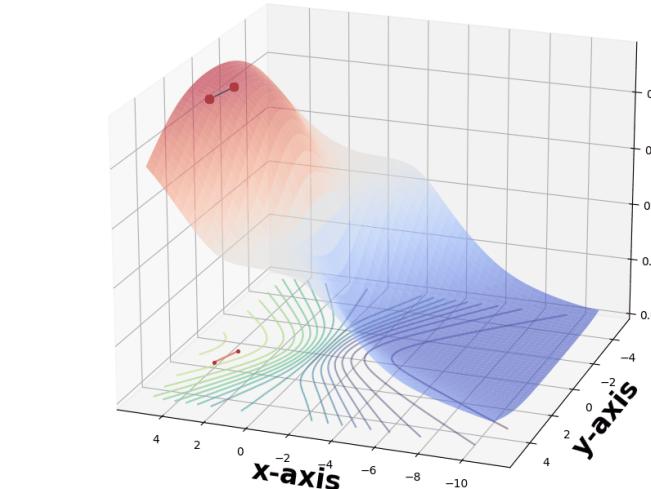


# Example: Logistic regression(MSE loss)

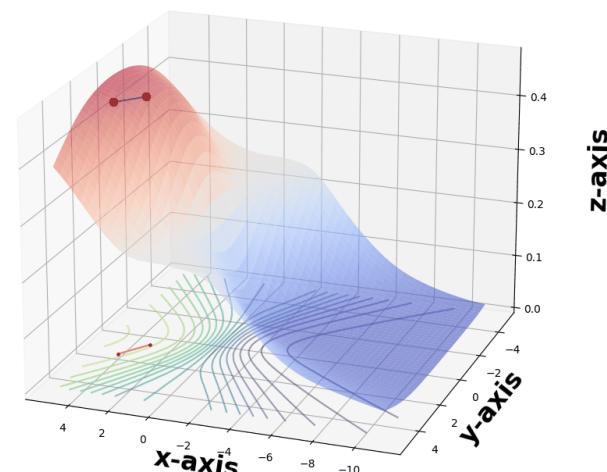
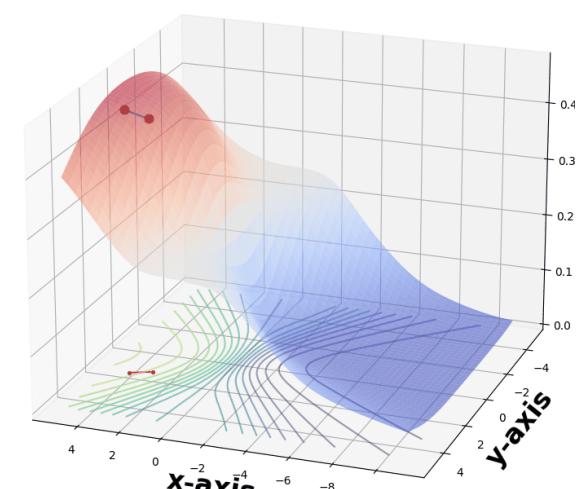
step\_size = 1 for all  
algorithms



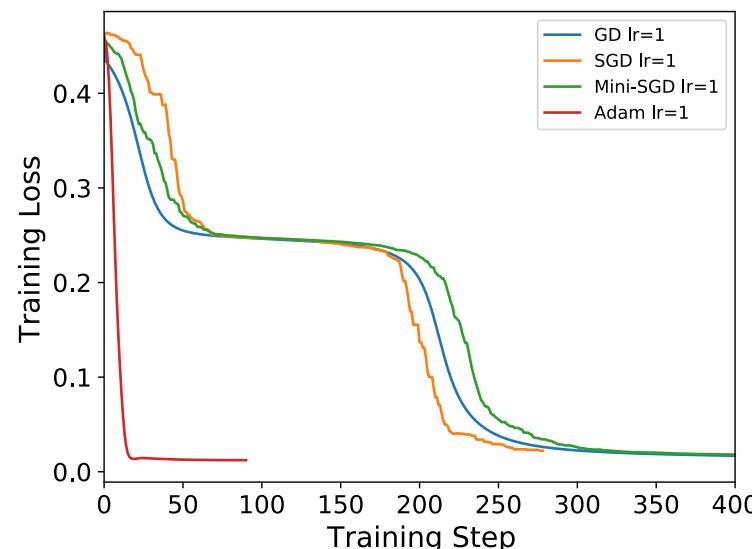
GD



SGD

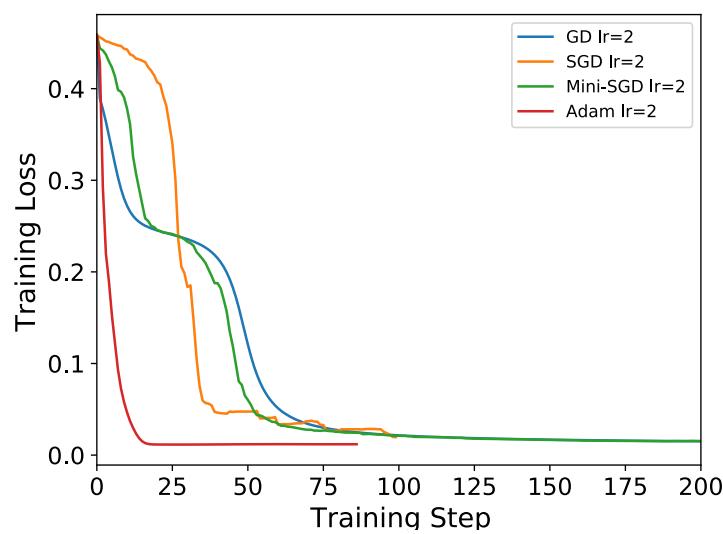
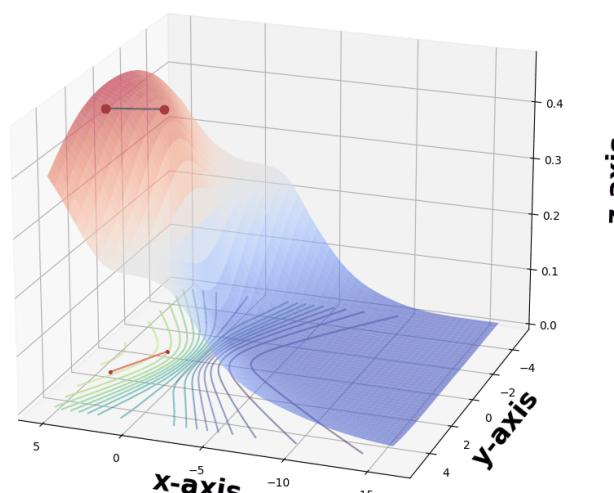
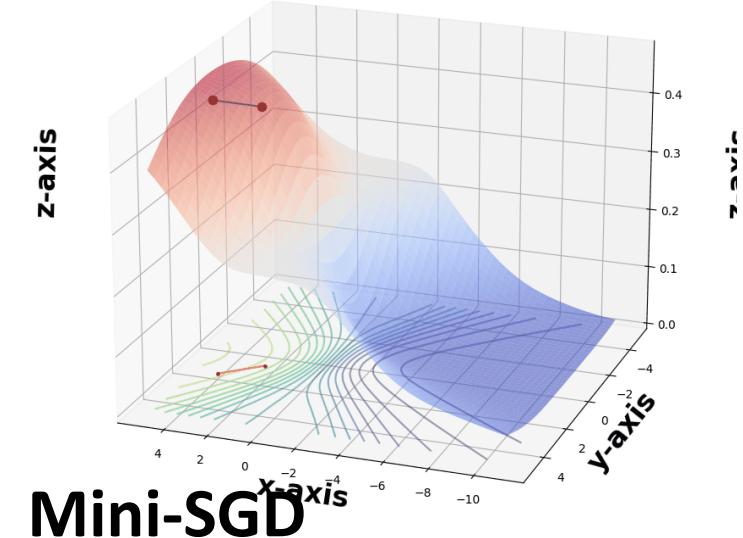
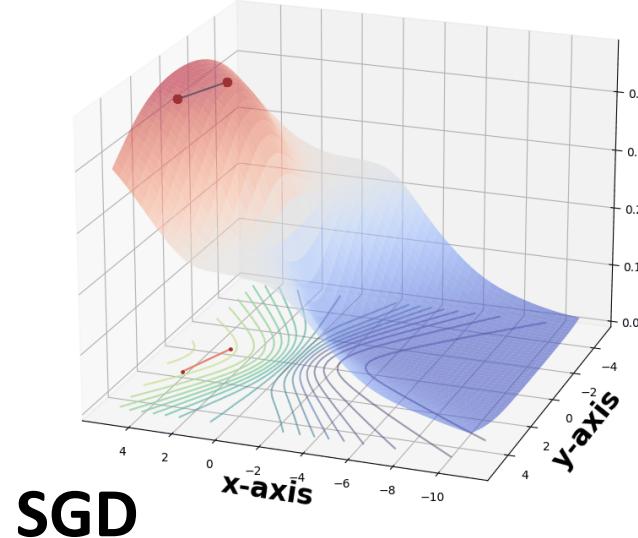
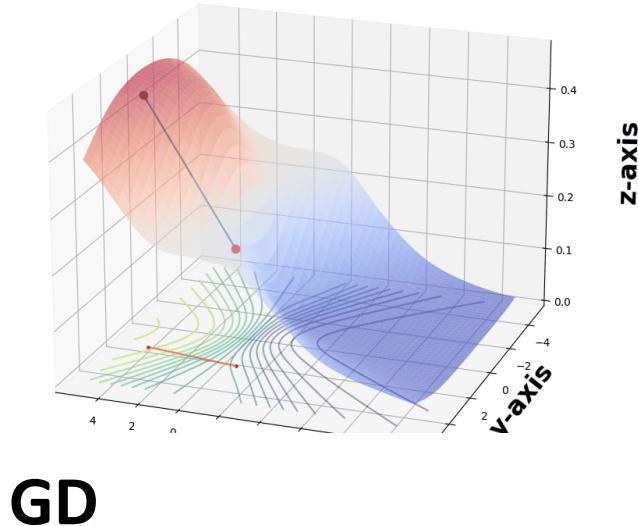


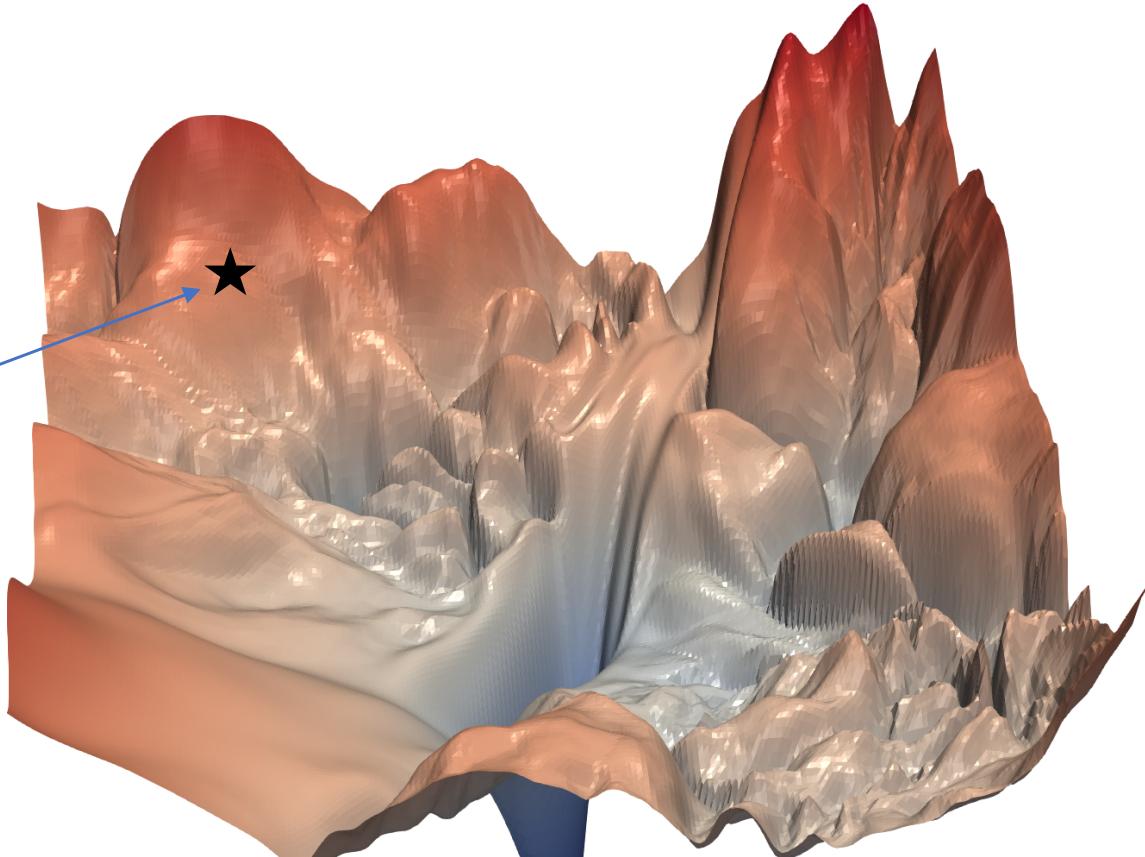
Adam



# Example: Logistic regression(MSE loss)

step\_size = 2 for all  
algorithms





We are here!



# Thank you!

<https://yalidu.github.io/>