

# **6.390** Intro to Machine Learning

Lecture 3: Gradient Descent Methods

Shen Shen

Sept 18, 2025

11am, Room 10-250

Interactive Slides and Lecture Recording

### Logistics

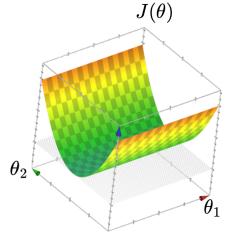
- 1. There is a student in our class who needs copies of class notes as an approved accommodation. If you're interested in serving as a paid note taker, please reach out to DAS, at 617-253-1674 or das-student@mit.edu.
- 2. **Midterm 1: October 8, 730pm-9pm**. It covers all the materials up to and including week 4 (linear classification). If you need to take the conflict or accommodation exam, please get in touch with us at 6.390-personal@mit.edu by **Sept 24**.
- 3. Heads-up: Midterm 2 is November 12, 730pm-9pm. Final is December 15, 9am-12pm.

More details on introML homepage

## Outline

- Gradient descent (GD)
  - The gradient vector
  - GD algorithm
  - Gradient decent properties
- Stochastic gradient descent (SGD)
  - SGD algorithm and setup
  - SGD vs. GD

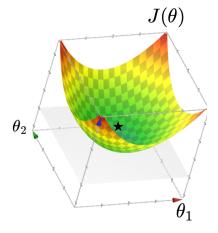
#### Recall:



When *X* is not full column rank

- $J(\theta)$  has a "flat" bottom, like a half pipe
- This formula is not well-defined
- Infinitely many optimal hyperplanes

• No way yet to obtain an optimal parameter



Typically, X is full column rank

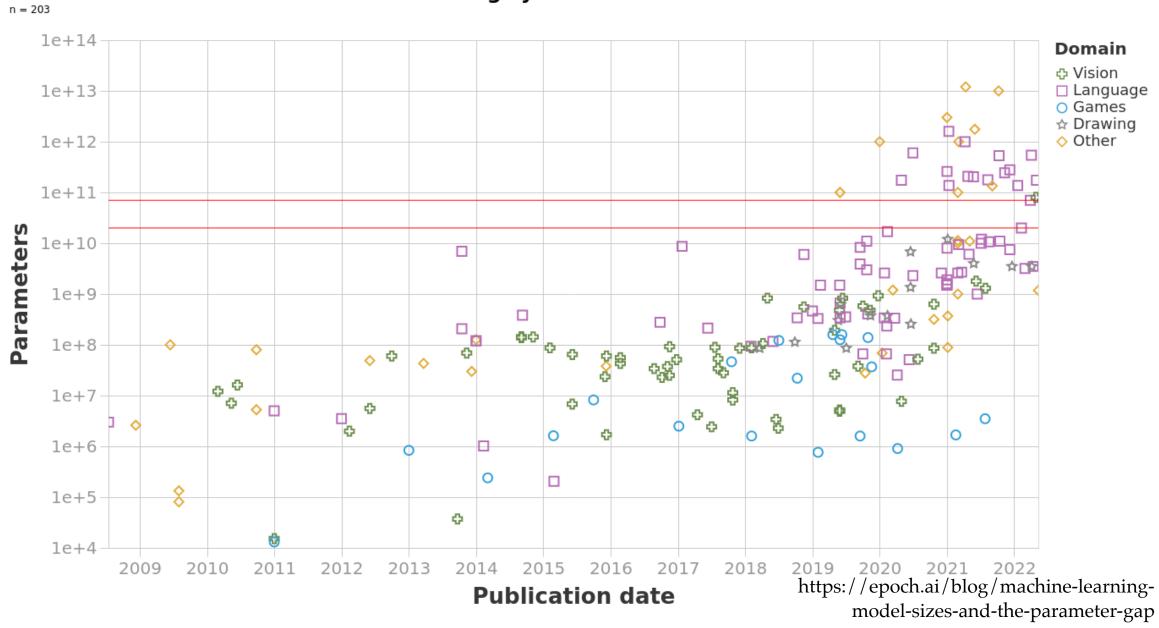
•  $J(\theta)$  "curves up" everywhere

$$ullet \; heta^* = \left( X^ op X 
ight)^{-1} X^ op Y$$

•  $\theta^*$  gives the unique optimal hyperplane

•  $\theta^*$  can be costly to compute (lab2, Q2.7)

### Parameters of milestone Machine Learning systems over time



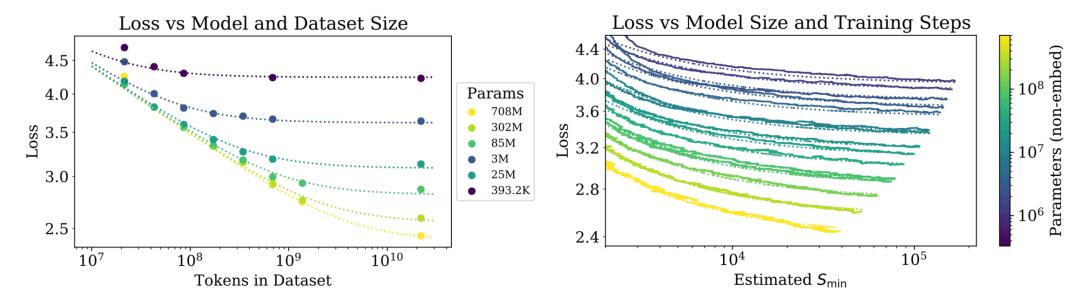
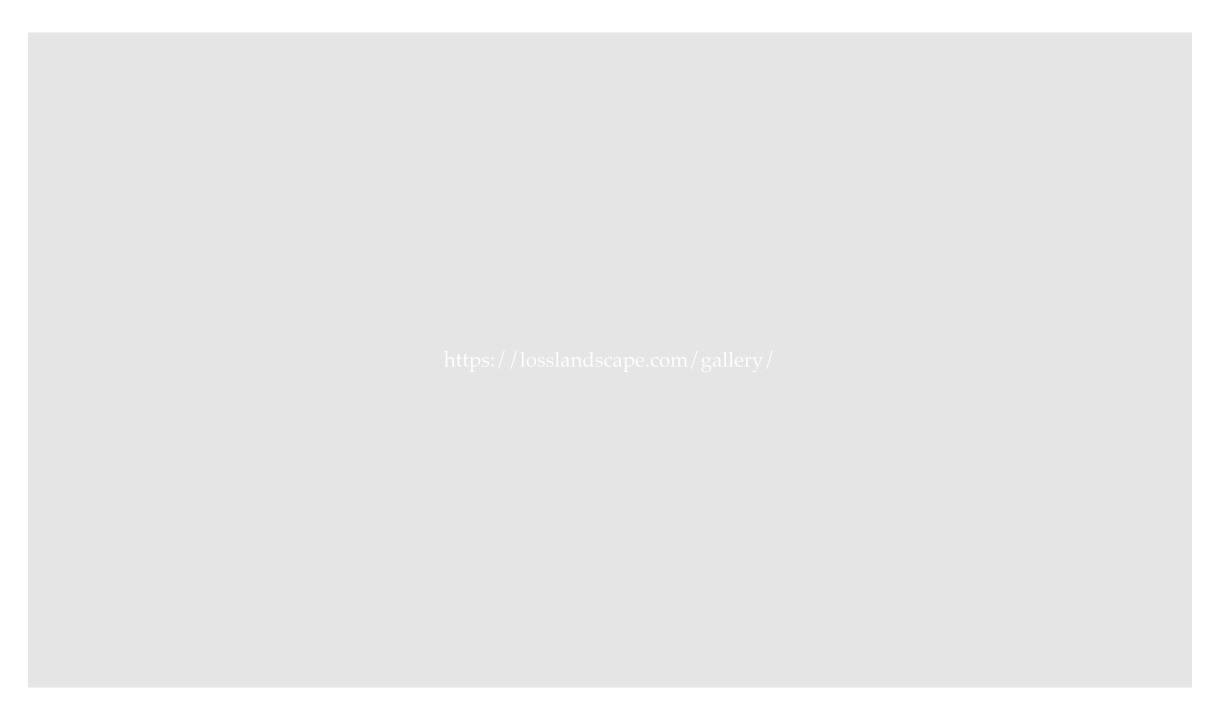


Figure 4 Left: The early-stopped test loss L(N, D) varies predictably with the dataset size D and model size N according to Equation (1.5). **Right**: After an initial transient period, learning curves for all model sizes N can be fit with Equation (1.6), which is parameterized in terms of  $S_{\min}$ , the number of steps when training at large batch size (details in Section 5.1).



In the real world,

- the number of parameters is huge
- the number of training data points is huge
- hypothesis class is typically highly nonlinear
- loss function is rarely as simple as squared error

Need a more **efficient** and **general** algorithm to train => gradient descent methods

## Outline

- Gradient descent algorithm (GD)
  - The gradient vector
  - GD algorithm
  - Gradient decent properties
- Stochastic gradient descent (SGD)
  - SGD algorithm and setup
  - SGD vs. GD

For  $f:\mathbb{R}^m \to \mathbb{R}$ , its *gradient*  $\nabla f:\mathbb{R}^m \to \mathbb{R}^m$  is defined at the point  $p=(x_1,\ldots,x_m)$  as:

$$abla f(p) = \left[ egin{array}{c} rac{\partial f}{\partial x_1}(p) \ dots \ rac{\partial f}{\partial x_m}(p) \end{array} 
ight]$$

- 1. The gradient generalizes the concept of a derivative to multiple dimensions.
- 2. By construction, the gradient's dimensionality always matches the function input.

3. The gradient can be symbolic or numerical.

example: 
$$f(x, y, z) = x^2 + y^3 + z$$

$$abla f(p) = \left[egin{array}{c} rac{\partial f}{\partial x_1}(p) \ dots \ rac{\partial f}{\partial x_m}(p) \end{array}
ight]$$

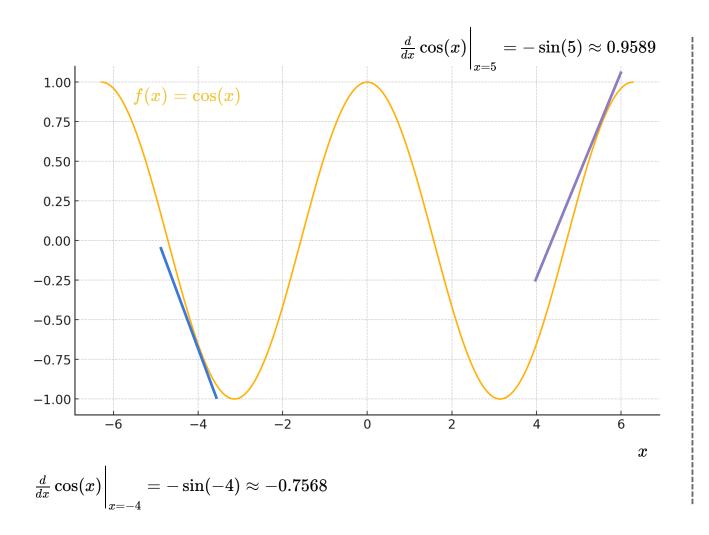
its *symbolic* gradient:

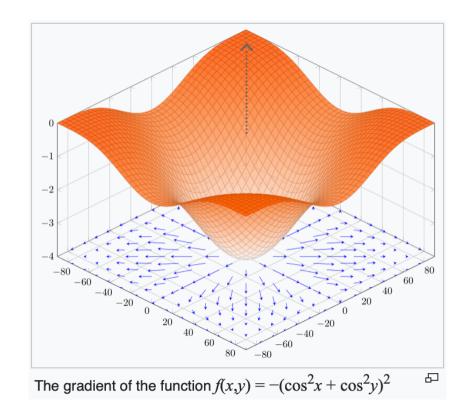
$$abla f(x,y,z) = egin{bmatrix} 2x \ 3y^2 \ 1 \end{bmatrix}$$

evaluating the symbolic gradient at a point gives a numerical gradient:

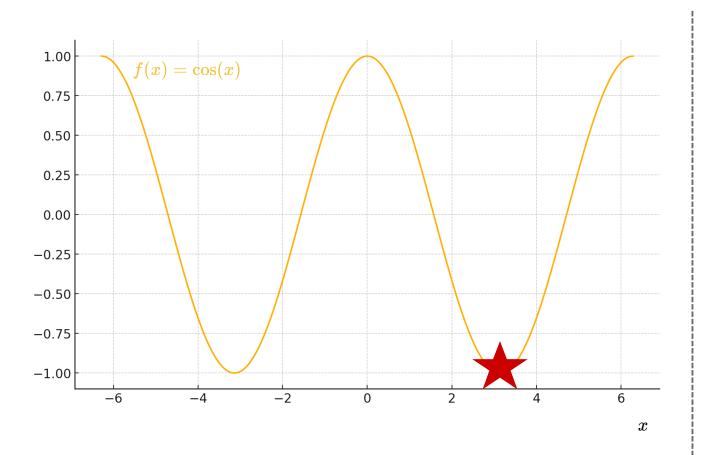
$$\left. 
abla f(3,2,1) = 
abla f(x,y,z) 
ight|_{(x,y,z)=(3,2,1)} = egin{bmatrix} 6 \ 12 \ 1 \end{bmatrix}$$

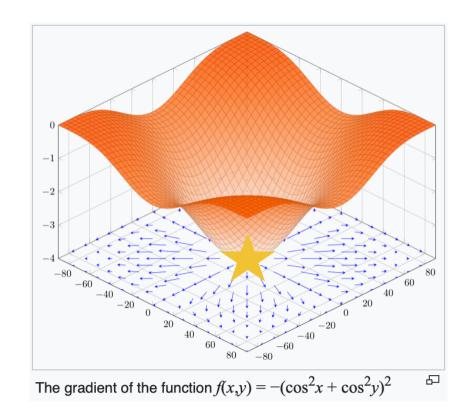
### 4. The gradient points in the direction of the (steepest) *increase* in the function value.





### 5. The gradient at the function minimizer is *necessarily* zero.





For  $f:\mathbb{R}^m \to \mathbb{R}$ , its *gradient*  $\nabla f:\mathbb{R}^m \to \mathbb{R}^m$  is defined at the point  $p=(x_1,\ldots,x_m)$  as:

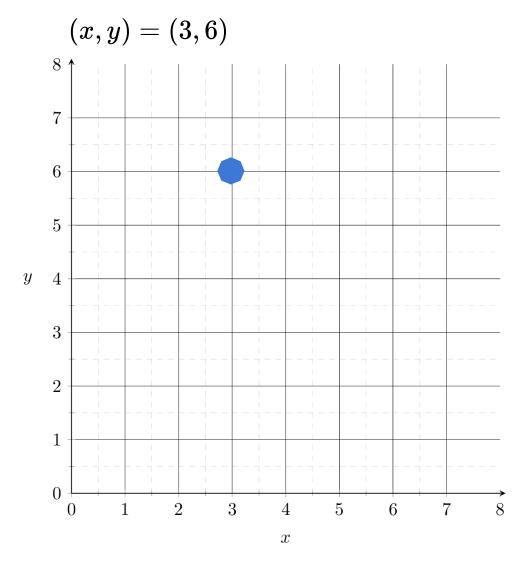
$$abla f(p) = \left[ egin{array}{c} rac{\partial f}{\partial x_1}(p) \ dots \ rac{\partial f}{\partial x_m}(p) \end{array} 
ight]$$

- 1. The gradient generalizes the concept of a derivative to multiple dimensions.
- 2. By construction, the gradient's dimensionality always matches the function input.
- 3. The gradient can be symbolic or numerical.
- 4. The gradient points in the direction of the (steepest) *increase* in the function value.
- 5. The gradient at the function minimizer is *necessarily* zero.

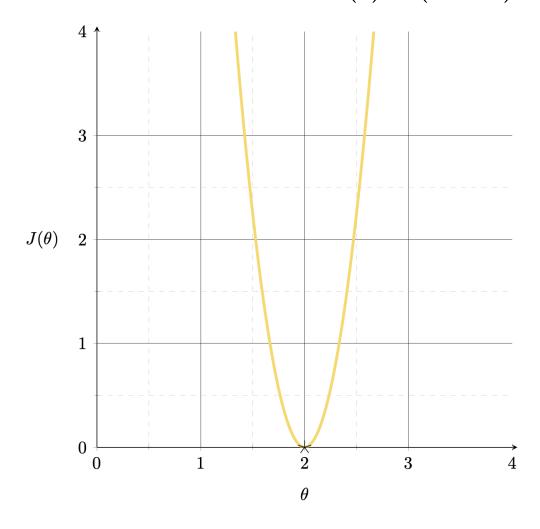
## Outline

- Gradient descent algorithm (GD)
  - The gradient vector
  - GD algorithm
  - Gradient decent properties
- Stochastic gradient descent (SGD)
  - SGD algorithm and setup
  - SGD vs. GD

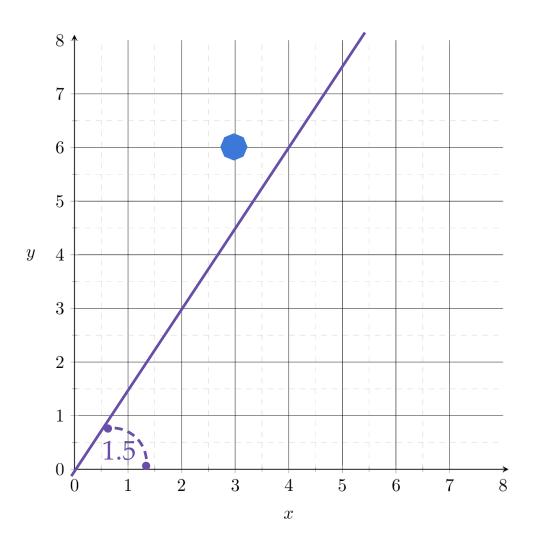
### A single training data point



Want to fit a line (without offset) to minimize the MSE:  $J(\theta) = (3\theta - 6)^2$ 

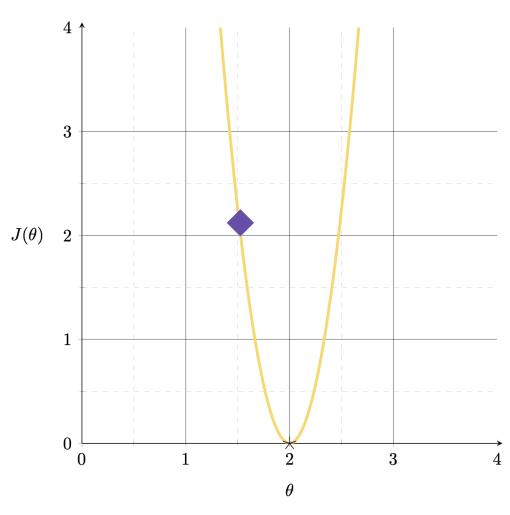


### Suppose we fit a line y = 1.5x



### MSE could get better.

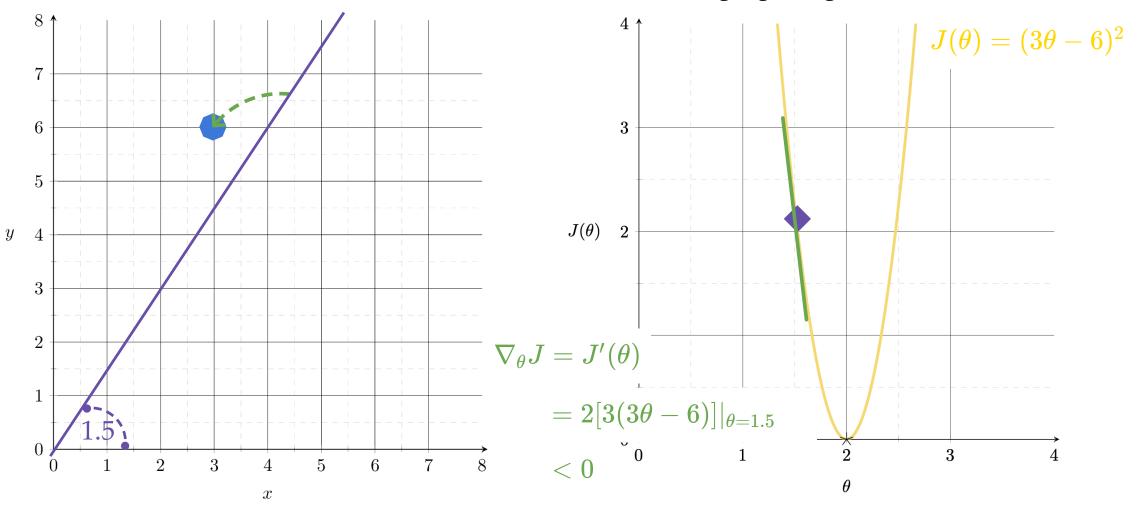
### How to formalize this?



### Suppose we fit a line y = 1.5x

### MSE could get better. How to?

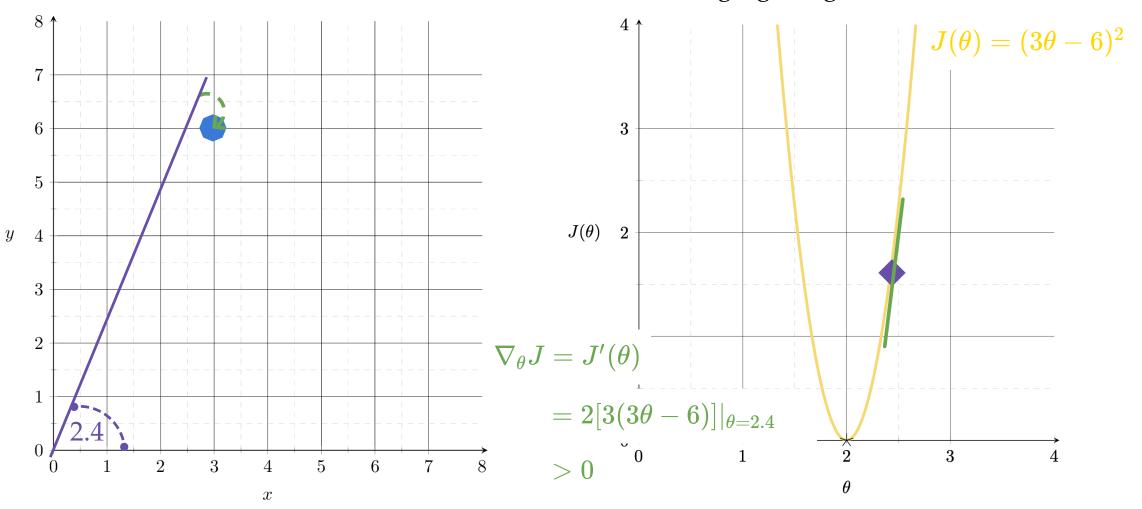
Leveraging the gradient.



### Suppose we fit a line y = 2.4x

MSE could get better. How to?

Leveraging the gradient.



### initial guess

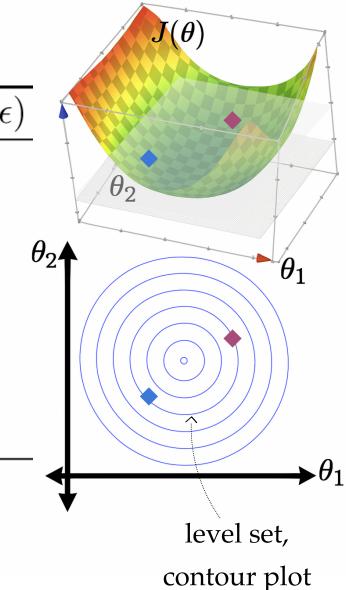
hyperparameters

of parameters learning rate

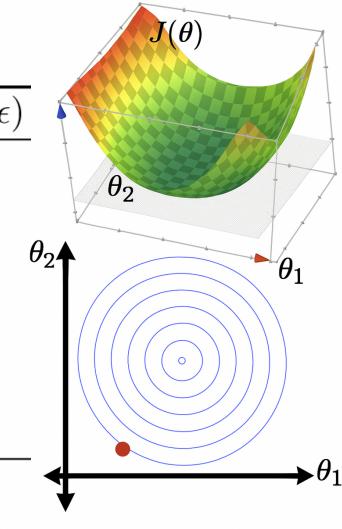
precision

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \ \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \ \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$

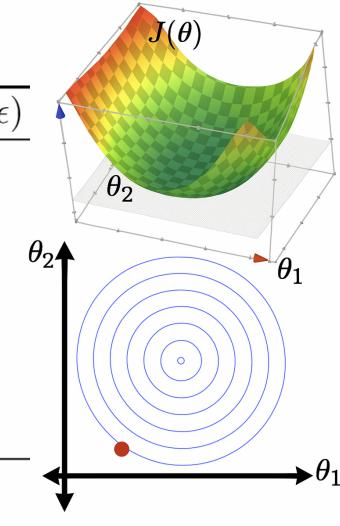


- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \ \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$

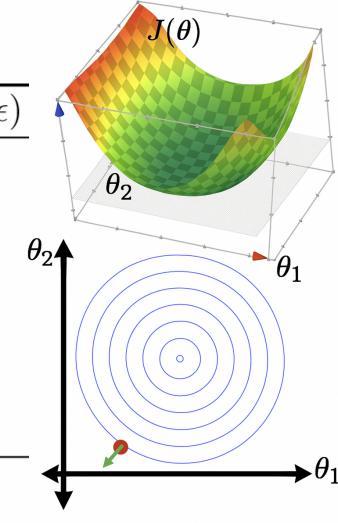


counter

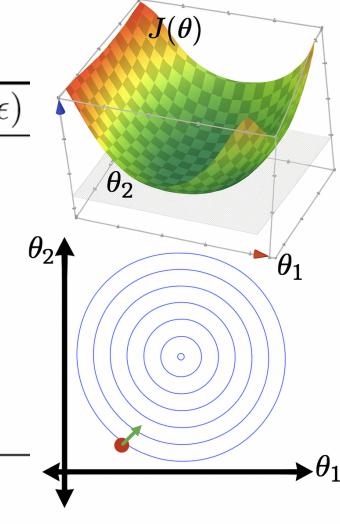
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t = 0iteration
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



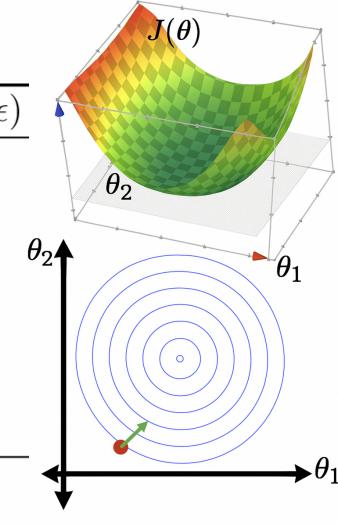
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6: **until**  $\left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



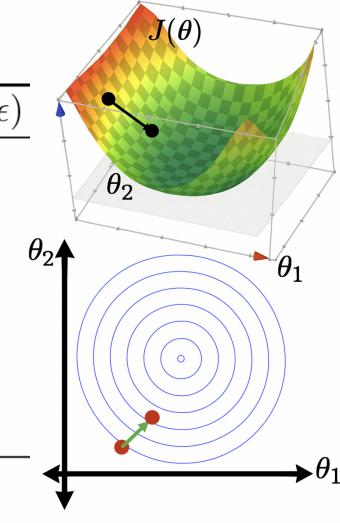
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



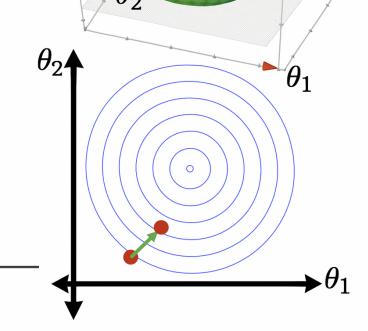
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



• What does this 3d vector \ represent? anything in the psuedocode?

## **Algorithm 1** Gradient Descent $(\theta_{\text{init}}, \eta, J, \nabla_{\theta} J, \epsilon)$

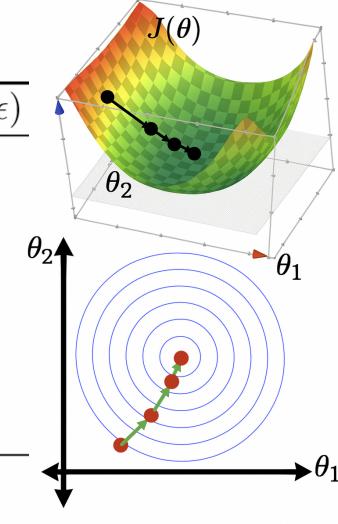
- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6:  $\mathbf{until} \left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



 $J(\theta)$ 

• What does this 2d vector / represent? anything in the psuedocode?

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$ 6: **until**  $\left| J(\theta^{(t)}) J(\theta^{(t-1)}) \right| < \epsilon$
- 7: **return**  $\theta^{(t)}$



- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat

4: 
$$t = t + 1$$

5: 
$$\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$$

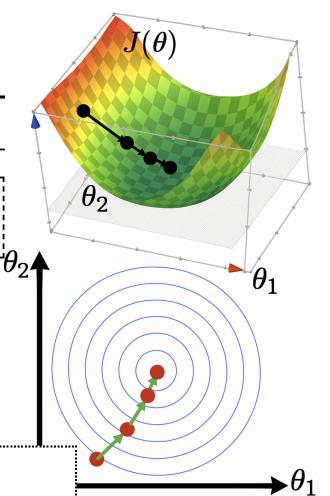
5: 
$$\theta^{(t)} = \theta^{(t-1)} - \eta \nabla_{\theta} J(\theta^{(t-1)})$$
6: 
$$\mathbf{until} \left| J(\theta^{(t)}) - J(\theta^{(t-1)}) \right| < \epsilon$$

7: **return**  $\theta^{(t)}$ 

objective improvement is nearly zero.

Other possible stopping criterion for line 6:

- Small parameter change:  $\|\theta^{(t)} \theta^{(t-1)}\| < \epsilon$ , or
- Small gradient norm:  $\|\nabla_{\theta}J(\theta^{(t-1)})\| < \epsilon$



## Outline

- Gradient descent algorithm (GD)
  - The gradient vector
  - GD algorithm
  - Gradient decent properties
- Stochastic gradient descent (SGD)
  - SGD algorithm and setup
  - SGD vs. GD

When minimizing a function, we aim for a global minimizer.

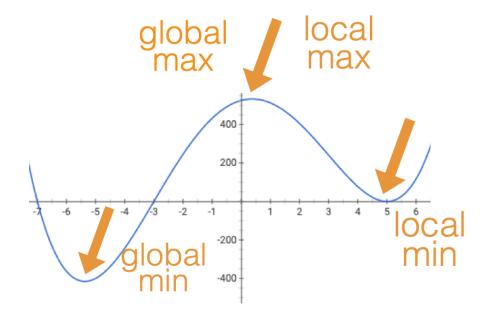
 $\Rightarrow$ 

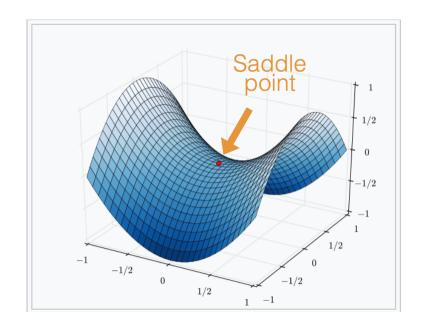
At a global minimizer

the gradient vector is zero

gradient descent can achieve this (to arbitrary precision)





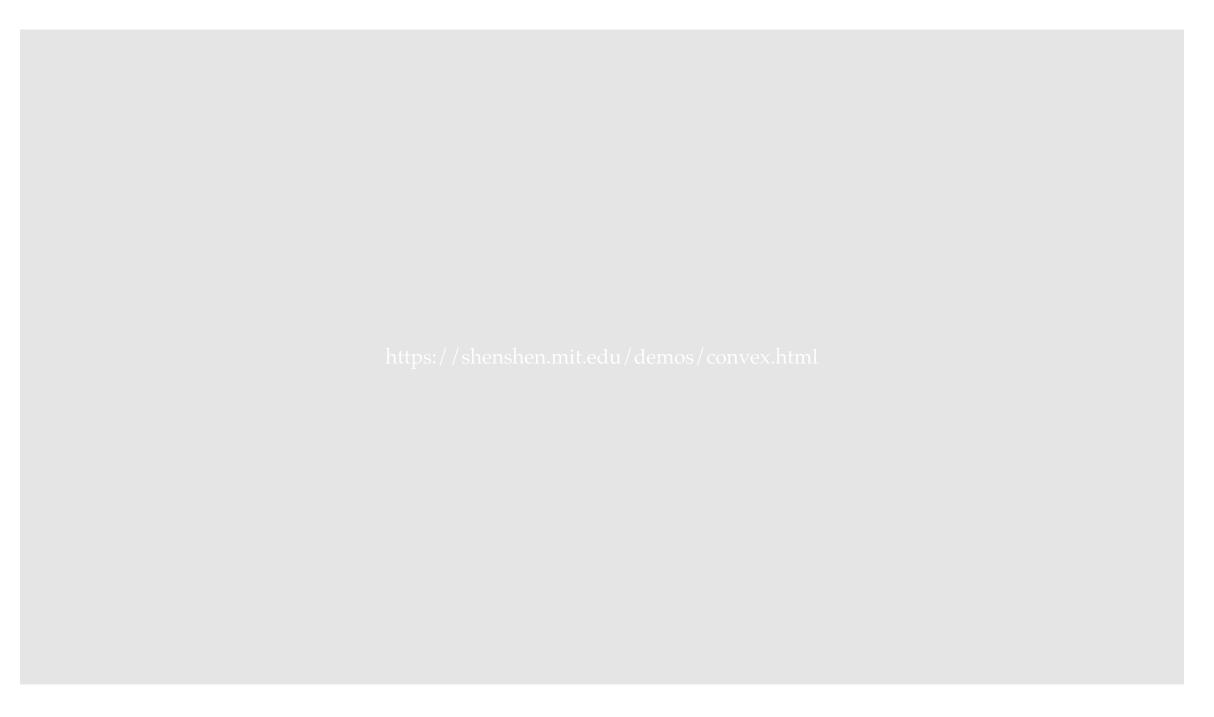


When minimizing a function, we aim for a global minimizer.

At a global minimizer  $\leftarrow$  { the gradient vector is zero the objective function is convex

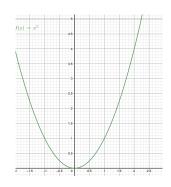
A function *f* is *convex* if any line segment connecting two points of the graph of *f* lies above or on the graph.

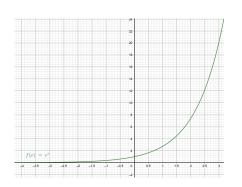
- f is concave if -f is convex.
- Convex functions are the largest well-understood class of functions where optimization theory guarantees convergence and efficiency

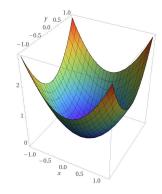


### Some examples

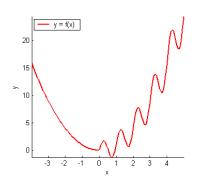
### Convex functions

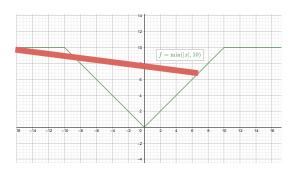


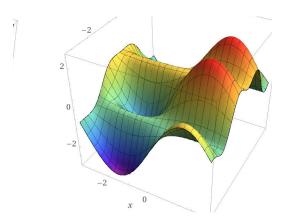




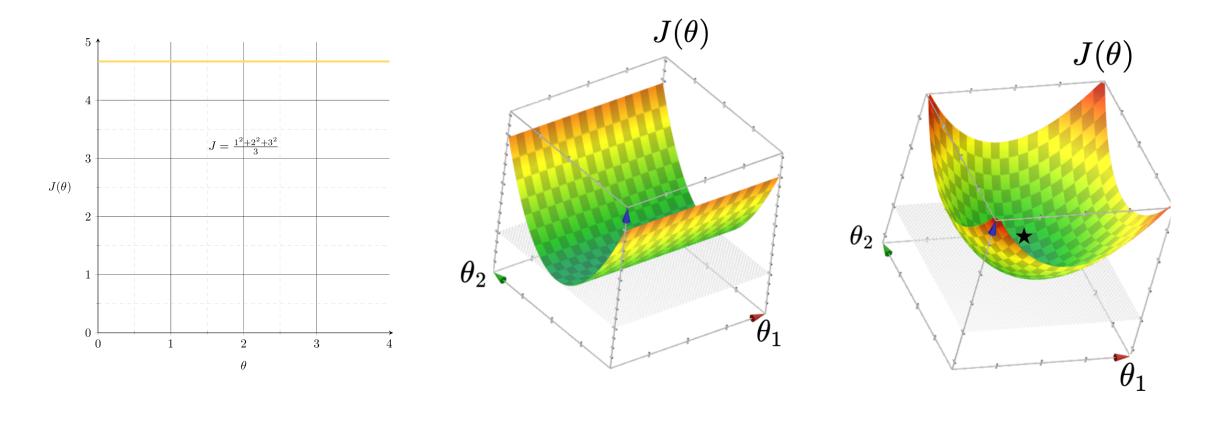
### Non-convex functions







• MSE:  $J(\theta) = \frac{1}{n}(X\theta - Y)^{\top}(X\theta - Y)$  is always convex



convexity is why we can claim the point whose gradient is zero is a global minimizer.

### case (c) training data set again

$$(x_1, x_2) = (6, 9), y = 9$$



$$(x_1,x_2)=(\underbrace{4}_{_{_{\!\!x_{_{\! 2}}}}},6),y=8$$

ttps://shenshen.mit.e

$$(x_1,x_2^{ ilde{\mathsf{y}}})=(2,3),y=7$$

• Ridge objective with  $\lambda > 0$  is always (strongly) convex

convexity is why we can claim the point whose gradient is zero is a global minimizer.

- Assumptions:
  - *f* is sufficiently "smooth"
  - f is convex
  - *f* has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of f.

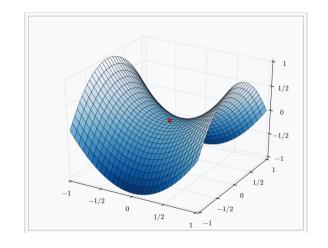
- Assumptions:
  - f is sufficiently "smooth"
  - f is convex
  - *f* has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small

if violated, may not have gradient, can't run gradient descent

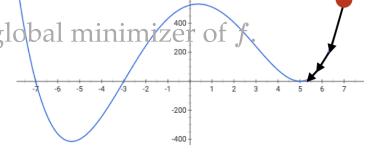
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of f.

- Assumptions:
  - f is sufficiently "smooth"
  - f is convex
  - *f* has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of

if violated, may get stuck at a saddle point

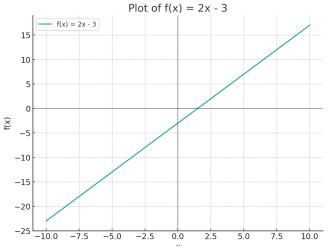


or a local minimum



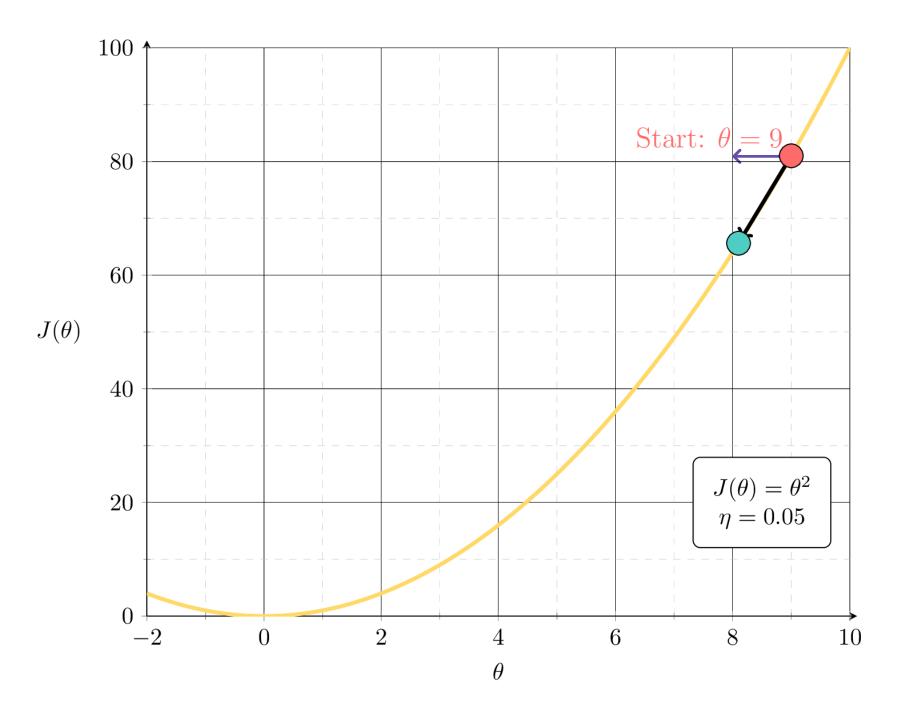
- Assumptions:
  - *f* is sufficiently "smooth"
  - f is convex
  - *f* has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of f.

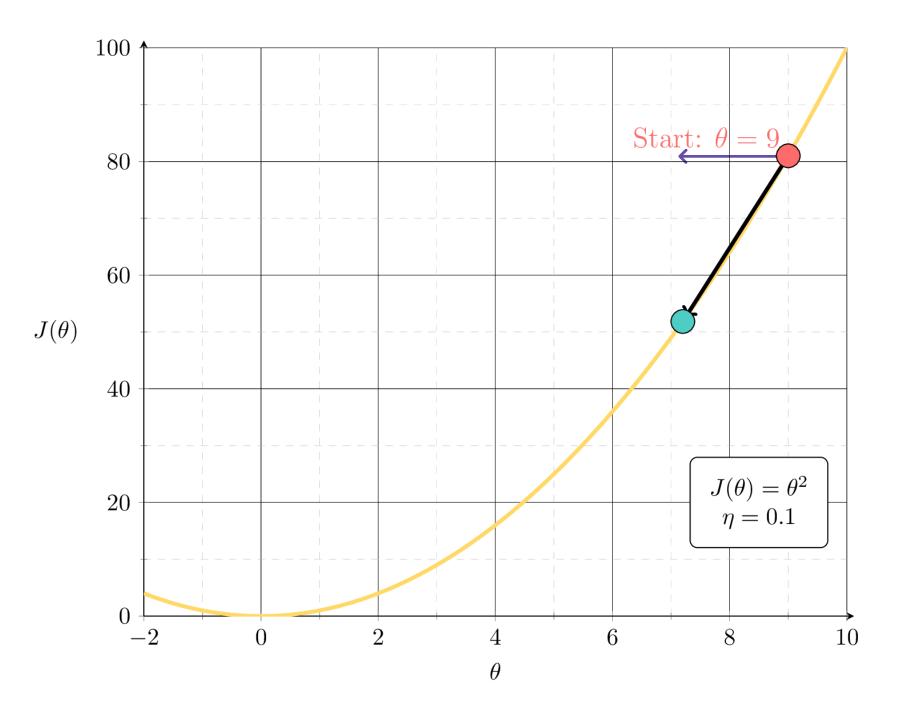
if violated:
may not terminate/no
minimum to converge to

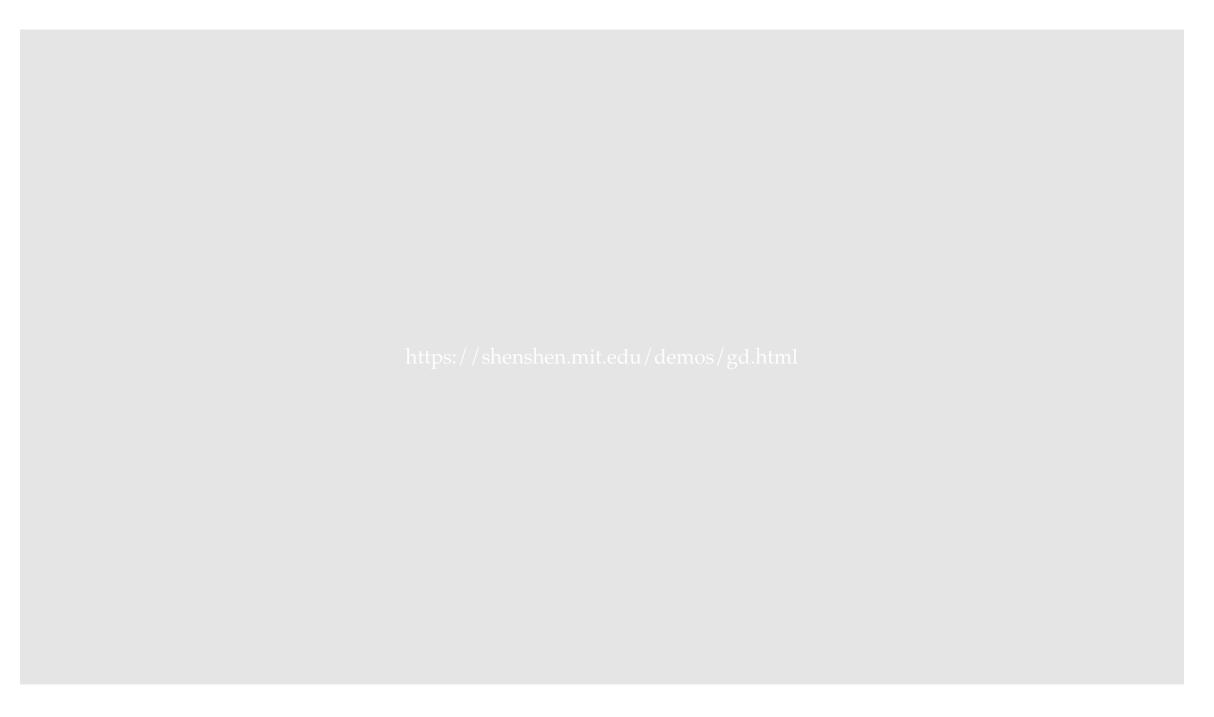


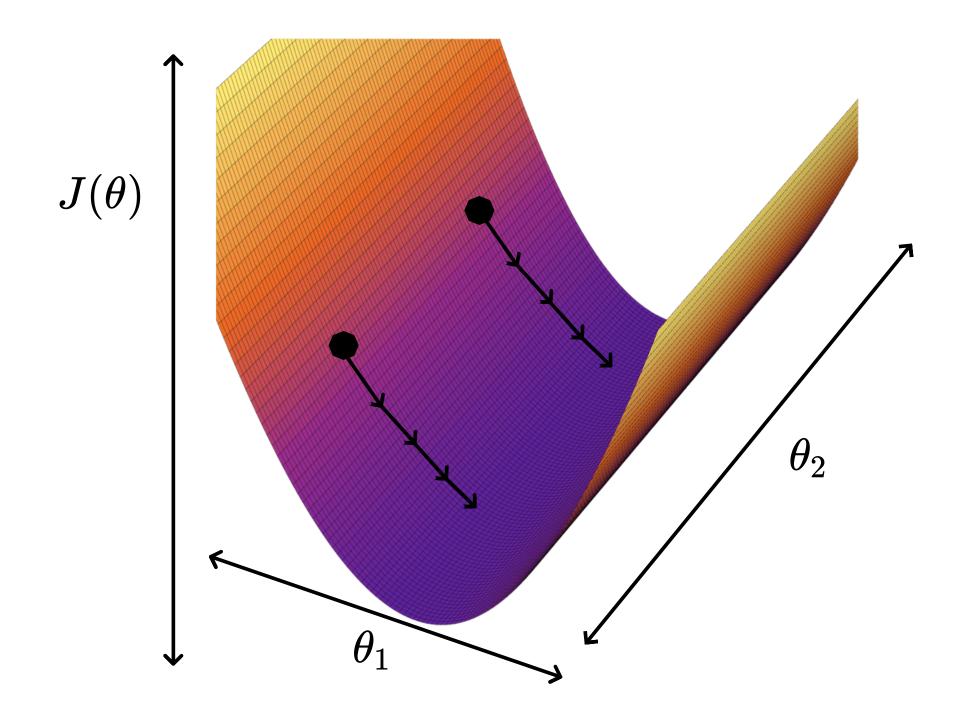
- Assumptions:
  - *f* is sufficiently "smooth"
  - f is convex
  - *f* has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small
- Conclusion:
  - Gradient descent converges arbitrarily close to a global minimizer of f.

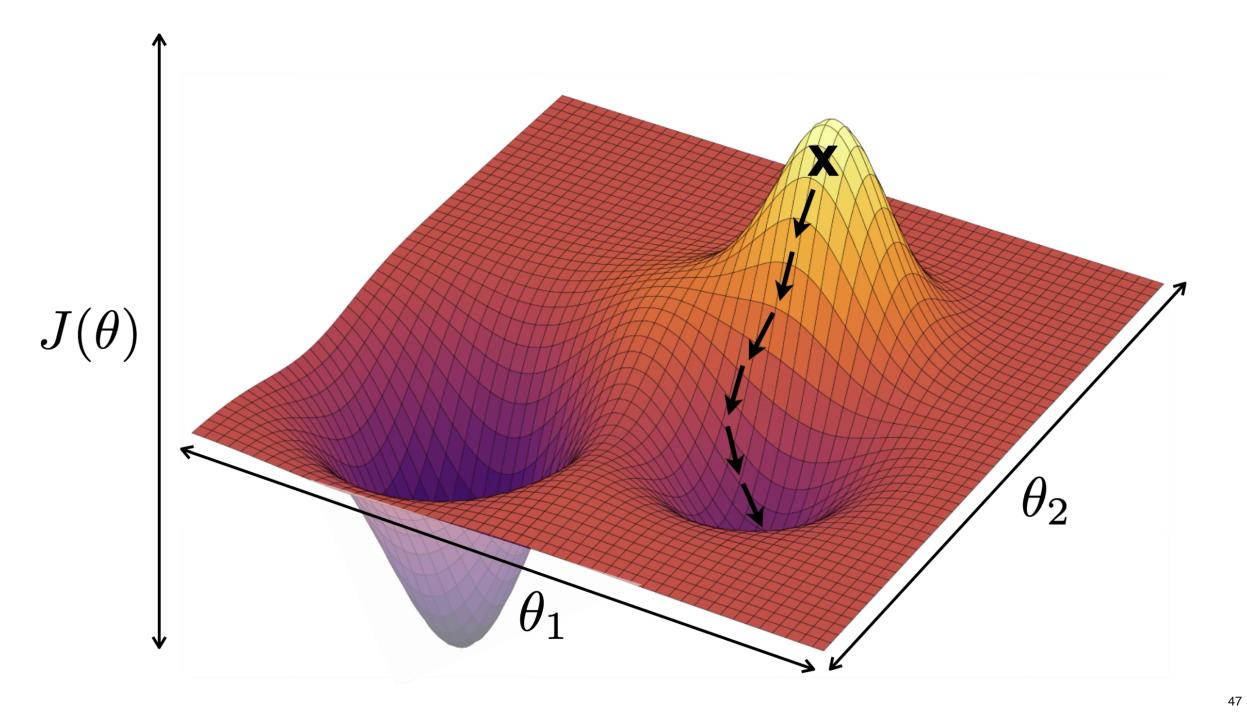
if violated: see demo on next slide, also lab/hw









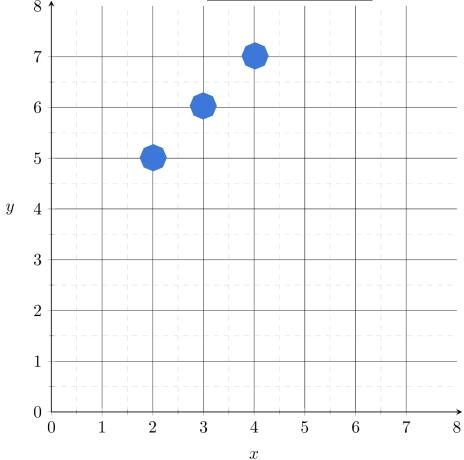


# Outline

- Gradient descent algorithm (GD)
  - The gradient vector
  - GD algorithm
  - Gradient decent properties
- Stochastic gradient descent (SGD)
  - SGD algorithm and setup
  - SGD vs. GD

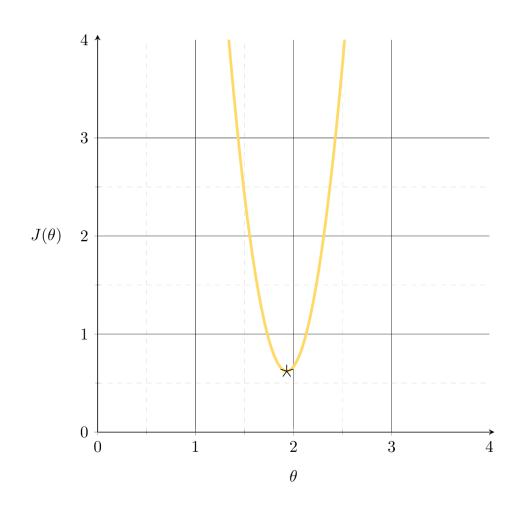
training data

	x	y
p1	2	5
p2	3	6
р3	4	7



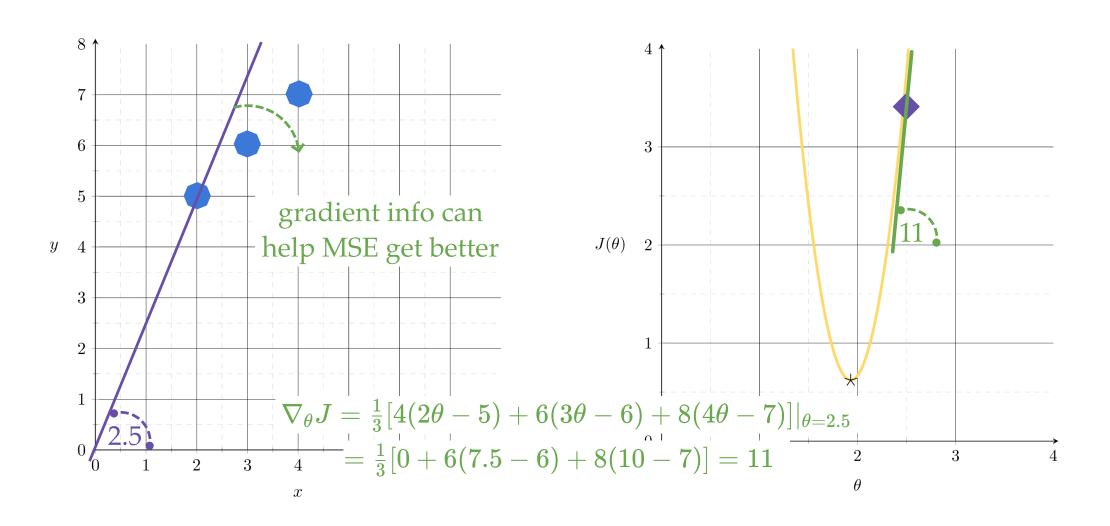
Fit a line (without offset) to the dataset, the MSE:

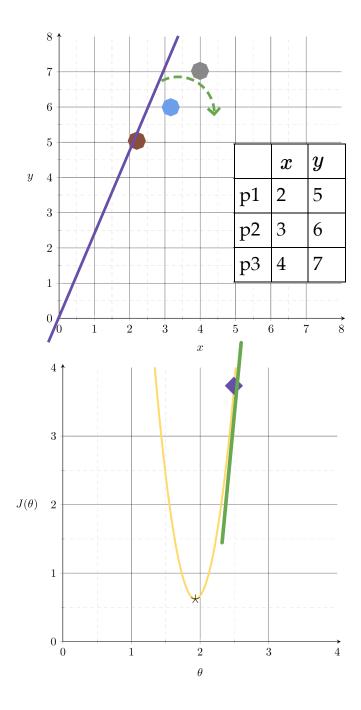
$$J( heta) = rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight]$$



#### Suppose we fit a line y = 2.5x

$$J( heta) = rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight]$$



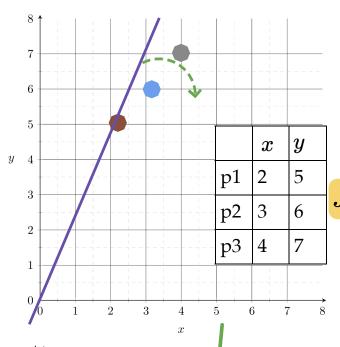


$$egin{align} egin{align} oldsymbol{J}( heta) &= rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight] \ &= rac{1}{3} \left[ egin{align} J_1 &+ J_2 &+ J_3 \end{array} 
ight] \end{aligned}$$

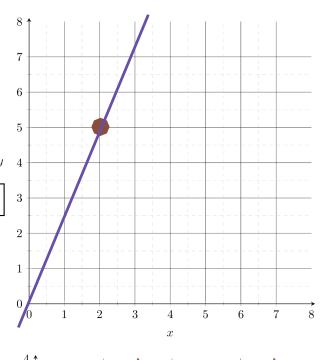
• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J = rac{1}{3}[4(2 heta-5)+6(3 heta-6)+8(4 heta-7)]$$

$$= \frac{1}{3} \left[ \begin{array}{ccc} \nabla_{\theta} J_1 & + & \nabla_{\theta} J_2 & + & \nabla_{\theta} J_3 \end{array} \right]$$

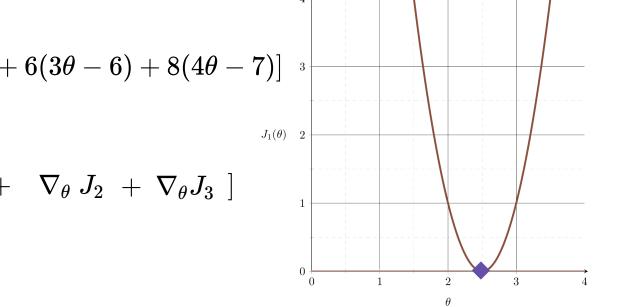


$$egin{align} oldsymbol{J( heta)} &= rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight] {}_3 \ &= rac{1}{3} \left[ egin{align} J_1 &+ & J_2 &+ J_3 
ight] \end{array}$$

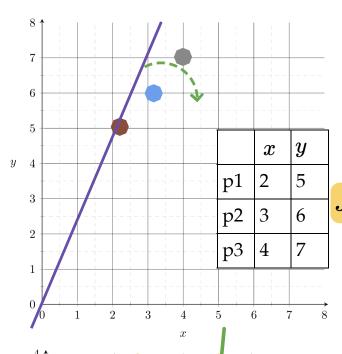


• and its gradient w.r.t.  $\theta$ :

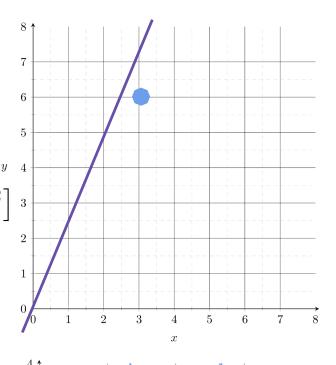
$$oldsymbol{
abla}_{ heta} J = rac{1}{3} [4(2 heta - 5) + 6(3 heta - 6) + 8(4 heta - 7)]^{-3}$$



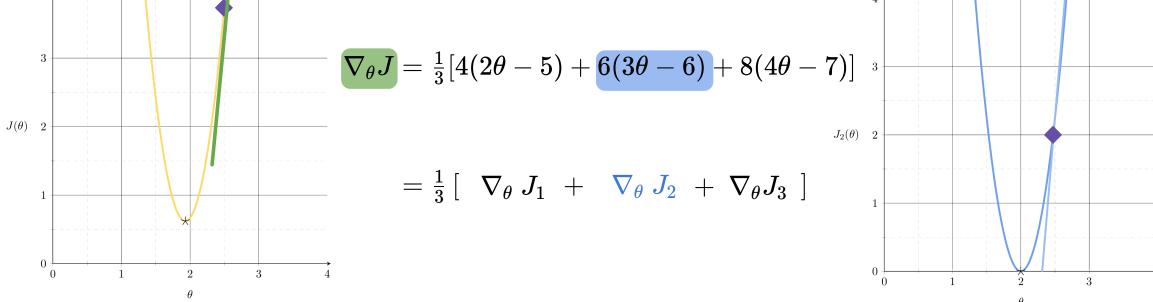
 $J(\theta) = 2$ 

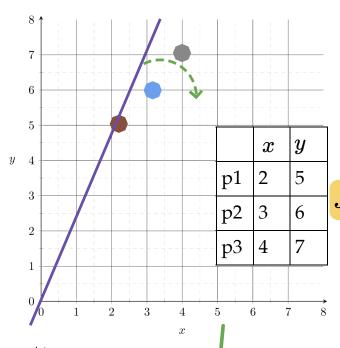


$$egin{align} m{J( heta)} &= rac{1}{3} \left[ (2 heta - 5)^2 + m{(3 heta - 6)^2} + (4 heta - 7)^2 
ight]{}^3 \ &= rac{1}{3} \left[ m{J_1} + m{J_2} + m{J_3} 
ight] \end{array}$$

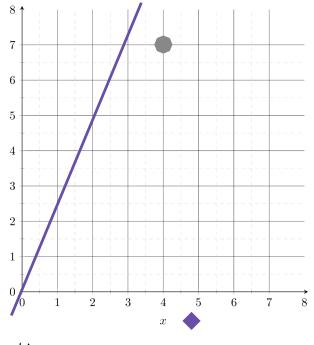


• and its gradient w.r.t.  $\theta$ :



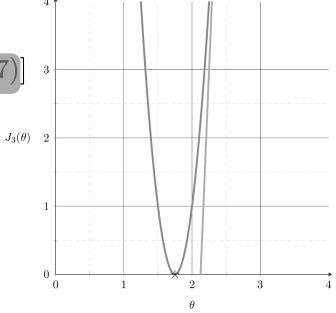


$$egin{align} egin{align} oldsymbol{J( heta)} &= rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight] \ &= rac{1}{3} \left[ egin{align} J_1 &+ & J_2 &+ J_3 
ight] \end{aligned}$$



• and its gradient w.r.t.  $\theta$ :

$$oxed{
abla_{ heta} J} = rac{1}{3}[4(2 heta - 5) + 6(3 heta - 6) + 8(4 heta - 7)]$$



 $J(\theta) = 2$ 

### Gradient of an ML objective

Using our example data set,

• the MSE of a linear hypothesis:

$$J( heta) = rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight]$$

• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J = rac{1}{3}[4(2 heta-5)+6(3 heta-6)+8(4 heta-7)]$$

Using any dataset,

• the MSE of a linear hypothesis:

$$J( heta) = rac{1}{n} \sum_{i=1}^n \left( heta^ op x^{(i)} - y^{(i)} 
ight)^2$$

• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J( heta) = rac{1}{n}\sum_{i=1}^n 2\left( heta^ op x^{(i)} - y^{(i)}
ight)x^{(i)}$$

### Gradient of an ML objective

• the MSE of a linear hypothesis:

$$J( heta) = rac{1}{n} \sum_{i=1}^n \left( heta^ op x^{(i)} - y^{(i)} 
ight)^2$$

• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J( heta) = rac{1}{n}\sum_{i=1}^n 2\left( heta^ op x^{(i)} - y^{(i)}
ight)x^{(i)}$$

In general,

• An ML objective function is a finite sum

$$J( heta) = rac{1}{n} \sum_{i=1}^n J_i( heta)$$

• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J( heta) = rac{1}{n}\sum_{i=1}^n 2\left( heta^ op x^{(i)} - y^{(i)}
ight)x^{(i)} egin{aligned} 
abla_{ heta}J( heta) &= 
abla(rac{1}{n}\sum_{i=1}^n J_i( heta)) &= rac{1}{n}\sum_{i=1}^n 
abla_{ heta}J_i( heta) \end{aligned}$$



(gradient of the sum) = (sum of the gradient)

### Gradient of an ML objective

#### In general,

• An ML objective function is a finite sum

$$J( heta) = rac{1}{n} \sum_{i=1}^n J_i( heta)$$

loss incurred on a single  $i^{
m th}$  data point

• and its gradient w.r.t.  $\theta$ :

$$abla_{ heta}J( heta) = 
abla_{ heta}(rac{1}{n}\sum_{i=1}^n J_i( heta)) = rac{1}{n}\sum_{i=1}^n 
abla_{ heta_i}( heta)$$

need to add n of these,  $\operatorname{each} \nabla_{\theta} J_i(\theta) \in \mathbb{R}^d$ 

gradient info from a single  $i^{
m th}$  data point's loss

Costly in practice!

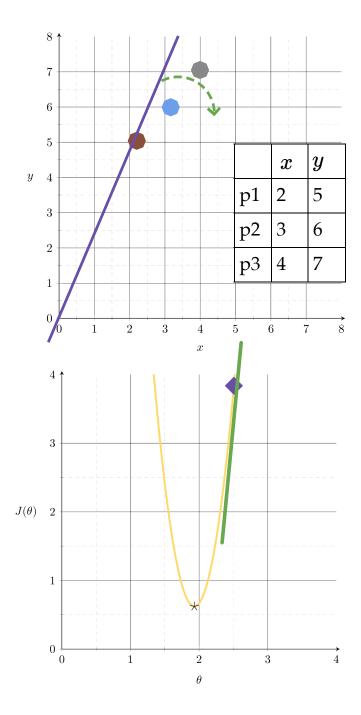
### **Algorithm 1** Gradient Descent $(\theta_{\text{init}}, \eta, J, \nabla_{\theta} J, \epsilon)$

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J(\theta^{(t-1)})$
- 6: until  $|J(\theta^{(t)}) J(\overline{\theta^{(t-1)}})| < \epsilon$
- 7: **return**  $\theta^{(t)}$

#### **Algorithm 2** Stochastic Gradient Descent $(\theta_{\text{init}}, \eta, J, \nabla_{\theta} J, \epsilon)$

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t=0
- 3: repeat
- 4: t = t + 1
- 5:  $i = \text{randomly sample from } \{1, \dots, n\}$ 6:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J_i(\theta^{(t-1)})$
- 7: until  $|J(\theta^{(t)}) J(\overline{\theta^{(t-1)}})| < \epsilon$
- 8: return  $\theta^{(t)}$

$$oxed{
abla_{ heta} J( heta)} = rac{1}{n} \sum_{i=1}^n 
abla_{ heta} J_i( heta) lpha oxedown_{ heta} J_i( heta)$$

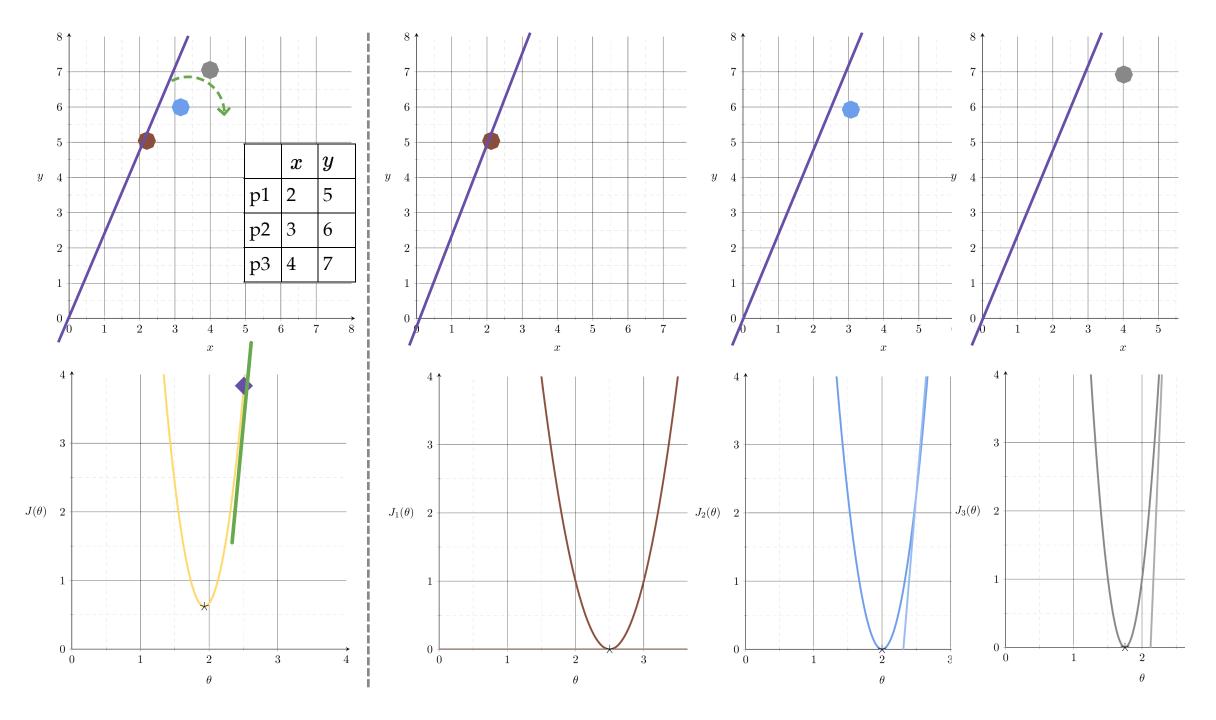


$$egin{align} egin{align} oldsymbol{J( heta)} &= rac{1}{3} \left[ (2 heta - 5)^2 + (3 heta - 6)^2 + (4 heta - 7)^2 
ight] \ &= rac{1}{3} \left[ egin{align} J_1 &+ & J_2 &+ J_3 \end{array} 
ight] \end{aligned}$$

• and its gradient w.r.t.  $\theta$ :

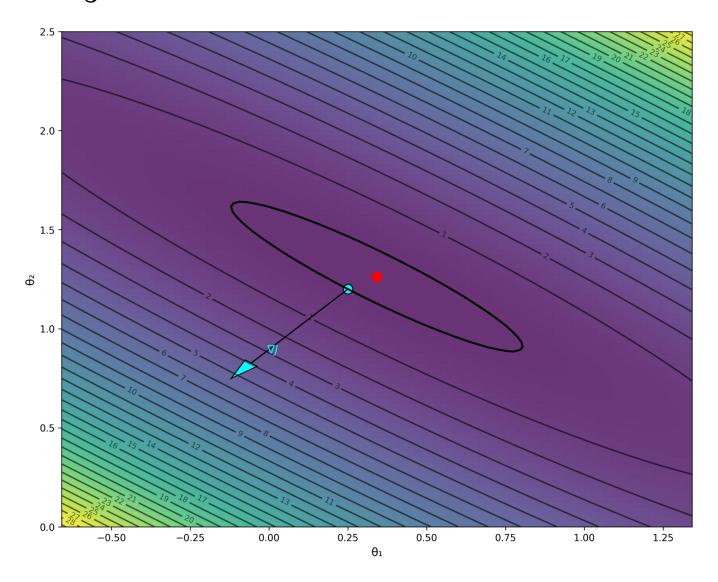
$$abla_{ heta} J = rac{1}{3} [4(2 heta - 5) + 6(3 heta - 6) + 8(4 heta - 7)]$$

$$= \frac{1}{3} \left[ \nabla_{\theta} J_1 + \nabla_{\theta} J_2 + \nabla_{\theta} J_3 \right]$$



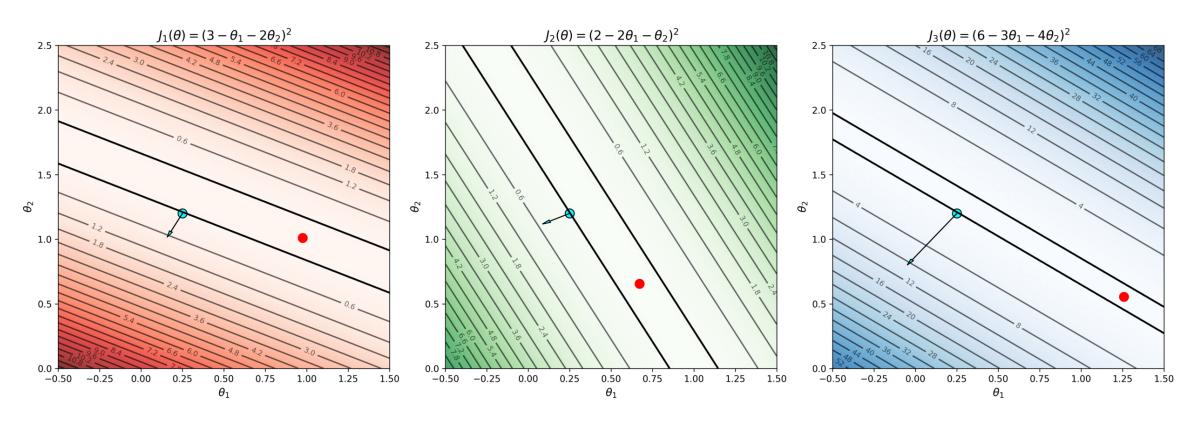
	<b>x</b> 1	<b>x2</b>	y
<b>p</b> 1	1	2	3
p2	2	1	2
р3	3	4	6

$$J( heta) = rac{1}{3} \left[ (3 - heta_1 - 2 heta_2)^2 + (2 - 2 heta_1 - heta_2)^2 + (6 - 3 heta_1 - 4 heta_2)^2 
ight]$$



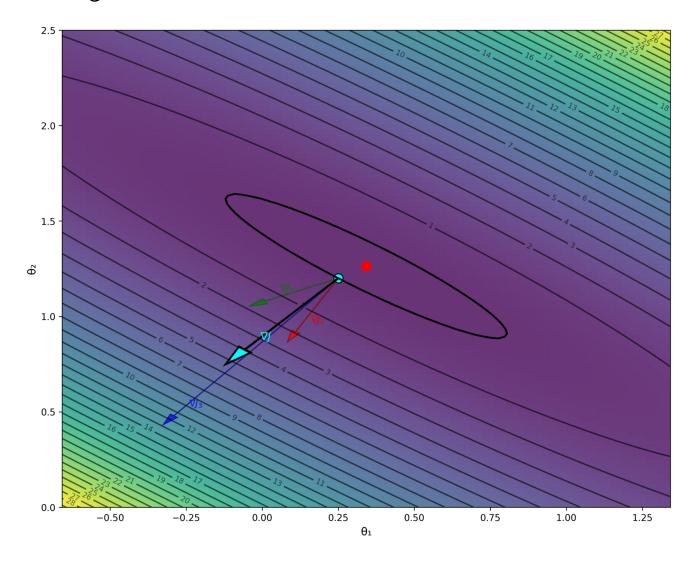
	<b>x</b> 1	<b>x2</b>	y
p1	1	2	3
p2	2	1	2
р3	3	4	6

$$J( heta) = rac{1}{3} \left[ (3 - heta_1 - 2 heta_2)^2 + (2 - 2 heta_1 - heta_2)^2 + (6 - 3 heta_1 - 4 heta_2)^2 
ight]$$



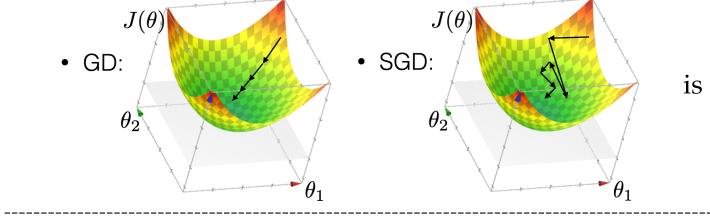
	<b>x</b> 1	<b>x2</b>	y
<b>p</b> 1	1	2	3
p2	2	1	2
р3	3	4	6

$$J( heta) = rac{1}{3} \left[ (3 - heta_1 - 2 heta_2)^2 + (2 - 2 heta_1 - heta_2)^2 + (6 - 3 heta_1 - 4 heta_2)^2 
ight]$$

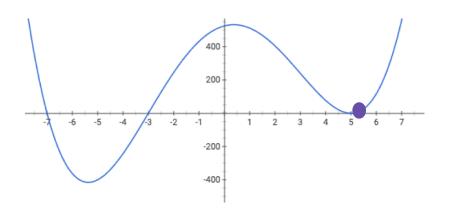


### Compared with GD, SGD:

$$abla_{ heta}J( heta)=rac{1}{n}\sum_{i=1}^{n}
abla_{ heta}J_{i}( heta)\ pprox 
abla_{ heta}J_{i}( heta) \qquad ext{is more efficient}$$



is much more "random"



may get us out of a local min

## Stochastic gradient descent performance

- Assumptions:
  - *f* is sufficiently "smooth"
  - f is convex
  - f has at least one global minimum
  - Run gradient descent for sufficient iterations
  - $\eta$  is sufficiently small and satisfies additional "scheduling" condition
- Conclusion:

$$\sum_{t=1}^{\infty} \eta(t) = \infty$$
 and  $\sum_{t=1}^{\infty} \eta(t)^2 < \infty$ 

• Stochastic gradient descent converges arbitrarily close to a global minimum of *f* with probability 1.

### **Algorithm 3** Mini-batch Gradient Descent( $\theta_{\text{init}}, \eta, b, J, \nabla_{\theta} J, \epsilon$ )

- 1: Initialize  $\theta^{(0)} = \theta_{\text{init}}$
- 2: Initialize t = 0

batch size

- 3: repeat
- 4: t = t + 1
- 5:  $B = \text{random mini-batch of size } b \text{ from } \{1, \dots, n\}$
- 6:  $\theta^{(t)} = \theta^{(t-1)} \eta \nabla_{\theta} J_B(\theta^{(t-1)})$
- 7: **until**  $|J(\theta^{(t)}) J(\overline{\theta^{(t-1)}})| < \epsilon$
- 8: **return**  $\theta^{(t)}$

**SGD** 

 $abla_{ heta}J_{i}( heta)$ 

mini-batch GD

$$rac{1}{b}\sum_{i=1}^b 
abla_ heta J_i( heta)$$

GD

$$rac{1}{n}\sum_{i=1}^n 
abla_{ heta} J_i( heta) = 
abla_{ heta} J( heta)$$

more accurate gradient estimate

stronger theoretical guarantee

### Summary

- Most ML methods can be formulated as optimization problems. We won't always be able to solve optimization problems analytically (in closed-form) nor efficiently.
- We can still use numerical algorithms to good effect. Lots of sophisticated ones available. Gradient descent is one of the simplest.
- The GD algorithm, iterative algorithm, keeps applying the parameter update rule.
- Under appropriate conditions (most notably, when objective function is convex, and when learning rate is small enough), GD can guarantee convergence to a global minimum.
- SGD is approximated GD, it uses a single data point to approximate the entire data set, it's more efficient, more random, and less guarantees.
- mini-batch GD is a middle ground between GD and SGD.

https://docs.google.com/forms/d/e/1FAIpQLScj9i83AI8TuhWDZXSjiWzX6gZpnPugjGsH-i3RdrBCtF-opg/viewform? embedded=true

We'd love to hear your thoughts.

Thanks!