

LECTURE 14

# Gradient Descent

An optimization method to numerically minimize loss functions.

**Data Science, Spring 2024 @ Knowledge Stream**

Sana Jabbar

# Optimization: Where Are We?

---

Lecture 14

- **Optimization: where are we?**
- Minimizing an arbitrary 1D function
- Gradient descent on a 1D model
- Gradient descent on high-dimensional models
- Batch, mini-batch, and stochastic gradient descent

Takeaways from the past few lectures:

- Choose a model
- Choose a loss function
- Optimize parameters – choose the values of  $\theta$  that minimize the model's loss

How have we optimized?

1. Use calculus to solve for  $\theta$

Take derivatives, set equal to 0, solve.

1. Use a geometric argument

Using orthogonality, derive the OLS solution  $\hat{\theta} = (\mathbb{X}^T \mathbb{X})^{-1} \mathbb{X}^T \mathbb{Y}$

## Where We're Going

---

We made some big assumptions with the calculus-based and geometric techniques.

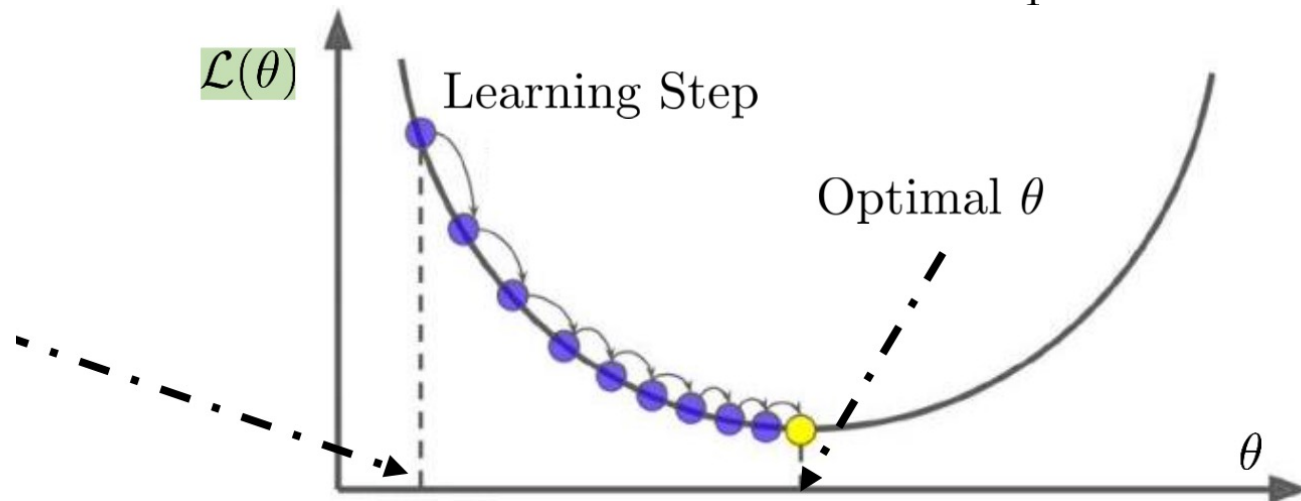
- Calculus: assumed that the loss function was differentiable at all points and that the algebra was manageable
- Geometric: OLS *only* applies when using a linear model with MSE loss

To design more complex models with different loss functions, we need a new optimization technique: **gradient descent**.

Big Idea: use an algorithm instead of solving for an exact answer

## Gradient Descent

- Gradient descent is an iterative algorithm in nature:
- Initially, chose the coefficients to be something reasonable (e.g., all zeros).
- Iteratively update the coefficients in the direction of steepest descent until convergence.
- Ensures that the new coefficients are better than the previous coefficients.



# Minimizing an Arbitrary 1D Function

---

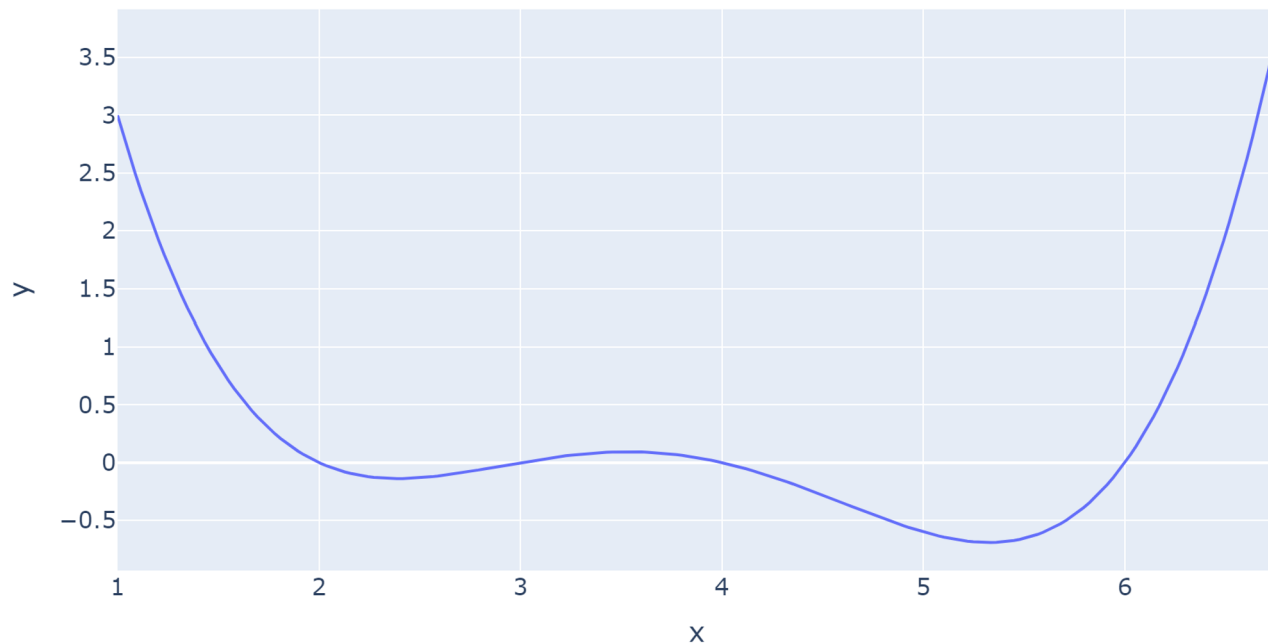
Lecture 14

- Optimization: where are we?
- **Minimizing an arbitrary 1D function**
- Gradient descent on a 1D model
- Gradient descent on high-dimensional models
- Batch, mini-batch, and stochastic gradient descent

## An Arbitrary Function

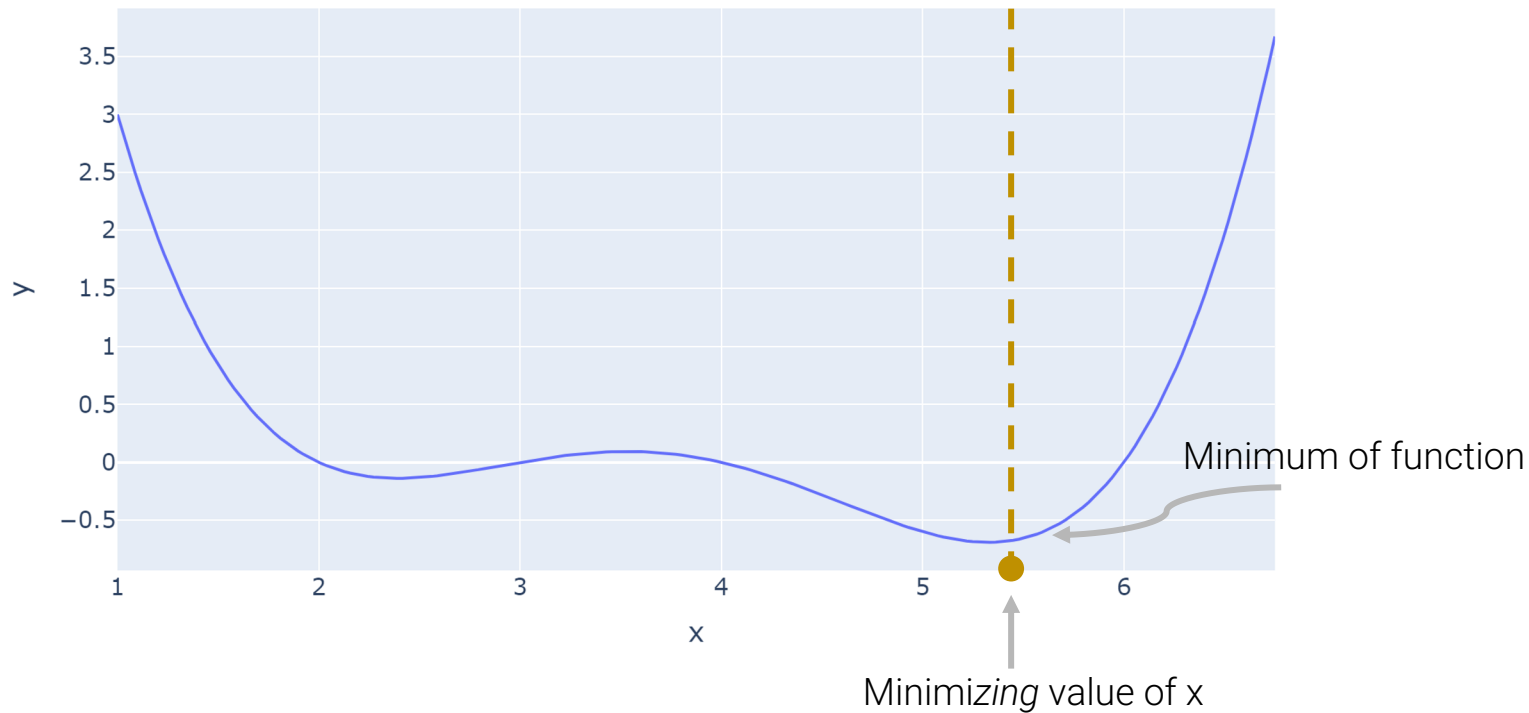
```
def arbitrary(x):  
    return (x**4 - 15*x**3 + 80*x**2 - 180*x + 144)/10
```

```
x = np.linspace(1, 6.75, 200)  
fig = px.line(y = arbitrary(x), x = x)
```



## An Arbitrary Function

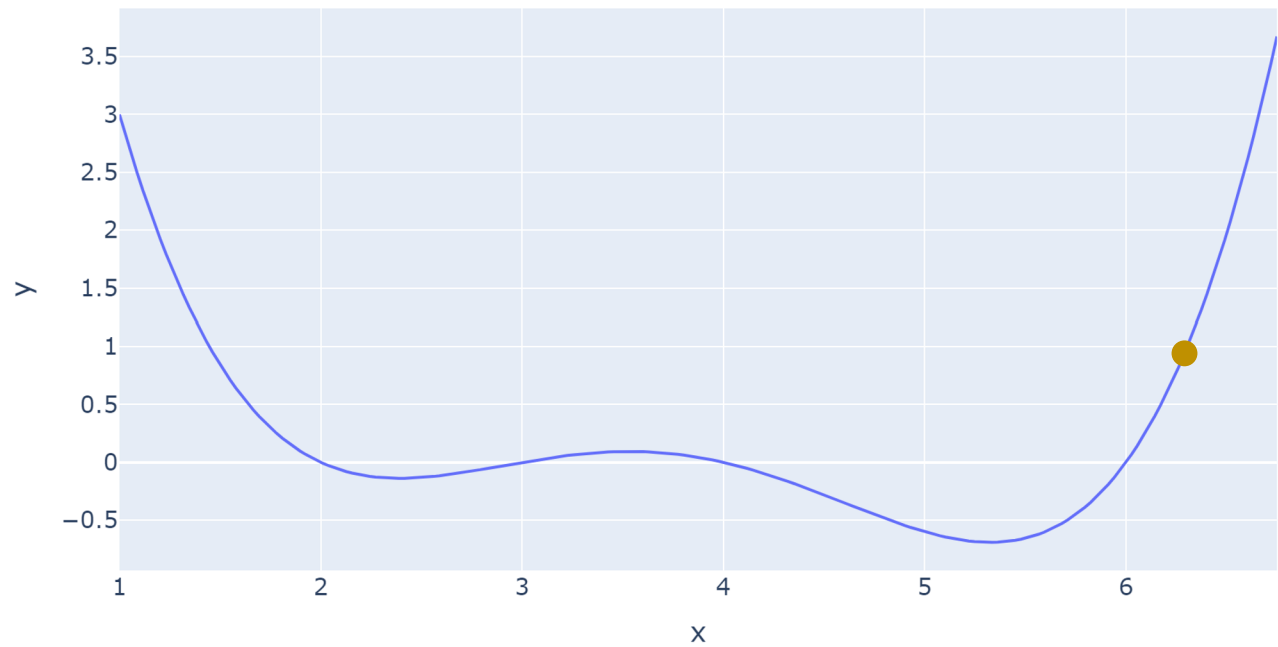
Our goal is to find the value of  $x$  that minimizes our function.





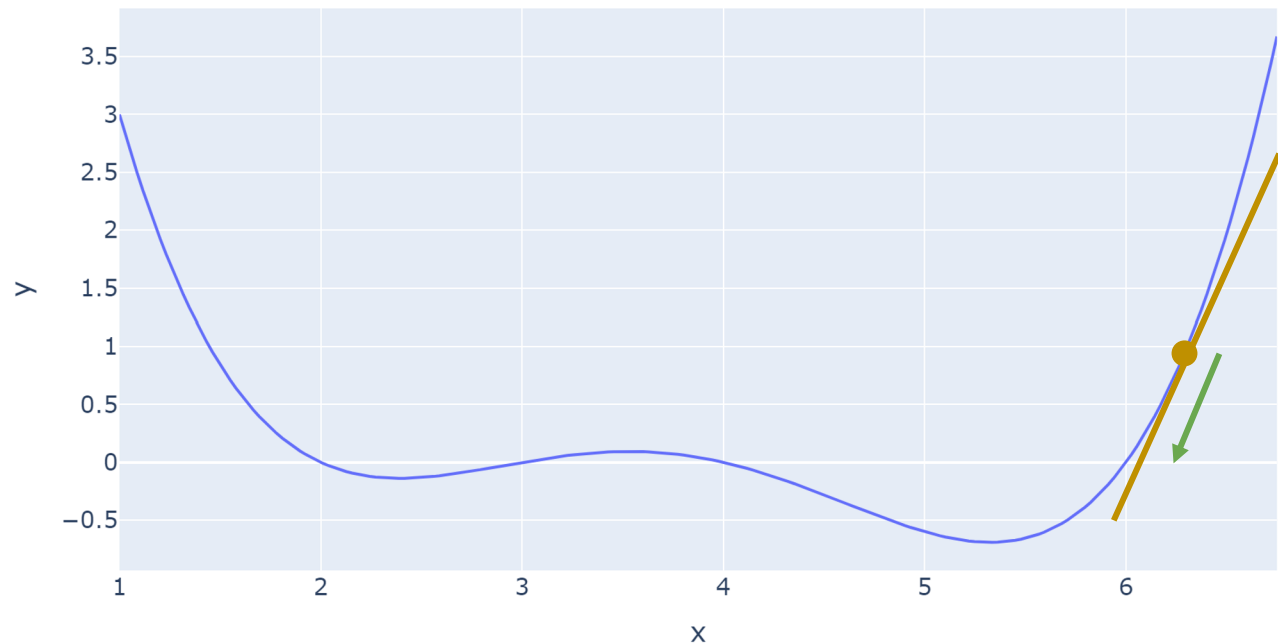
# Finding the Minimum

We could start with a random guess.



# Finding the Minimum

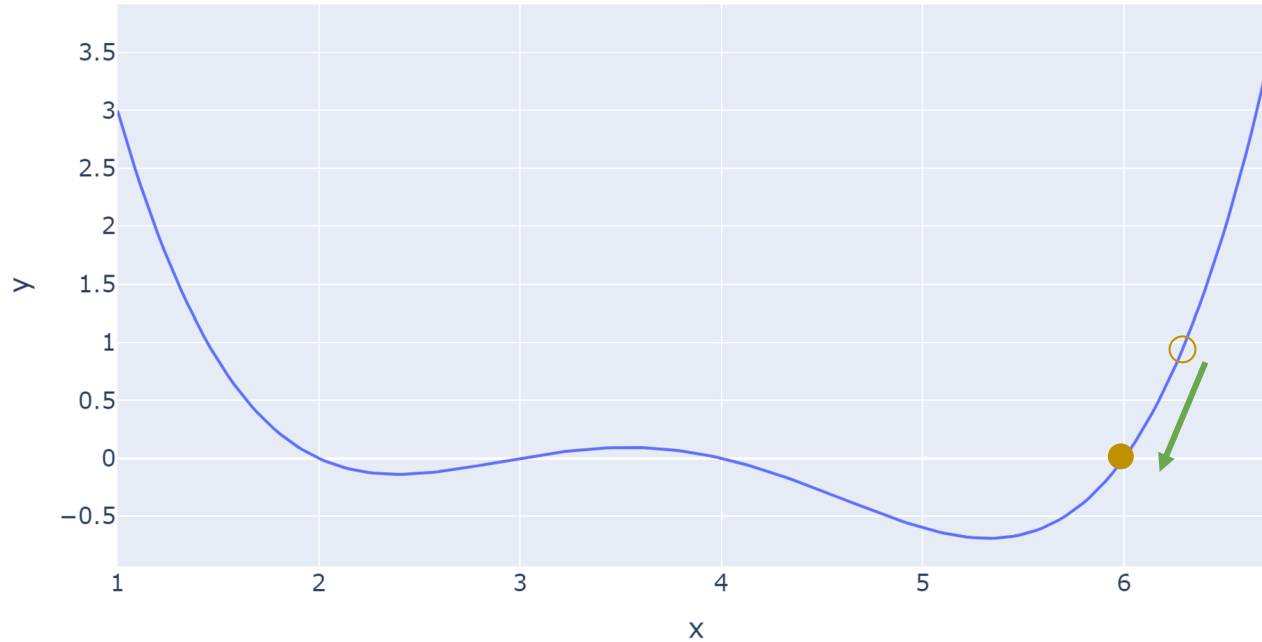
Where do we go next? We “step” downhill.



Follow the slope of the line down to the minimum.

## Finding the Minimum

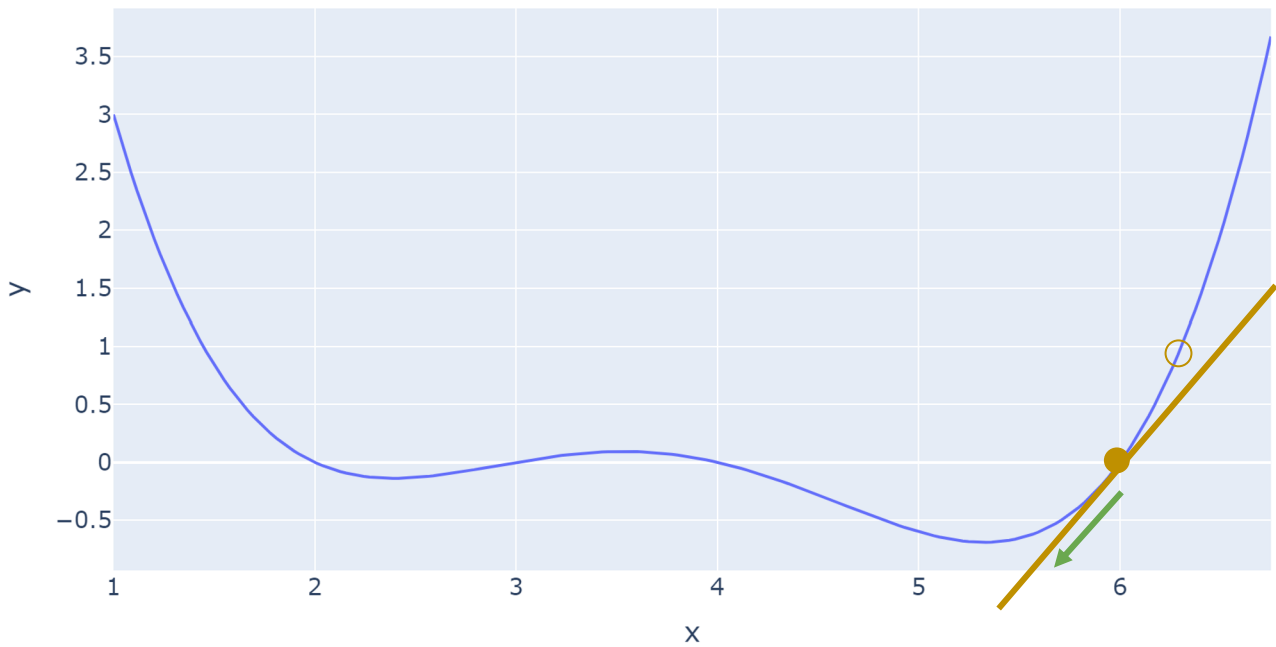
We arrive closer to the minimum.



Positive slope  $\rightarrow$  step to the left

# Finding the Minimum

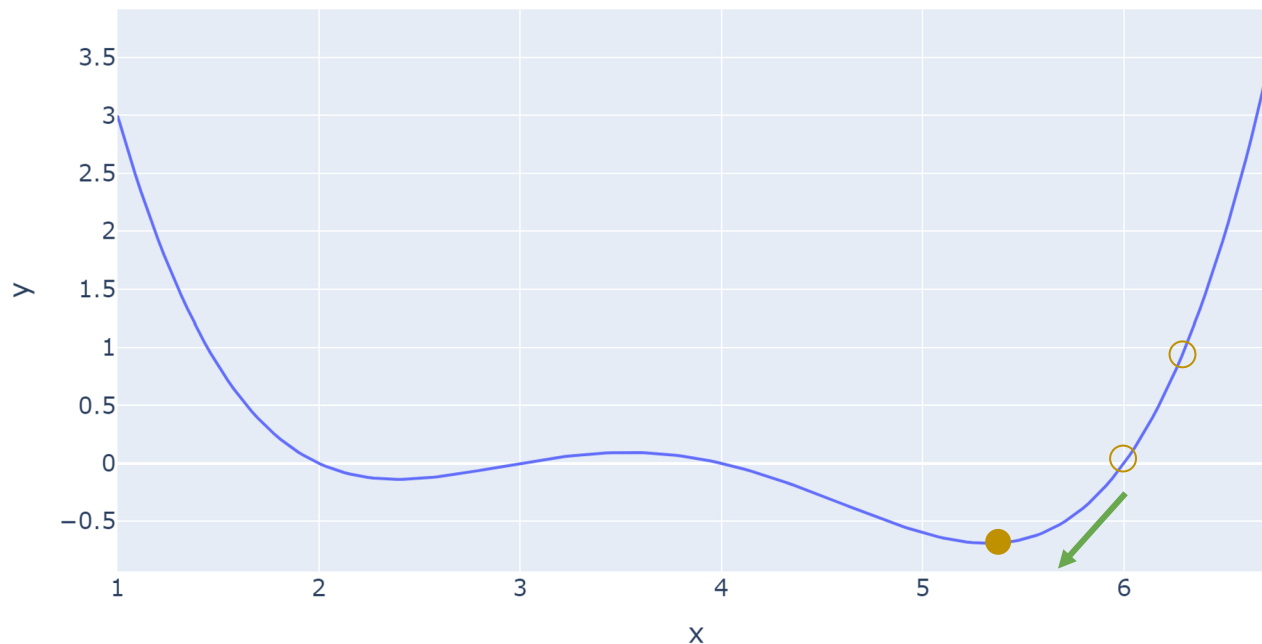
Do this again: follow the slope downwards towards the minimum



Positive slope → step to the left

## Finding the Minimum

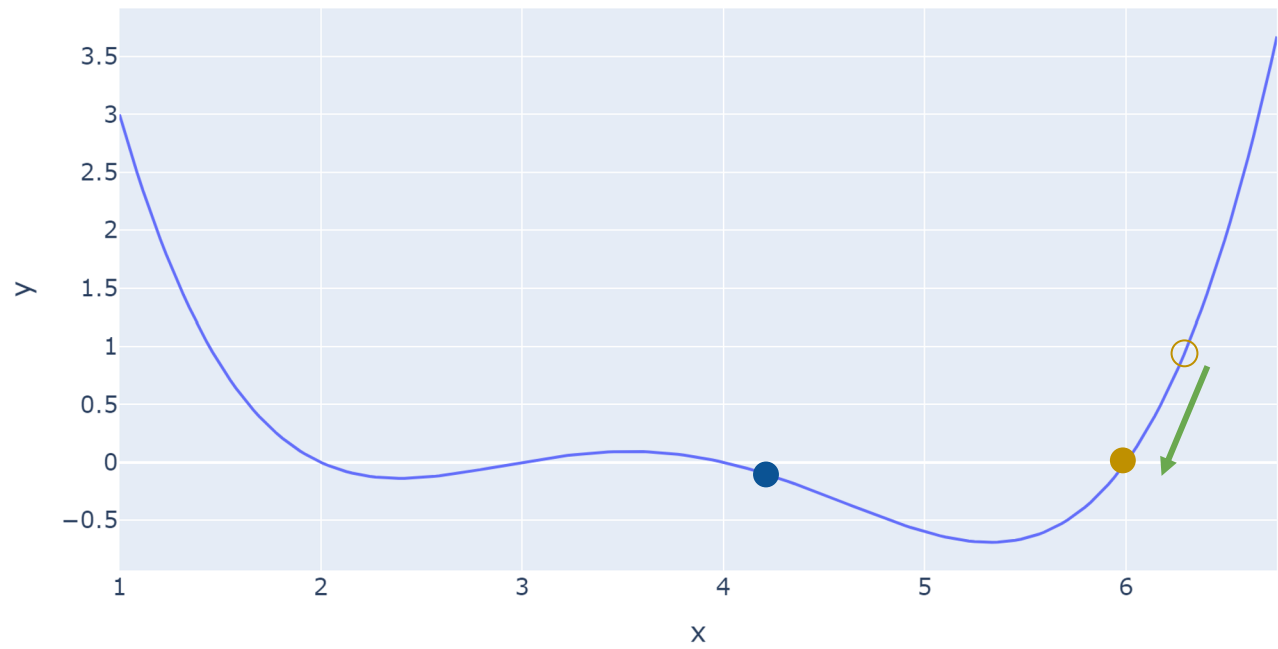
Do this again: follow the slope downwards towards the minimum



Positive slope  $\rightarrow$  step to the left

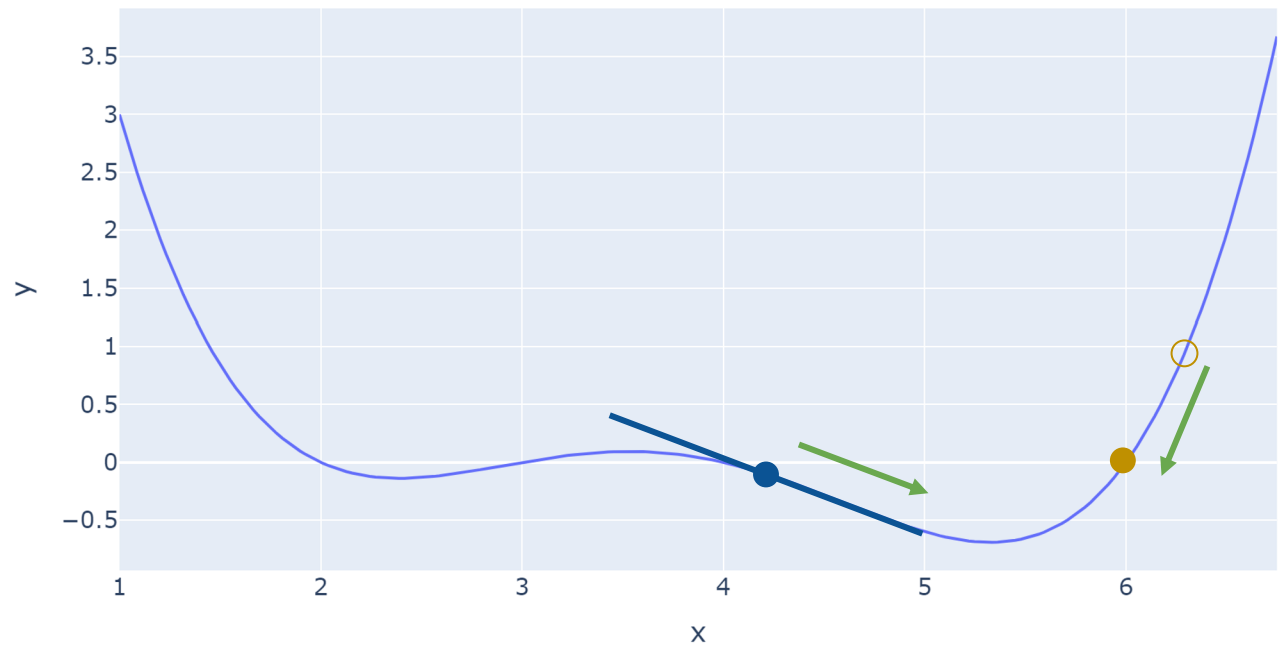
# Finding the Minimum

What if we had started elsewhere?



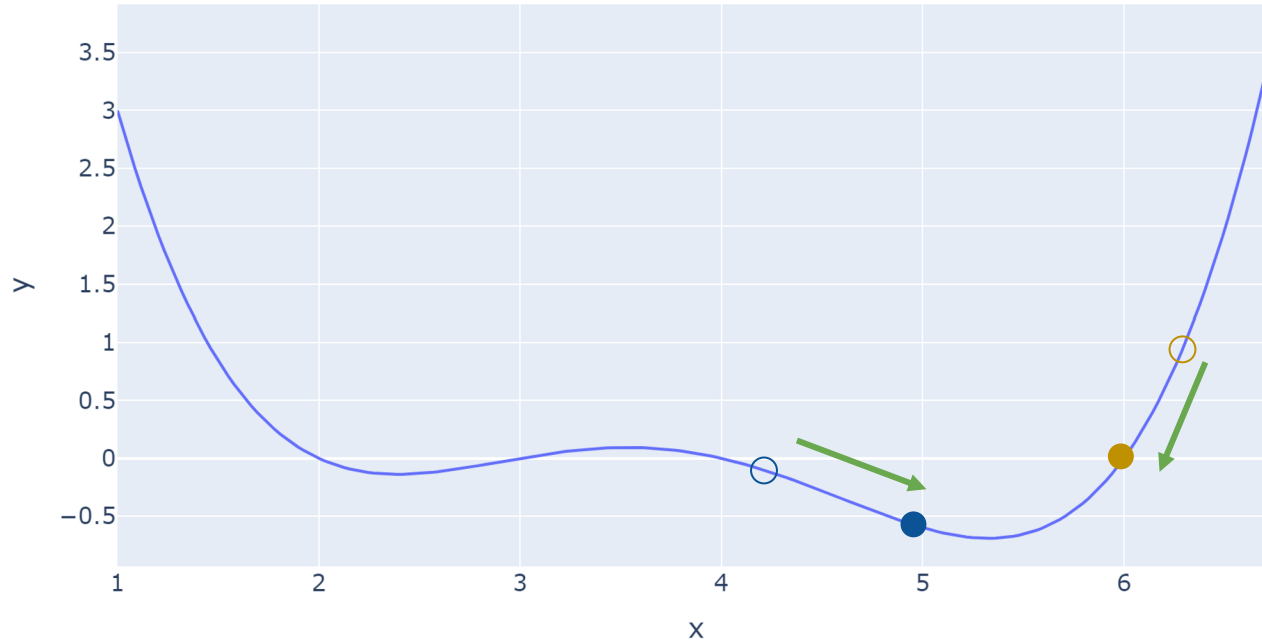
# Finding the Minimum

What if we had started elsewhere?



## Finding the Minimum

What if we had started elsewhere?

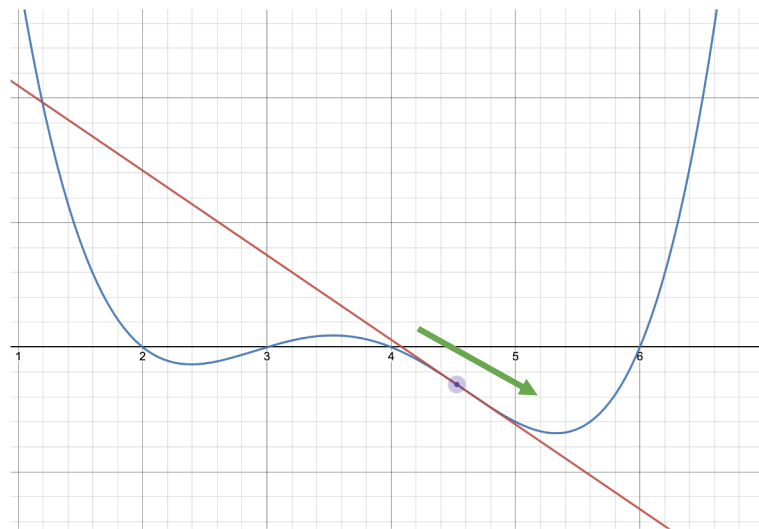


Negative slope  $\rightarrow$  step to the right

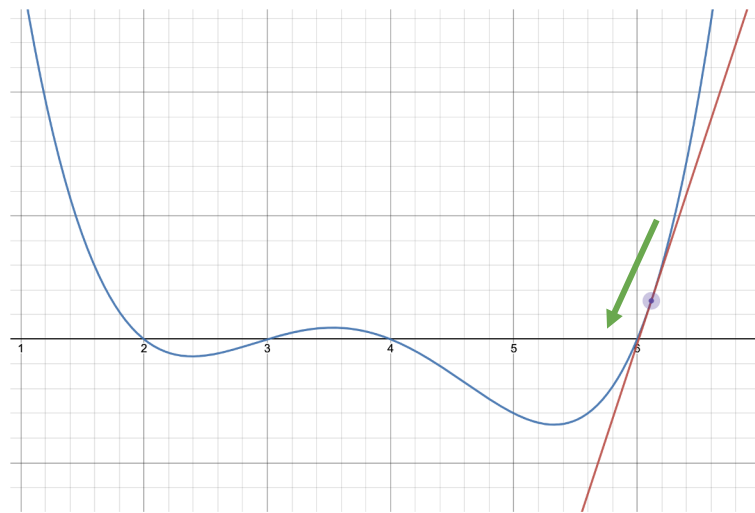


# Slopes Tell Us Where to Go

Negative slope  $\rightarrow$  step to the right  
Move in the *positive* direction



Positive slope  $\rightarrow$  step to the left  
Move in the *negative* direction



The derivative of the function at each point tells us the direction of our next guess.

## Slopes Tell Us Where to Go

---

The derivative of the function at each point tells us the direction of our next guess.

Negative slope  $\rightarrow$  step to the right

Move  $x$  in the *positive* direction

Positive slope  $\rightarrow$  step to the left

Move  $x$  in the *negative* direction

How can we use this?

## Slopes Tell Us Where to Go

---

The derivative of the function at each point tells us the direction of our next guess.

Negative slope  $\rightarrow$  step to the right

Move  $x$  in the *positive* direction

Positive slope  $\rightarrow$  step to the left

Move  $x$  in the *negative* direction

Our first attempt at making an algorithm: step in the opposite direction to the slope

$$x^{(t+1)} = x^{(t)} - \frac{d}{dx} f(x^{(t)})$$

Our next guess for the  
minimizing  $x$

...starts at the previous  
guess

...and moves opposite to the  
slope

## Introducing a Learning Rate

---

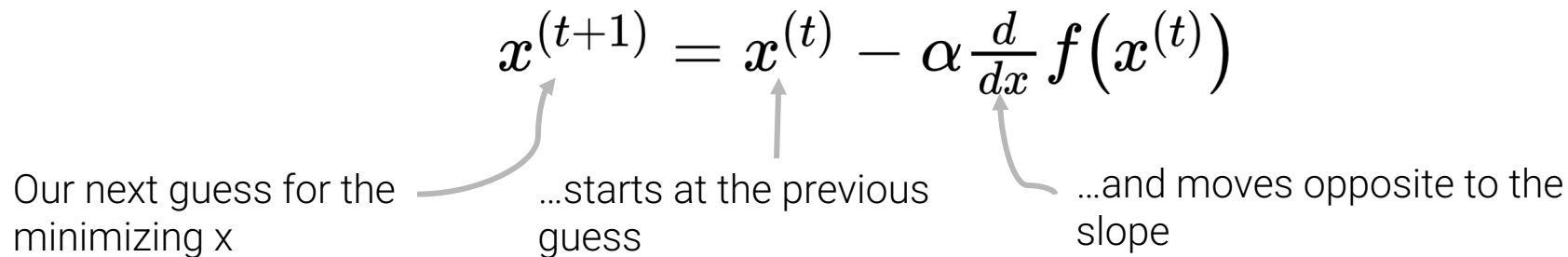
Problem: each step is too big, so we overshoot the minimizing  $x$

Solution: decrease the size of each step

Updated algorithm:  $\alpha$  represents a **learning rate** that we choose. It controls the size of each step.

$$x^{(t+1)} = x^{(t)} - \alpha \frac{d}{dx} f(x^{(t)})$$

Our next guess for the minimizing  $x$       ...starts at the previous guess      ...and moves opposite to the slope

The diagram shows the equation  $x^{(t+1)} = x^{(t)} - \alpha \frac{d}{dx} f(x^{(t)})$ . Three arrows originate from text labels below and point to specific parts of the equation: one from 'Our next guess for the minimizing x' to  $x^{(t+1)}$ , one from '...starts at the previous guess' to  $x^{(t)}$ , and one from '...and moves opposite to the slope' to  $-\alpha \frac{d}{dx} f(x^{(t)})$ .

Let's try  $\alpha = 0.3$

# Gradient Descent on a 1D Model

---

## Lecture 14

- Optimization: where are we?
- Minimizing an arbitrary 1D function
- **Gradient descent on a 1D model**
- Gradient descent on high-dimensional models
- Batch, mini-batch, and stochastic gradient descent

## From Arbitrary Functions to Loss Functions

In a modeling context, we aim to minimize a *loss function* by choosing the minimizing model *parameters*.

Terminology clarification:

- In applications, we usually care more about the average error across *all* datapoints

Going forward, we will take the “model’s loss” to mean the model’s average error across the dataset. This is sometimes also known as the empirical risk, cost function, or objective function.

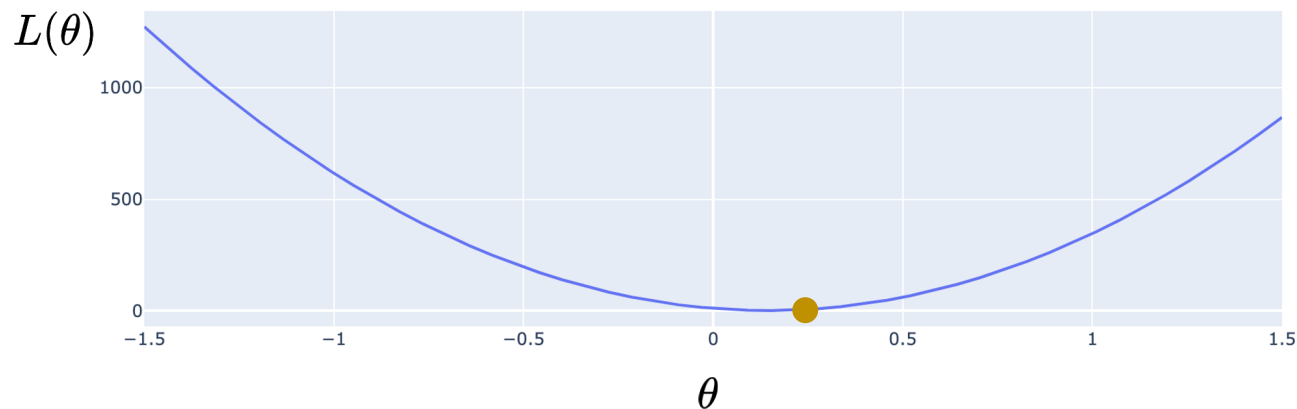
$$L(\theta) = R(\theta) = \frac{1}{n} \sum_{i=1}^n l(y, \hat{y})$$

## From Arbitrary Functions to Loss Functions

In a modeling context, we aim to minimize a *loss function* by choosing the minimizing model *parameters*.

Goal: choose the value of  $\theta$  that minimizes  $L(\theta)$ , the model's loss on the dataset

Our new framework:



Compute the model's loss for this choice of parameter

Test several values of the parameter  $\theta$

## From Arbitrary Functions to Loss Functions

---

Goal: choose the value of  $\theta$  that minimizes  $L(\theta)$ , the model's loss on the dataset

The **1D gradient descent** algorithm:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{d}{d\theta} L(\theta^{(t)})$$



## Gradient Descent on the tips Dataset

---

We want to predict the tip ( $y$ ) given the price of a meal ( $x$ ). To do this:

- Choose a model:  $\hat{y} = \theta_1 x$
- Choose a loss function:  $L(\theta) = MSE(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_1 x_i)^2$

## Gradient Descent on the tips Dataset

---

We want to predict the tip ( $y$ ) given the price of a meal ( $x$ ). To do this:

- Choose a model:  $\hat{y} = \theta_1 x$
- Choose a loss function:  $L(\theta) = MSE(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_1 x_i)^2$
- Optimize the model parameter – apply gradient descent

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{d}{d\theta} L(\theta^{(t)})$$

## Gradient Descent on the tips Dataset

---

- Optimize the model parameter – apply gradient descent

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{d}{d\theta} L(\theta^{(t)})$$

Our loss function

$$L(\theta) = MSE(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_1 x_i)^2$$

## Demo: Gradient Descent on the tips Dataset

- Optimize the model parameter – apply gradient descent

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{d}{d\theta} L(\theta^{(t)})$$

Our loss function

$$L(\theta) = MSE(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - \theta_1 x_i)^2$$

Take the derivative wrt  $\theta_1$

$$\frac{d}{d\theta_1} L(\theta_1^{(t)}) = \frac{-2}{n} \sum_{i=1}^n (y_i - \theta_1^{(t)} x_i) x_i$$

The gradient descent update rule

$$\theta_1^{(t+1)} = \theta_1^{(t)} - \alpha \frac{-2}{n} \sum_{i=1}^n (y_i - \theta_1^{(t)} x_i) x_i$$

# Gradient Descent on Multi- Dimensional Models

---

## Lecture 14

- Optimization: where are we?
- Minimizing an arbitrary 1D function
- Gradient descent on a 1D model
- **Gradient descent on multi-dimensional models**
- Batch, mini-batch, and stochastic gradient descent

## Models in 2D or Higher

---

Usually, models will have more than one parameter that needs to be optimized.

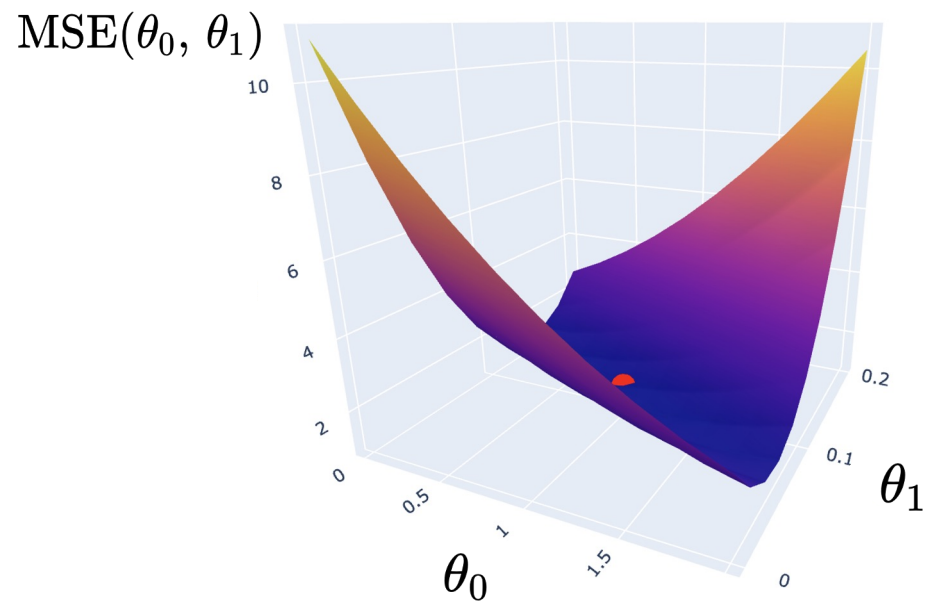
Simple linear regression:  $\hat{y} = \theta_0 + \theta_1 x$

Multiple linear regression:  $\hat{Y} = \theta_0 + \theta_1 X_{:,1} + \theta_2 X_{:,2} \dots + \theta_p X_{:,p}$

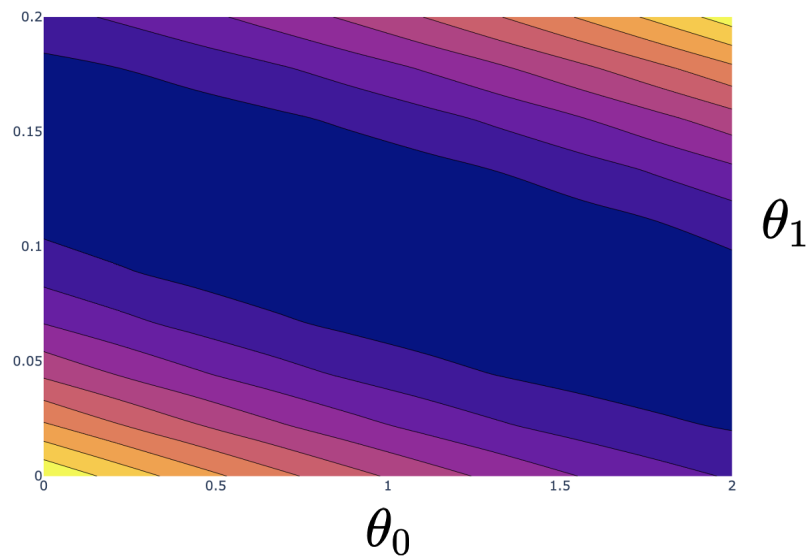
Idea: expand gradient descent so we can update our guesses for *all* model parameters, all in one go

With multiple parameters to optimize, we consider a **loss surface**

- What is the model's loss for a particular *combination* of possible parameter values?



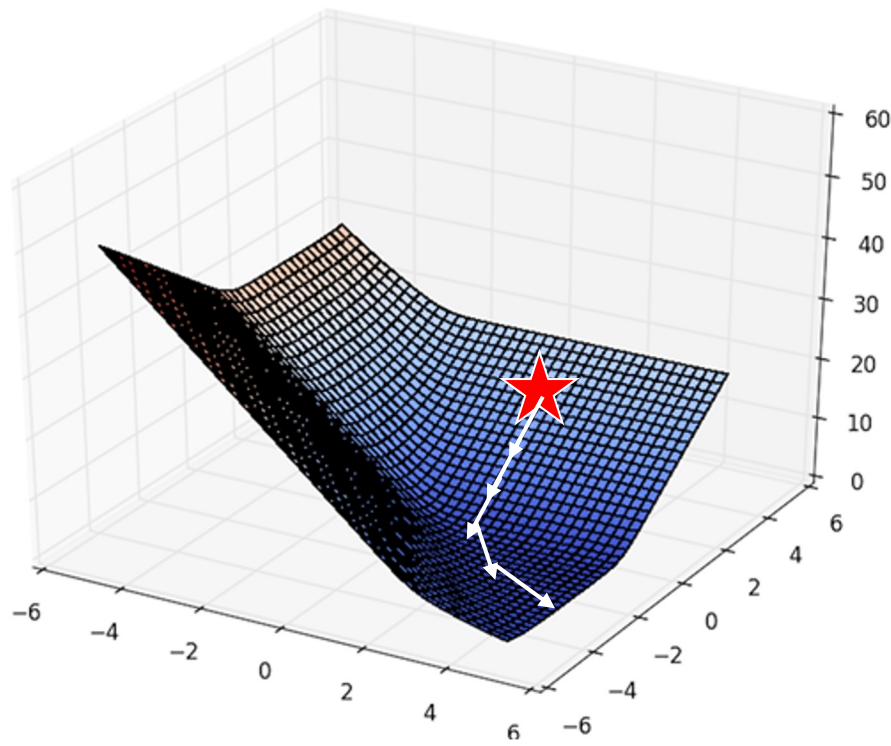
A bird's eye view from above:



# The Gradient Vector

As before, the derivative of the loss function tells us the best way towards the minimum value

On a 2D (or higher) surface, the best way to go down (gradient) is described by a *vector*



For the vector of parameter values  $\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$

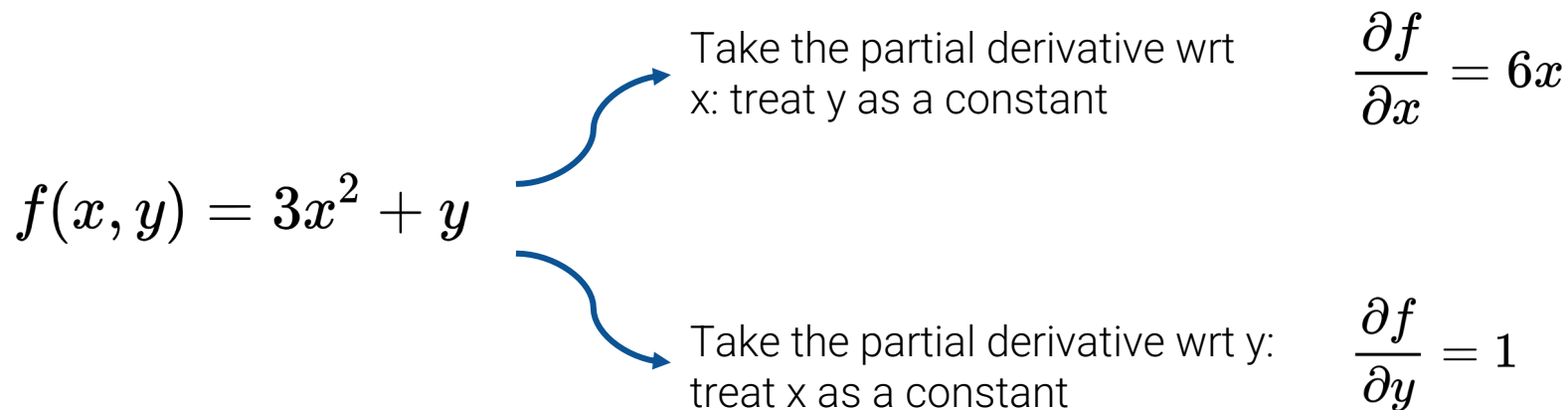
Take the *partial derivative* of loss with respect to each parameter  $\theta_i$



## A Math Aside: Partial Derivatives

For an equation with multiple variables, we take a **partial derivative** by differentiating with respect to just one variable at a time.

Intuitively: how does the function change if we vary one variable, while holding the others constant?



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

Take the partial derivative wrt  $x$ : treat  $y$  as a constant

$$\frac{\partial f}{\partial x} = 6x$$

Take the partial derivative wrt  $y$ : treat  $x$  as a constant

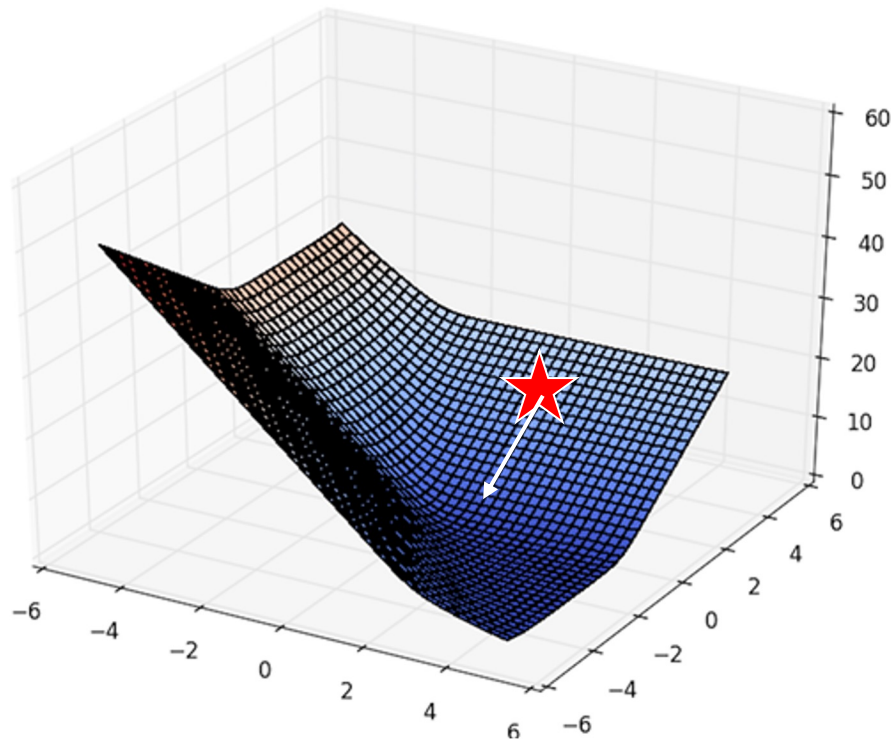
$$\frac{\partial f}{\partial y} = 1$$

This symbol means  
"partial derivative"

## The Gradient Vector

For the vector of parameter values  $\vec{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix}$

Take the *partial derivative* of loss with respect to each parameter:  $\frac{\partial L}{\partial \theta_0}$ ,  $\frac{\partial L}{\partial \theta_1}$



The **gradient vector** is

$$\nabla_{\vec{\theta}} L = \begin{bmatrix} \frac{\partial L}{\partial \theta_0} \\ \frac{\partial L}{\partial \theta_1} \end{bmatrix}$$

—  $\nabla_{\vec{\theta}} L$  always points in the downhill direction of the surface.

## Gradient Descent in Multiple Dimensions

---

Recall our 1D update rule:

$$\theta^{(t+1)} = \theta^{(t)} - \alpha \frac{d}{d\theta} L(\theta^{(t)})$$

Now, for models with multiple parameters, we work in terms of vectors:

$$\begin{bmatrix} \theta_0^{(t+1)} \\ \theta_1^{(t+1)} \\ \vdots \end{bmatrix} = \begin{bmatrix} \theta_0^{(t)} \\ \theta_1^{(t)} \\ \vdots \end{bmatrix} - \alpha \begin{bmatrix} \frac{\partial L}{\partial \theta_0} \\ \frac{\partial L}{\partial \theta_1} \\ \vdots \end{bmatrix}$$

Written in a more compact form:

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha \nabla_{\vec{\theta}} L(\vec{\theta}^{(t)})$$

## Gradient Descent Update Rule

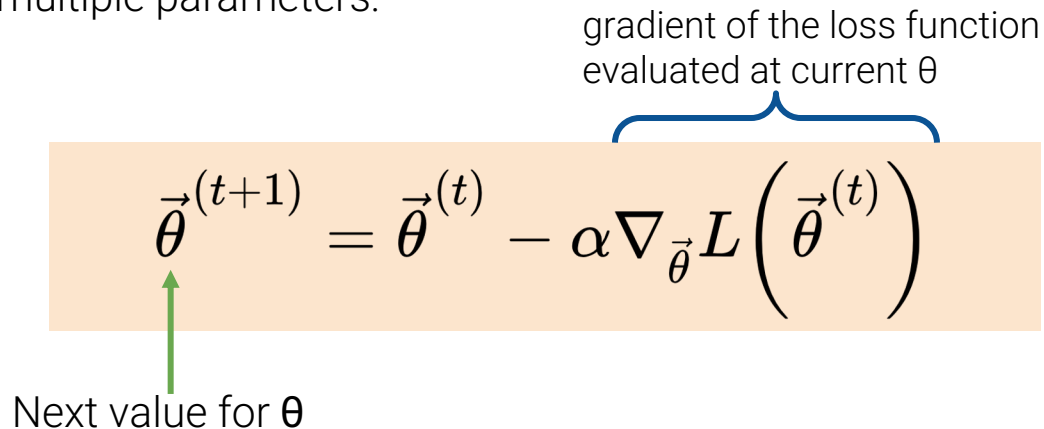
Gradient descent algorithm: nudge  $\theta$  in a negative gradient direction until  $\theta$  converges.

For a model with multiple parameters:

gradient of the loss function  
evaluated at current  $\theta$

$$\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha \nabla_{\vec{\theta}} L\left(\vec{\theta}^{(t)}\right)$$

Next value for  $\theta$

The diagram shows the gradient descent update rule equation:  $\vec{\theta}^{(t+1)} = \vec{\theta}^{(t)} - \alpha \nabla_{\vec{\theta}} L\left(\vec{\theta}^{(t)}\right)$ . The equation is enclosed in a light orange rectangular box. A blue curly bracket is positioned above the term  $\nabla_{\vec{\theta}} L\left(\vec{\theta}^{(t)}\right)$ , with the text "gradient of the loss function evaluated at current  $\theta$ " centered above it. A green arrow points upwards from the text "Next value for  $\theta$ " to the term  $\vec{\theta}^{(t+1)}$  on the left side of the equation.

# Batch, Mini-Batch, and Stochastic Gradient Descent

---

## Lecture 14

- Optimization: where are we?
- Minimizing an arbitrary 1D function
- Gradient descent on a 1D model
- Gradient descent on multi-dimensional models
- **Batch, mini-batch, and stochastic gradient descent**

We have just derived **batch gradient descent**.

- We used our *entire* dataset (as one big batch) to compute gradients
- Recall the derivative of MSE for our 1D model – involves working with *all*  $n$  datapoints

$$\frac{d}{d\theta_1} L(\theta_1^{(t)}) = \frac{-2}{n} \sum_{i=1}^n (y_i - \theta_1^{(t)} x_i) x_i$$

Using all datapoints is often impractical when our dataset is large.

Computing each gradient will take a long time; gradient descent will converge slowly because each individual update is slow.

## Mini-batch Gradient Descent

---

An alternative: use only a *subset* of the full dataset at each update.

Estimate the true gradient of the loss surface using just this subset of the data.

**Batch size:** the number of datapoints to use in each subset

In mini-batch GD:

- Compute the gradient on the first x% of the data. Update the parameter guesses.
- Compute the gradient on the next x% of the data. Update the parameter guesses.
- ...
- Compute the gradient on the last x% of the data. Update the parameter guesses.

**Training  
Epoch**

## Mini-batch Gradient Descent

---

In mini-batch GD:

- Compute the gradient on the first  $x\%$  of the data. Update the parameter guesses.
- Compute the gradient on the next  $x\%$  of the data. Update the parameter guesses.
- ...
- Compute the gradient on the last  $x\%$  of the data. Update the parameter guesses.

**Training  
Epoch**

In a single training epoch, we use every datapoint in the data once.

We then perform several training epochs until we are satisfied.



## Stochastic Gradient Descent

In the most extreme case, we may perform gradient descent with a batch size of just *one* datapoint – this is called **stochastic gradient descent**.

- Works surprisingly well in practice! Averaging across several epochs gives a similar result as directly computing the true gradient on all the data.

In stochastic GD:

- Compute the gradient on the first datapoint. Update the parameter guesses.
- Compute the gradient on the next datapoint. Update the parameter guesses.
- ...
- Compute the gradient on the last datapoint. Update the parameter guesses.

} **Training  
Epoch**