

WOMEN WHO

**SSID: WWCode**

**Password: password**





WOMEN WHO  
**CODE**  
MANILA



# ML and AI Study Group

Twitter: @wwcodemania  
FB: fb.com/wwcodemania

#WWCodeManila  
#YourProgrammingLanguage  
#StudyGroup



**Issa**

Research Fellow  
PCARI



WOMEN WHO

# **New Member's Introduction**





**I am <name>**

<your current profession>

<why did you join this study group?>



# OUR MISSION

Inspiring women to excel in technology careers.



# OUR VISION

A world where women are representative as technical executives, founders, VCs, board members and software engineers.



# STUDY GROUP

Study groups are events where women can come together and help each other learn and understand a specific programming language, technology, or anything related to coding or engineering.



# GUIDELINES

- If you have a question, just **ask**
- If you have an idea, **share it**
- **Make friends** and learn from your study groupmates
- **Do not** recruit or promote your business

WOMEN WHO

CODE

**SHOW & TELL**



# **LINEAR REGRESSION IN ONE VARIABLE**

(a.k.a Univariate Linear Regression)



# Agenda

1. Review
2. Hypothesis Function
3. Cost Function
4. Gradient Descent Algorithm

# Review

- Supervised vs. Unsupervised

# Review

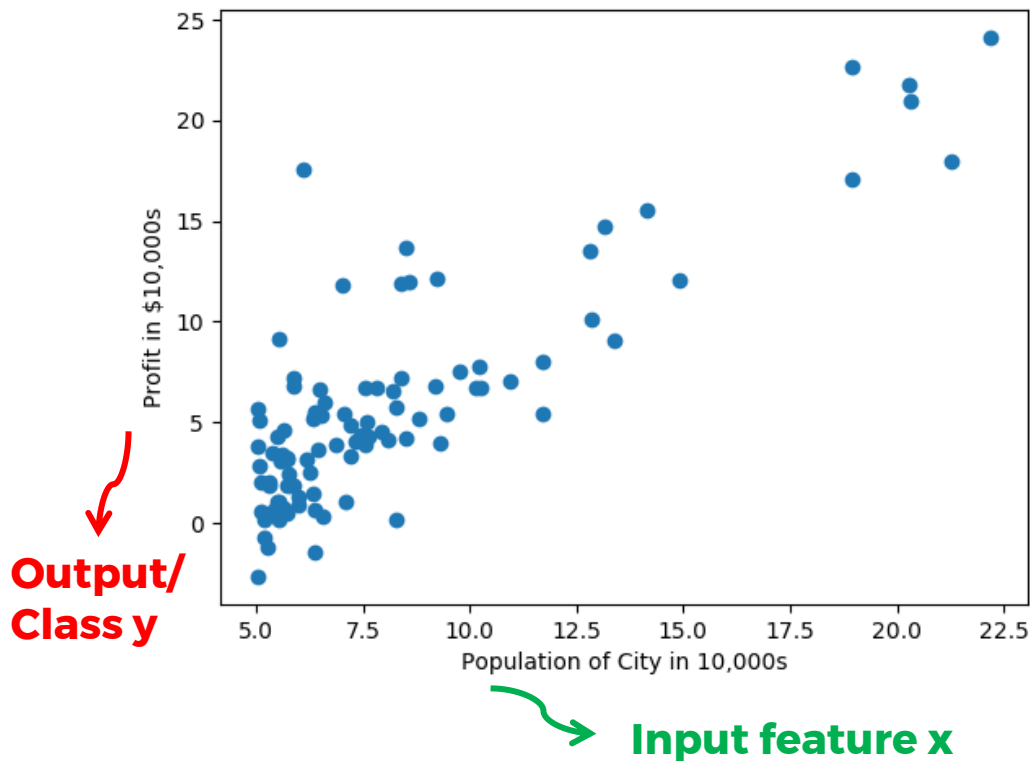
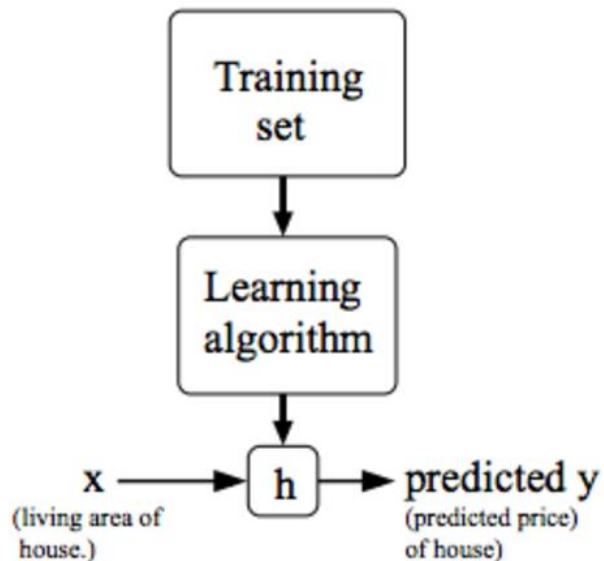
- Supervised vs. Unsupervised
- Supervised Learning:
  - **Classification** – output variables take discrete class labels
  - **Regression** – output variable takes continuous values

# Review

- Supervised vs. Unsupervised
- Supervised Learning:
  - **Classification** – output variables take discrete class labels
  - **Regression** – output variable takes continuous values
- **Univariate Linear Regression** – predicts a single output  $y$  from a single input value  $x$

# Linear Regression in One Variable

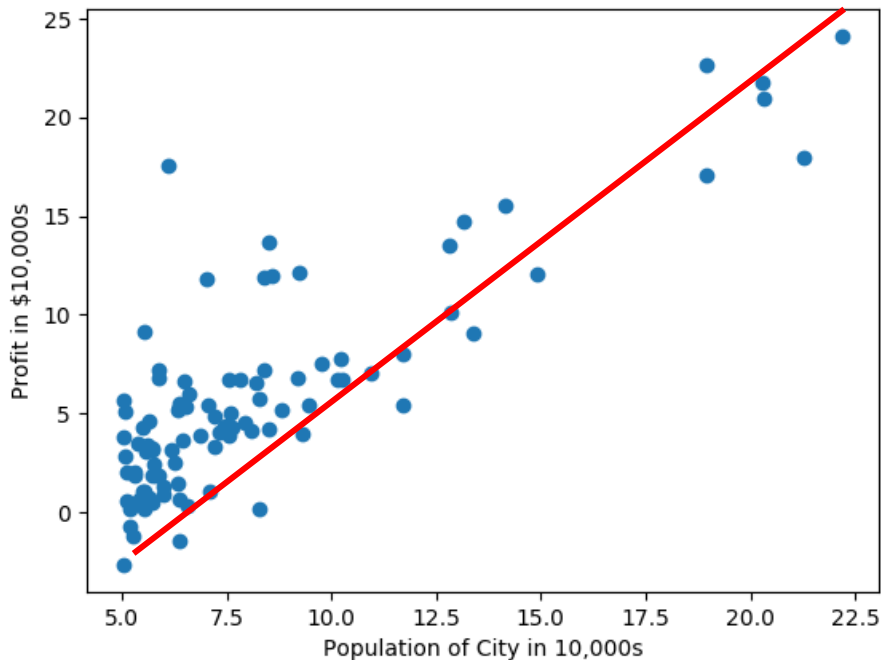
e.g. **Predicting Profit of a City**





# Linear Regression

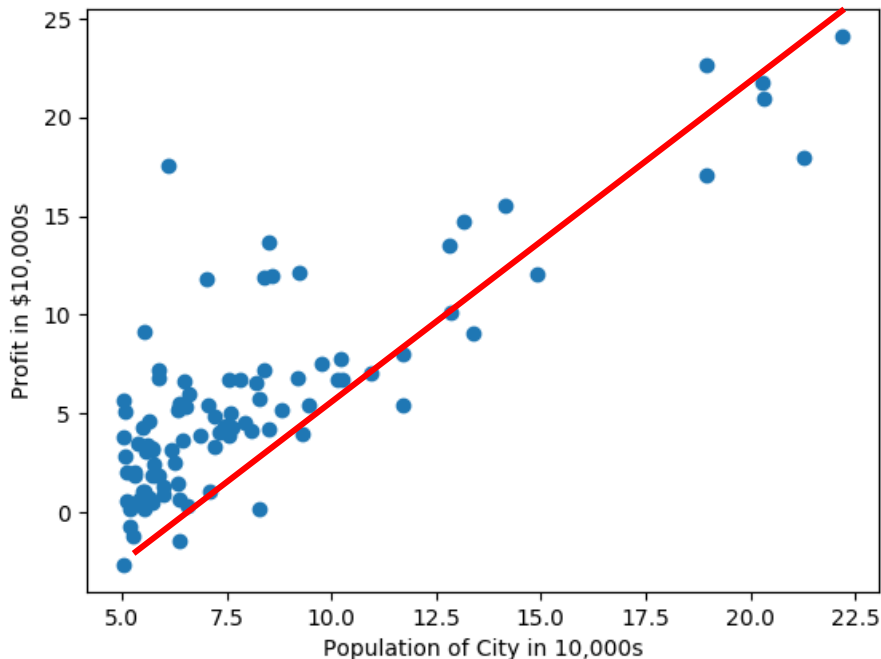
We want to **fit** a line through the data.



# Linear Regression

We want to **fit** a line through the data.

And make it so that the line **generalizes** the data well.

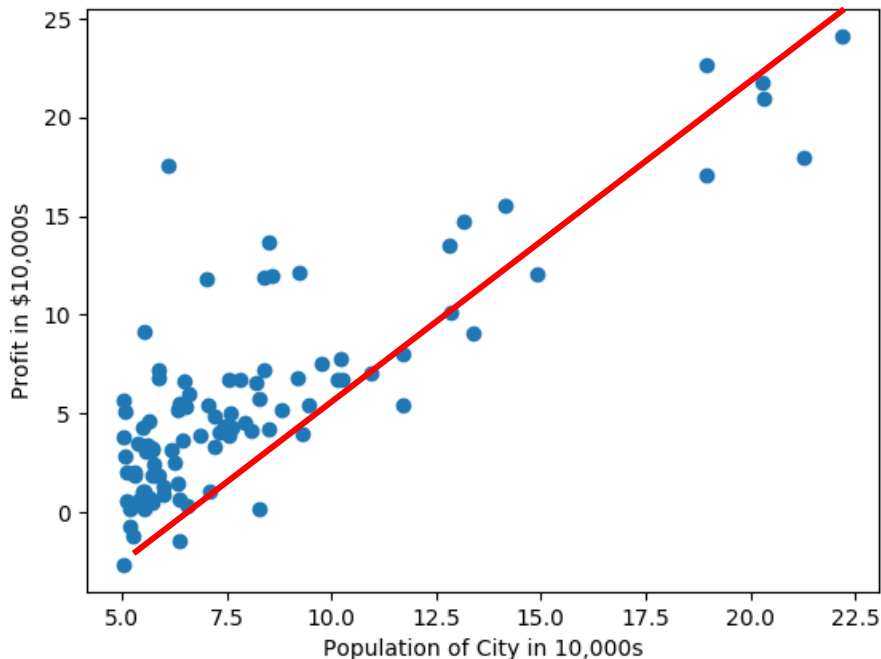


# Linear Regression

We want to **fit** a line through the data.

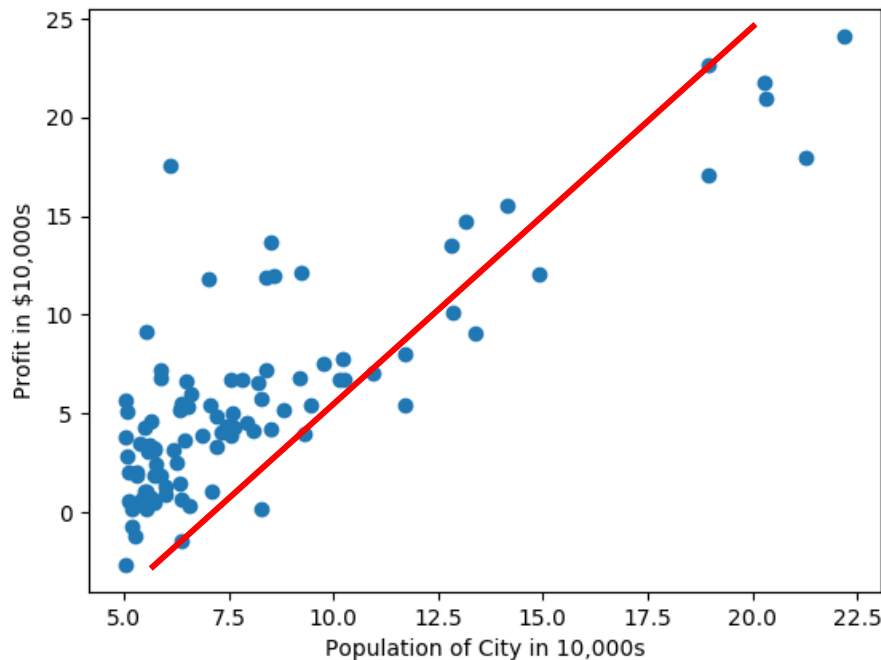
And make it so that the line **generalizes** the data well.

e.g. If Imus City has a population of about 21k, then its predicted profit is...?



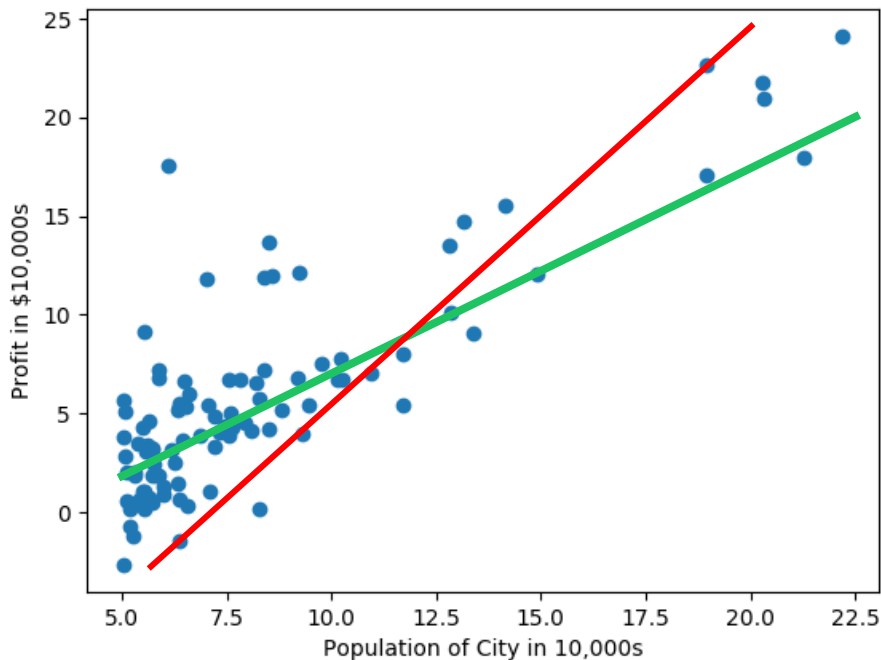
# Linear Regression

How do we come up with the “right” line?



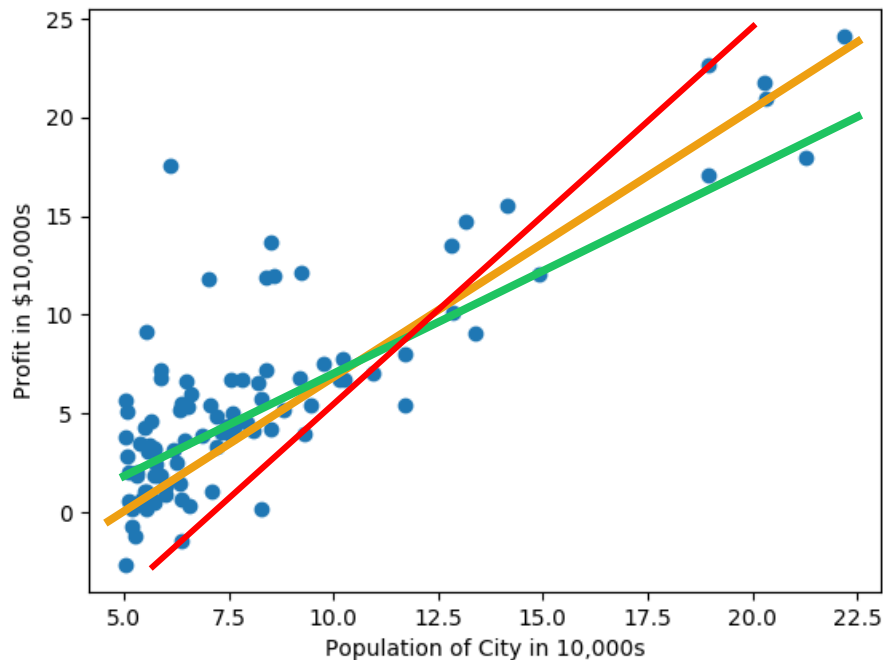
# Linear Regression

How do we come up with the “right” line?



# Linear Regression

How do we come up with the “right” line?



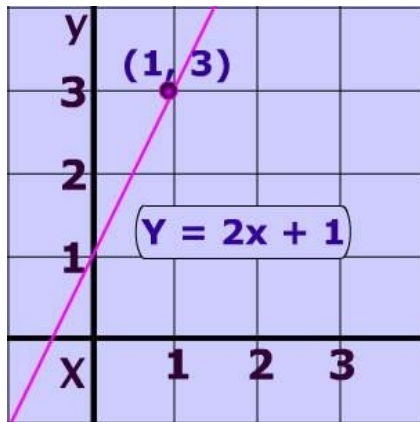
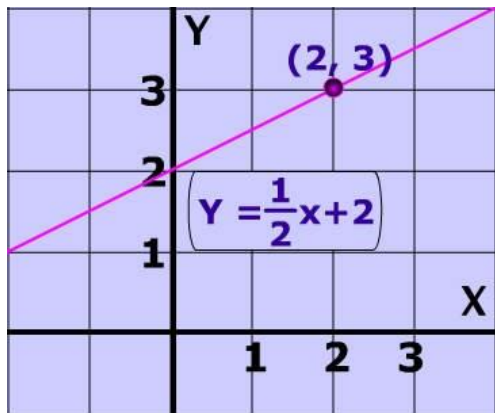
# Equation of a Line

- Recall the equation of a line?

# Equation of a Line

- Recall the equation of a line?

$$y = mx + b$$



## Exercise:

Draw the lines:

$$y = \frac{3}{2}x + 1$$

$$y = -2x - 5$$



# Equation of a Line

- Recall the equation of a line?

$$y = mx + b$$

- More generally,

$$y = \theta_1 x + \theta_0$$

where  $\theta_0, \theta_1$  are called “weights”.

$\theta_0$  has a special name called the “bias”

# Hypothesis Function

- We need to find the right values for  $\theta_0$  and  $\theta_1$ .
- How do we get the right values?

# Hypothesis Function

- We need to find the right values for  $\theta_0$  and  $\theta_1$ .
- How do we get the right values?
- We start with a guess (or *hypothesis*) using random weights

# Hypothesis Function

- We need to find the right values for  $\theta_0$  and  $\theta_1$ .
- How do we get the right values?
- We start with a guess (or *hypothesis*) using random weights
- We define our **hypothesis function** as

$$h(x) = \theta_1 x + \theta_0$$

where  $\theta_0$  and  $\theta_1$  are set to some random values (e.g.  $\theta_0 = 0, \theta_1 = 0$ )

# Hypothesis Function

- Let  $x$  be the input features and  $y$  be the true output values.
- Suppose we come up with two hypotheses:

Input $x$	Output $y$	$h_{\theta}(x)$	$h_{\theta}(x)$
		$\theta_0 = 2, \theta_1 = 3$	$\theta_0 = 4, \theta_1 = 3$
0	4	?	?
1	8	?	?
2	9	?	?
3	13	?	?

# Hypothesis Function

- Let  $x$  be the input features and  $y$  be the true output values.
- Suppose we come up with two hypotheses:

Input $x$	Output $y$	$h_{\theta}(x)$	$h_{\theta}(x)$
		$\theta_0 = 2, \theta_1 = 3$	$\theta_0 = 4, \theta_1 = 3$
0	4	2	?
1	8	5	?
2	9	8	?
3	13	11	?

# How good is our hypothesis?

- We need a way to measure how close our hypothesis  $h_{\theta}(x)$  is to the true output values  $y$ .
- i.e. we want a function that measures how good or bad our hypothesis function performs

# How good is our hypothesis?

- We need a way to measure how close our hypothesis  $h_{\theta}(x)$  is to the true output values  $y$ .
- i.e. we want a function that measures how good or bad our hypothesis function performs
- We will call such a function the **Cost Function**, which will measure the average error of the  $h_{\theta}(x)$ .



# Cost Function

$$\text{Cost } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(h(x_i) - y_i)^2}_{\text{error - we want this to be close to zero}}$$

*“Mean Squared Error (MSE)”*

average over  $m$   
examples

error - we want this to be close to zero

sum over all  $m$  examples

Note: The error is squared and the mean is halved as a convenience for computing the gradient descent later.

# Minimizing the Cost Function

- We want to minimize the cost function, or the mean squared error (MSE).

# Minimizing the Cost Function

- We want to minimize the cost function, or the mean squared error (MSE).
- In other words, we want to find the  $\theta_0$  and  $\theta_1$  that minimizes  $J(\theta_0, \theta_1)$

# Minimizing the Cost Function

Input x	Output y
0	0
1	2
2	4
3	6

Suppose  $\theta_0 = 0$ .

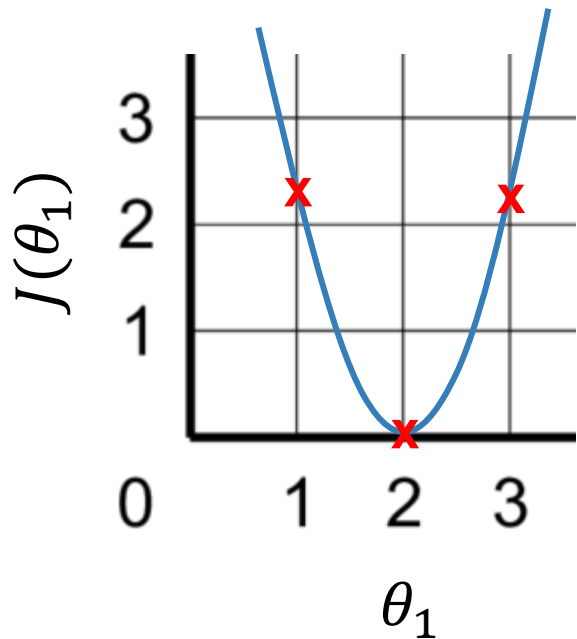
1. Plot  $h_{\theta}(x)$  when:

- $\theta_1 = 1$
- $\theta_1 = 2$
- $\theta_1 = 3$

2. Solve for:

- $J(\theta_1 = 1)$
- $J(\theta_1 = 2)$
- $J(\theta_1 = 3)$

# $J(\theta)$ vs $\theta$ plot

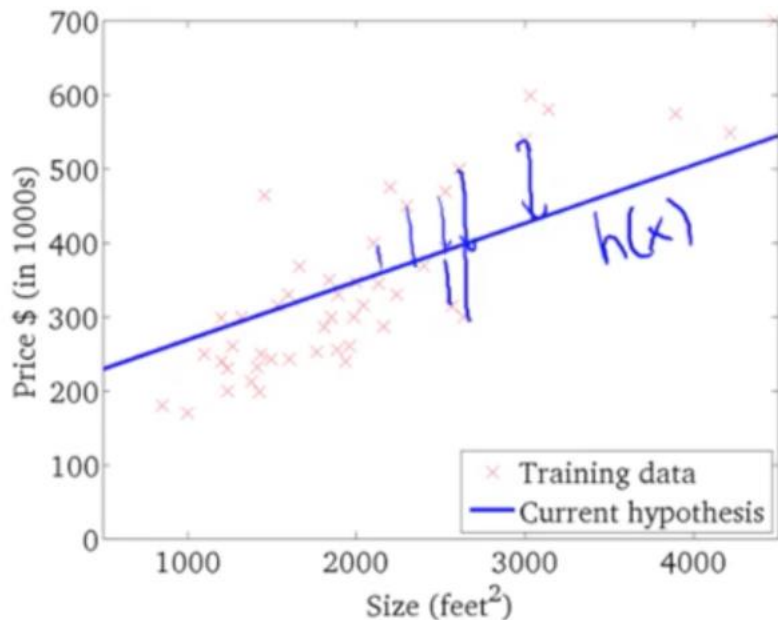


- A parabola makes sense, since  $J$  is a quadratic function.
- The value of  $\theta_1$  that minimizes  $J$  is 2.

# Minimizing J

$$h_{\theta}(x)$$

(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



In general,

the values of the weights  
 $\theta_1, \theta_0$  that minimize  $J(\theta_1, \theta_0)$



minimizes the distance  
between  $h(x)$  and every data  
point

# Gradient Descent Algorithm

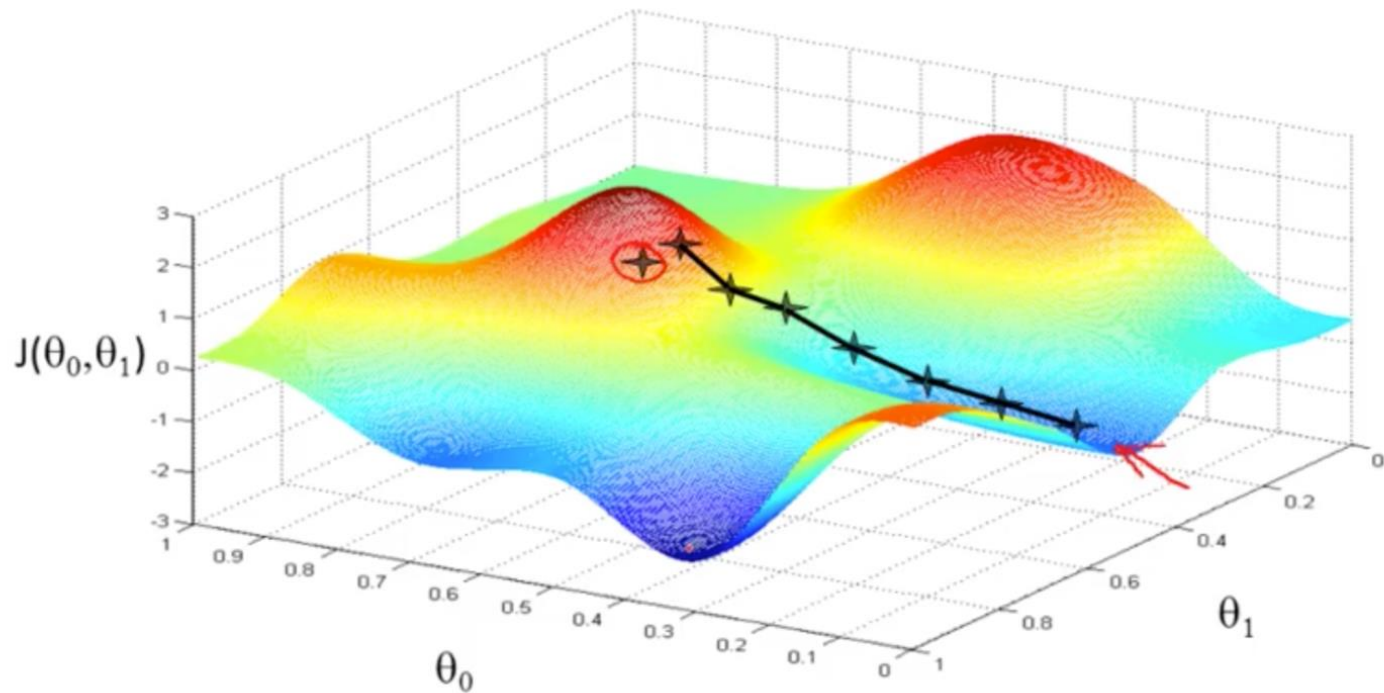
- Have some function  $J(\theta_0, \theta_1)$
- Want to find  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$

# Gradient Descent Algorithm

- Have some function  $J(\theta_0, \theta_1)$
- Want to find  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$
- Outline of Gradient Descent Algorithm
  - Start with some  $\theta_0, \theta_1$
  - Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at the minimum.



# Gradient Descent Algorithm



# Gradient Descent Algorithm

- We want to take steps down the cost function in the direction of the steepest descent

- The **direction** is given by the **derivative of the cost function**:

$$\frac{\partial}{\partial \theta_j} J(\theta_j) \quad (\text{for } j = 0 \text{ and } j = 1)$$

- The **size** of each step is given by some **learning rate**  $\alpha$ .

# Gradient Descent Algorithm

repeat until convergence:

$$\theta_j := \theta_j - \alpha \underbrace{\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)}_{\text{Derivative of the cost function}} \quad (\text{for } j = 0 \text{ and } j = 1)$$

**Learning rate**  
(gives the step size  
of the descent)

**Derivative of the cost function**  
(gives the direction of the descent)

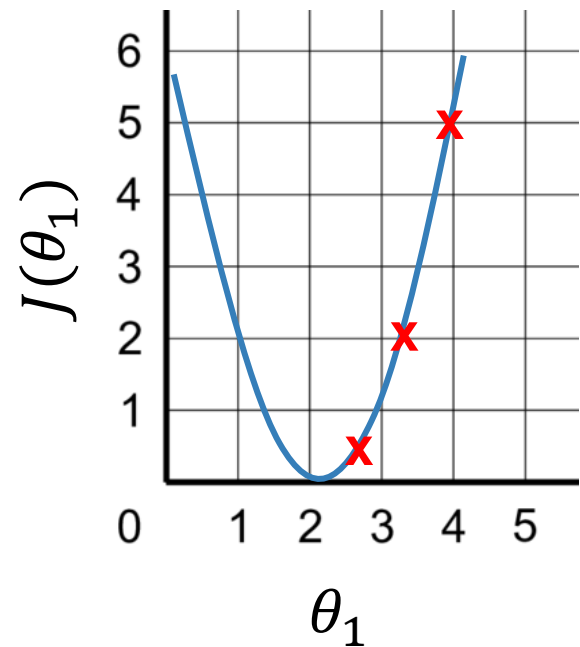
- $a := b$  (this means assignment)
- $a = b$  (truth assertion)

WOMEN WHO

**Let's talk about derivatives for  
a sec.**

CODE<sup>®</sup>

# Derivatives

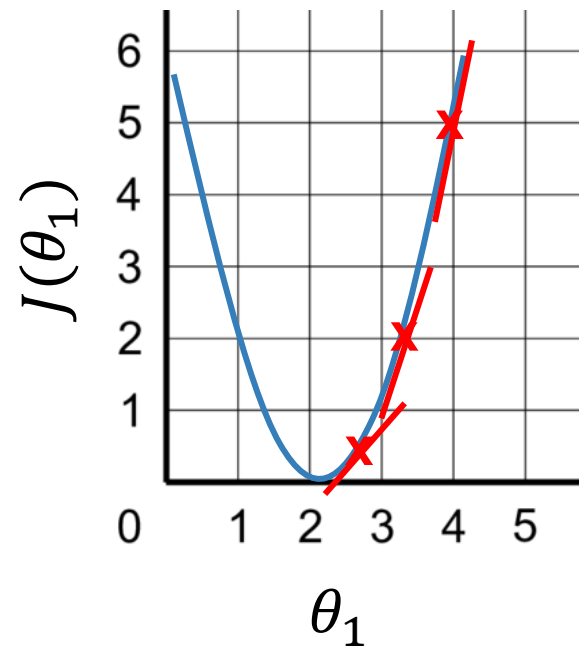


To get to the minimum value, we take the **derivative** of our cost function.

$$\frac{\partial}{\partial \theta_1} J(\theta_1)$$

**Derivative** – the slope of the tangent line at a point

# Derivatives

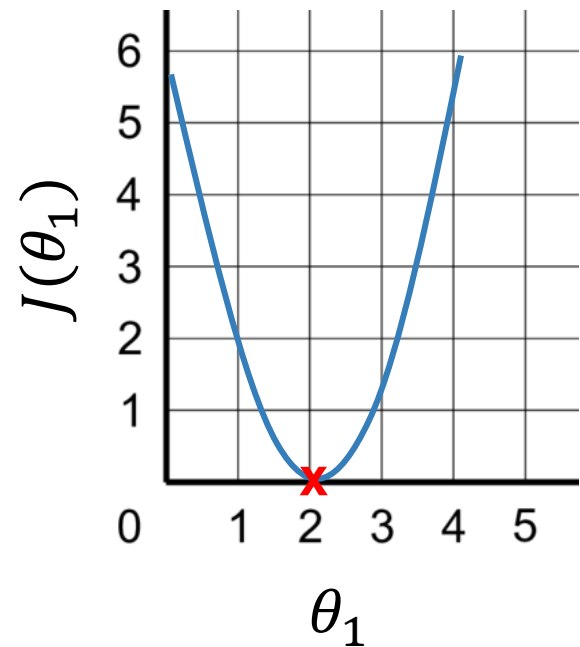


To get to the minimum value, we take the **derivative** of our cost function.

$$\frac{\partial}{\partial \theta_1} J(\theta_1)$$

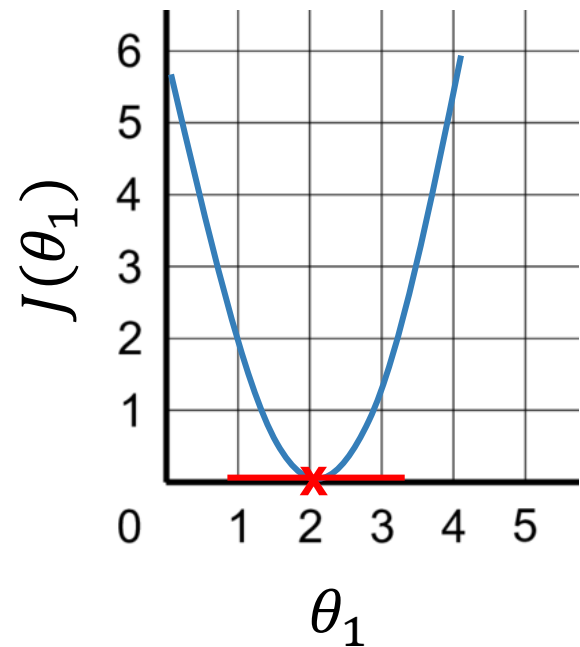
**Derivative** – the slope of the tangent line at a point

# Derivatives



What is the slope at the minimum of  $J$ ?

# Derivatives

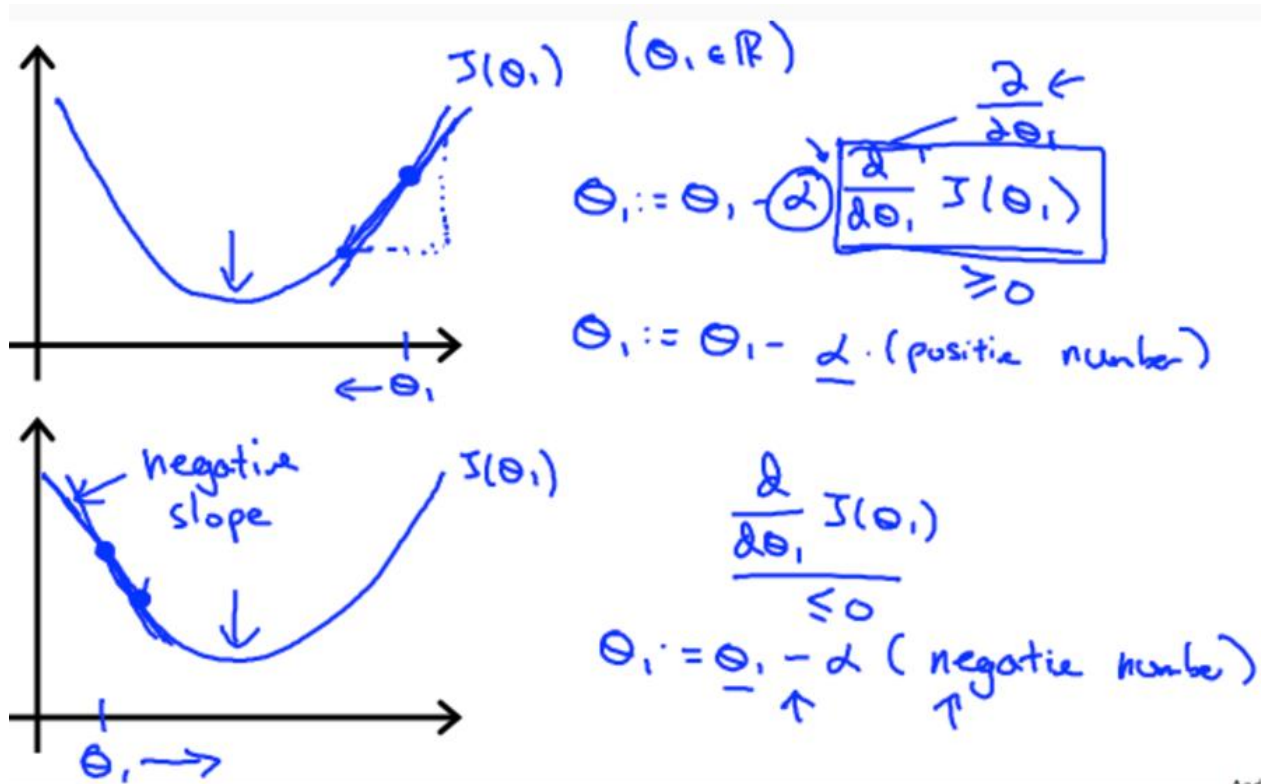


What is the slope at the minimum of  $J$ ?

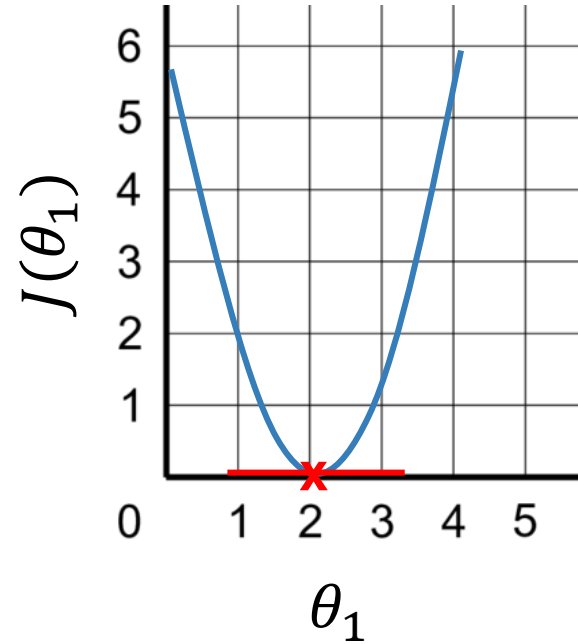
The slope given by  $\frac{\partial}{\partial \theta_1} J(\theta_1)$  will always be zero!



# Derivative Intuition



# Covergence



We say that our algorithm has converged when

$$\frac{\partial}{\partial \theta_1} J(\theta_1) = 0$$

$$\theta_1 := \theta_1 - \alpha(0)$$

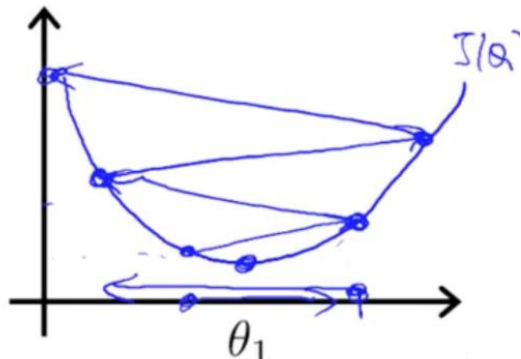
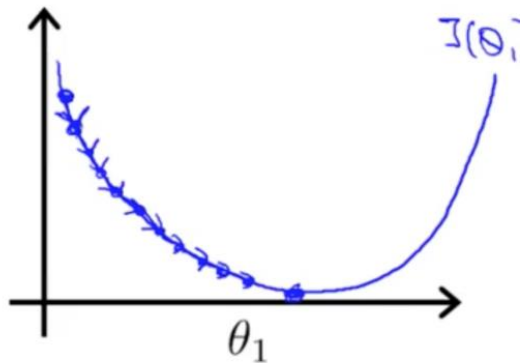
$\theta_1$  ceases to update to a new value.

# Learning Rate Intuition

$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If  $\alpha$  is too small, gradient descent can be slow.

If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



# Gradient Descent Algorithm

repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$     (for  $j = 0$  and  $j = 1$ )  
}

---

Correct: Simultaneous update

$\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$

$\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$

$\theta_0 := \text{temp0}$

$\theta_1 := \text{temp1}$

# Gradient Descent Algorithm

```
repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$     (for  $j = 0$  and  $j = 1$ )  
}
```

---

Correct: Simultaneous update

```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_0 :=$  temp0  
 $\theta_1 :=$  temp1
```

Incorrect:

```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 $\theta_0 :=$  temp0  
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_1 :=$  temp1
```

# Gradient Descent Algorithm

repeat until convergence: {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m ((h_{\theta}(x_i) - y_i)x_i)$$

}

WOMEN WHO

# **Partner/Group/Individual Exercise**

CODE®

WOMEN WHO

# **Partner/Group/Individual Presentation**





WOMEN WHO

CODE

# References



# T.I.L.

**SHARE IT!**  
**In front!**

On Twitter: @wwcodemanila  
Or FB: fb.com/wwcodemanila

Don't forget to tag WWCodeManila so we can retweet or share it.

# Feedback Form

<https://goo.gl/YzSqcS>

Please don't rate the event on meetup.

Not helpful. It is best to just tell your concerns via the feedback form. We are building a community not a Yelp restaurant.

WOMEN WHO

**THANK YOU :)**

CODE®