

Introduction to Machine Learning

SCP8084699 - LT Informatica

Linear Regression, Gradient Descent
Prof. Lamberto Ballan

Summary

- Inductive (learning) Bias
- Bias-Variance Tradeoff

- Linear Regression
 - ▶ (Univariate) Linear Regression model
 - ▶ Cost Function - Intuition
 - ▶ Examples

Today

Linear Regression

- **Example:** housing prices in Portland (OR)

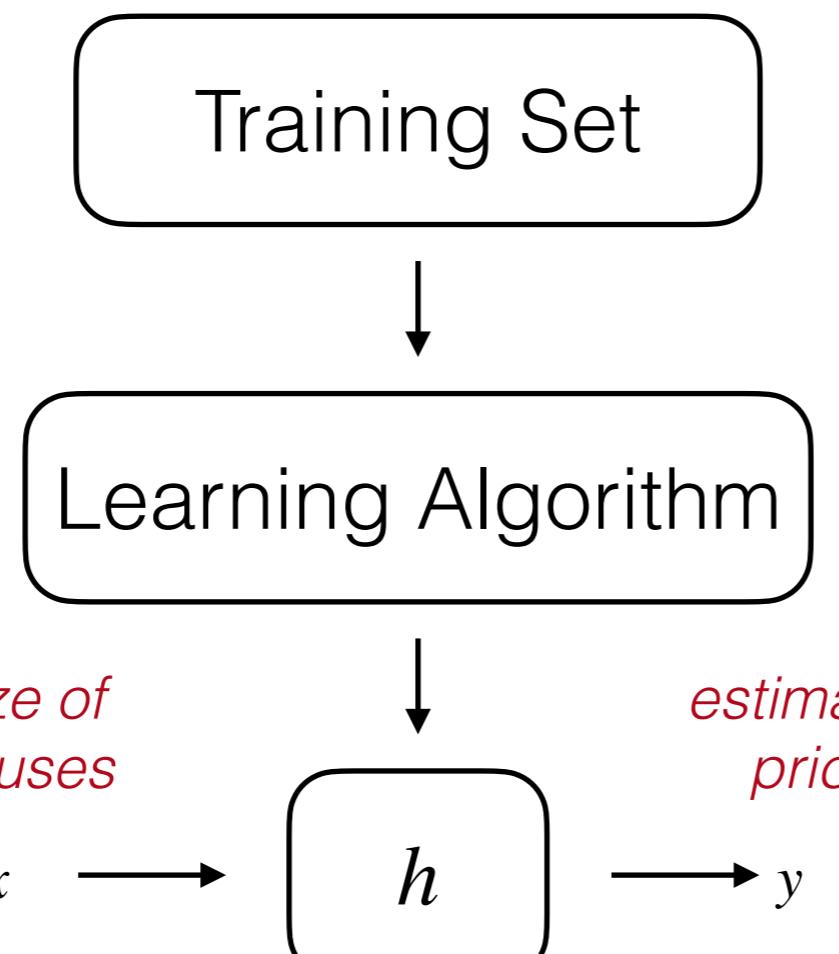
- Notation:
 - m = Number of training examples
 - x 's = “input” variable / features
 - y 's = “output” variable / target variable

Size in feet ² (x)	Price (\$) in 1K's (y)
2104	460
1416	232
1534	315
852	178
...	...

Linear Regression

$\{(x^{(i)}, y^{(i)})\}$

Size in feet ² (x)	Price (\$) in 1K's (y)
2104	460
1416	232
1534	315
852	178
...	...



How do we represent h ?

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Univariate linear regression model
(i.e. one variable)

Linear Regression

- Training set:

Size in feet ² (x)	Price (\$) in 1K's (y)
2104	460
1416	232
1534	315
852	178
...	...

- Hypothesis: $h_{\theta}(x) = \theta_0 + \theta_1 x$

θ_i 's parameters

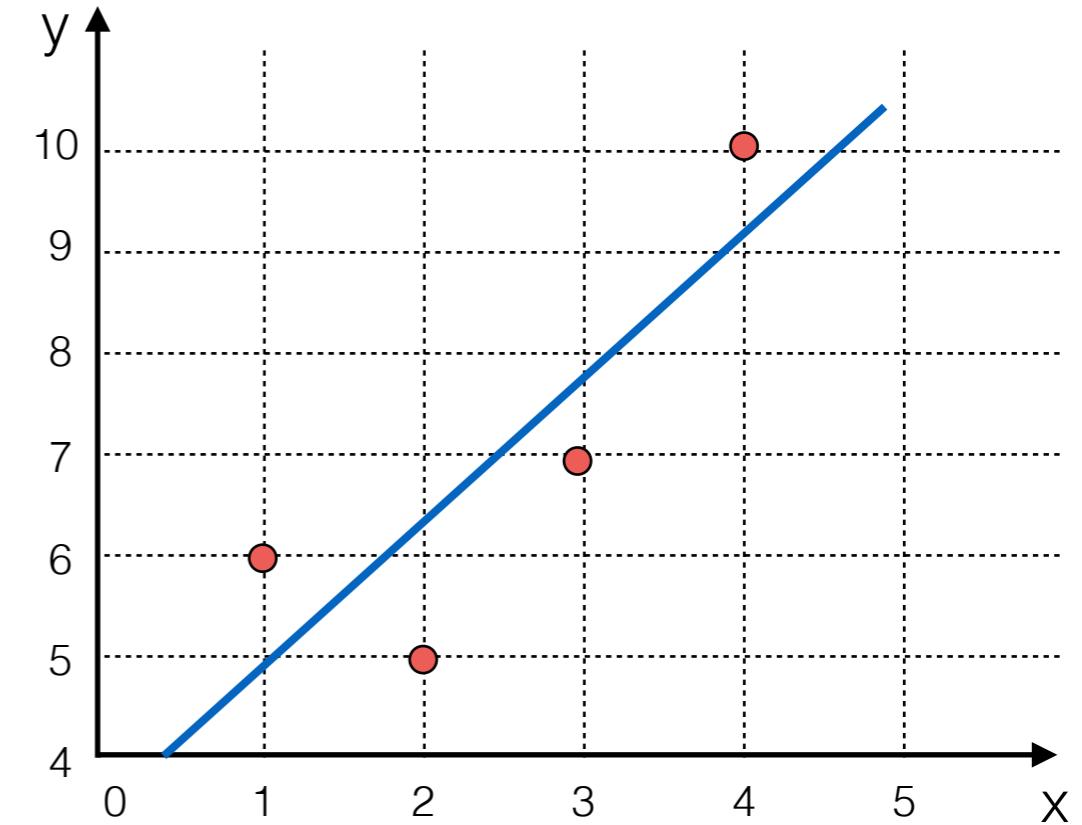
- How to chose the parameters?

Linear Regression

- Let's start with a simple example:

x	y
1	6
2	5
3	7
4	10

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



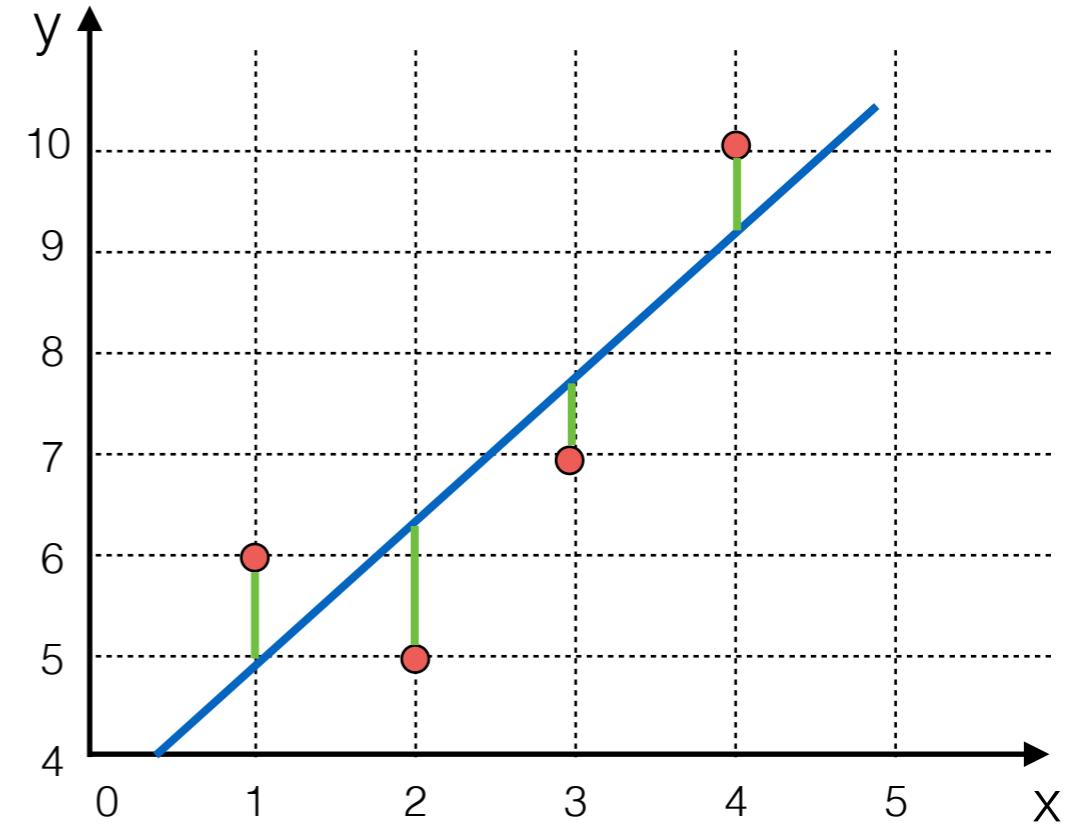
- Idea: choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples $\{(x^{(i)}, y^{(i)})\}$

Linear Regression

m training examples

x	y
1	6
2	5
3	7
4	10

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



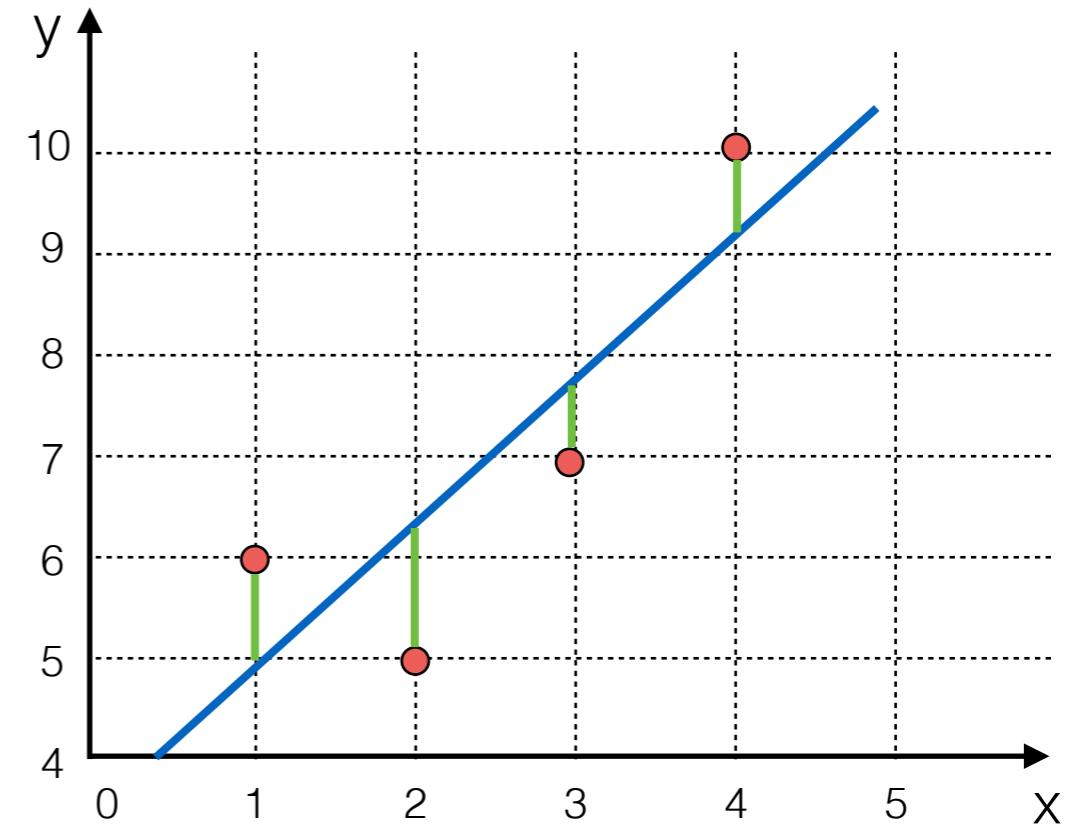
- Idea: choose θ_0, θ_1 so that $h_{\theta}(x)$ is close to y for our training examples $\{(x^{(i)}, y^{(i)})\}$
- Least squares: sum of squared distances (residuals)

Cost Function - Intuition

m training examples

x	y
1	6
2	5
3	7
4	10

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



- Residuals: $\{h_{\theta}(x^{(i)}) - y^{(i)}\}$
- Cost function: $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$
- Objective: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Cost Function - Intuition

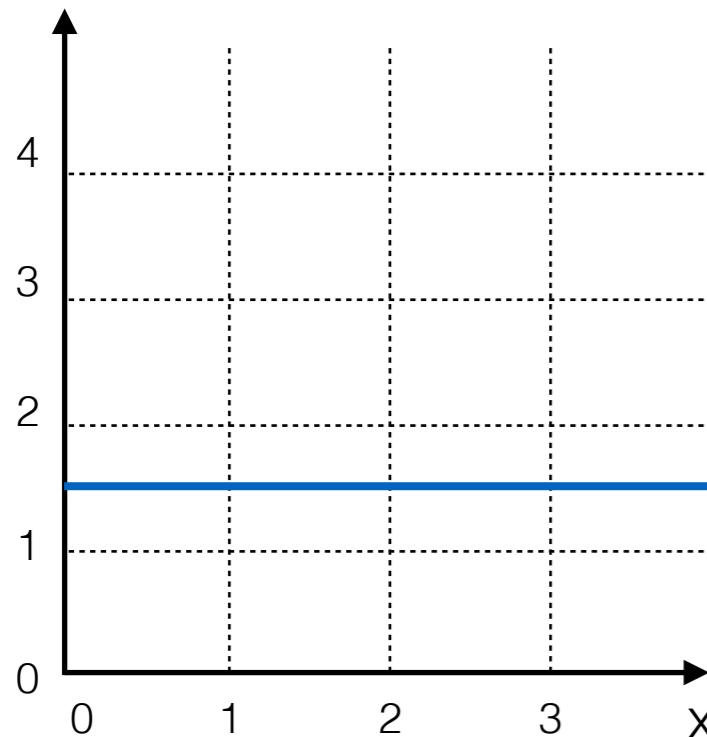
- How to choose θ_i 's? Let's see a few examples...

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

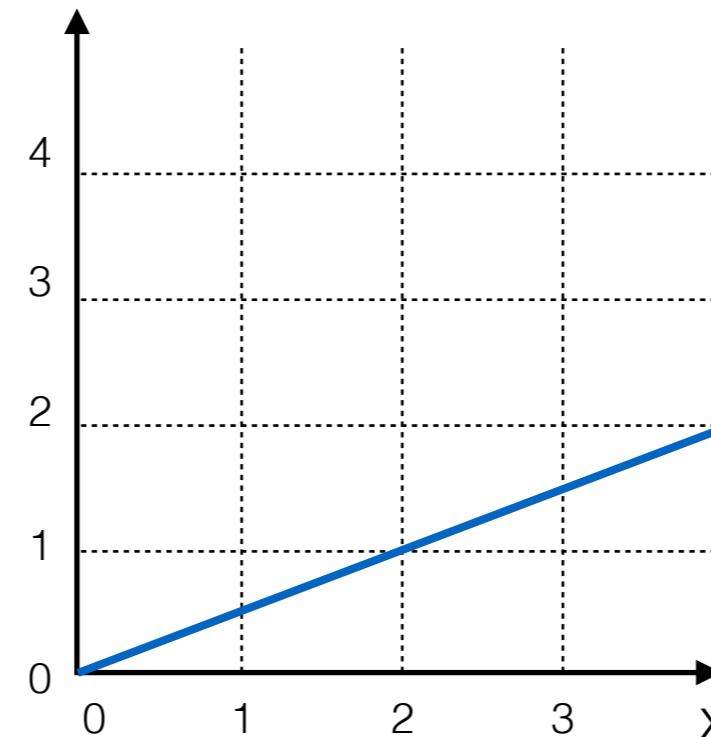
hypothesis

$$\theta_0, \theta_1$$

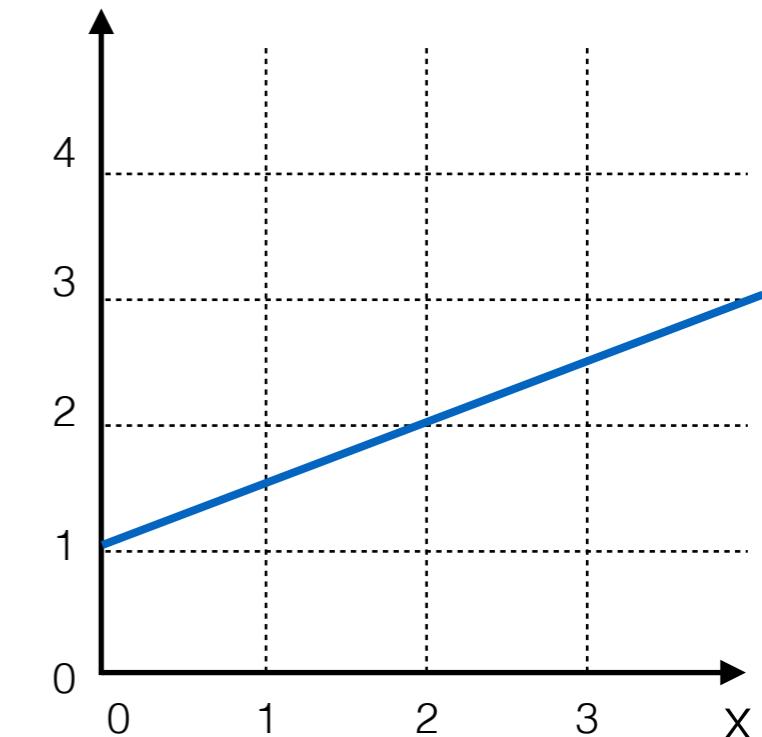
parameters



$$\theta_0 = 3/2, \theta_1 = 0$$



$$\theta_0 = 0, \theta_1 = 1/2$$



$$\theta_0 = 1, \theta_1 = 1/2$$

Cost Function - Intuition

- Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Simplified version

$$h_{\theta}(x) = \theta_1 x$$

- Parameters: θ_0, θ_1

$$\theta_1$$

- Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

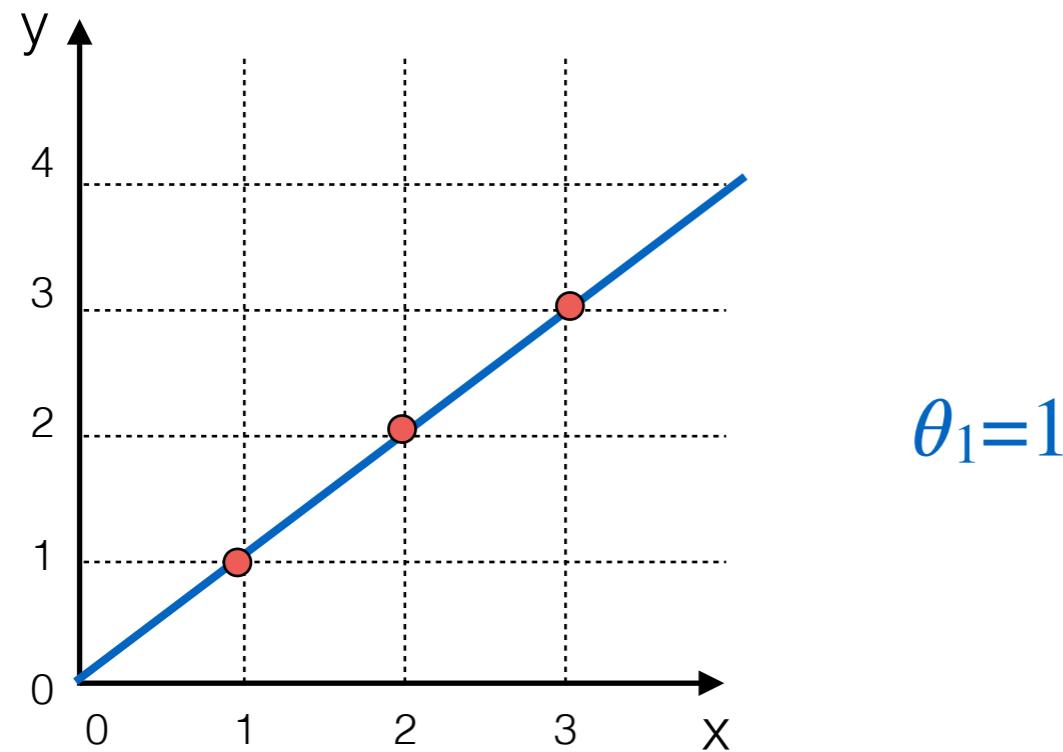
- Objective: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

$$\underset{\theta_1}{\text{minimize}} J(\theta_1)$$

Cost Function - Intuition

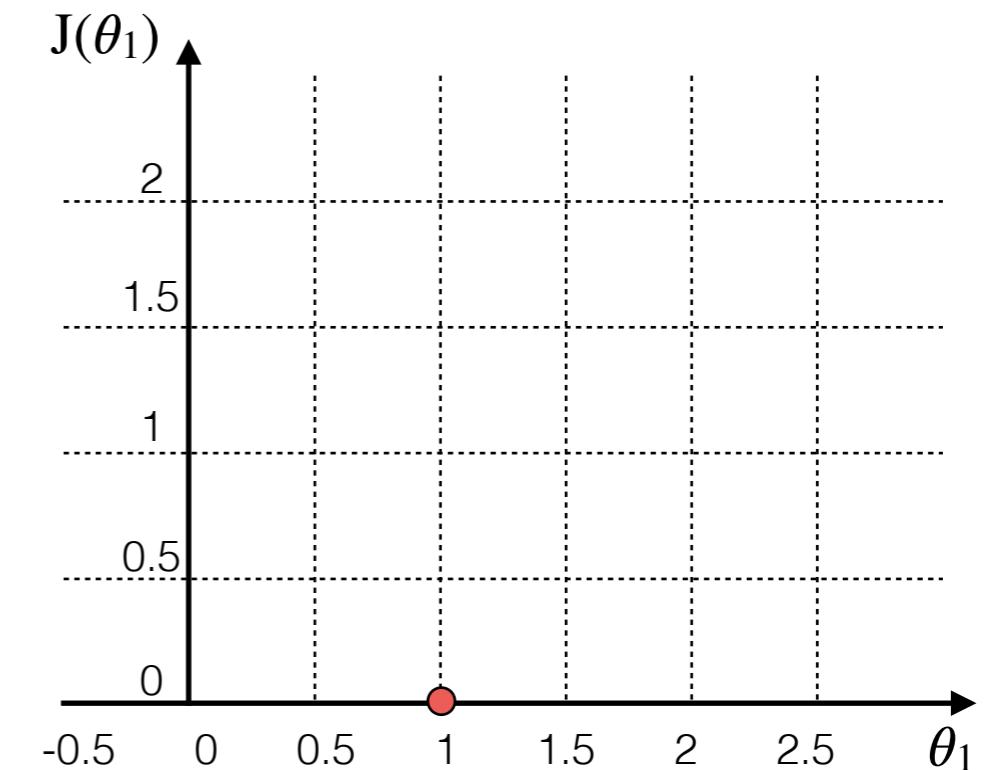
$$h_{\theta}(x) = \theta_1 x$$

(for fixed θ_1 this is a function of x)



$$J(\theta_1)$$

(function of the parameter θ_1)



$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

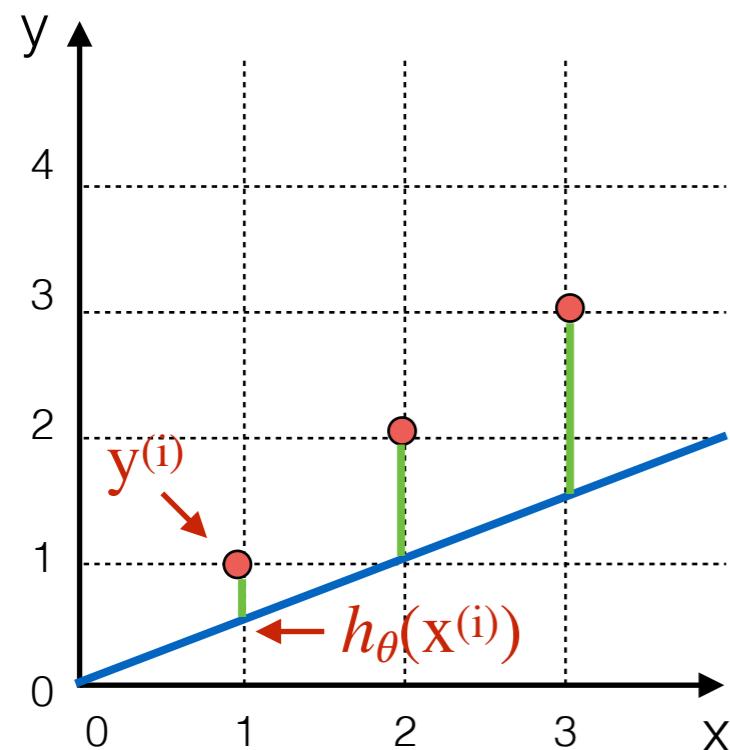
$\theta_1=1$

$$= \frac{1}{2m} (0+0+0)^2$$

Cost Function - Intuition

$$h_{\theta}(x) = \theta_1 x$$

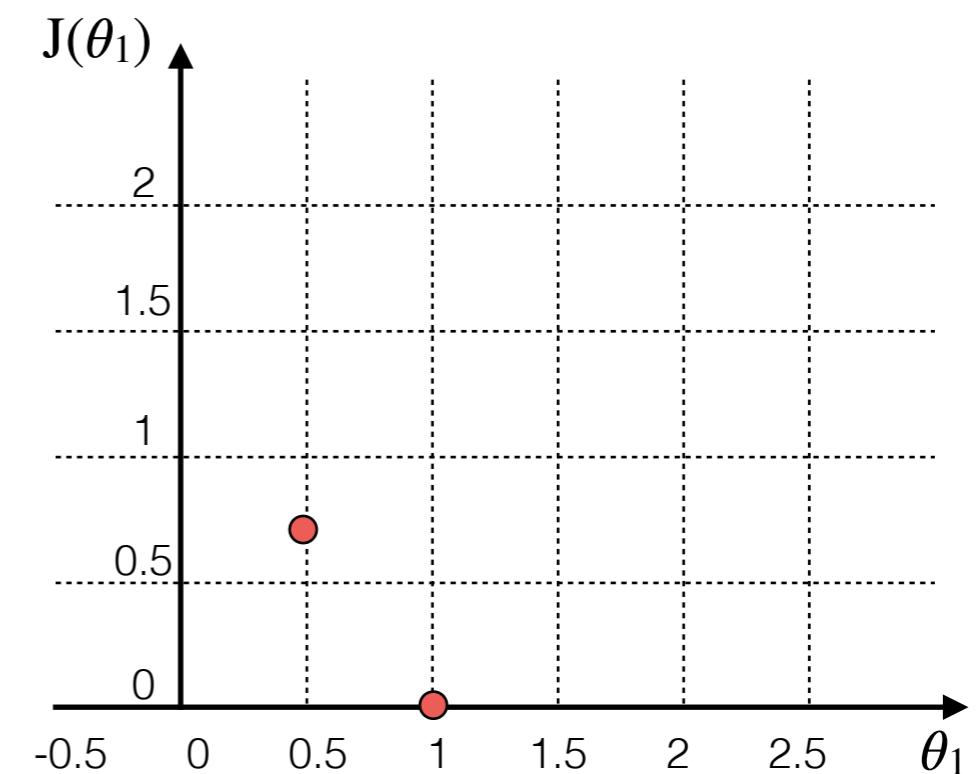
(for fixed θ_1 this is a function of x)



$$\theta_1 = 1/2 ?$$

$$J(\theta_1)$$

(function of the parameter θ_1)

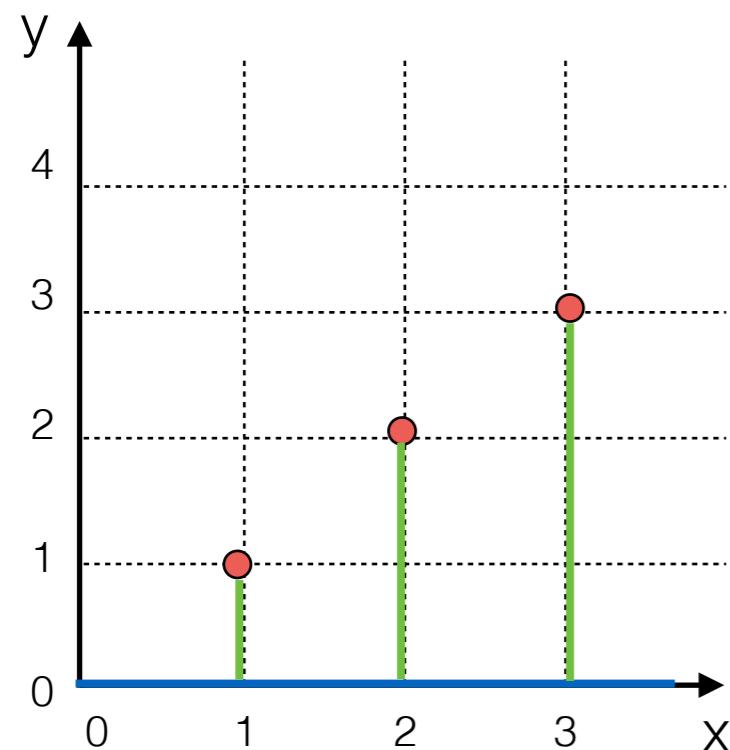


$$J(0.5) = \frac{1}{2m} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] = \frac{1}{6} [0.25 + 1 + 2.25] \approx 0.67$$

Cost Function - Intuition

$$h_{\theta}(x) = \theta_1 x$$

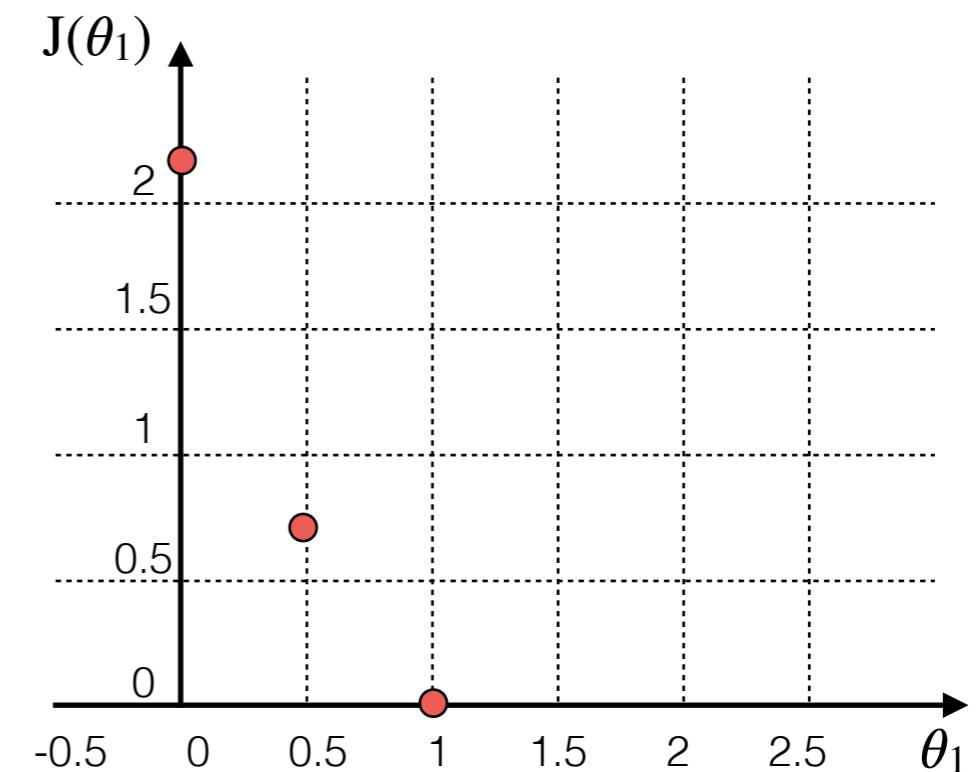
(for fixed θ_1 this is a function of x)



$\theta_1=0$?

$$J(\theta_1)$$

(function of the parameter θ_1)

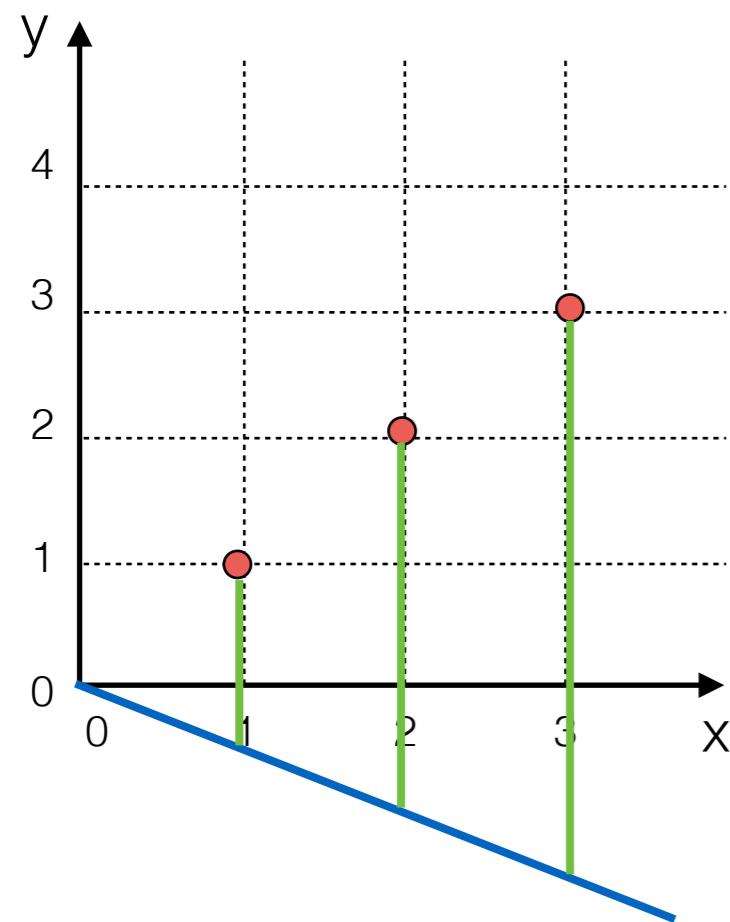


$$J(0) = \frac{1}{2m} [(0-1)^2 + (0-2)^2 + (0-3)^2] = \frac{1}{6} [1+4+9] \approx 2.33$$

Cost Function - Intuition

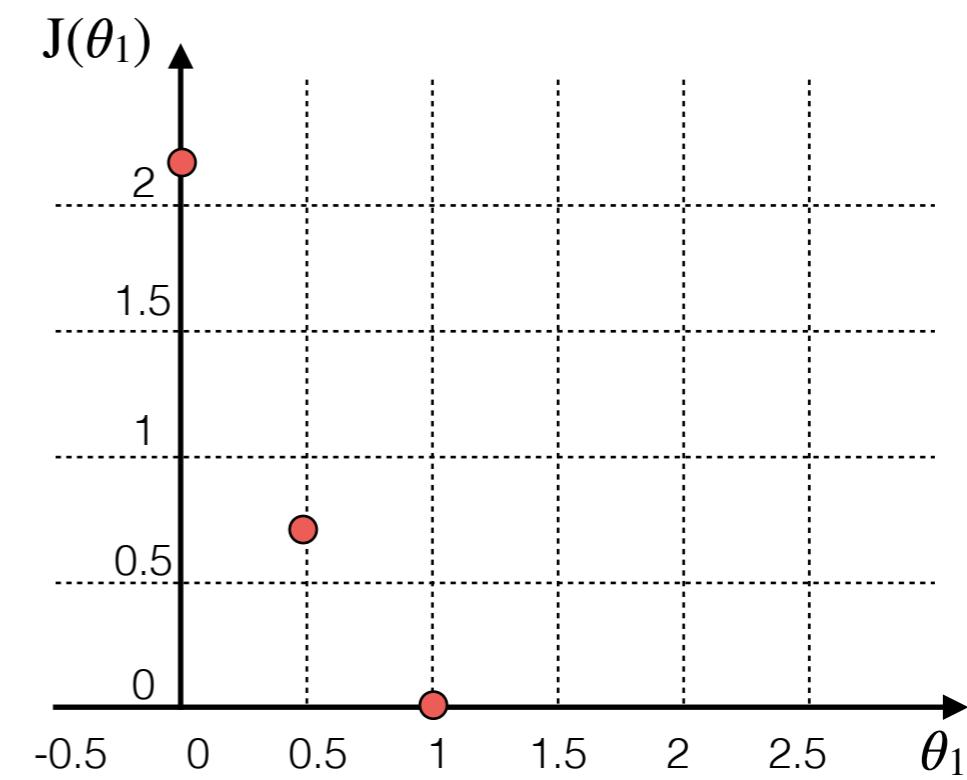
$$h_{\theta}(x) = \theta_1 x$$

(for fixed θ_1 this is a function of x)



$$\bullet J(\theta_1)$$

(function of the parameter θ_1)

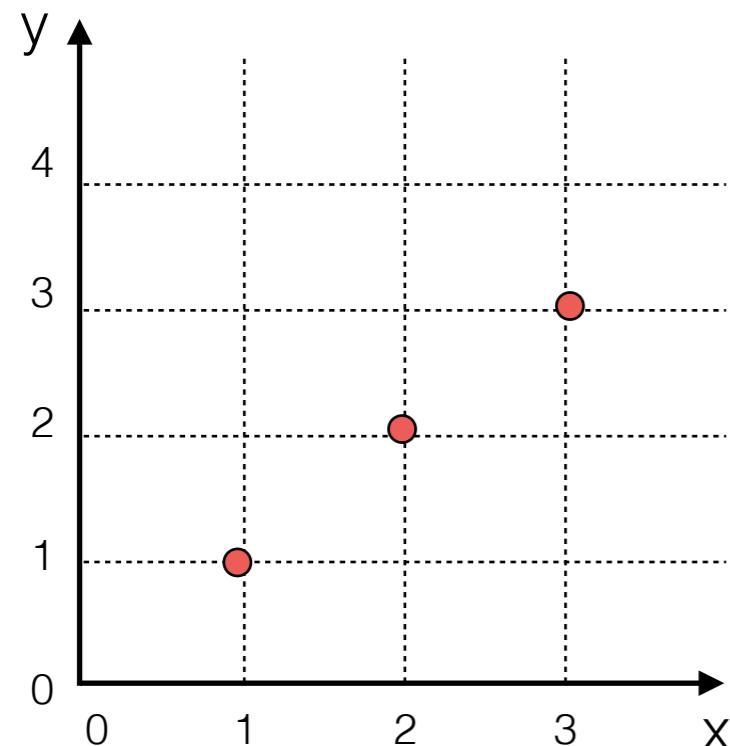


$$J(-0.5) = \frac{1}{2m} [(-0.5-1)^2 + (-1-2)^2 + (-1.5-3)^2] \approx 5.25$$

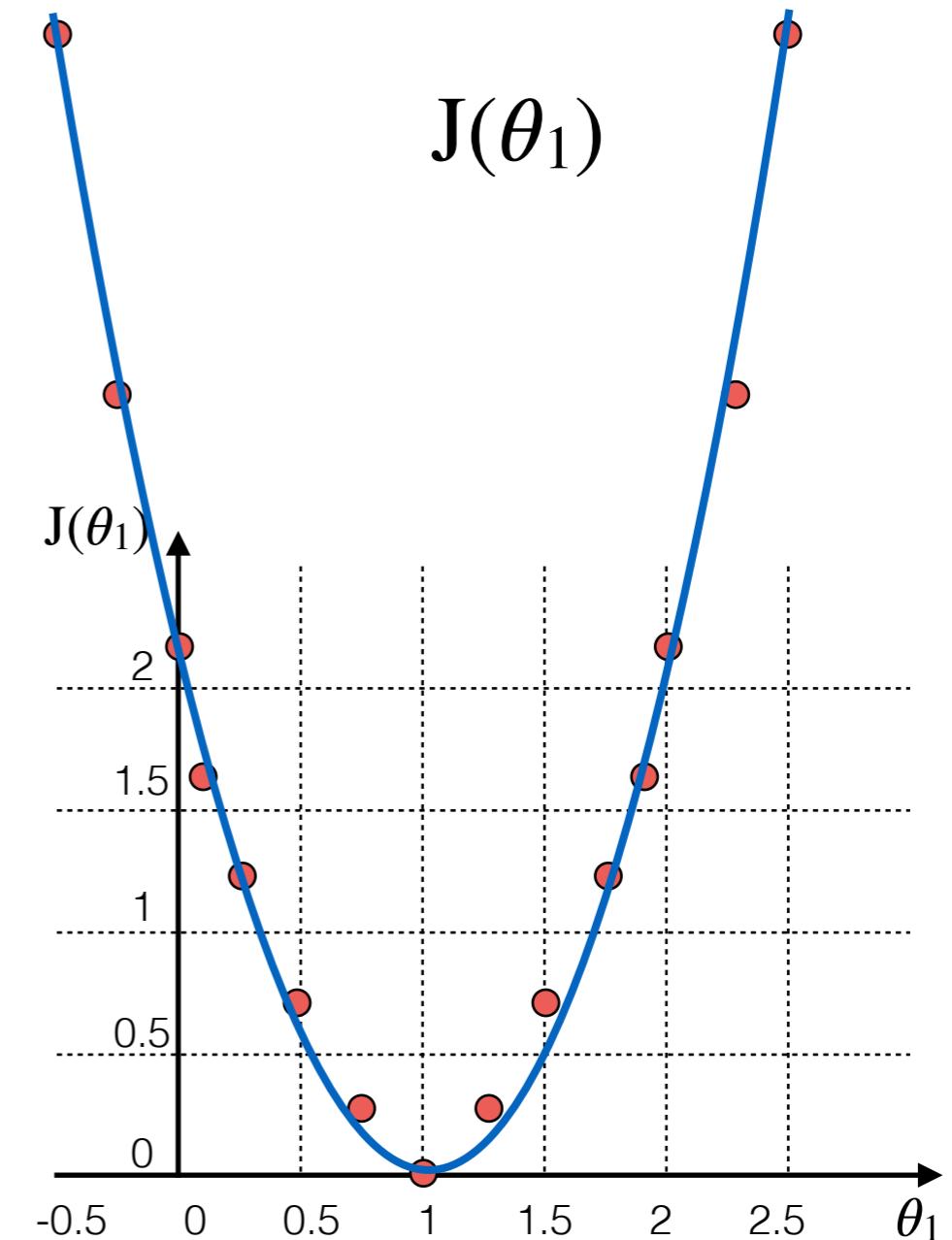
Cost Function - Intuition

$$h_{\theta}(x) = \theta_1 x$$

(for fixed θ_1 this is a function of x)



$$\theta_1 = ?$$

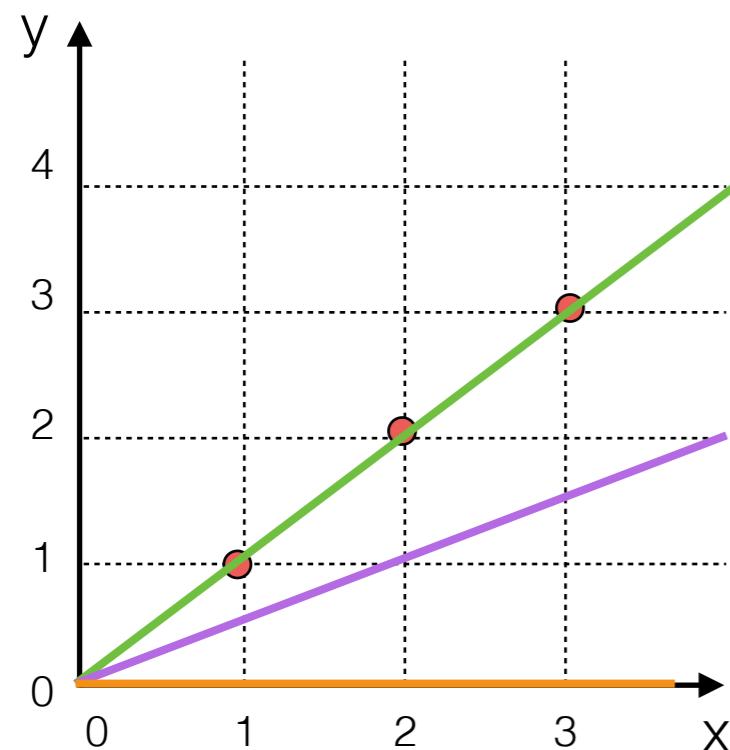


So, for different values, $J(\theta_1) = \dots$

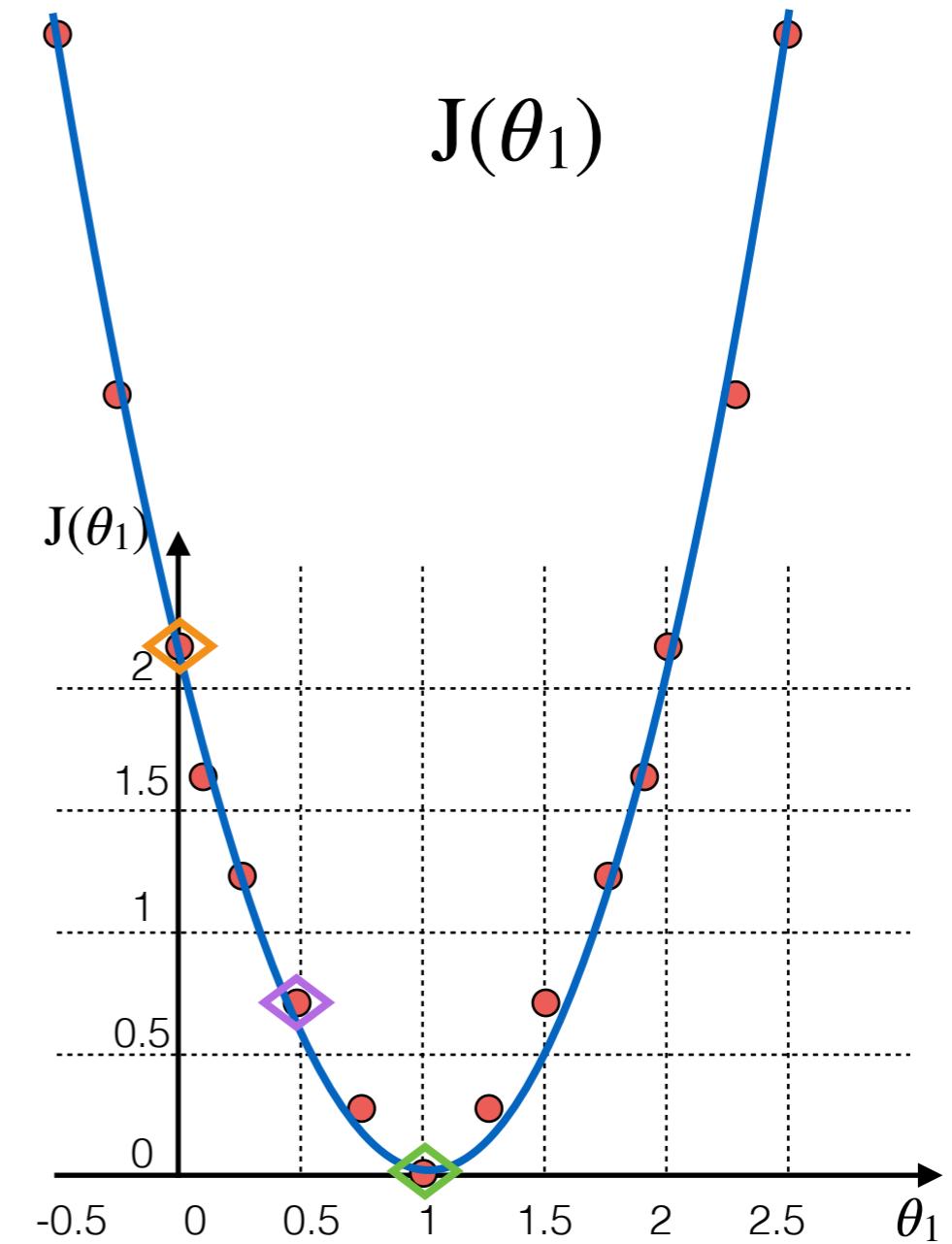
Cost Function - Intuition

$$h_{\theta}(x) = \theta_1 x$$

(for fixed θ_1 this is a function of x)



$\theta_1=1$
 $\theta_1=1/2$
 $\theta_1=0$

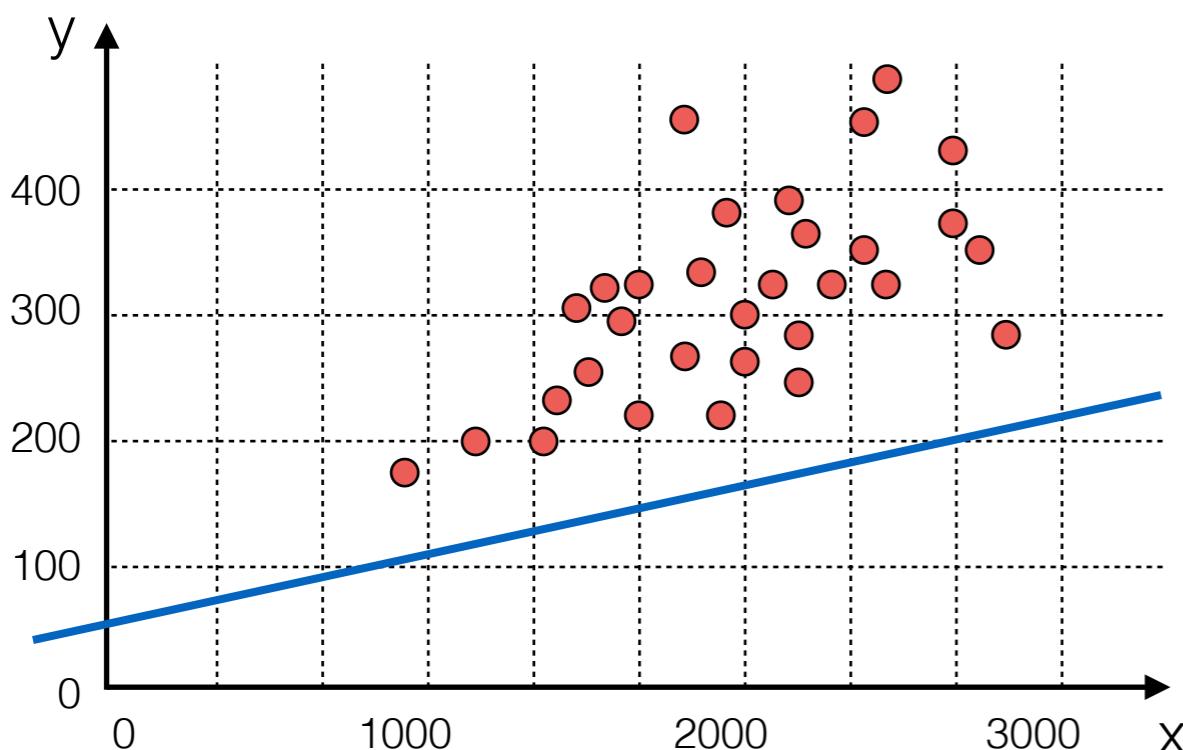


- To recap, each value of θ_1 corresponds to a different hypothesis (i.e. a different line through the data)

Cost Function - Intuition

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

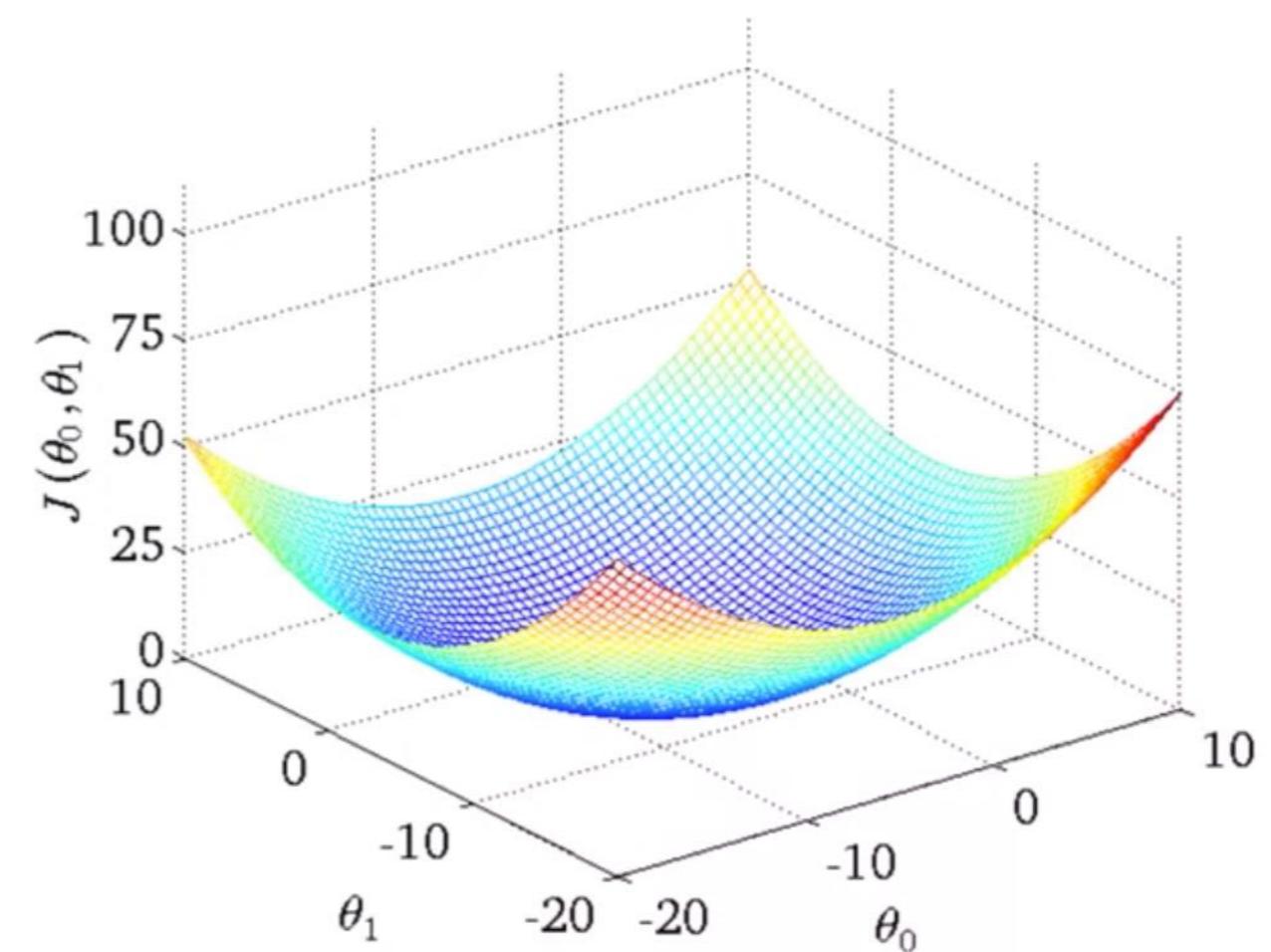
(for fixed θ_0, θ_1 this is a function of x)



$$\theta_0 = 50, \theta_1 = 0.06$$

$$J(\theta_0, \theta_1)$$

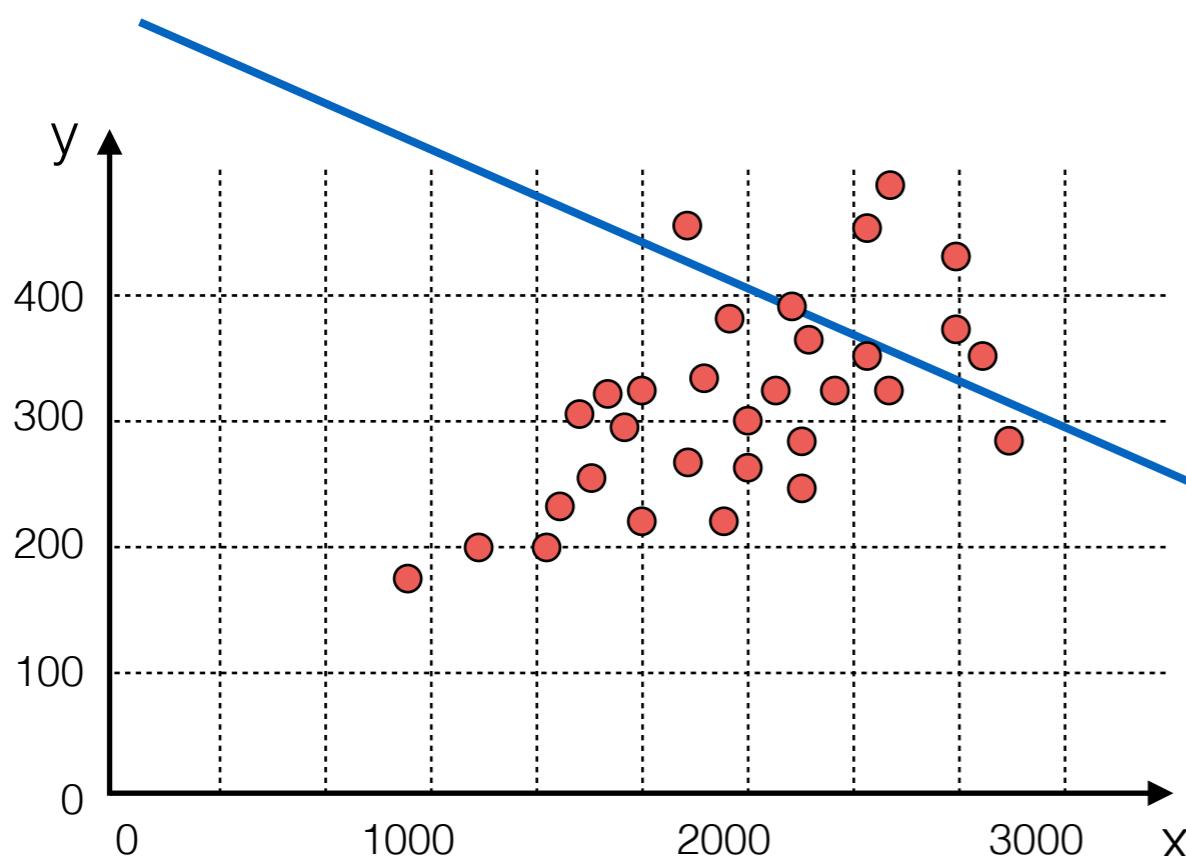
(function of the parameters θ_0, θ_1)



Cost Function - Intuition

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

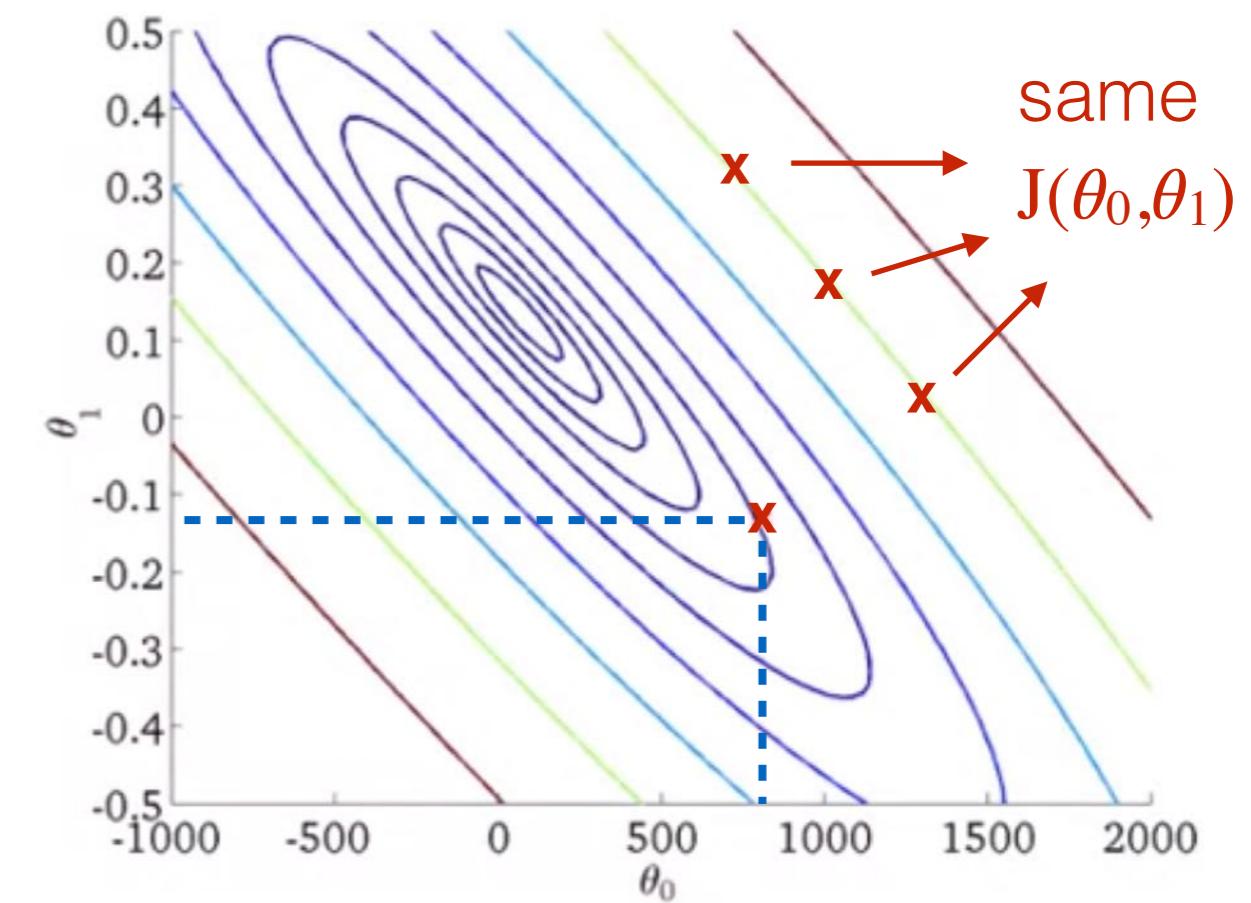
(for fixed θ_0, θ_1 this is a function of x)



- Training data
- Current hypothesis

$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)

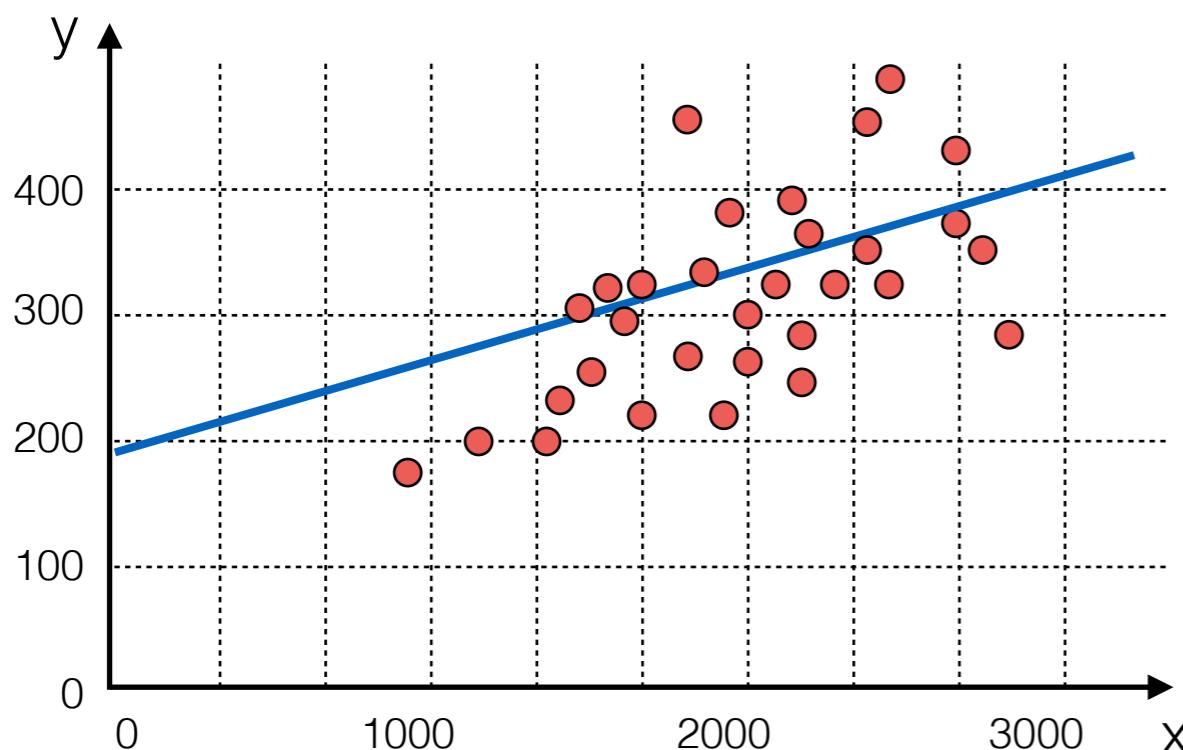


Contour plot

Cost Function - Intuition

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

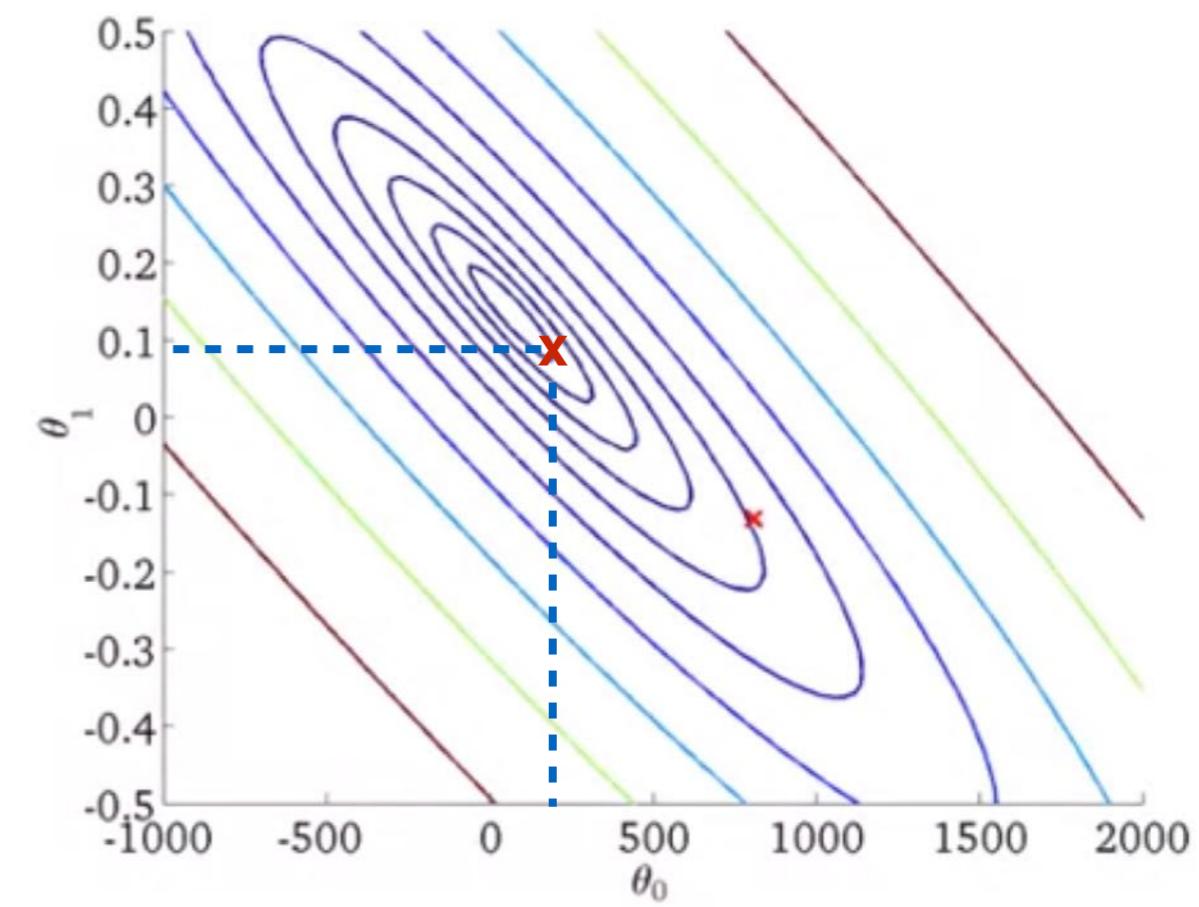
(for fixed θ_0, θ_1 this is a function of x)



- Training data
- Current hypothesis

$$J(\theta_0, \theta_1)$$

(function of the parameters θ_0, θ_1)



Contour plot

Parameter Learning

- What we really want is an efficient algorithm for automatically finding the value of θ_0 and θ_1
- In other words we need an efficient algorithm that minimizes the cost function J

Gradient Descent

- You have a cost function $J(\theta_0, \dots, \theta_n)$
- Objective: minimize $\underset{\theta_0, \dots, \theta_n}{J(\theta_0, \dots, \theta_n)}$
- Outline:
 - ▶ Start with some configuration of parameters θ_i 's
 - ▶ Keep changing θ_i 's parameters to reduce the cost function J until we (hopefully) end up at a minimum

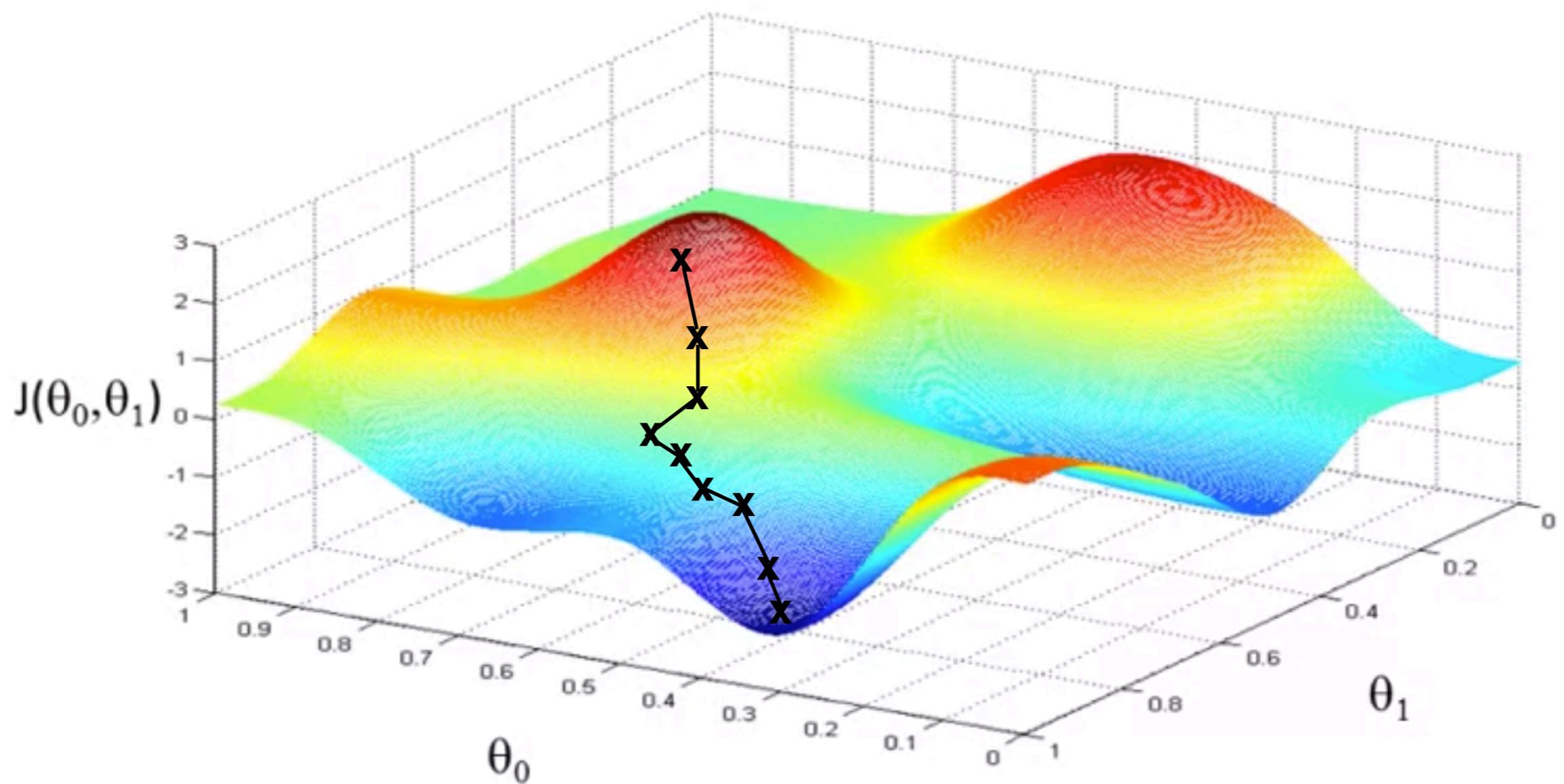
Gradient Descent - Intuition

- “Fog in the mountains” analogy:
 - You are stuck in the mountains and visibility is low
 - Look at the steepness of the hill at you current position and proceed in the direction with the steepest descent



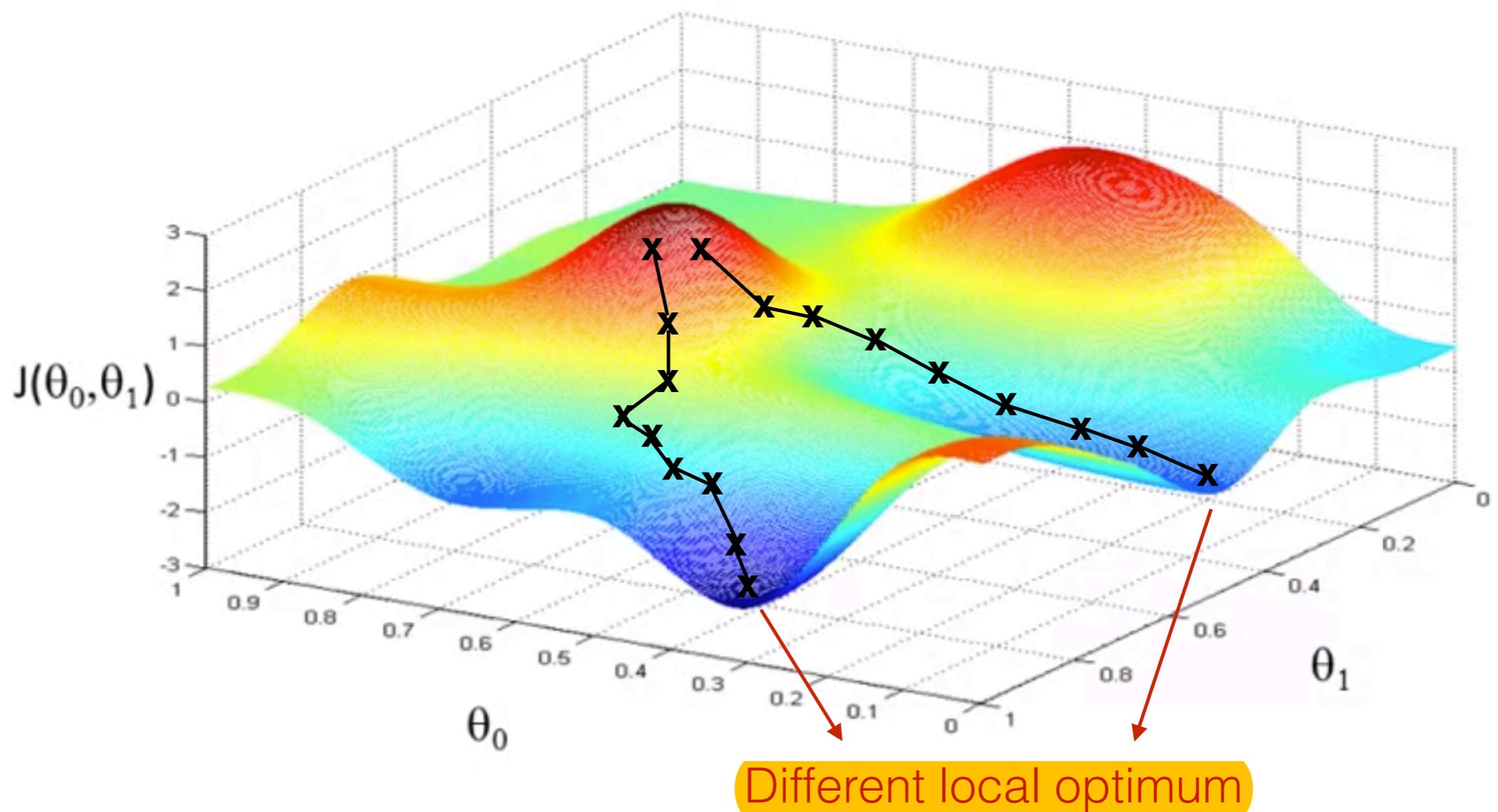
Gradient Descent - Intuition

- Let's start with a visualization



Gradient Descent - Intuition

- Let's start with a visualization



Gradient Descent

- Given $J(\theta_0, \dots, \theta_n)$, we take a small step in the direction of the negative gradient, so that:

$$\theta_{j+1} = \theta_j - \eta \nabla J(\theta_j) \quad \text{where} \quad \theta = \{\theta_0, \dots, \theta_n\}$$

- The parameter $\eta > 0$ is known as the *learning rate*
- After each update, the gradient is re-evaluated for the new weight vector θ_{j+1} , and the process repeated
- Note: the cost function is computed on the entire training set in order to evaluate ∇J (i.e. *batch* method)

Gradient Descent

- A simple implementation (using our example):

```
repeat until convergence {
```

$$\theta_j := \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

```
}
```

Note: simultaneously update θ_0 and θ_1 !

Correct:

$$\text{tmp0} := \theta_0 - \eta \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1);$$

$$\text{tmp1} := \theta_1 - \eta \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1);$$

$$\theta_0 := \text{tmp0}; \quad \theta_1 := \text{tmp1};$$

Incorrect:

$$\theta_0 := \theta_0 - \eta \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1);$$

$$\theta_1 := \theta_1 - \eta \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1);$$

Gradient Descent

- A simple implementation (using our example):

```
repeat until convergence{
```

$$\theta_j := \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{for } j = 0 \text{ and } j = 1)$$

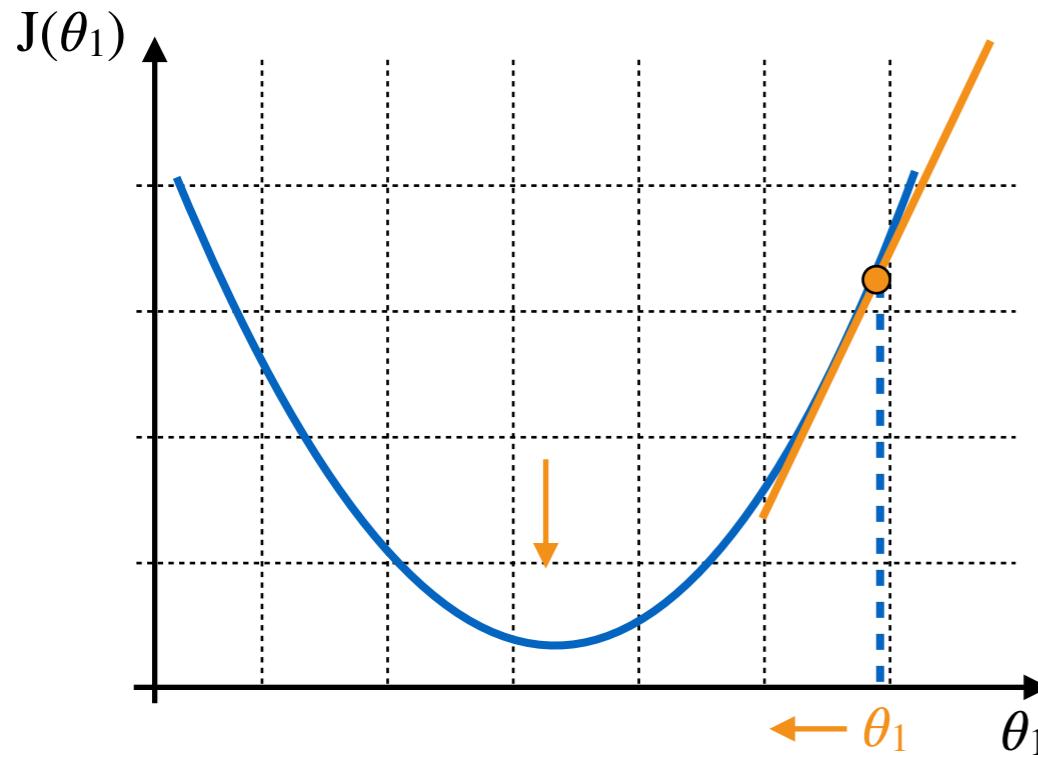
}

learning rate *derivative term*

Simplified version (like we did in a previous lecture):

$$\underset{\theta_1}{\text{minimize}} J(\theta_1) \quad \text{where} \quad \theta_1 \in R$$

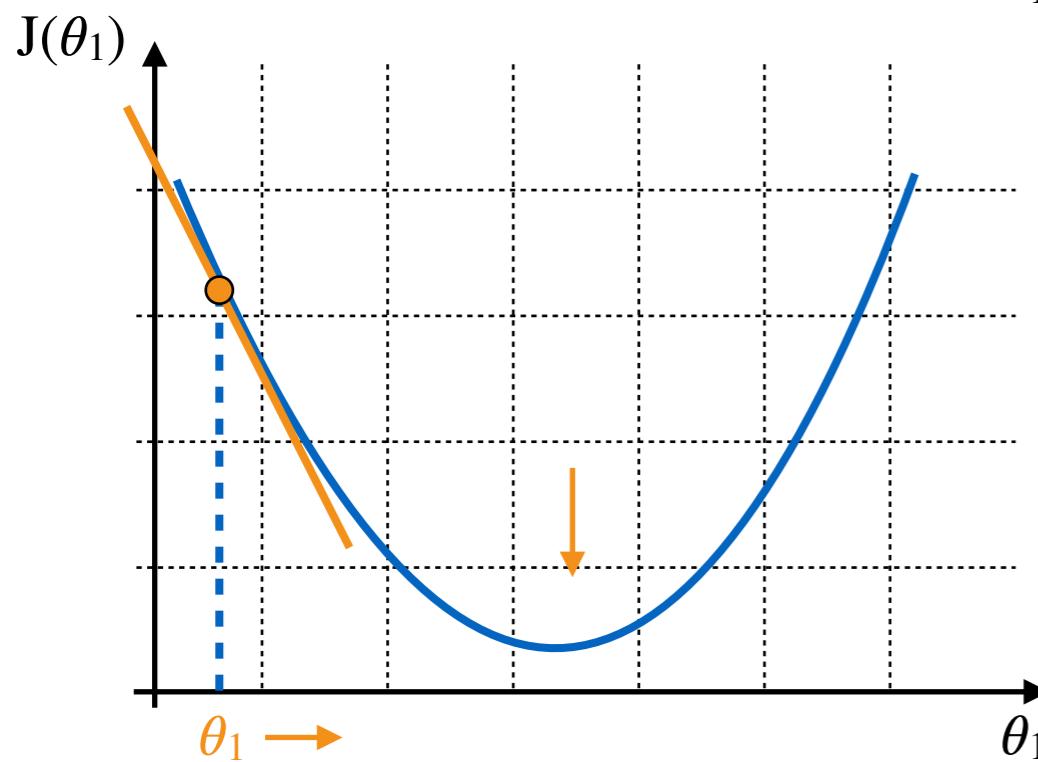
Gradient Descent



$$\theta_1 := \theta_1 - \eta \frac{d}{d\theta_1} J(\theta_1)$$

≥ 0

then we are decreasing θ_1
(in the “right direction”)



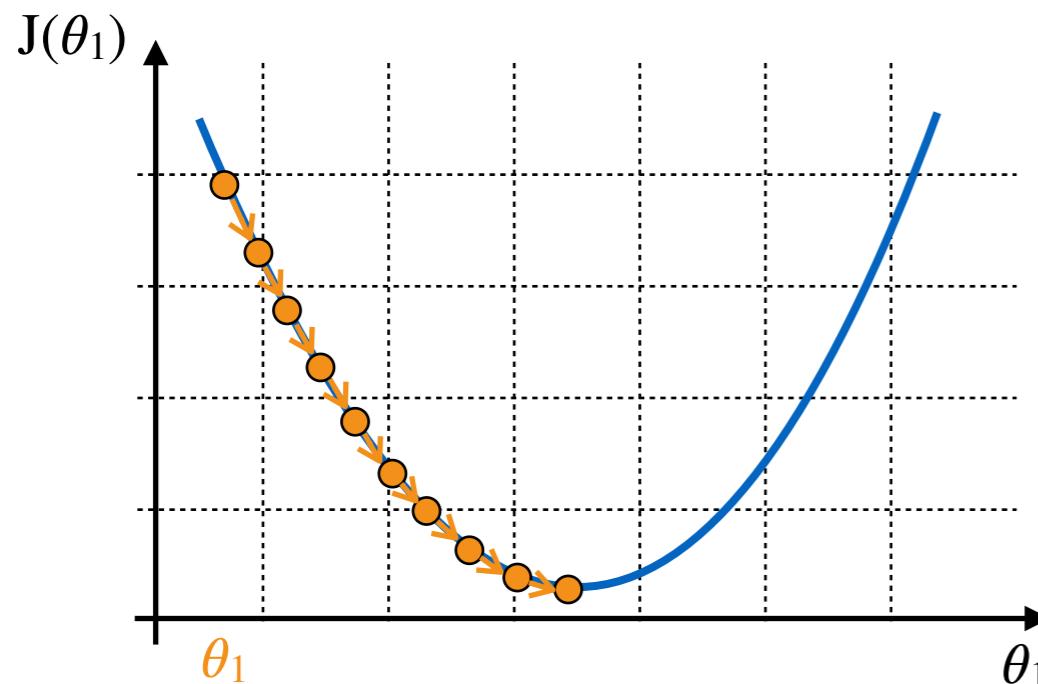
La parte della Derivata indica la DIREZIONE

$$\theta_1 := \theta_1 - \eta \frac{d}{d\theta_1} J(\theta_1)$$

≤ 0

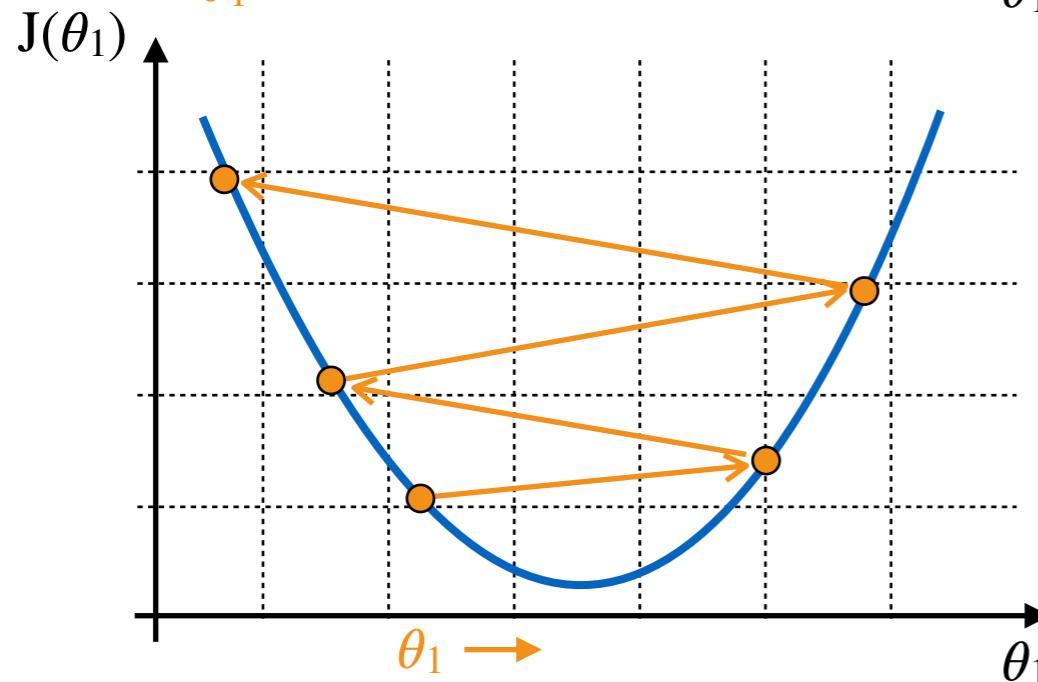
Gradient Descent

- Let's take a look now at the learning rate term η



$$\theta_1 := \theta_1 - \eta \frac{d}{d\theta_1} J(\theta_1)$$

- If η is too small, gradient descent can be slow

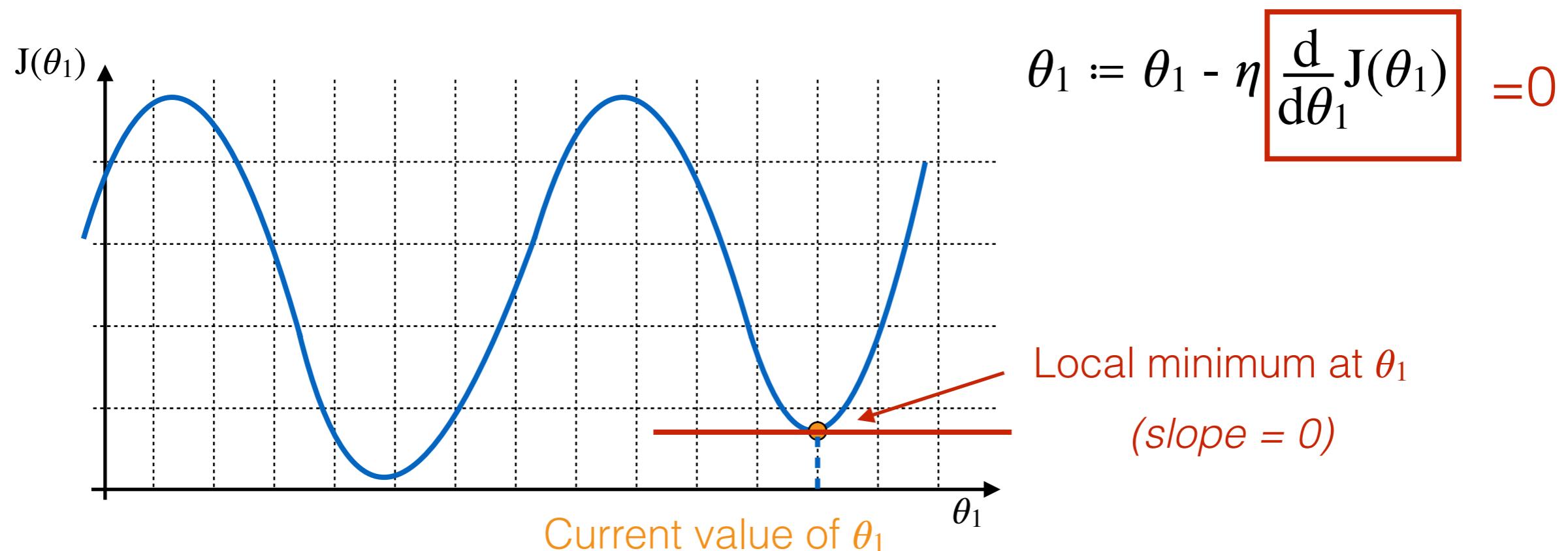


Il parametro controlla l'ampiezza del passo iterativo

- If η is too large, it can overshoot the minimum
(it may fail to converge or even diverge)

Gradient Descent

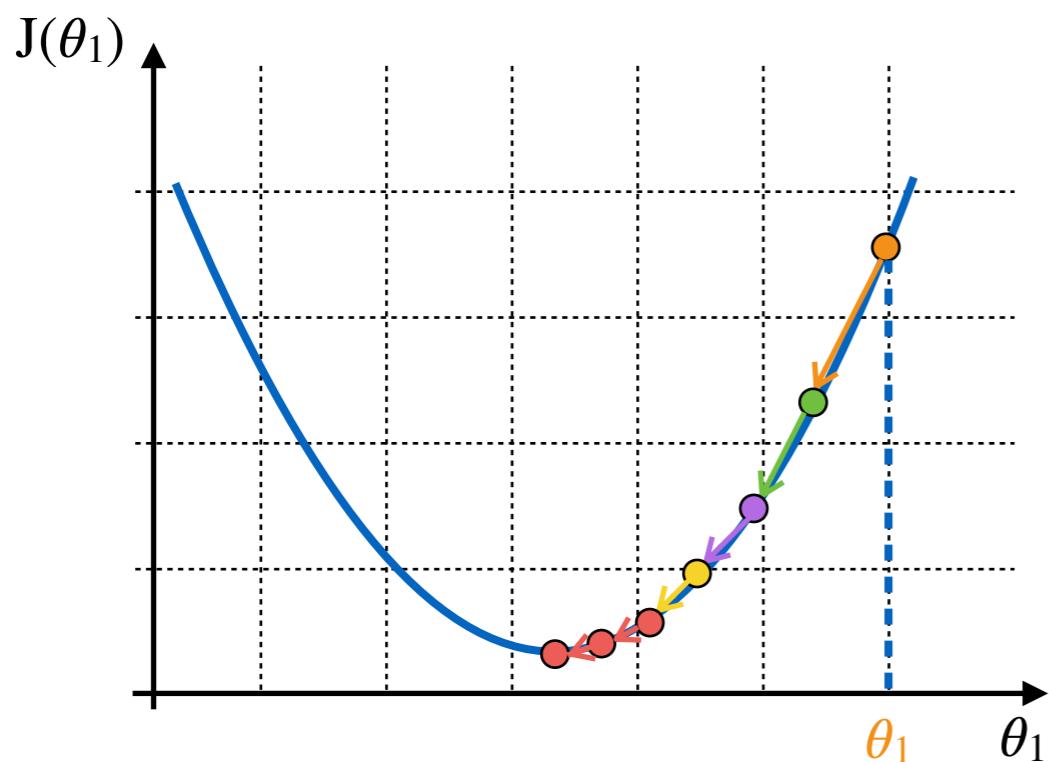
- Q: What if your parameter θ_1 is already at a local minimum?



- ▶ A: one step of gradient descent does not change θ_1

Gradient Descent

- Gradient descent can converge to a local minimum even with a fixed learning rate η



$$\theta_1 := \theta_1 - \eta \frac{d}{d\theta_1} J(\theta_1)$$

- As we approach a local minimum, gradient descent will “naturally” take smaller steps

Linear Regression & Gradient Descent

Linear Regression Model

- Hypothesis:

$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

- Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- Objective: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Gradient Descent

repeat until convergence {

$$\theta_j := \theta_j - \eta \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j = 0$ and $j = 1$)

}

Contact

- **Office:** Torre Archimede, room 6CD3
- **Office hours (ricevimento):** Friday 9:00-11:00

✉ lamberto.ballan@unipd.it
⬆ <http://www.lambertoballan.net>
⬆ <http://vimp.math.unipd.it>
{@} twitter.com/lambertoballan