



# AI 330: Machine Learning

## Fall 2023

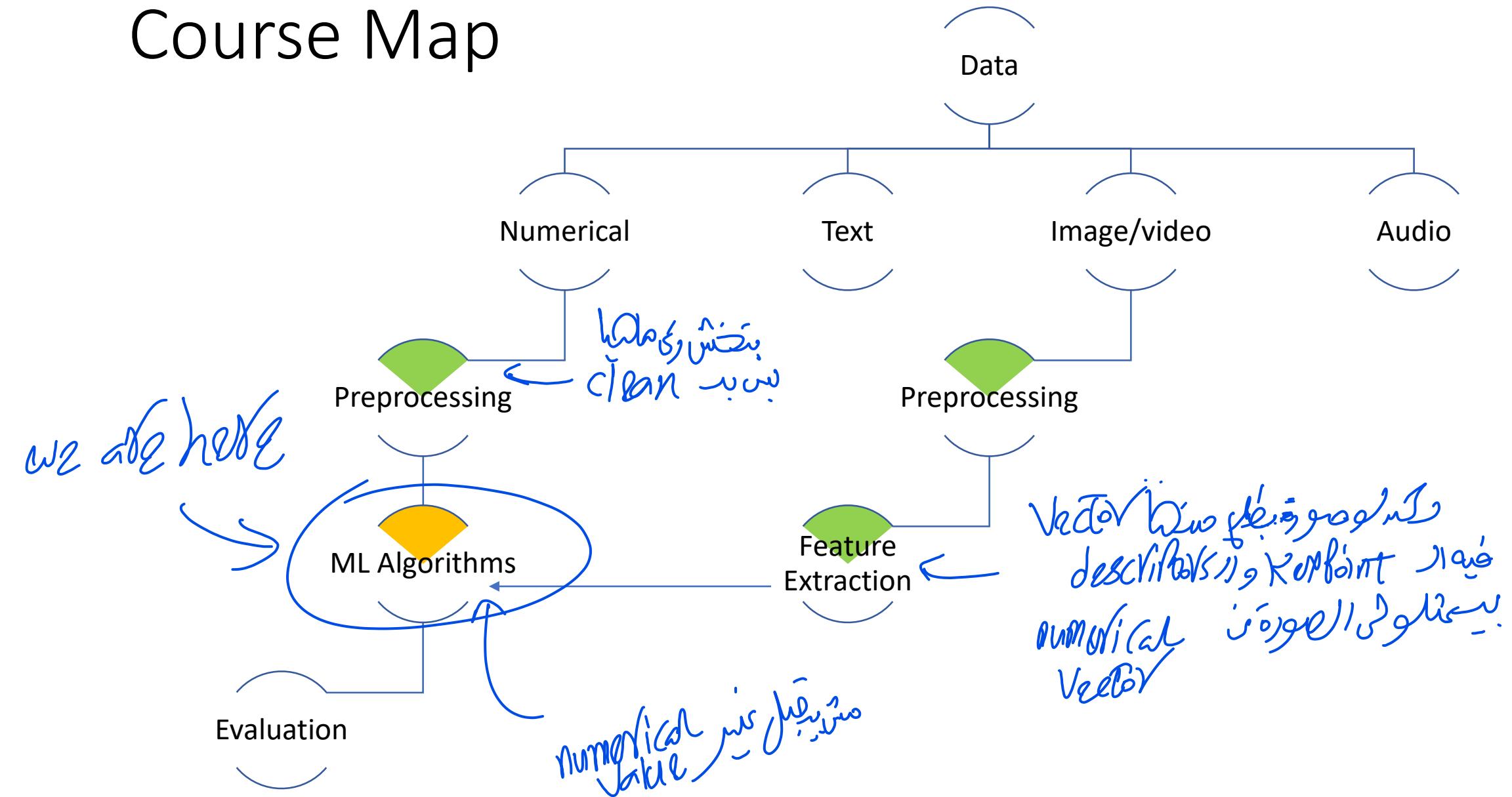
**Dr. Wessam EL-Behaidy**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

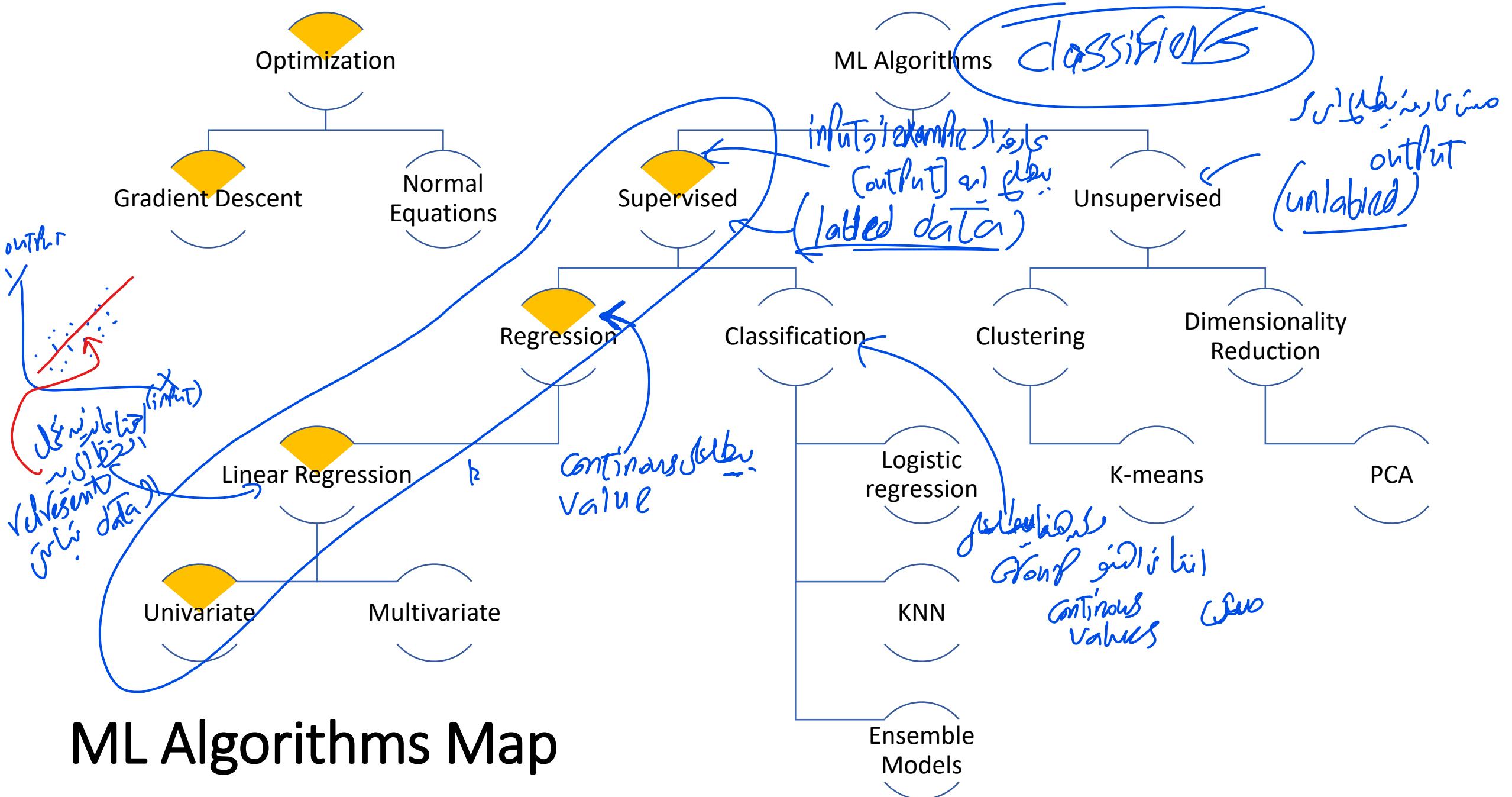
**Dr. Ensaif Hussein**

Associate Professor, Computer Science Department,  
Faculty of Computers and Artificial Intelligence,  
Helwan University.

# Course Map



# ML Algorithms Map



# Lecture 4

---

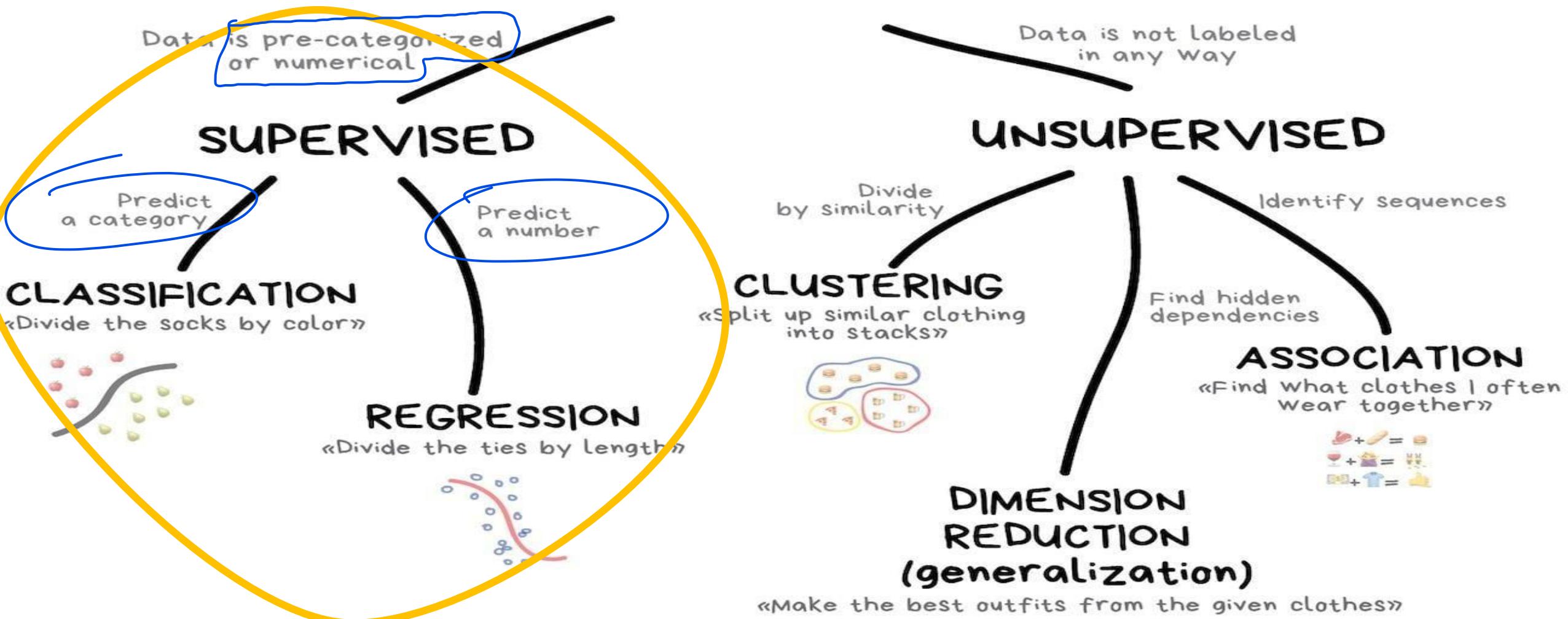
## Univariate Linear Regression & Gradient Descent

Hypothesis – Cost function – Optimization (Gradient Descent)

Slides of:

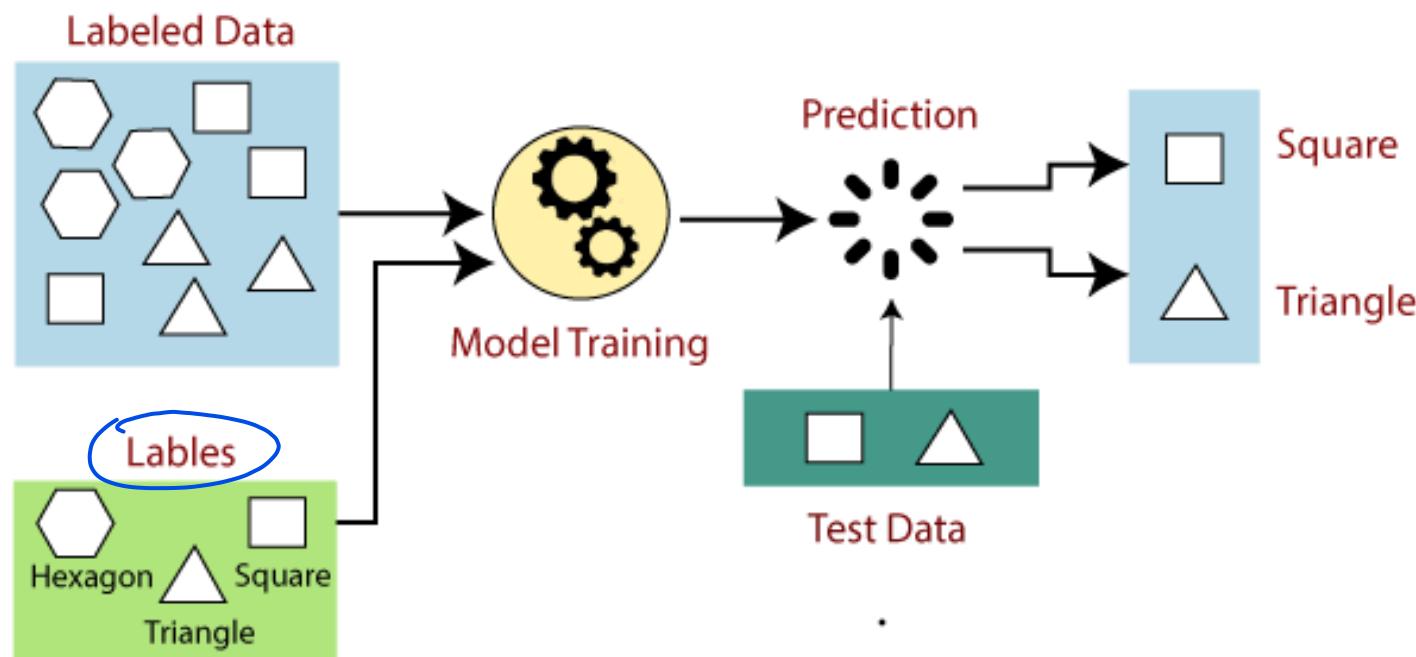
Machine Learning Specialization <https://www.coursera.org/specializations/machine-learning-introduction>  
at Stanford University (**Andrew Ng**)

# CLASSICAL MACHINE LEARNING



# Supervised Learning

- In Supervised learning problems, We are given a (**labelled** or **annotated**) dataset that we already know what our correct output should look like (**true label**), having the idea that there is a relationship between the input and the output.



# Supervised Learning

- Supervised learning problems are categorized into:

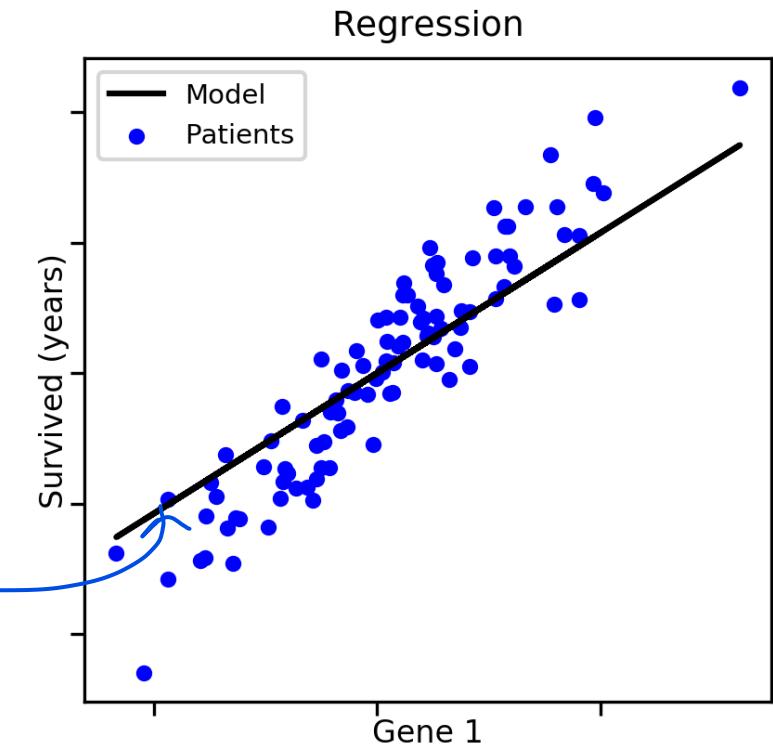
- Regression Problem. → Predict a number
- Classification Problem. → In a category

# Regression Problem

- In a **regression problem**,  
• we are trying to predict results within a **continuous output**, meaning that we are trying to **map input variables** to some **continuous function**.

*Predict a continuous output*

*will be below  
predict value*



<https://aldro61.github.io/microbiome-summer-school-2017/sections/basics/>

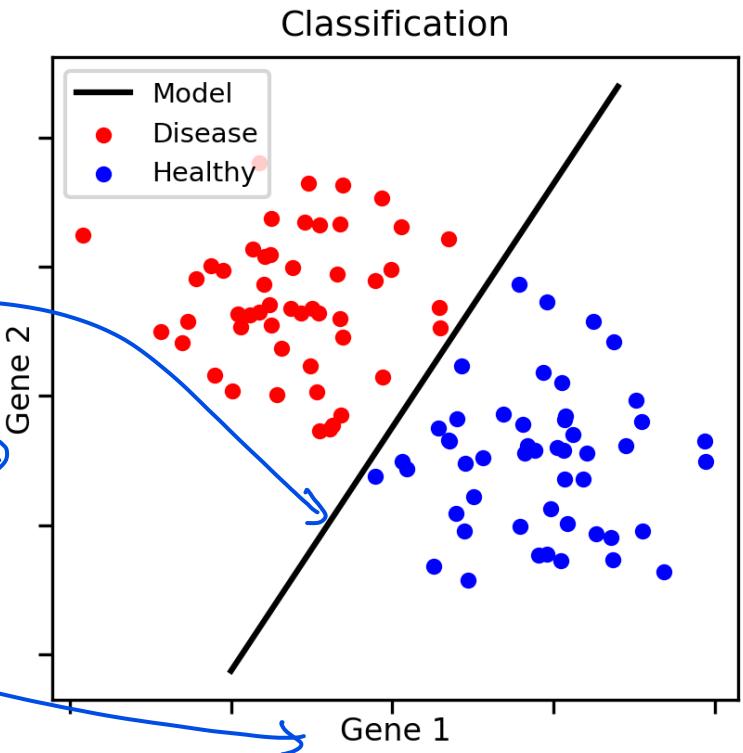
# Classification Problem

- In a **classification problem**,
  - we are trying to predict results in a **discrete output**.

کاربردی که نتایج را در گروه های مختلفی پیش بینی می کند.

Group  
class

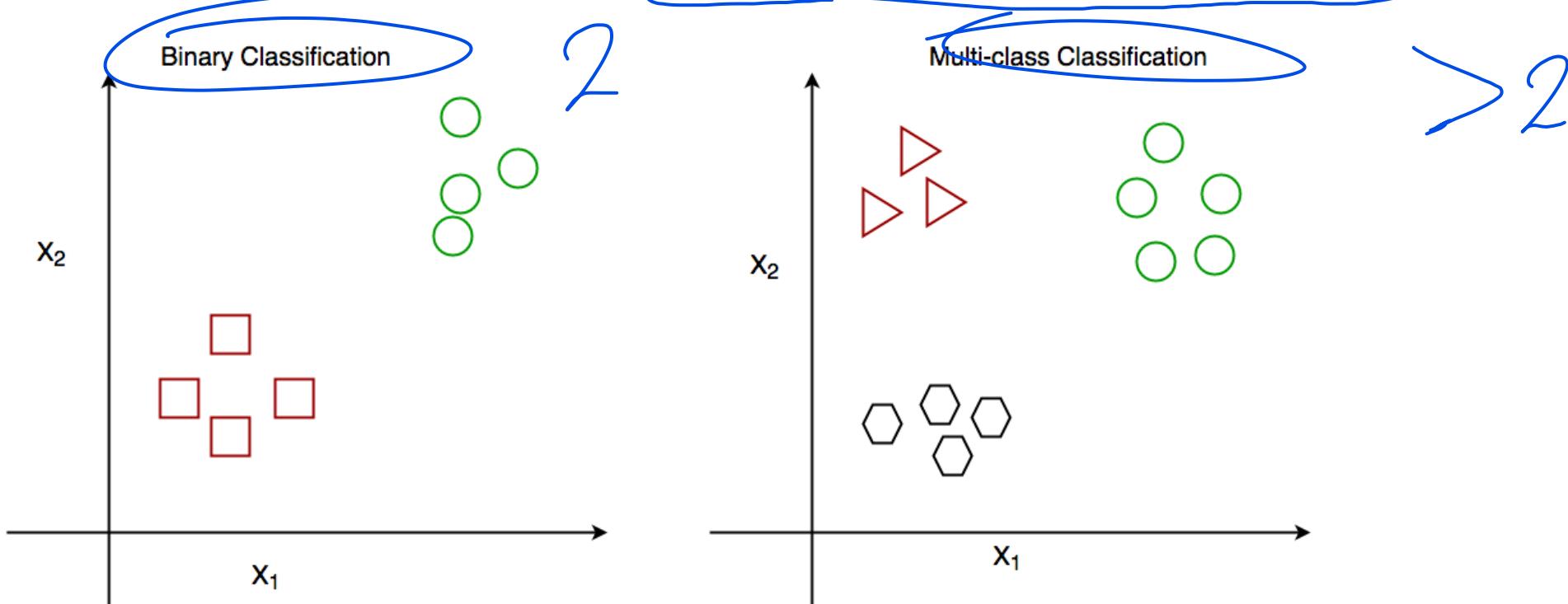
Predict a discrete output



<https://aldro61.github.io/microbiome-summer-school-2017/sections/basics/>

# Classification Problem

- Classification problems can be:
  - Binary classification has only 2 outputs
  - Multi-class classification has more than 2 outputs (3 or more).



# Regression or Classification?

- Given data about the size of houses on the real estate market, try to predict their price. R
- Given a picture of Male/Female, We have to predict his/her age on the basis of given picture. R
- Given a picture of Male/Female, We have to predict Whether She is of High school, College, Graduate age. C
- Banks have to decide whether or not to give a loan to someone on the basis of his credit history. C

Continuous  
Values

discrete  
Values

Groups

R

R

C

C

the first classifier (algorithm)

## Regression Problem:

### Univariate Linear Regression

Linear Regression with One Variable

Prediction is done using one variable

# Linear Regression with One Variable

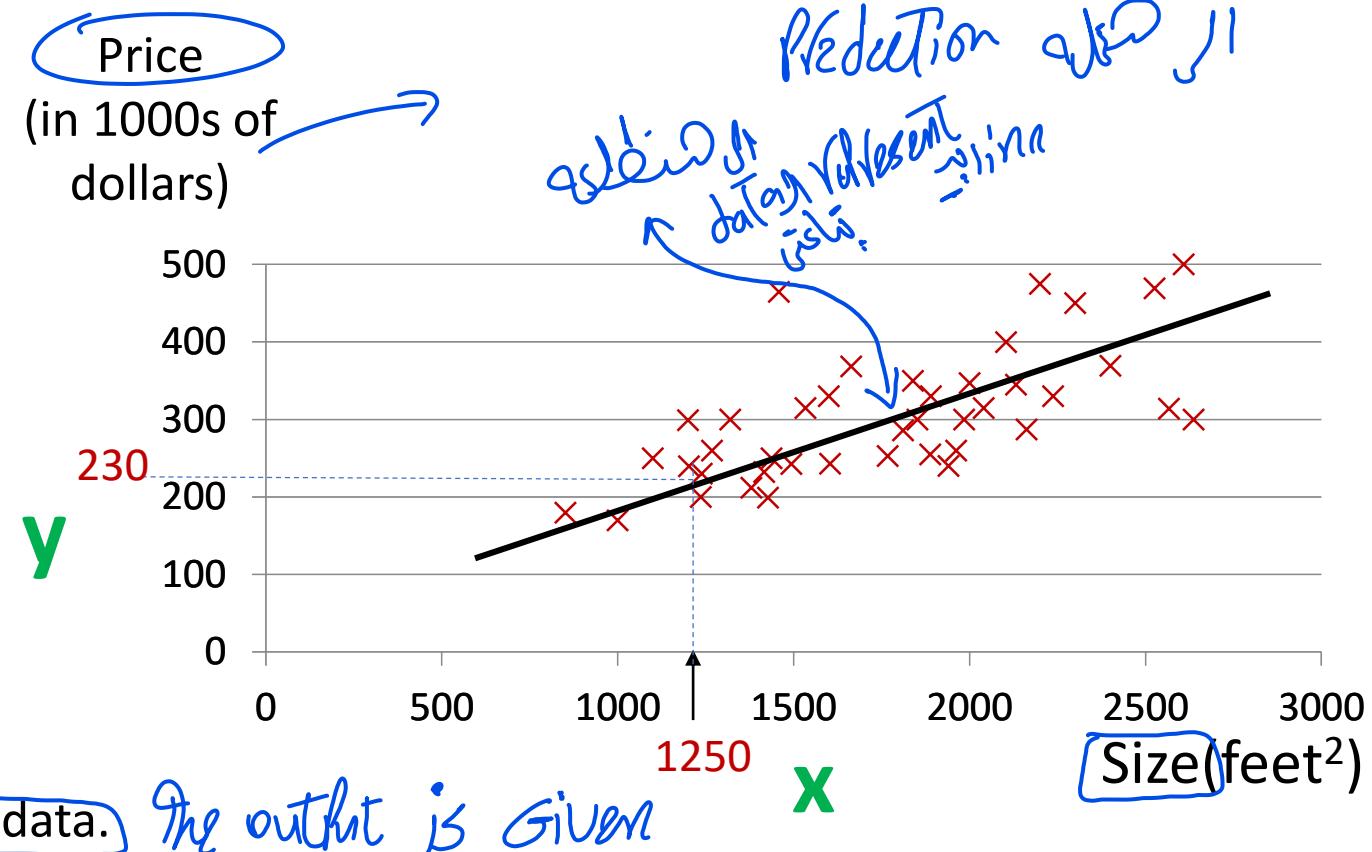
- Linear regression with **one variable** is also known as "**univariate**"
- We want to predict a **single output** value  $y$  from a **single input** value  $x$ .

Housing Prices (Portland, Oregon)

Living area (feet <sup>2</sup> )	Price (1000\$)
2104	400
1600	330
2400	369
1416	232
3000	540
:	:

## Supervised Learning

Given the "right answer" for each example in the data.



The output is given

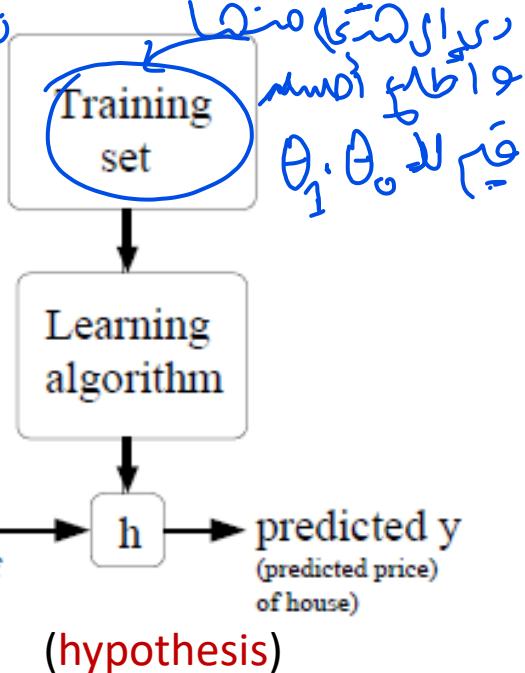
# The Hypothesis Function

$\Theta_1 \rightarrow \Theta_0$  の変換

$$\hat{Y} = mX + b$$

مقدار میانگین  $(\bar{Y})$  را مقدار میانگین  $(m)$  نویسند.

مقدار میانگین  $(\bar{X})$  را مقدار میانگین  $(b)$  نویسند.



- Our hypothesis function has the general form:

form:  
 $\hat{y} = h_{\theta}(x)$  =  $\theta_0 + \theta_1 x$

(detection) Targets (outfit)  
(predict)

features (input)

The diagram shows the equation  $\hat{y} = h_{\theta}(x)$  with handwritten annotations. The term  $\hat{y}$  is circled in blue and labeled 'Targets (outfit)'. The entire equation  $\hat{y} = h_{\theta}(x)$  is highlighted in yellow and labeled '(detection) predict'. The right side of the equation,  $= \theta_0 + \theta_1 x$ , is also highlighted in yellow and labeled '(detection) predict'. The terms  $\theta_0$  and  $\theta_1 x$  are circled in grey and labeled 'features (input)'.

- like the equation of a straight line
  - we are trying to create a function called  $h_{\theta}$  that is trying to map our input data (the  $x$ 's) to our output data (the  $y$ 's).

## Notation:

$x$ 's = “input” variable / features

**y's = “output” variable / “target” variable**

## Parameters/ weights

input $x$	output $y$
0	4
1	7
2	7
3	8

# Example

- Suppose we have the following set of training data:

- Now we can make a random guess about our  $h_\theta$

For example:  $\theta_0, \theta_1 = 2$

and the output of  $h_\theta(1)$  is  $4$ . So  $\theta_0, \theta_1$  are initialized as:

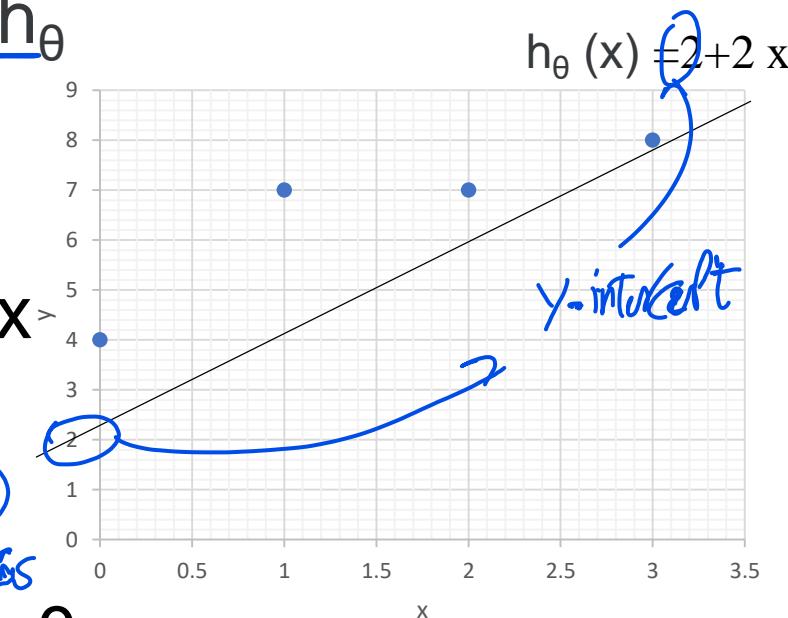
- The hypothesis function becomes  $h_\theta(x) = 2 + 2x$
- So for input of  $x=1$  to our hypothesis  $h_\theta(1)$ ,

$\hat{y} = 4$ . This is off by 3.

$\theta_0, \theta_1$  are initialized as:

Note that we will be trying out various values  $\theta_0, \theta_1$

- to try to find values which provide the best possible "fit" or the most representative "straight line" through the data points mapped on the x-y plane.



Loss Function

## Cost Function

[Squared error function]

w assigned techniques

also loss function

Linear regression → squared error function  
also assigned techniques

- We can measure the accuracy of our hypothesis function by using a cost function.
- This takes an average of all the results of the hypothesis with inputs from x's compared to the actual output y's.

just examples

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

Squared error function

loss function minimize J<sup>i,j</sup>

Objective function (Our goal) :  $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Choose  $\theta_0, \theta_1$  so that  $h_\theta(x)$  is close to output y for our training examples (x,y)

linear regression // learning technique

Notation

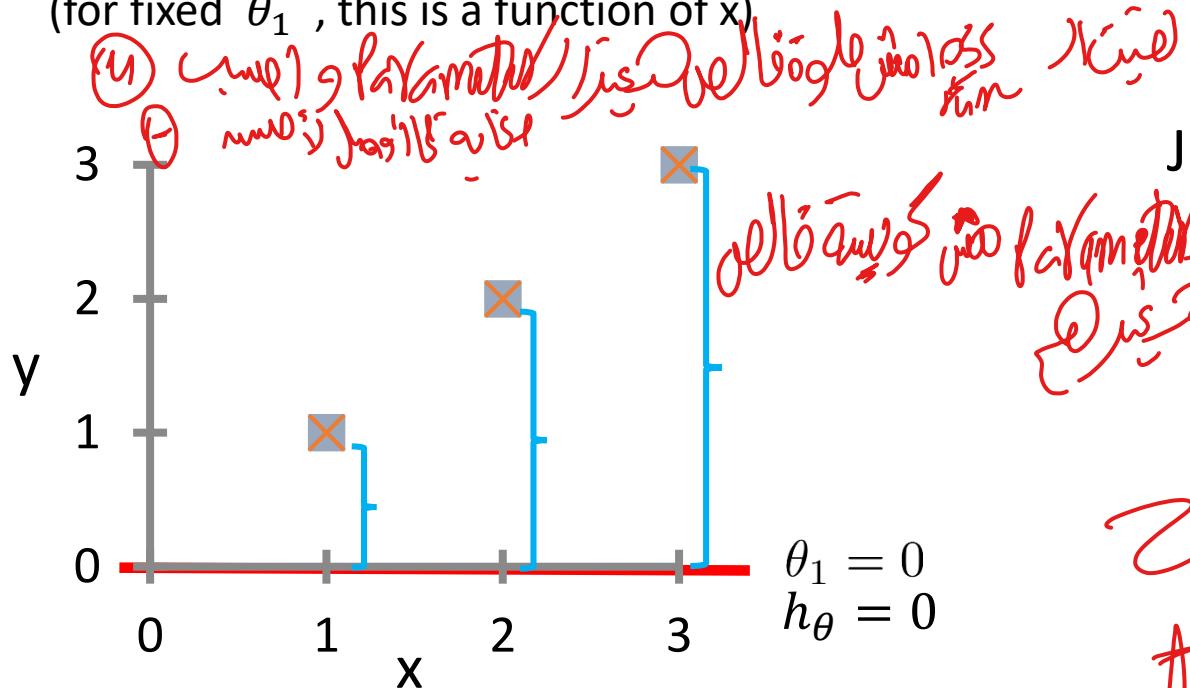
m = Number of training examples

$(x^{(i)}, y^{(i)})$  =  $i^{th}$  training example

# Simplified Example

- ①  $z \rightarrow \theta_1 x \rightarrow \theta_1$  افترض  $\theta_1$  لـ  $h_{\theta}$
- ②  $\theta_1 = 0$  لـ  $h_{\theta}$  Hypothesis
- ③  $h_{\theta}(x)$  squared error loss function

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



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

input $x$	output $y$
1	1
2	2
3	3

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

Where  $m$  (# of samples)= 3

$$\begin{aligned} J(0) &= \frac{1}{6} ((0-1)^2 + (0-2)^2 + (0-3)^2) \\ &= \frac{1}{6} (1+4+9) \approx 2.3 \end{aligned}$$

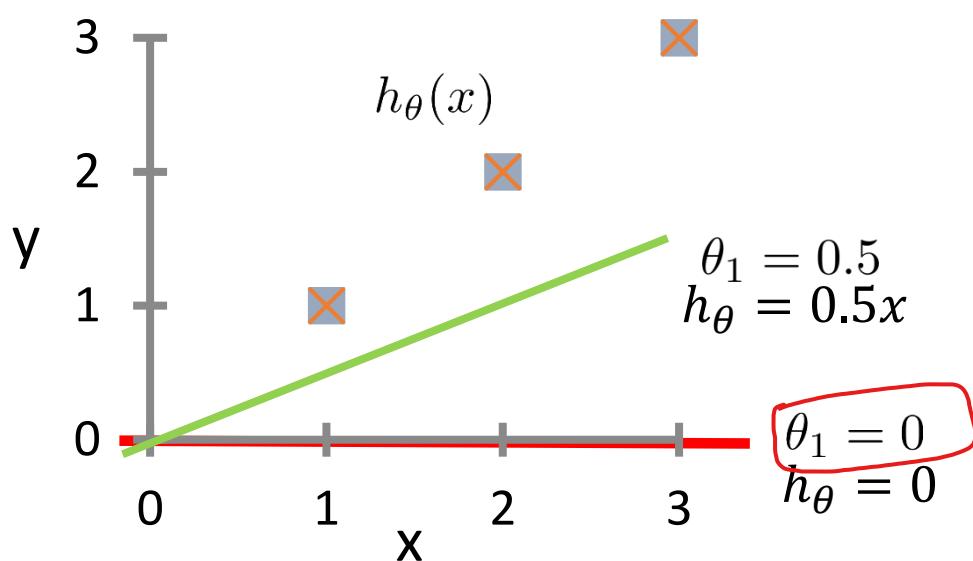
زنون )  
 $J(1)$

# Simplified Example

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

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

$\theta_1 = \frac{1}{2}$



$$h_{\theta} = \frac{1}{2}x$$

input $x$	output $y$
1	1
2	2
3	3

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

*square error*

Where m (# of samples)= 3

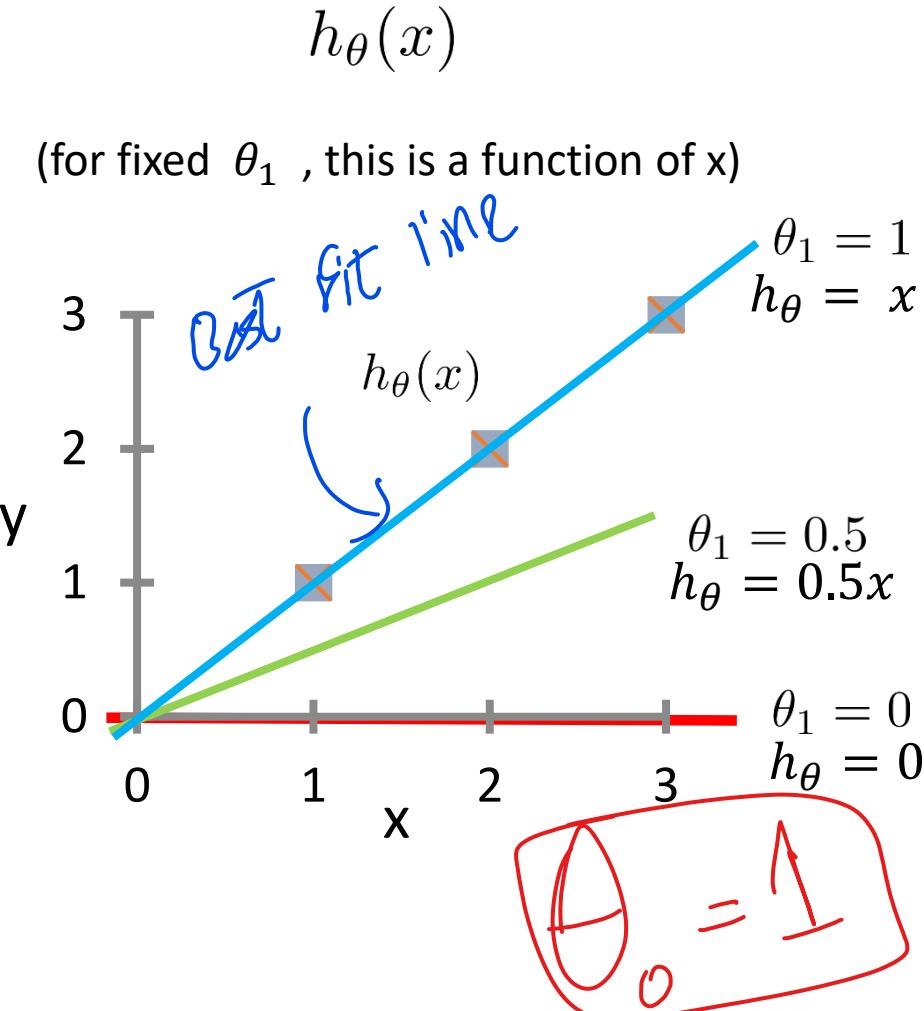
$$\begin{aligned} J(0) &= \frac{1}{6} ((0-1)^2 + (0-2)^2 + (0-3)^2) \\ &= \frac{1}{6} (1+4+9) \approx 2.3 \\ J(0.5) &= \frac{1}{6} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2) \\ &= \frac{1}{6} (0.25+1+2.25) \approx 0.58 \end{aligned}$$

$$\begin{aligned} J(0.5) &= \frac{1}{6} \left[ \left( \frac{1}{2} \times 1 - 1 \right)^2 + \left( \frac{1}{2} \times 2 - 2 \right)^2 + \left( \frac{1}{2} \times 3 - 3 \right)^2 \right] = \frac{1}{6} ( \frac{1}{4} + 1 + 2.25 ) \\ &= \frac{3.5}{6} = 0.58 \end{aligned}$$

Andrew Ng 18

# Simplified hypothesis

input $x$	output $y$
1	1
2	2
3	3



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

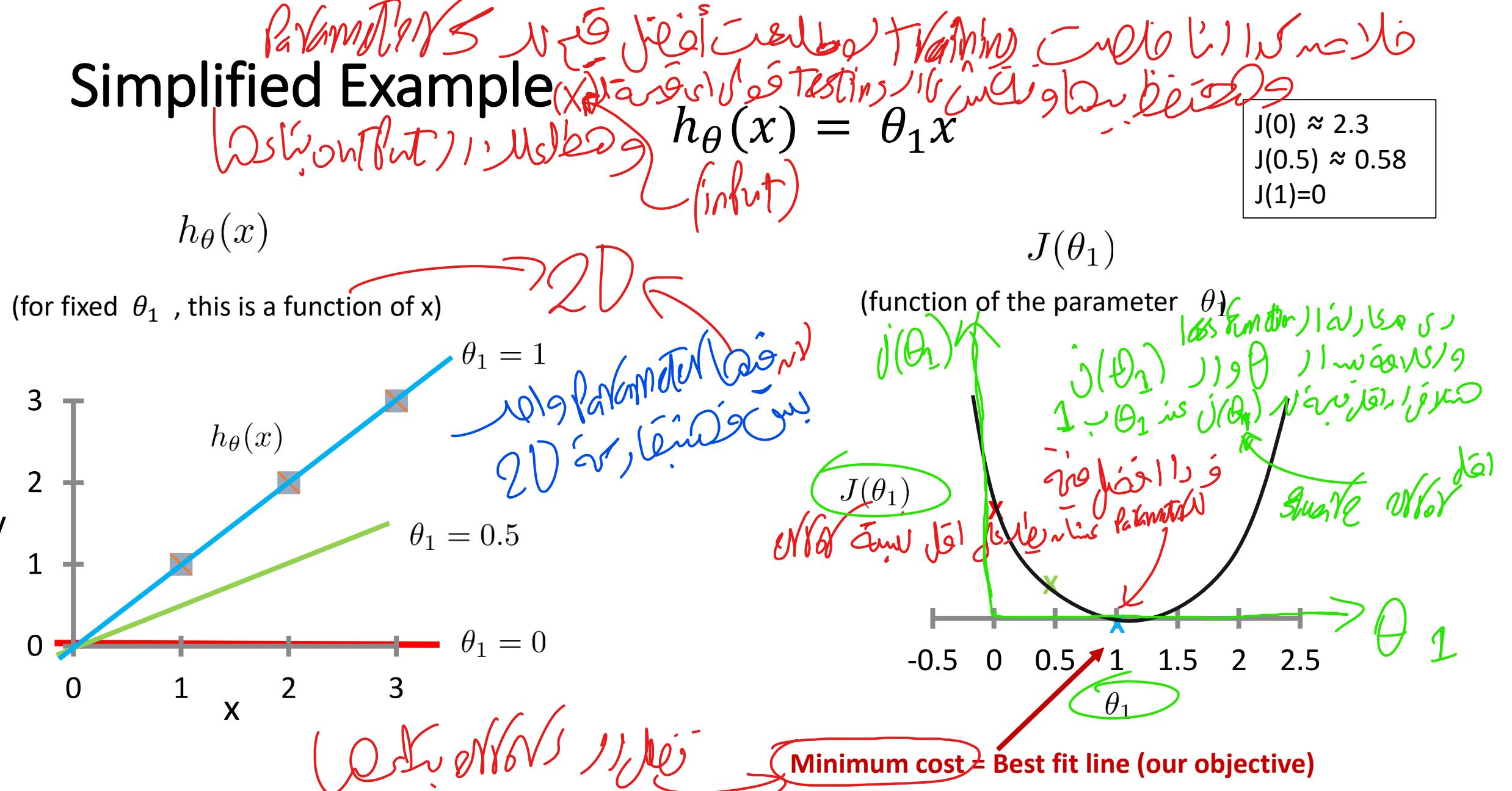
Where m (# of samples)= 3

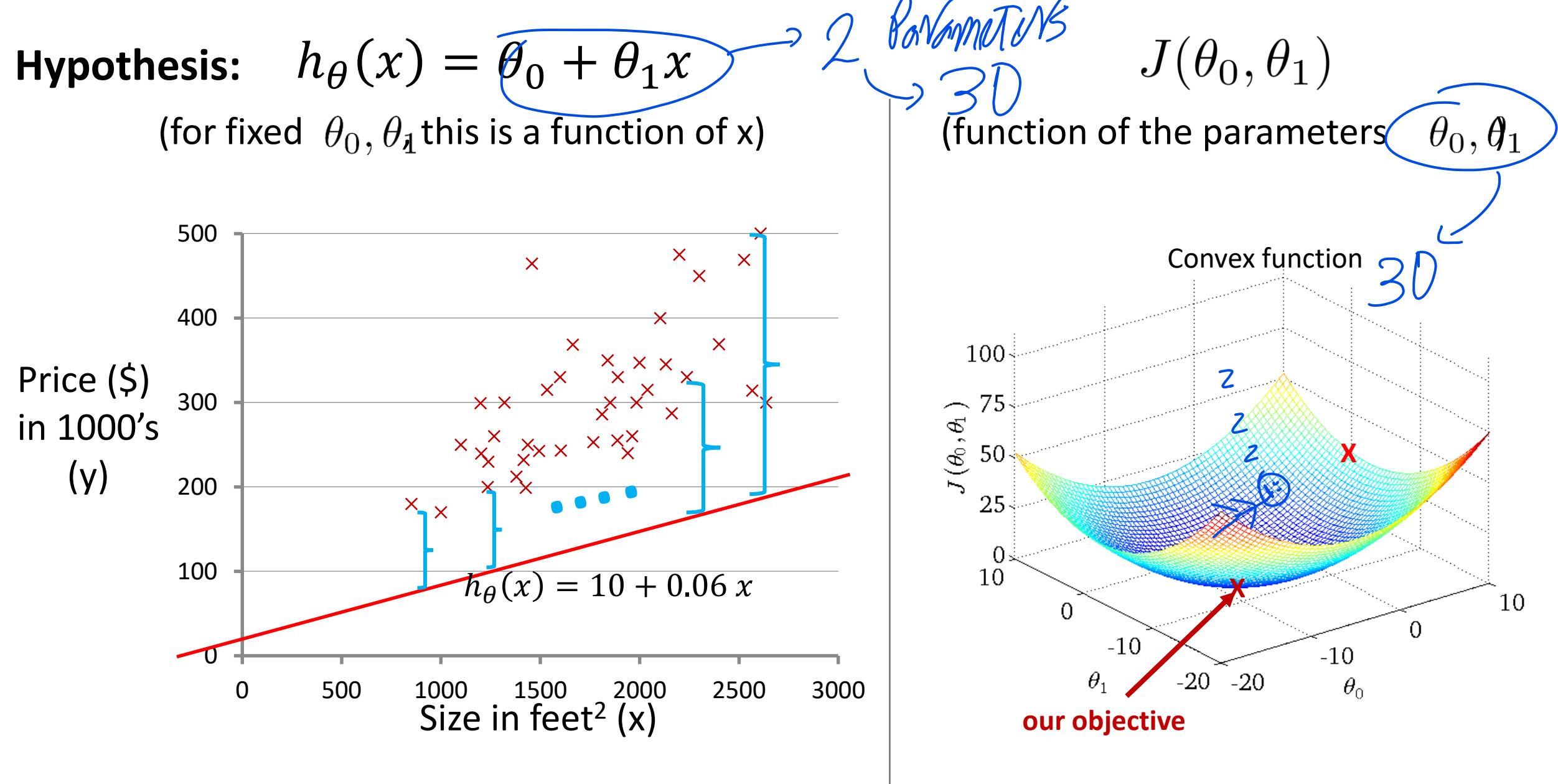
$$\begin{aligned} J(0) &= \frac{1}{6} ((0-1)^2 + (0-2)^2 + (0-3)^2) \\ &= \frac{1}{6} (1+4+9) \approx 2.3 \end{aligned}$$

$$\begin{aligned} J(0.5) &= \frac{1}{6} ((0.5-1)^2 + (1-2)^2 + (1.5-3)^2) \\ &= \frac{1}{6} (0.25+1+2.25) \approx 0.58 \end{aligned}$$

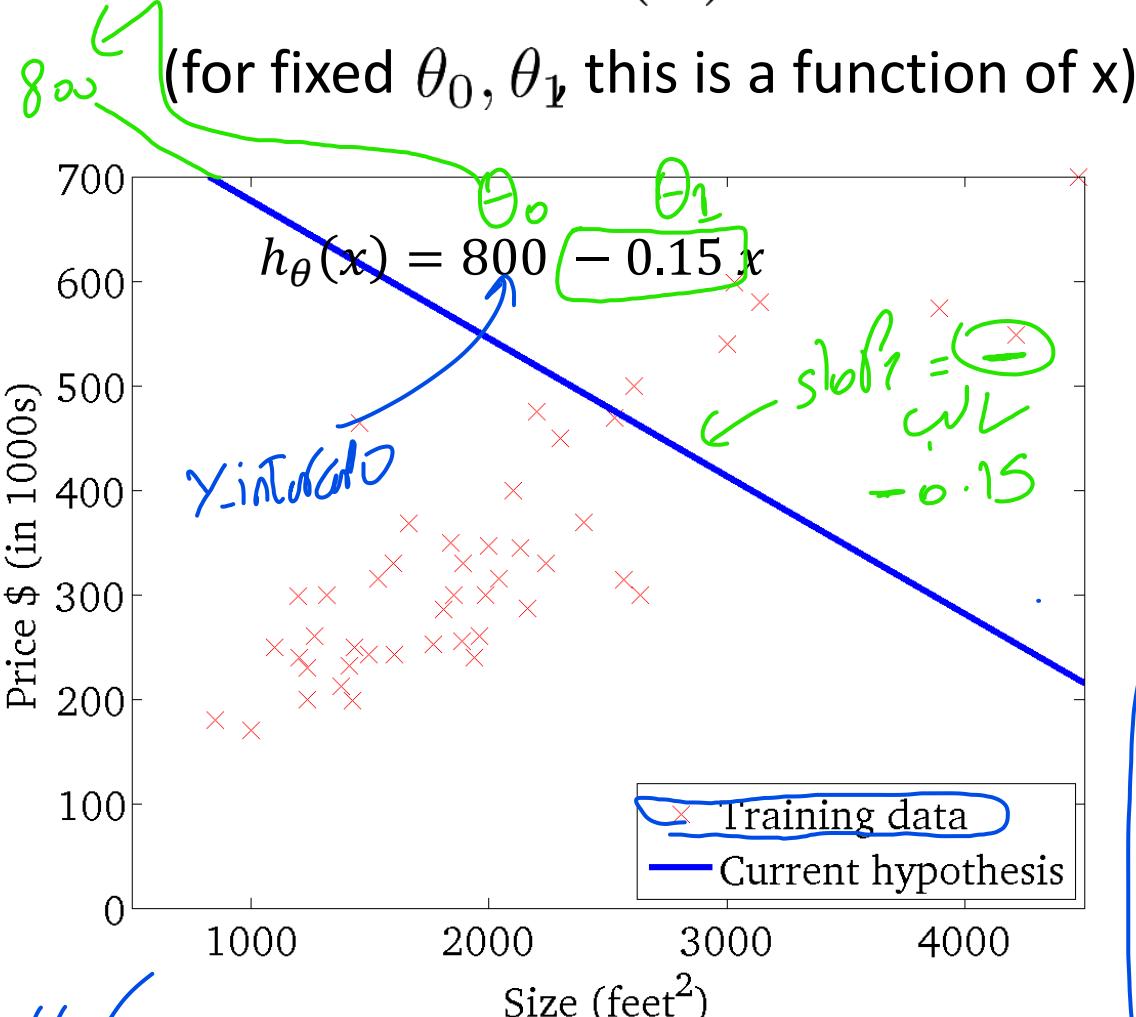
$$\begin{aligned} J(1) &= \frac{1}{6} ((1-1)^2 + (2-2)^2 + (3-3)^2) \\ &= \frac{1}{6} (0+0+0) = 0 \end{aligned}$$

# Simplified Example



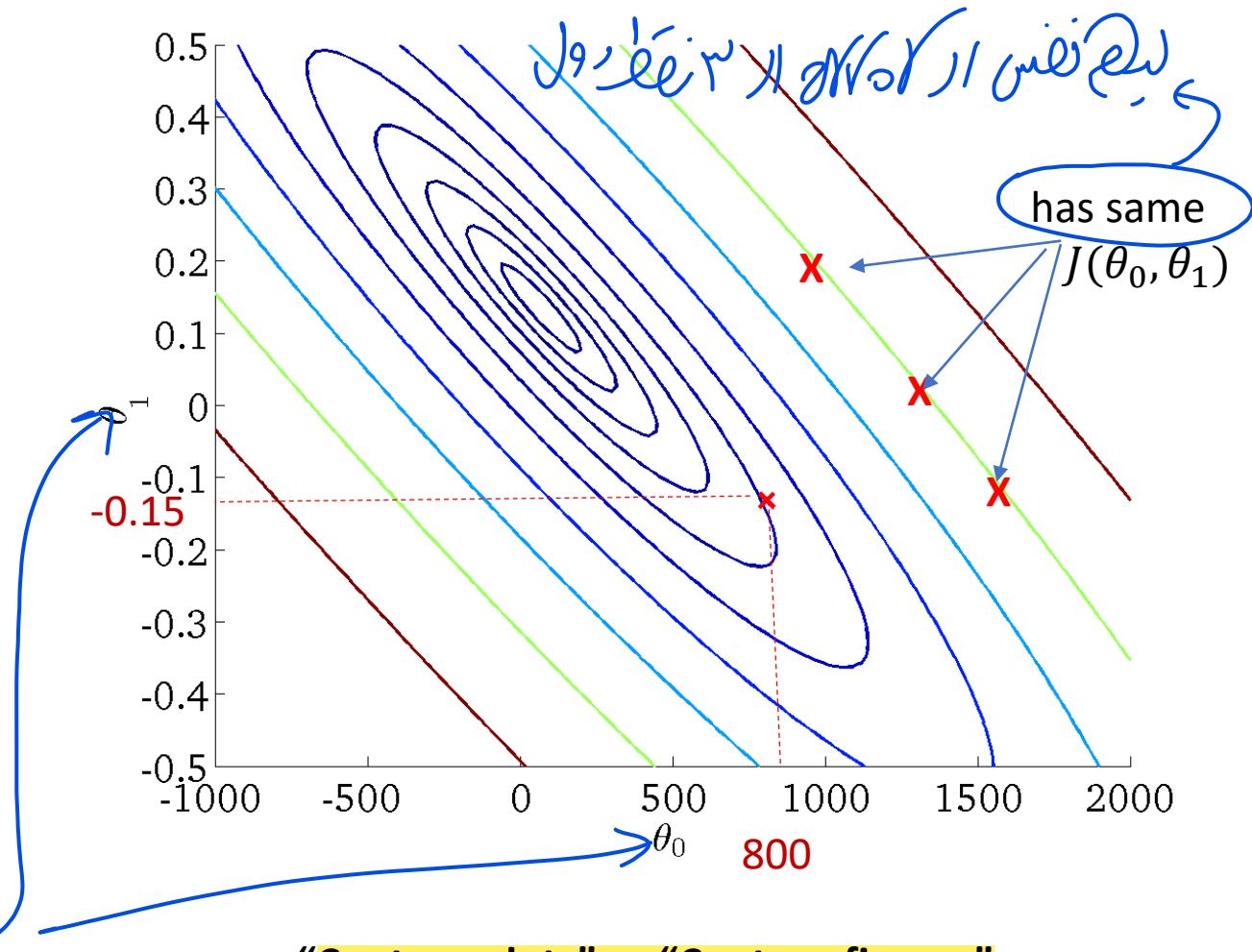


$$h_{\theta}(x)$$



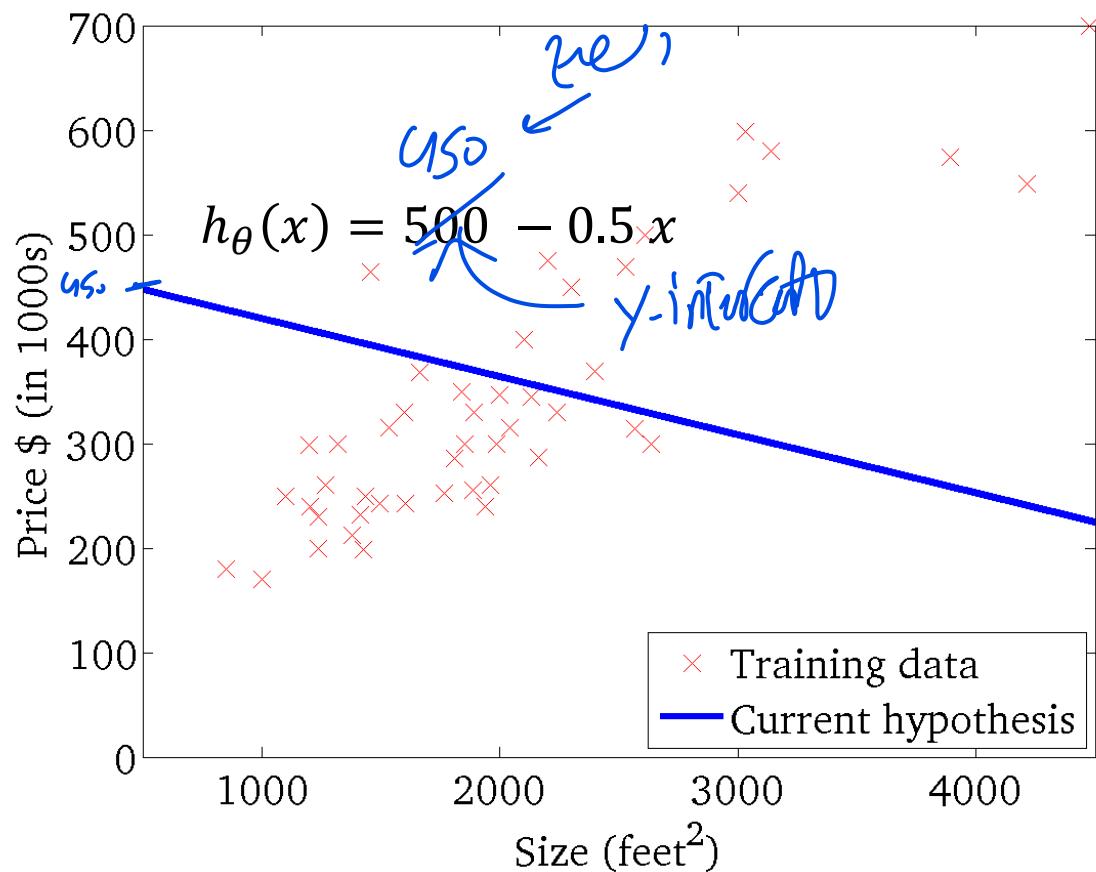
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



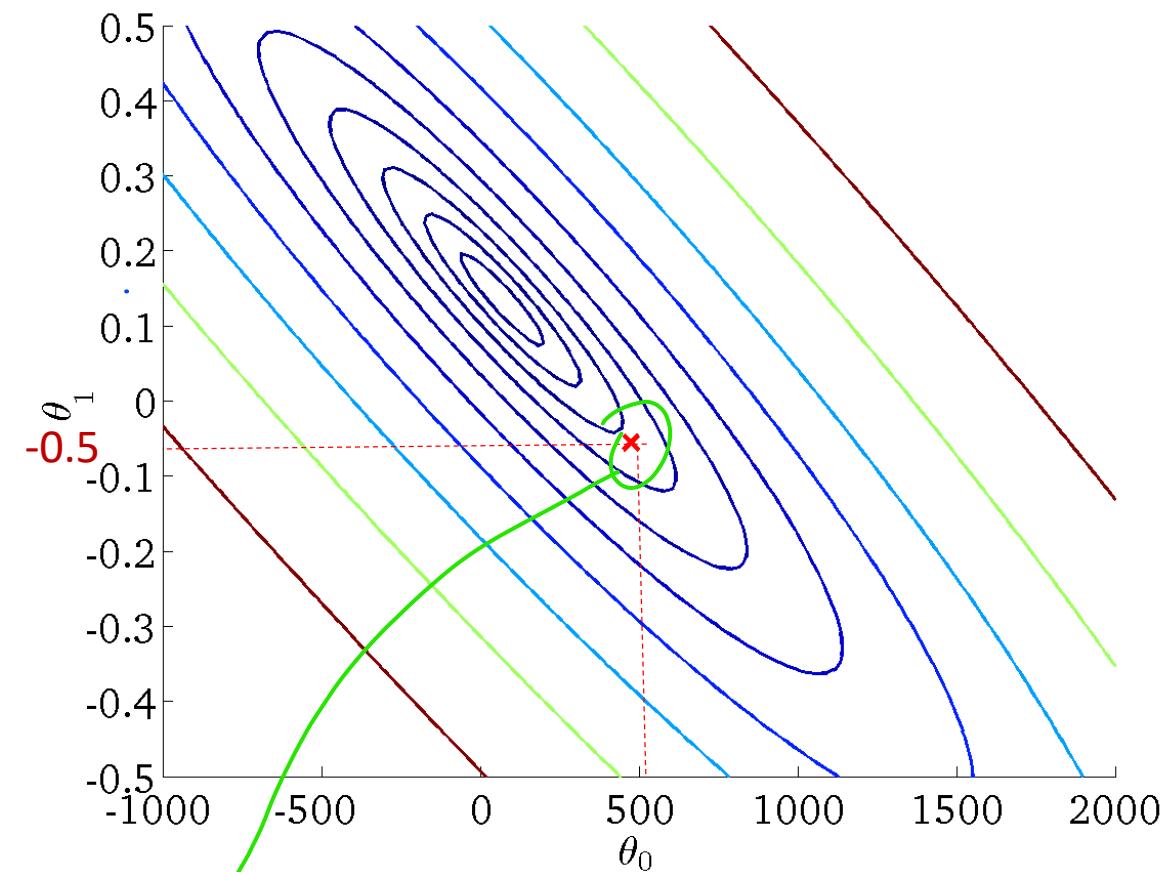
$$h_{\theta}(x)$$

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

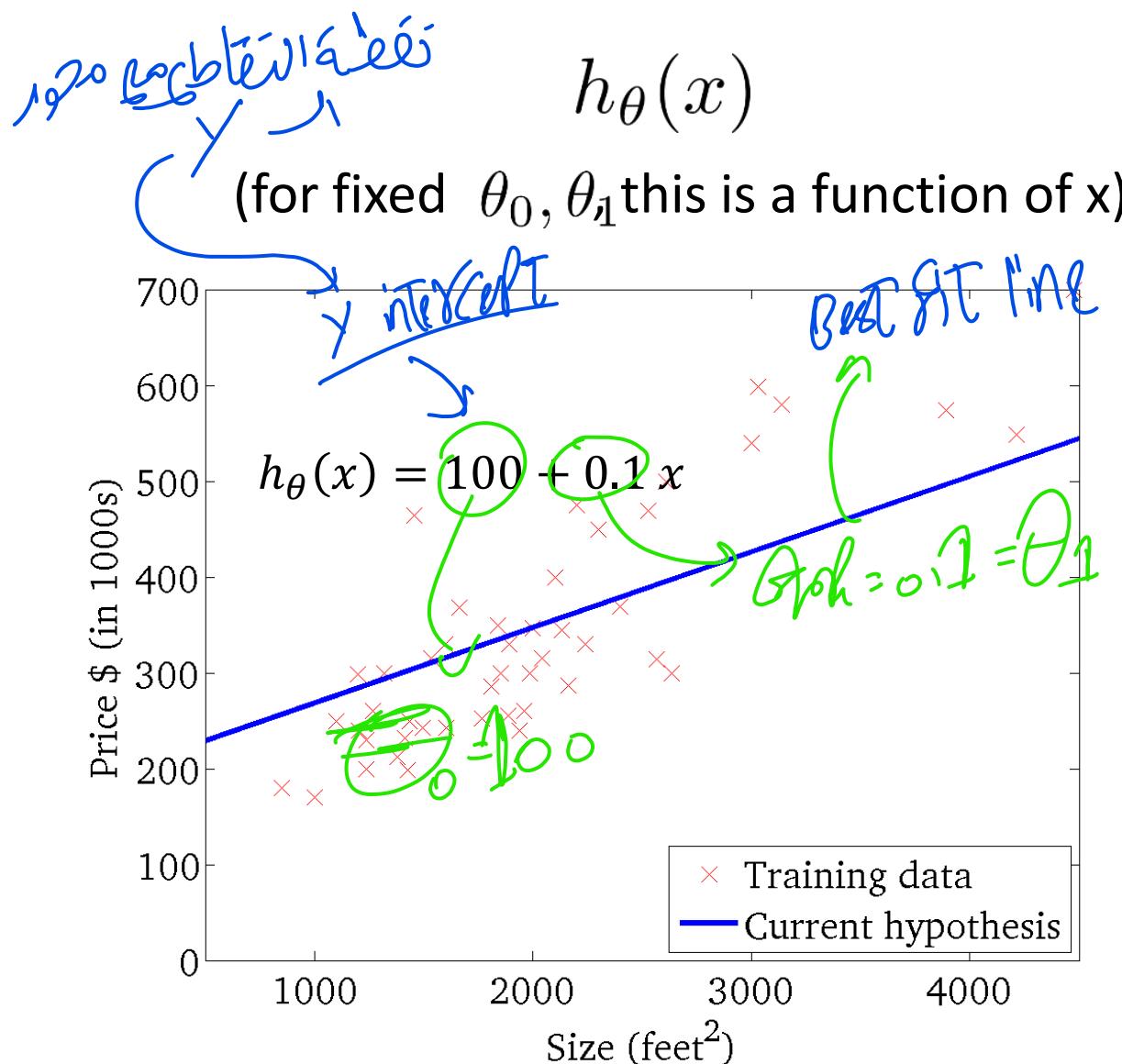


$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )

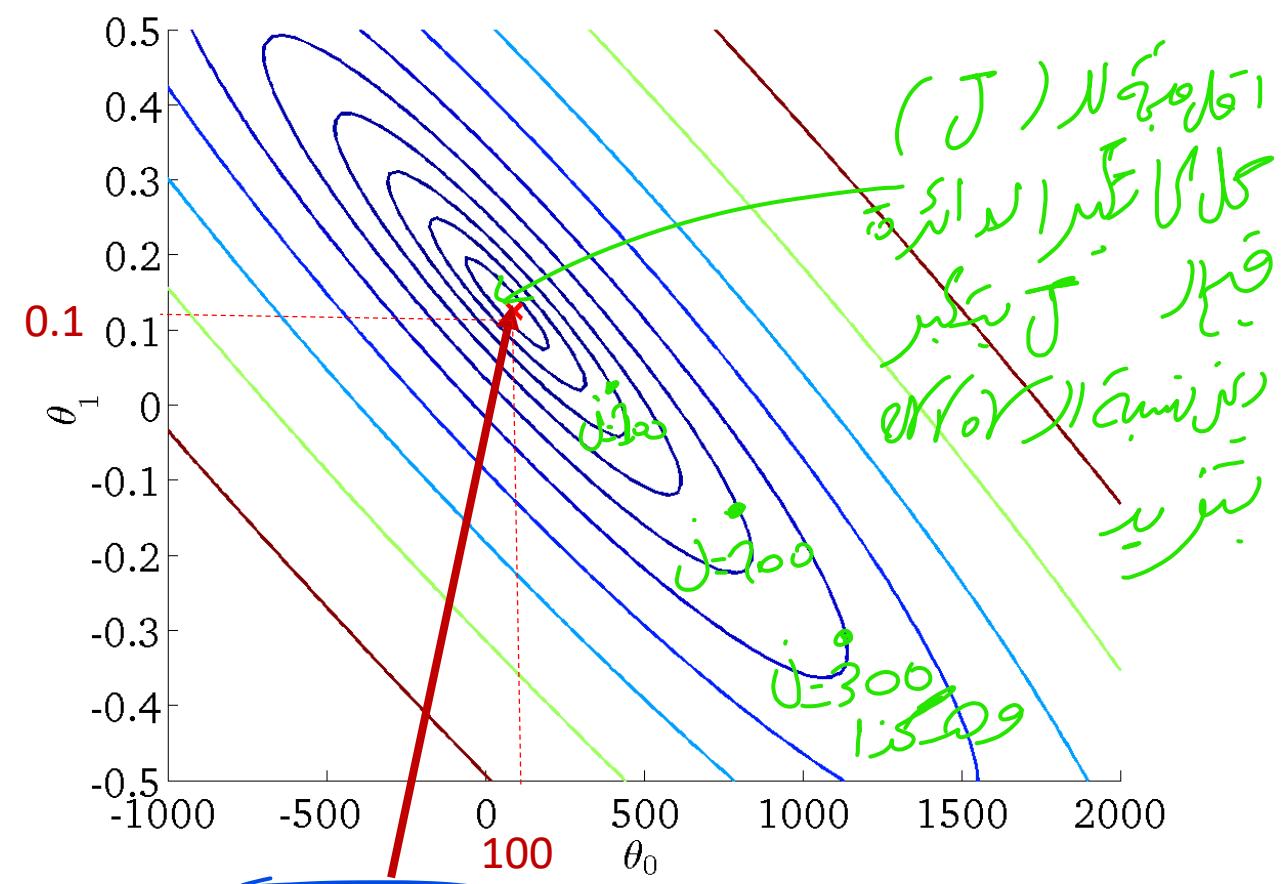


$(J)$  ) ریز نهاد و  $\theta_1$  )  $\theta_0$  نماید



$J(\theta_0, \theta_1)$

(function of the parameters  $\theta_0, \theta_1$ )



هدفنا هو العثور على القيم المثلثة

# Quick Summary until now

Hypothesis:

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

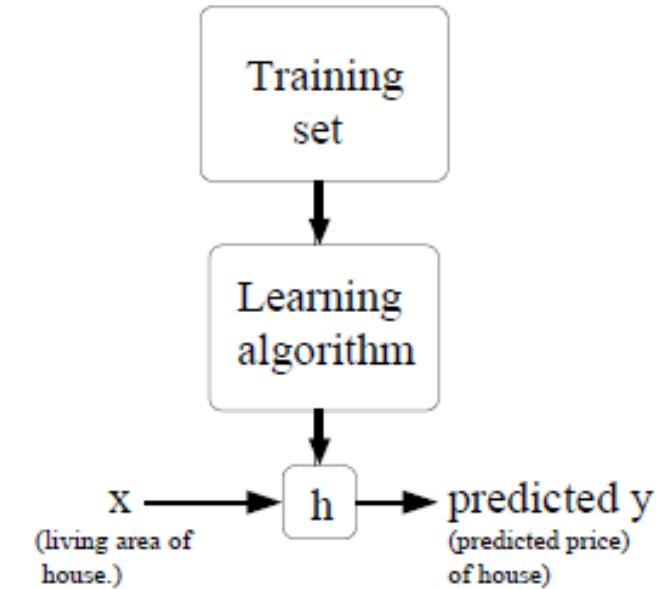
Cost Function:

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

Goal:

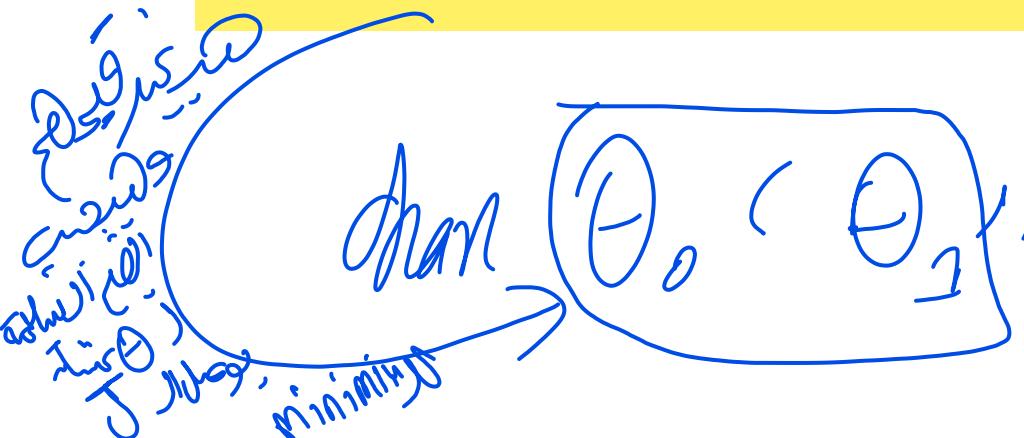
$$\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$$

Our strategy until now: Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at a minimum!!



# Optimization Technique:

## Gradient Descent



كايبر نورس ميجي  
Cost loss function  
automatically  
Min

# Optimization

- **Optimization** is the process of finding the set of parameters/weights  $\theta_0, \theta_1$  that minimize the **cost function**.
- **Strategy 1:** A first very bad idea solution: **RANDOM SEARCH**
  - ▲ is to simply try out many different random parameters and keep track of what works best. (iterative refinement.)
  - ▲ start with random weights and iteratively refine them over time to get lower cost

جذب العين

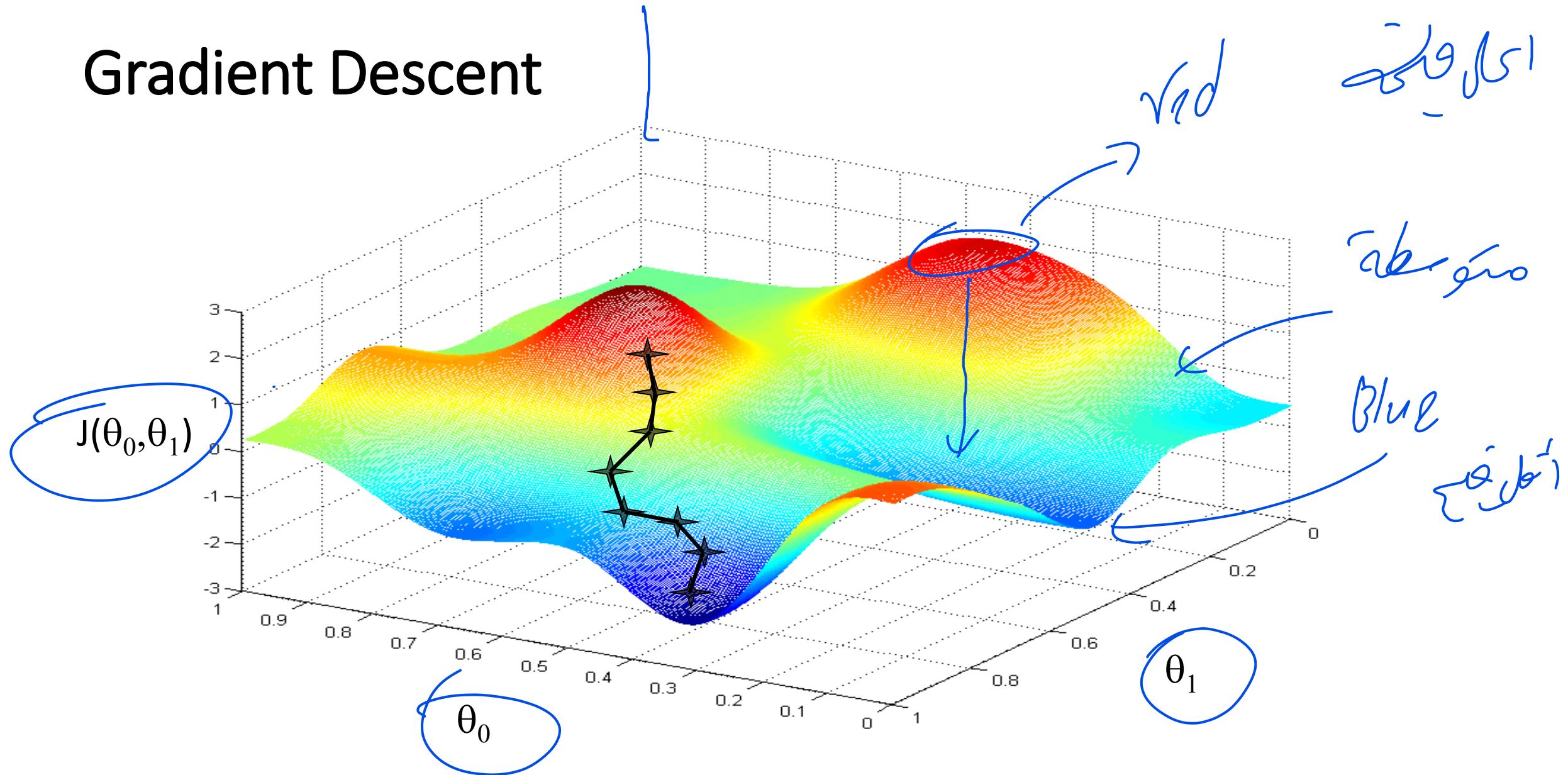
# Optimization

- **Strategy 2: Following the Gradient**

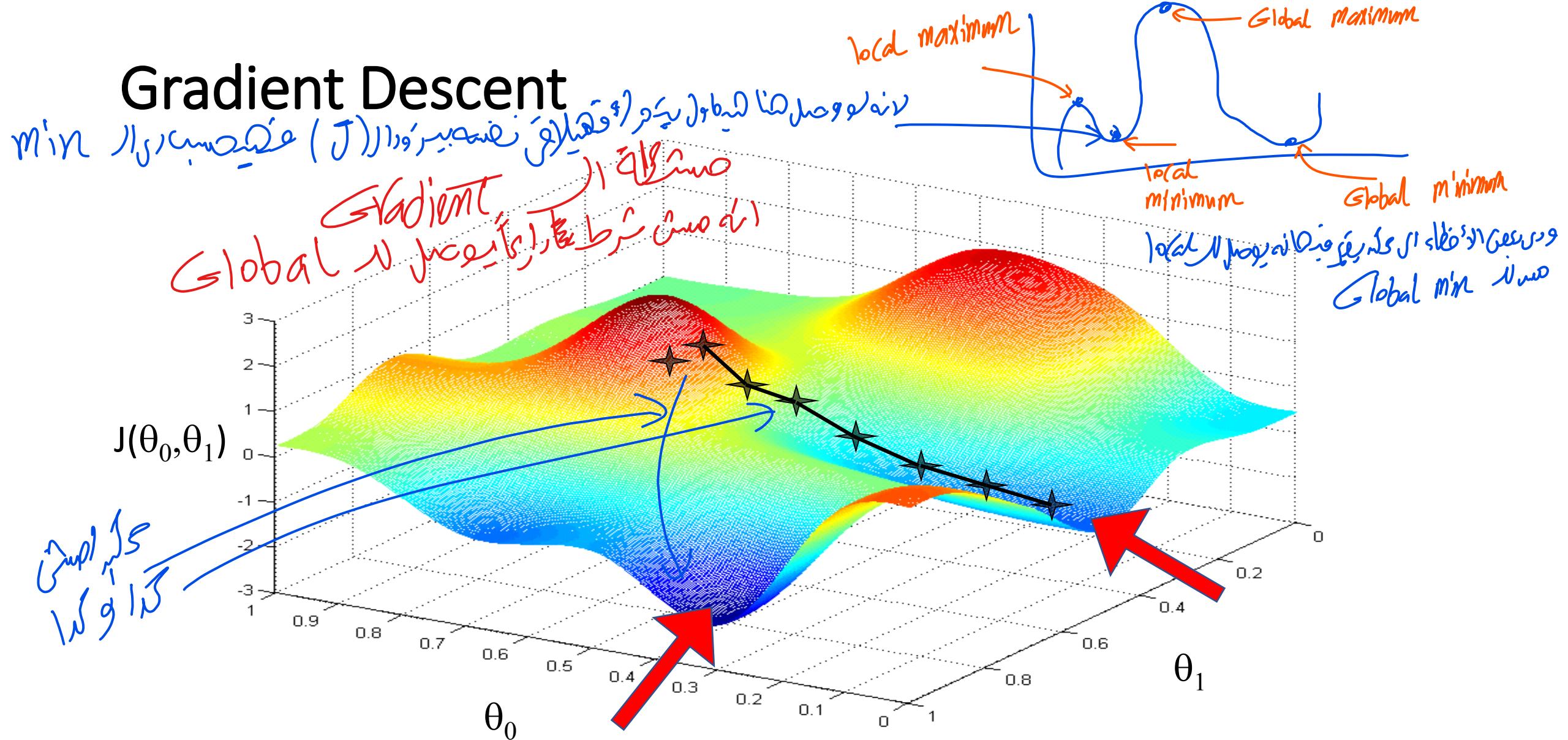
- ▲ Compute the best direction along which we should change our parameter (weight) vector that is mathematically guaranteed to be the direction of the steepest descend.
- ▲ This direction will be related to the **gradient** of the cost function.

Value of  $\theta_0, \theta_1$  to give minimum squared error function ( $J$ )

# Gradient Descent



# Gradient Descent

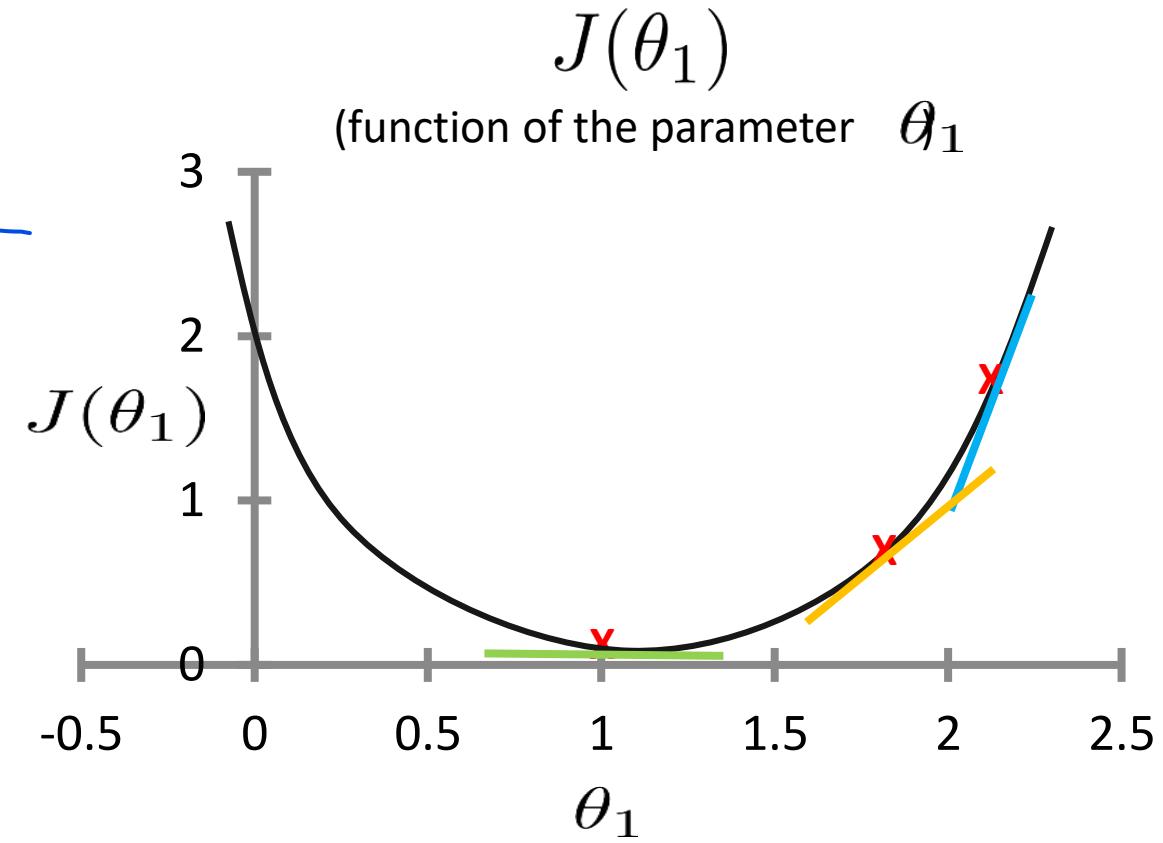


# Gradient Descent

- The way we do this is by taking the derivative (the tangential line to a function) of our cost function.
- The slope of the tangent is the derivative at that point and it will give us a direction to move towards.
- We will know that we have succeeded when our cost function is at the very bottom of the pits in our graph, i.e. when its value is the minimum.

stop if derivative is 0

(a)  $\theta_1$  vs  $J(\theta_1)$   $\rightarrow$   $J(\theta_1)$  error in terms of parameters  $\theta_1$



# Gradient Descent

- We make steps down the cost function in the direction with the steepest descent, and the size of each step is determined by the parameter  $\alpha$ , which is called the **learning rate**.  $J(\theta_1)$
  - The gradient descent algorithm is:  $\frac{\partial J(\theta)}{\partial \theta_1}$  (function of the parameter  $\theta_1$ )

- The gradient descent algorithm is:

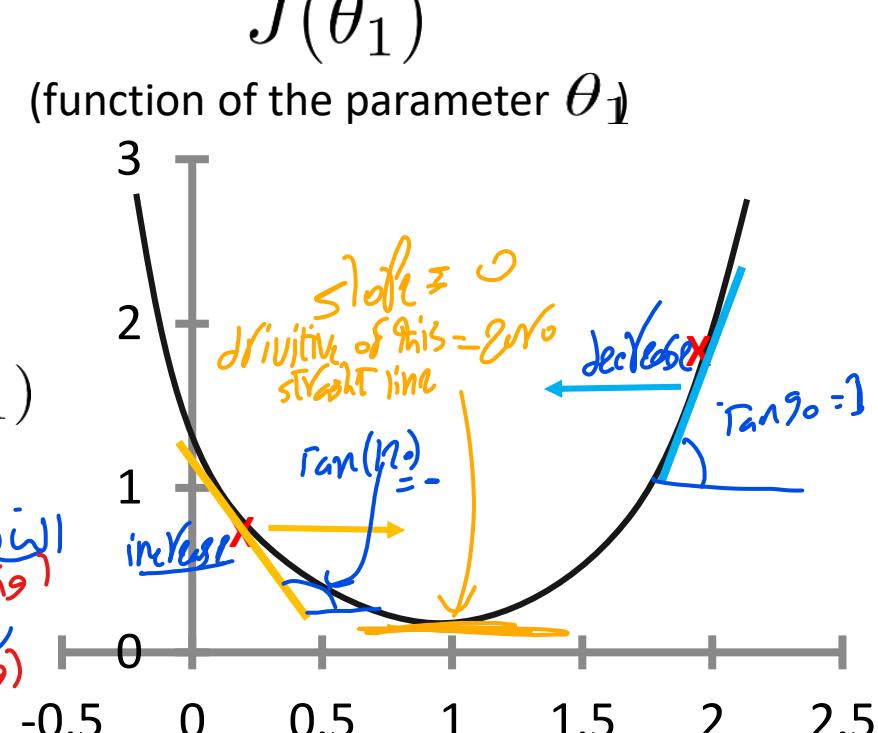
```

repeat until convergence {
    θj := θj - α * ∂J(θ0, θ1) / ∂θj
    for j = 0 and j = 1
        if (θj - previous θj) * direction < 0
            increase step size
        else
            decrease step size
}

```

*Note: The handwritten notes provide additional context and annotations:*

- Partial derivative of  $J$ :**  $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$
- Initial Normalization:**  $\theta_0$  and  $\theta_1$  are normalized.
- Slope Method:** The slope method is used to determine the learning rate.
- Learning rate (step size):** The value of the learning rate is shown as  $\alpha$ .
- Direction:** The direction of update is indicated by the sign of the derivative.
- Increase:** If the direction is negative, the step size is increased.
- Decrease:** If the direction is positive, the step size is decreased.
- Positive slope:** A positive slope is indicated by a green circle with a plus sign.
- Negative slope:** A negative slope is indicated by a green circle with a minus sign.
- Value:** The current value of the parameters  $\theta_0$  and  $\theta_1$  is shown.
- Iteration:** The iteration number  $j$  is shown.
- Convergence:** The algorithm repeats until convergence.



**Positive slope** (positive number)  $\rightarrow \theta_1$  will decrease  
**Negative slope** (negative number)  $\rightarrow \theta_1$  will increase

## Learning Rate

- The gradient tells us the direction, but it does not tell us how far along this direction we should step.
- The **learning rate (step size)** determines how big the step would be on each iteration. It determines how fast or slow we will move towards the optimal weights.

# Learning Rate

النسبة المئوية

- If learning rate is large, it may fail to converge and overshoot the minimum.
- If learning rate is very small, it would take long time to converge and become computationally expensive.
- The most commonly used rates are :

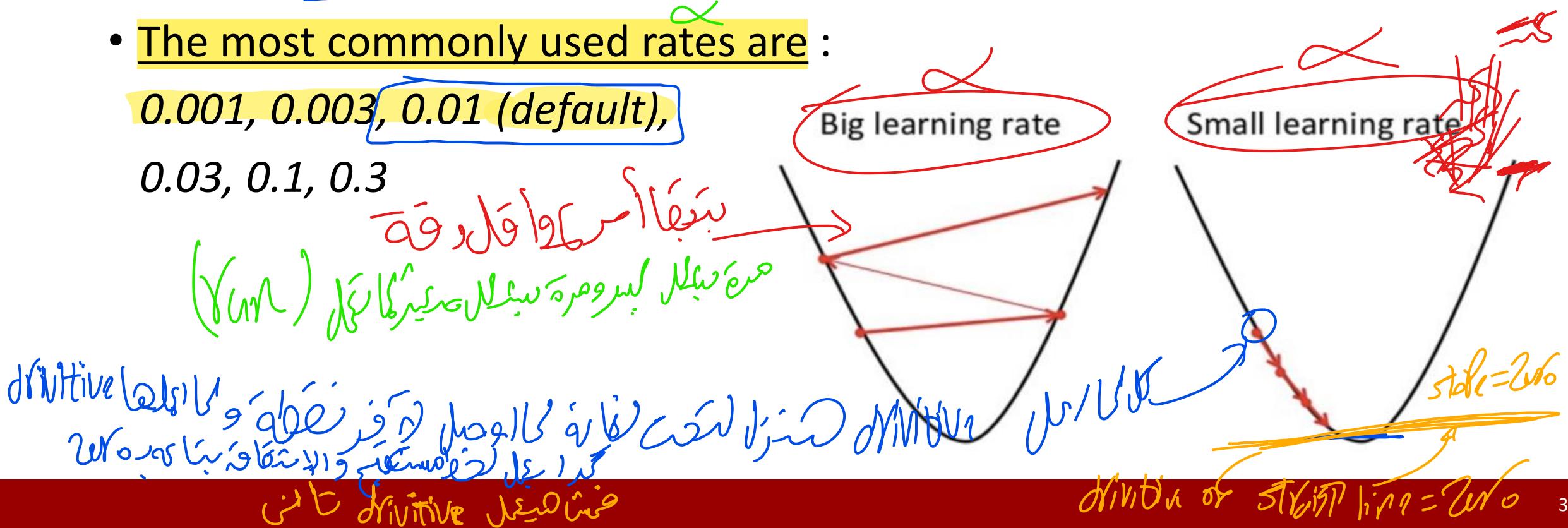
0.001, 0.003, 0.01 (default),

0.03, 0.1, 0.3

النسبة المئوية  
(%) لبيانات التعلم

Big learning rate

Small learning rate



gradient line  
وهي قدرة  
النهاية  
الآن  
خشن

# Gradient Descent for Linear Regression

analyse

When specifically applied to the case of **linear regression**, a new form of the gradient descent equation can be derived. We can substitute our actual cost function and our actual hypothesis function and modify the equation to:

- 1 - Start by initializing the parameters  $\theta_0, \theta_1$  randomly
- 2 - Repeat until convergence { // until error is small
  - Predicted values with linear regression hypothesis.
  - Calculate the cost function.
  - If cost is large, update parameters using GD

Should be done simultaneously

in parallel

}  $\theta_1$  is  $g(\theta_0)$

and  $\theta_0$  is  $g(\theta_1)$

and  $\theta_1$  is  $g(\theta_0)$

$$\text{Cost Function} \quad \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) :$$
$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) :$$

different dimensions

same

dimensions

hypothesis's function loss

using  $\theta_1, \theta_0$  to find minimum

loss(function) Partial derivative

new value

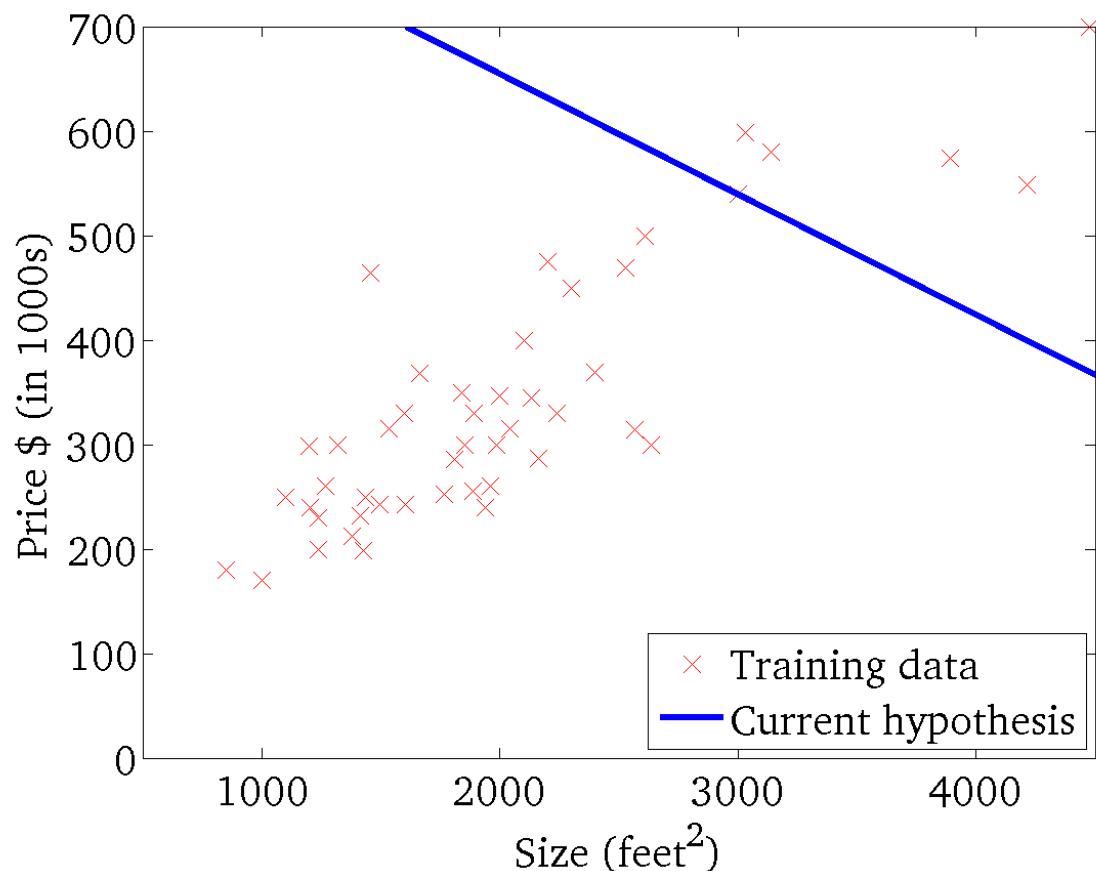
gradient descent

loop continue

convergence

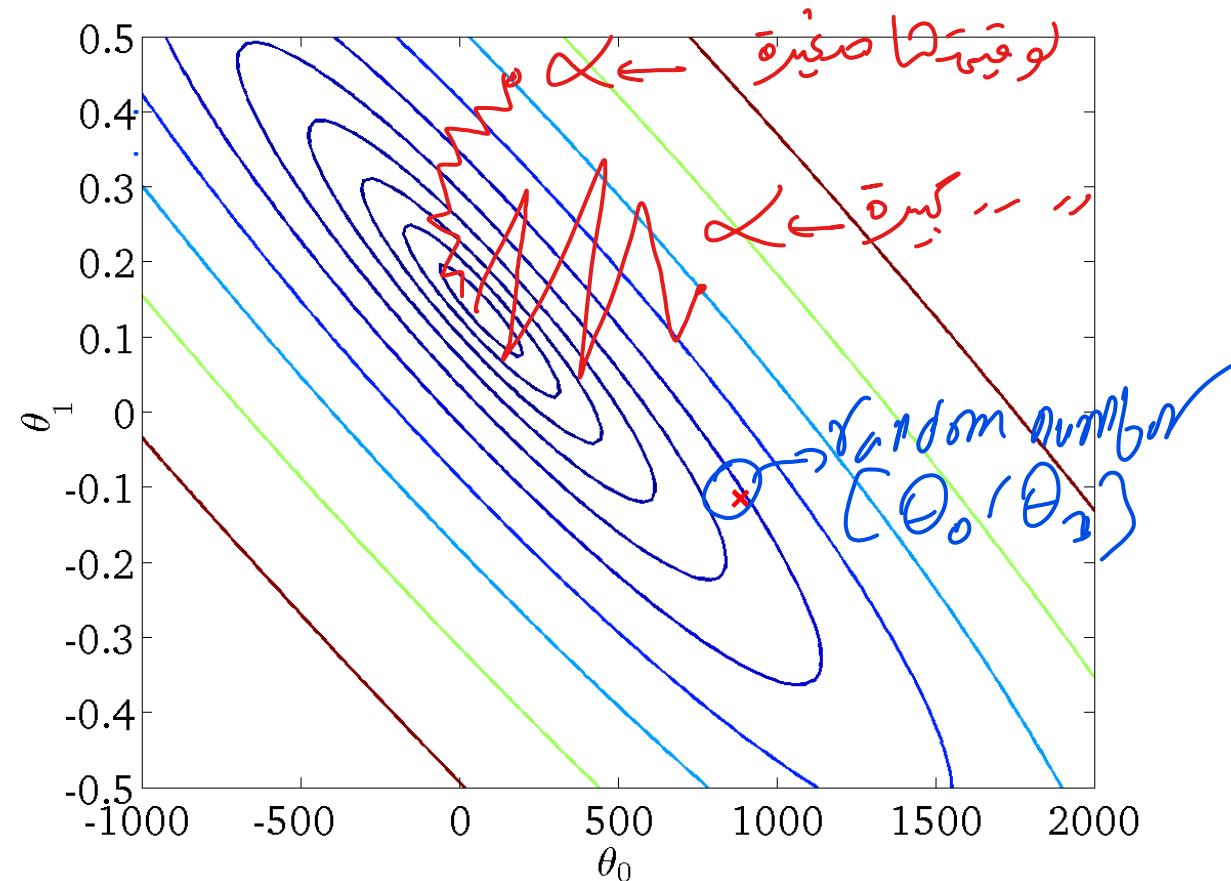
$$h_{\theta}(x)$$

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



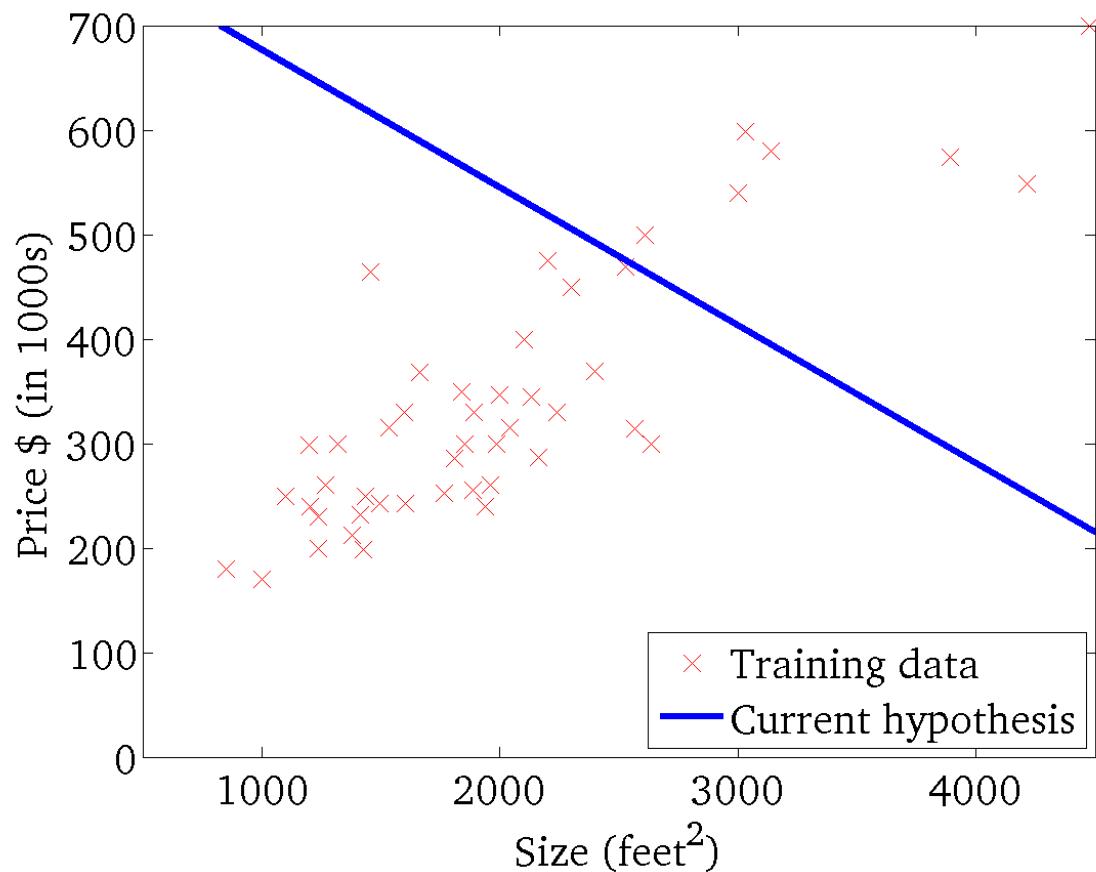
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



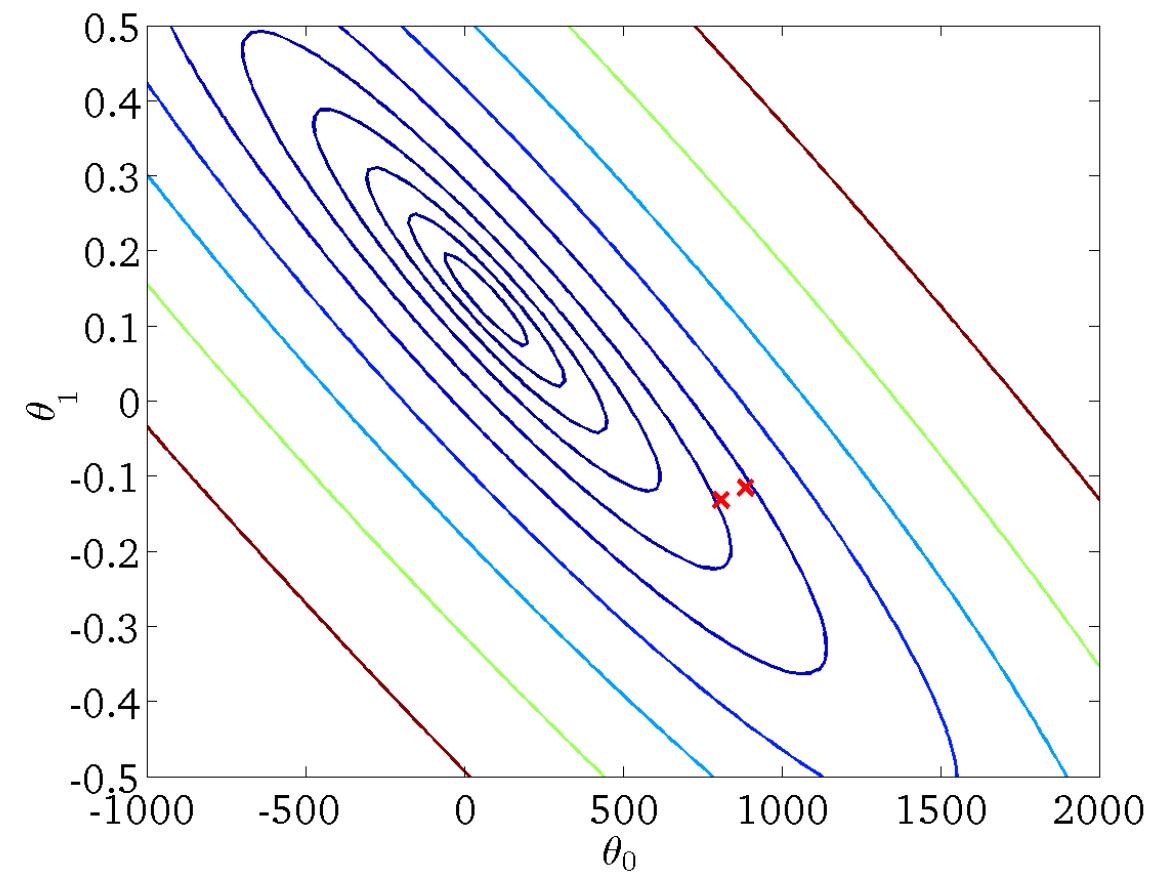
$$h_{\theta}(x)$$

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



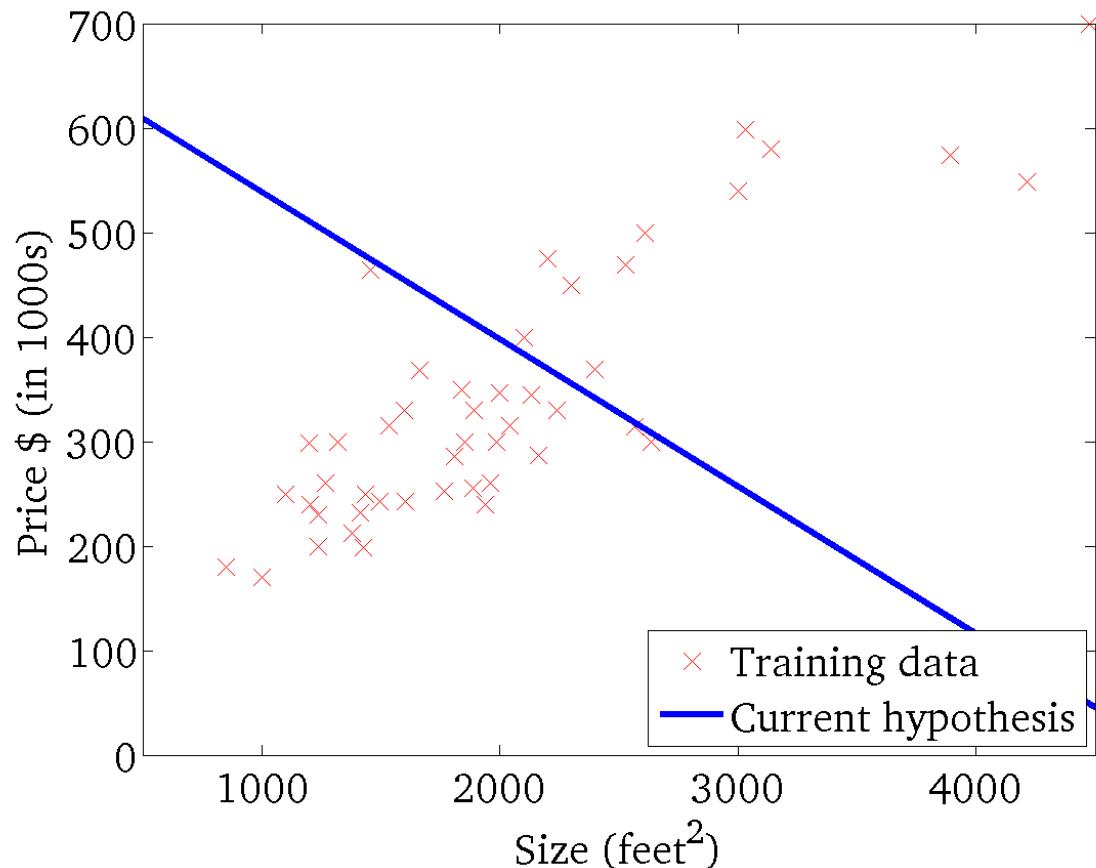
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



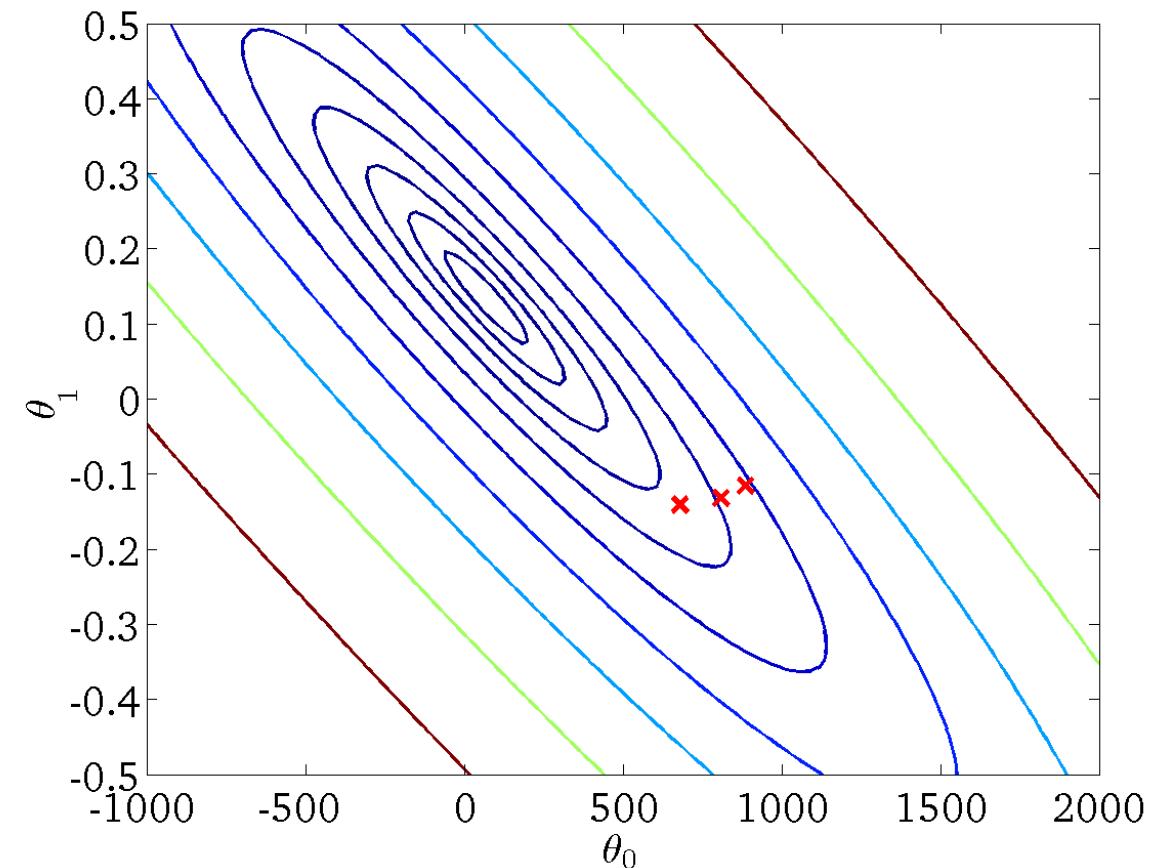
$$h_{\theta}(x)$$

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



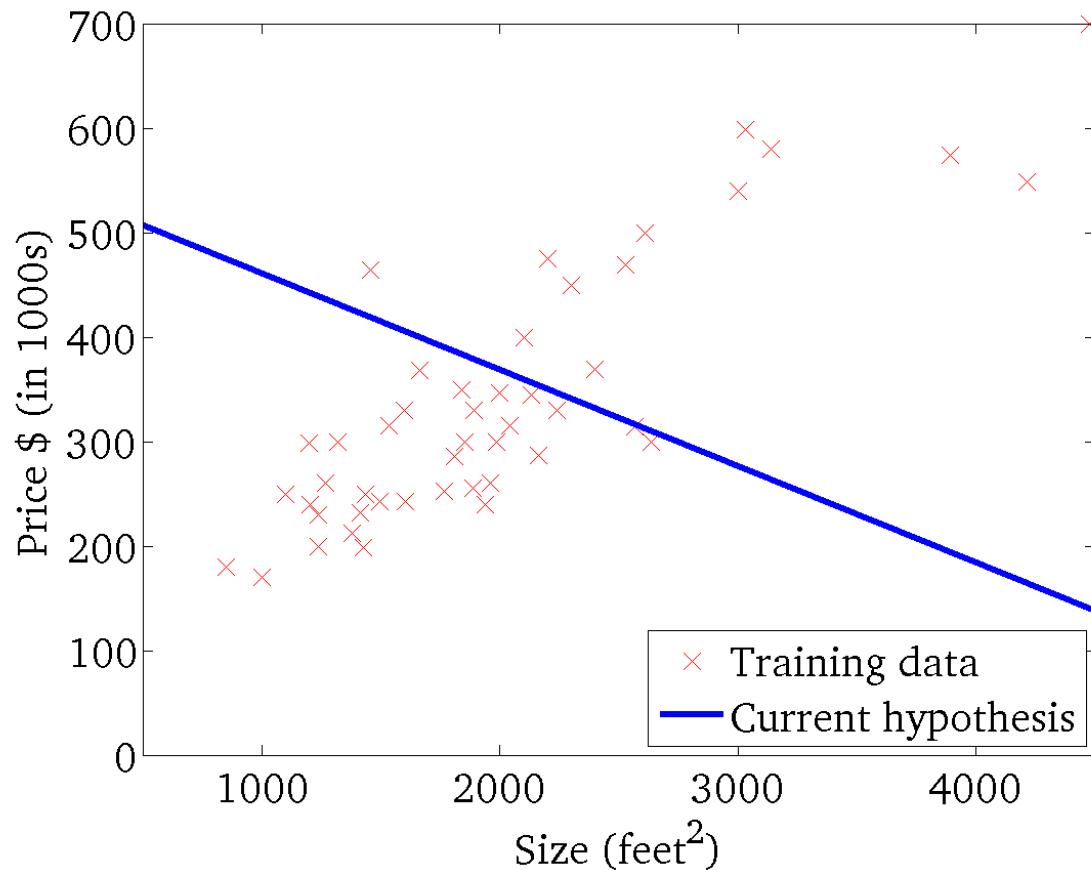
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



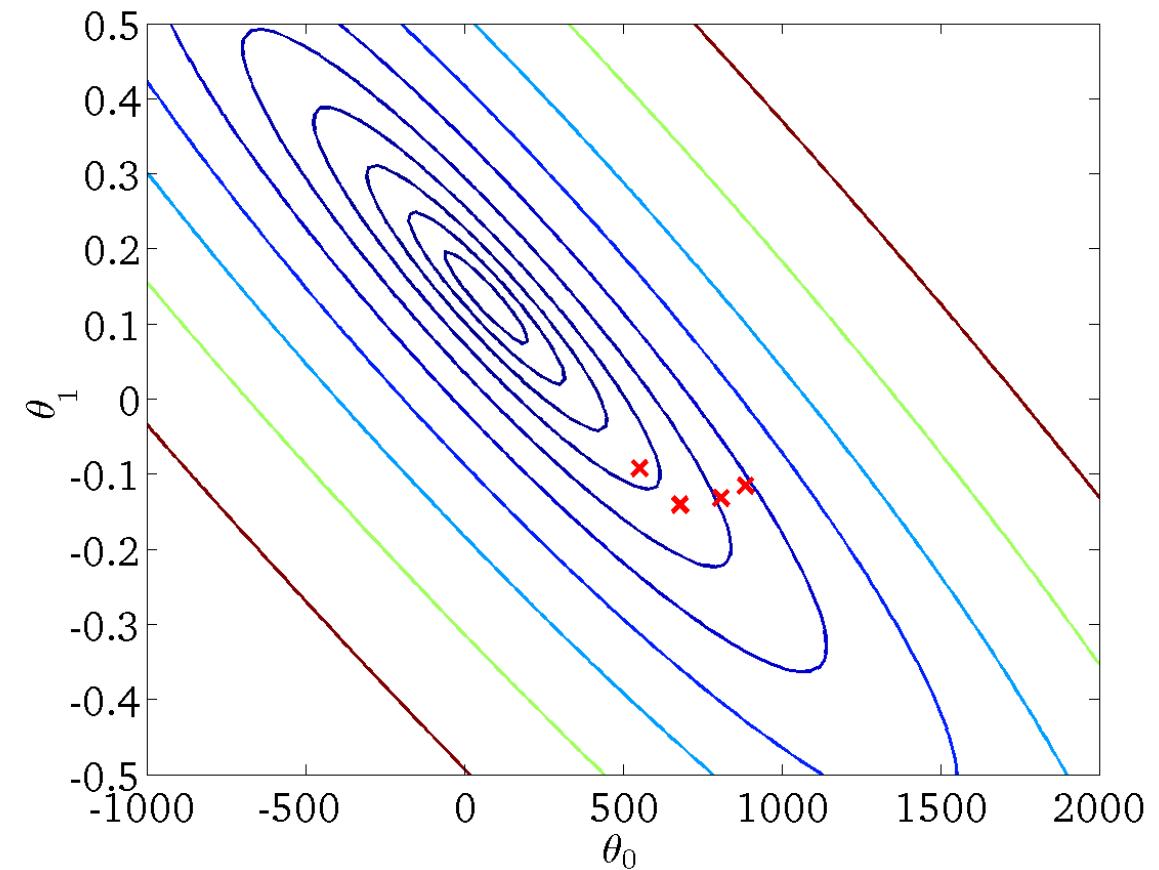
$$h_{\theta}(x)$$

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



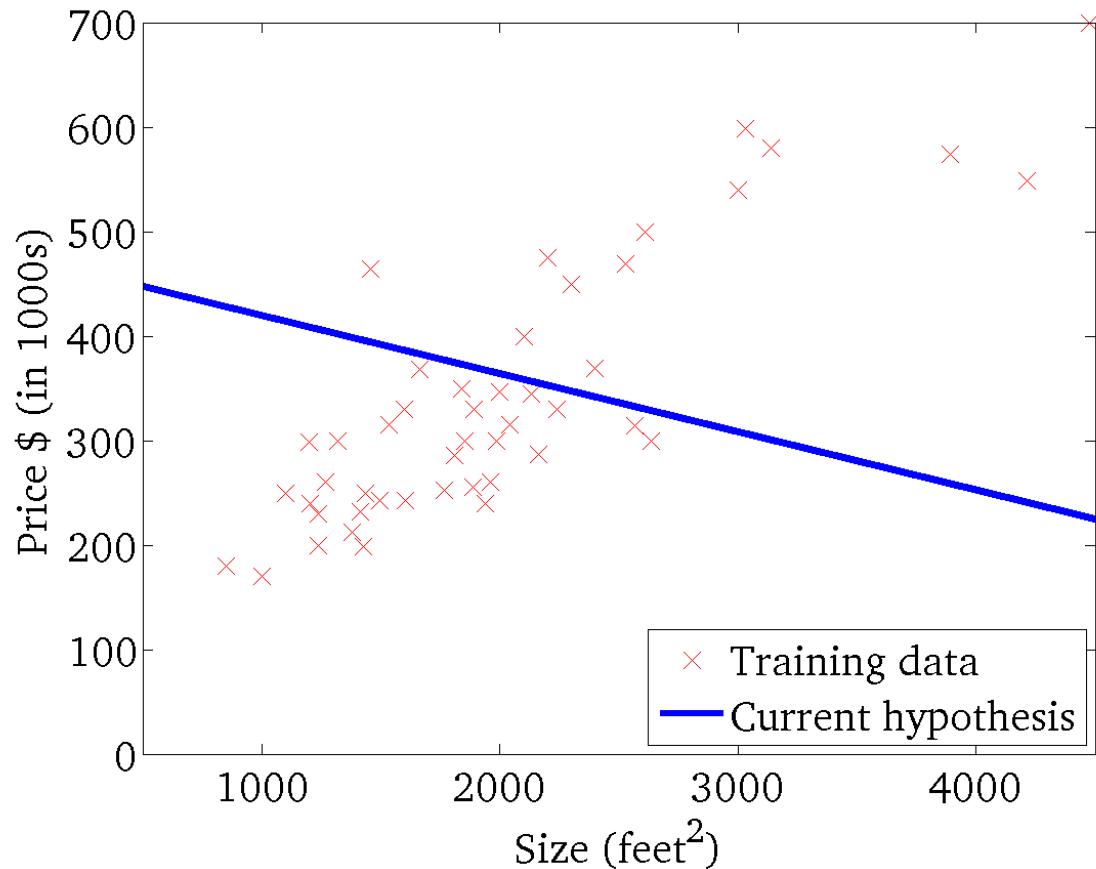
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



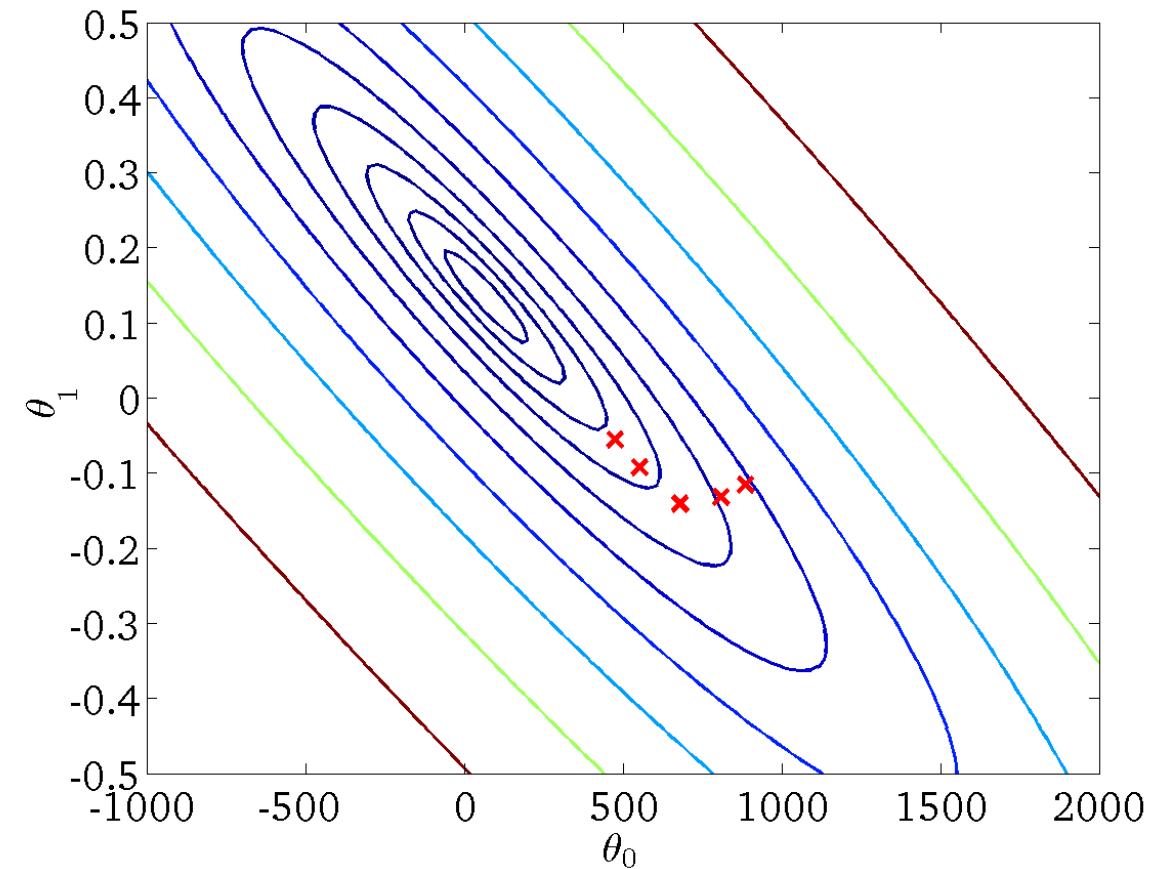
$$h_{\theta}(x)$$

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



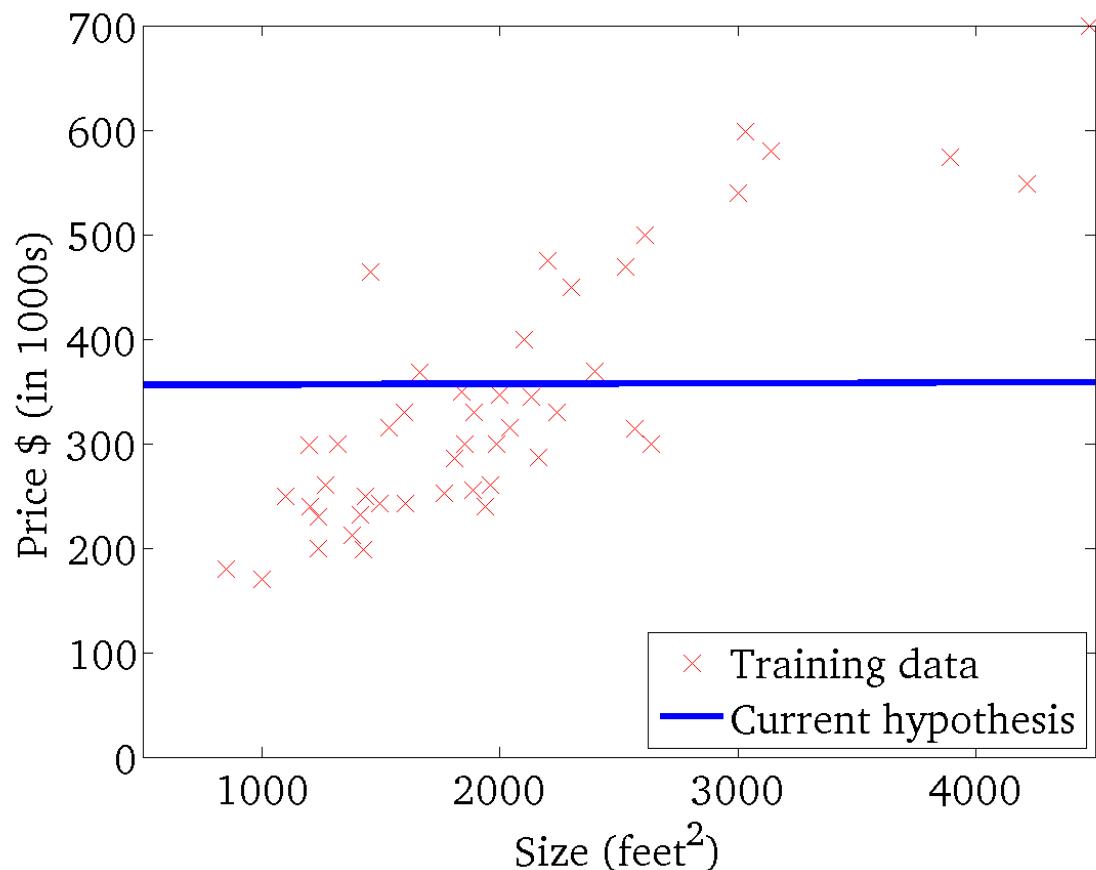
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



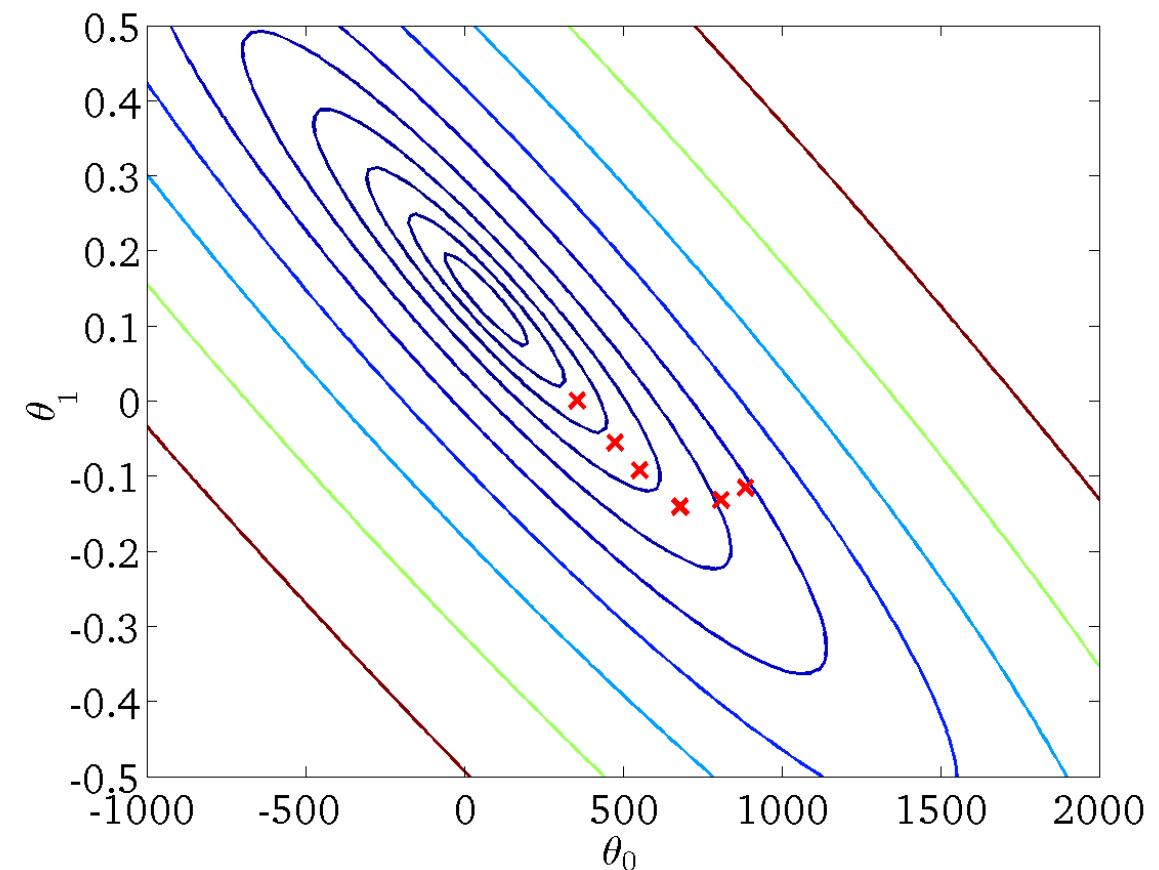
$$h_{\theta}(x)$$

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



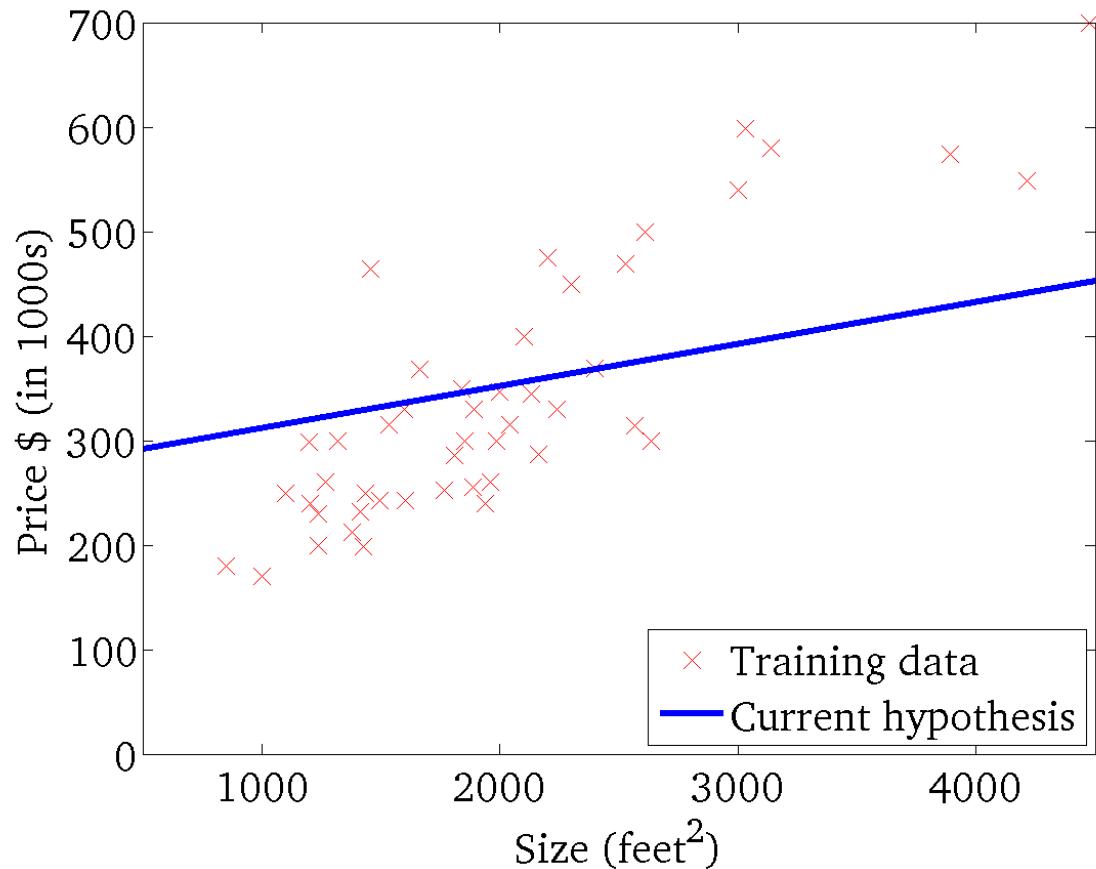
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



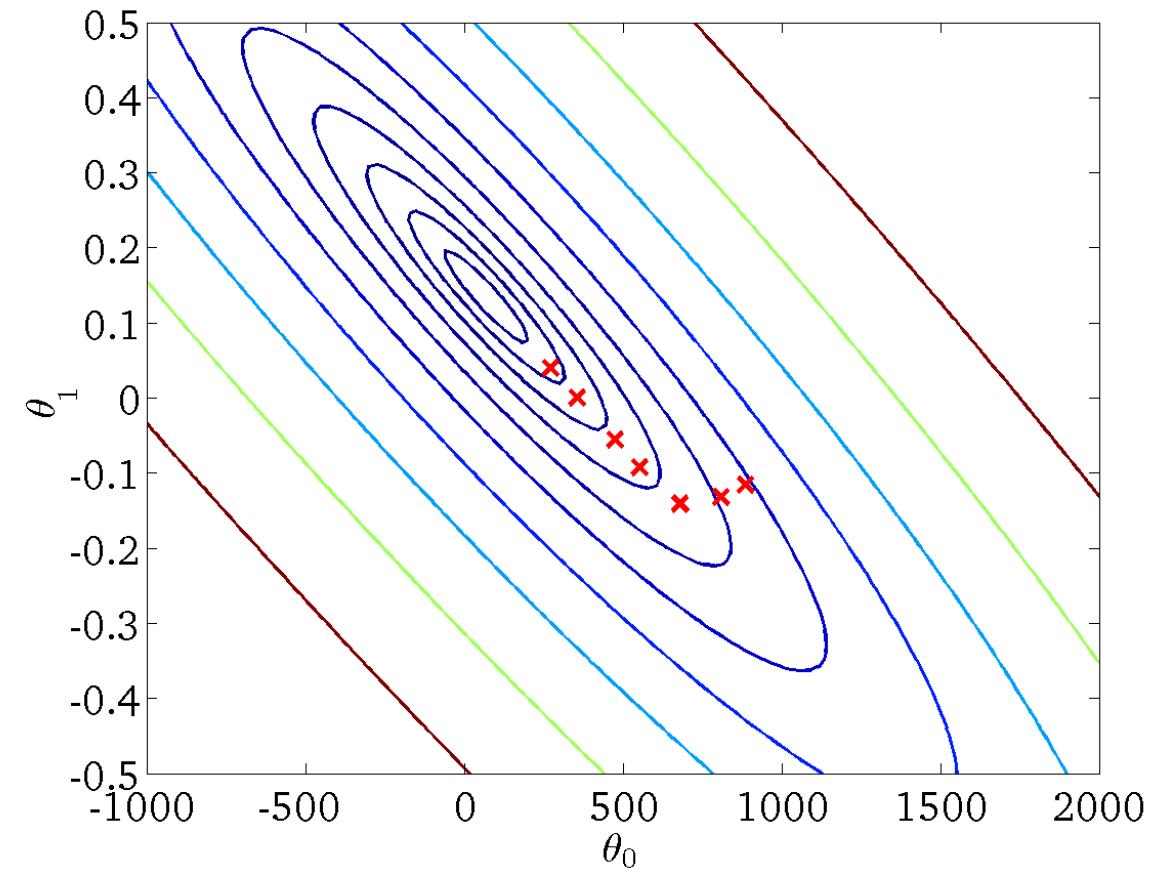
$$h_{\theta}(x)$$

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



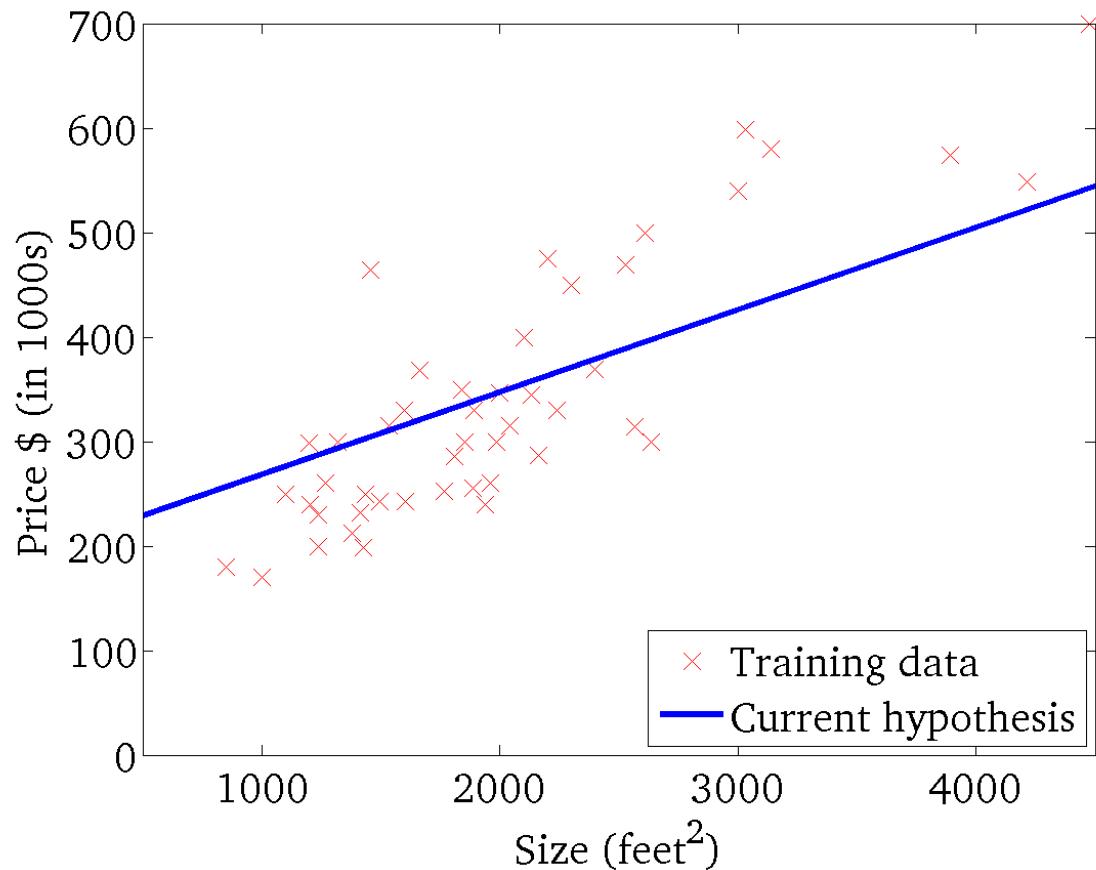
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



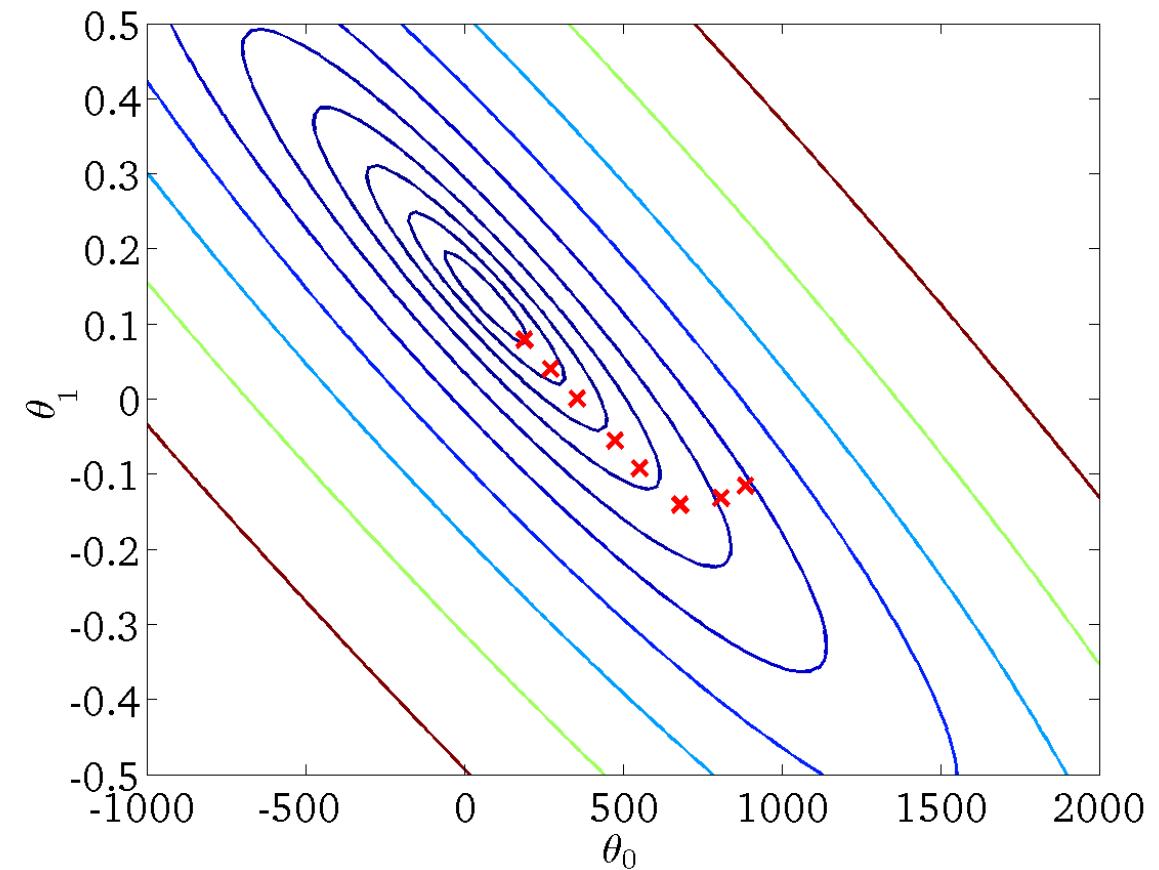
$$h_{\theta}(x)$$

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



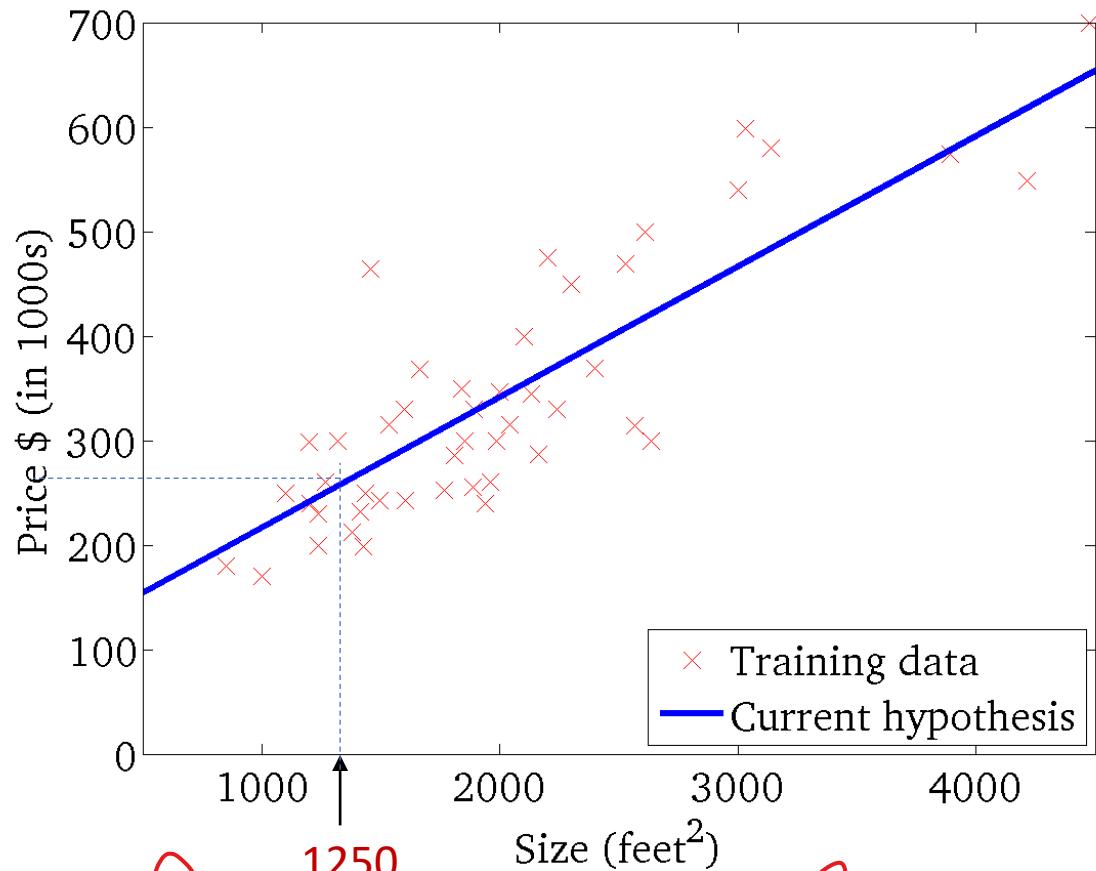
$$J(\theta_0, \theta_1)$$

(function of the parameters  $\theta_0, \theta_1$ )



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

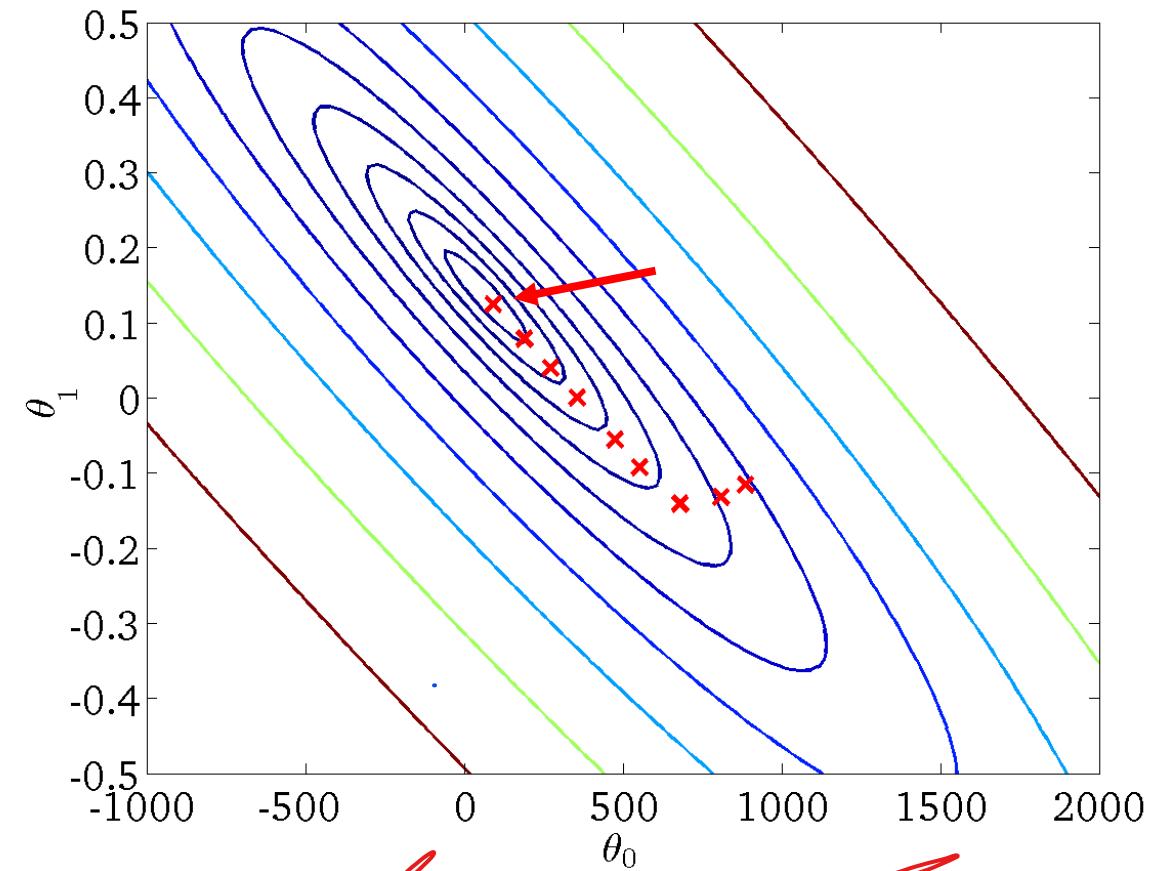
(for fixed  $\theta_0, \theta_1$  this is a function of  $x$ )



$\theta_0$ ,  $\theta_1$  Training

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

(function of the parameters  $\theta_0, \theta_1$ )



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

# Gradient Descent variants

- There are three variants of gradient descent based on the amount of data used to calculate the gradient:
  1. Batch gradient descent
  2. Stochastic gradient descent
  3. Mini-batch gradient descent

# Batch Gradient Descent

- Batch Gradient Descent, aka **Vanilla gradient descent**, calculates the error for each observation in the dataset but performs an update only after all observations have been evaluated.
- One cycle through the entire training dataset is called a **training epoch**. Therefore, it is often said that batch gradient descent performs model updates at the end of each training epoch.
- Batch gradient descent is not often used, because it represents a huge consumption of computational resources, as the entire dataset needs to remain in memory.

# Stochastic Gradient Descent (SGD)

$\theta_1 \rightarrow \theta_0$ . Just 2nd example  $N^1$

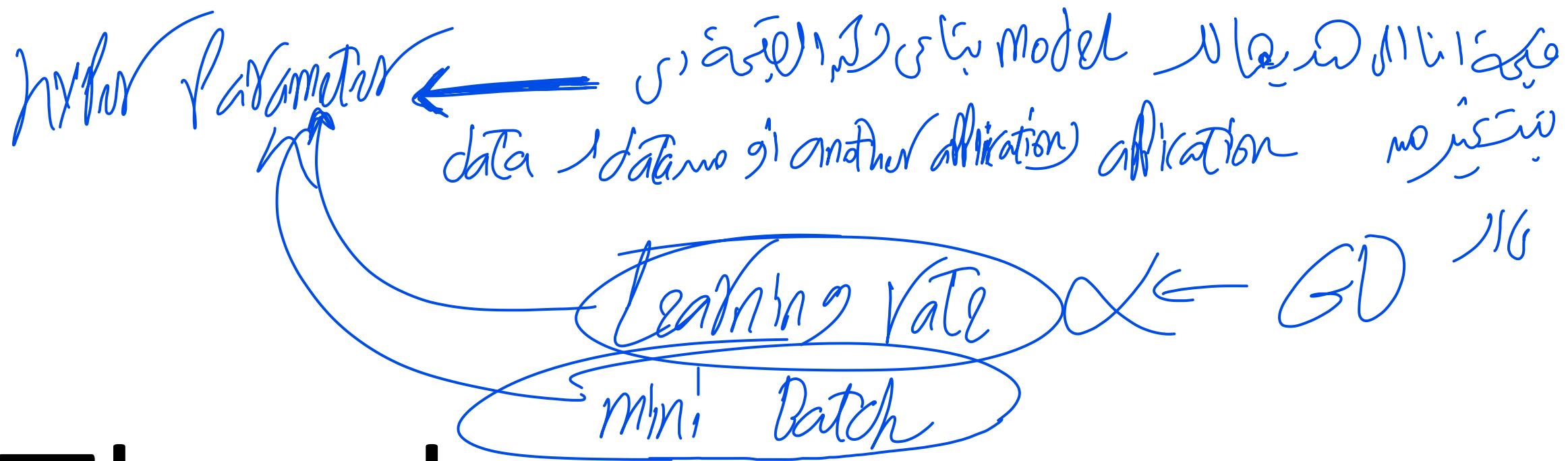
- Stochastic gradient descent, often abbreviated **SGD**, is a variation of the gradient descent algorithm that calculates the error and updates the model for each example in the training dataset.
- SGD is usually faster than batch gradient descent, but its frequent updates cause a higher variance in the error rate, that can sometimes jump around instead of decreasing.
- The noisy update process can allow the model to avoid local minima (e.g. premature convergence).

# Mini-Batch Gradient Descent

جهاز حاسوب مخصوص  
experiments كل 300، example 100 update 0.1%

- Mini-batch gradient descent seeks to find a balance between the robustness of stochastic gradient descent and the efficiency of batch gradient descent.
- It is the most common implementation of gradient descent used in the field of deep learning.  
*batches جزءاً من data المجموع*
- It splits the training dataset into small batches that are used to calculate model error and update model coefficients.  
*one full iteration update uses whole Batch to 15*
- “Batch size” commonly used as power of 2: 32, 64, 128, 256, and so on.  
*batch حجم المجموعة*





# Thanks