

# Linear regression

**Lesson #1:** Linear regression with one variable

**Lesson #2:** Linear regression with multiple variables

**Duration:** 5 hrs

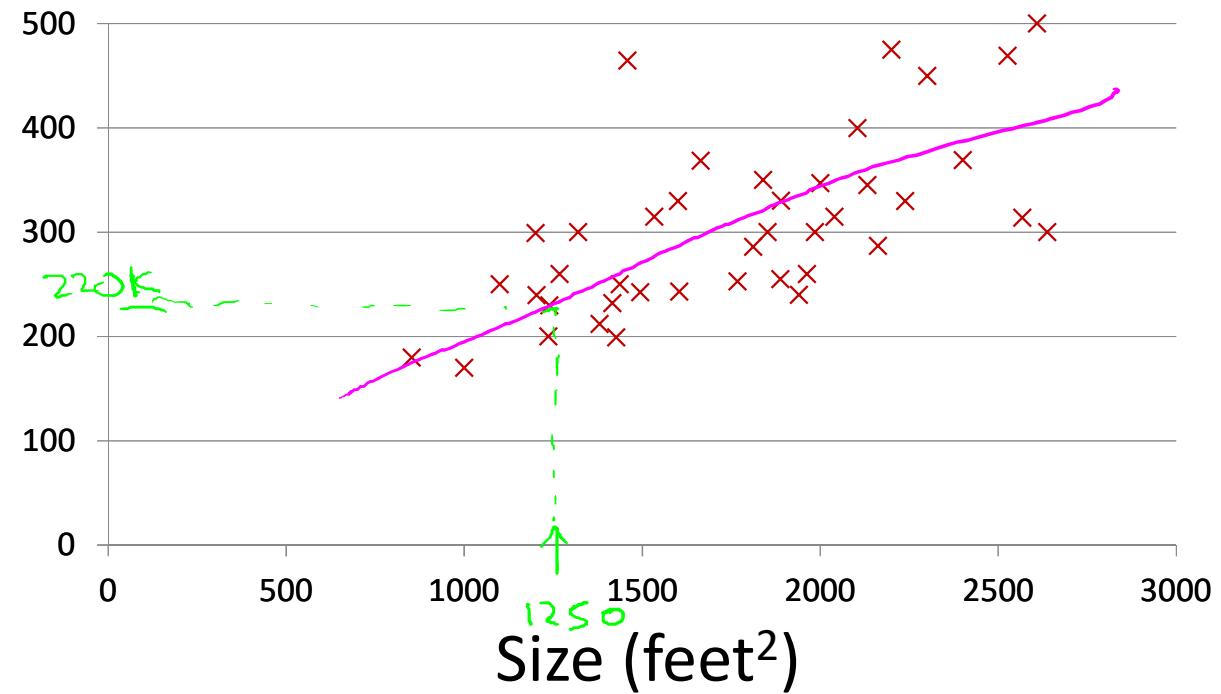
*Source: Machine Learning, Andrew Ng, coursera.org*

# Lesson #1: Linear regression with one variable

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Cost function
  - Gradient descent
  - Gradient descent for linear regression

# Housing Prices (Portland, OR)

Price  
(in 1000s  
of dollars)



## Supervised Learning

Given the “right answer” for each example in the data.

## Regression Problem

Predict real-valued output  
Classification: discrete-valued output

## Training set of housing prices (Portland, OR)

| Size in feet <sup>2</sup> (x) | Price (\$) in 1000's (y) |
|-------------------------------|--------------------------|
| 2104                          | 460                      |
| 1416                          | 232                      |
| 1534                          | 315                      |
| 852                           | 178                      |
| ...                           | ...                      |

Notation:

$m$  = Number of training examples

$$x^{(1)} = 2104$$

$x$ 's = “input” variable / features

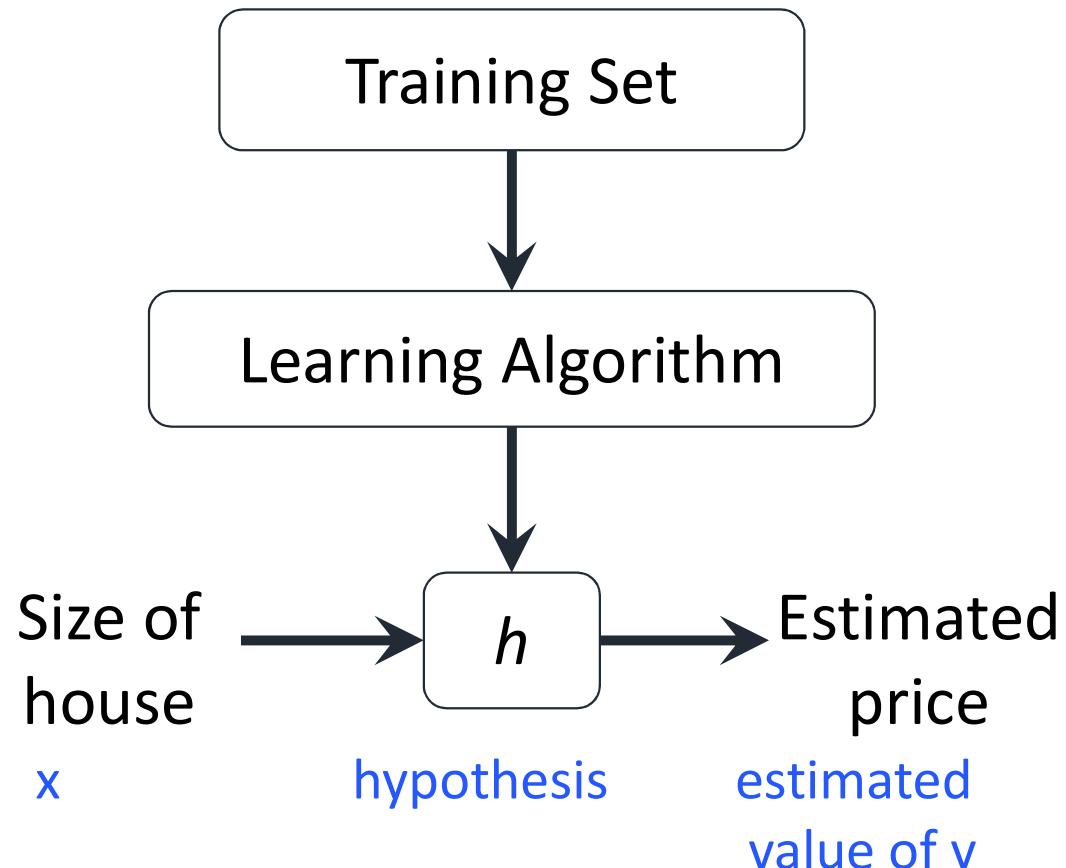
$$x^{(2)} = 1416$$

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

$$y^{(1)} = 460$$

$(x, y)$  – one training example

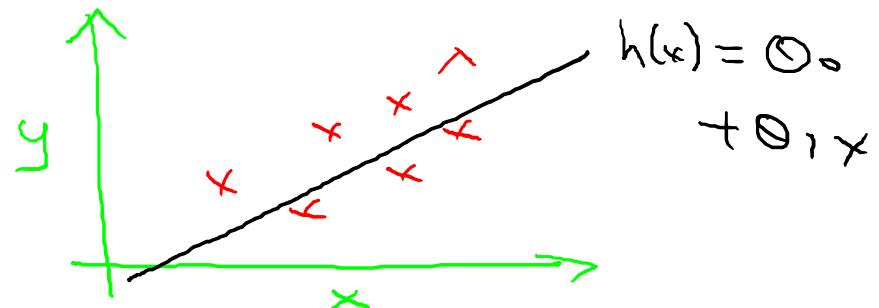
$(x^{(i)}, y^{(i)})$  –  $i^{\text{th}}$  training example



## How do we represent $h$ ?

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

Shorthand:  $h(x)$



Linear regression with one variable (x).  
Univariate linear regression.

# Lesson #1: Linear regression with one variable

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Cost function
  - Gradient descent
  - Gradient descent for linear regression

## Training Set

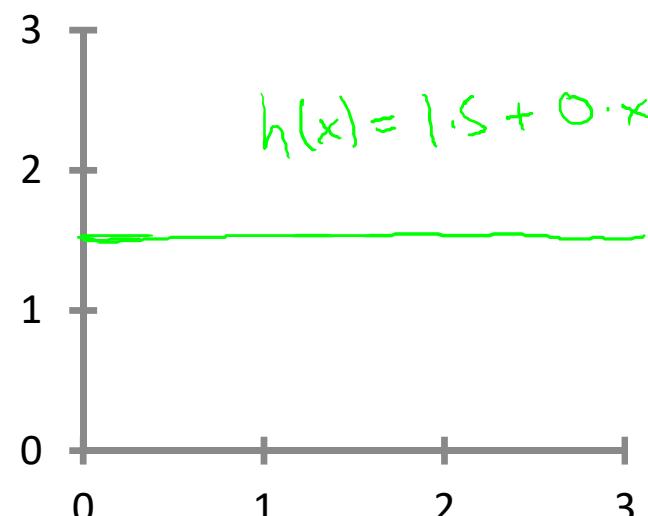
| Size in feet <sup>2</sup> (x) | Price (\$) in 1000's (y) |
|-------------------------------|--------------------------|
| 2104                          | 460                      |
| 1416                          | 232                      |
| 1534                          | 315                      |
| 852                           | 178                      |
| ...                           | ...                      |

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

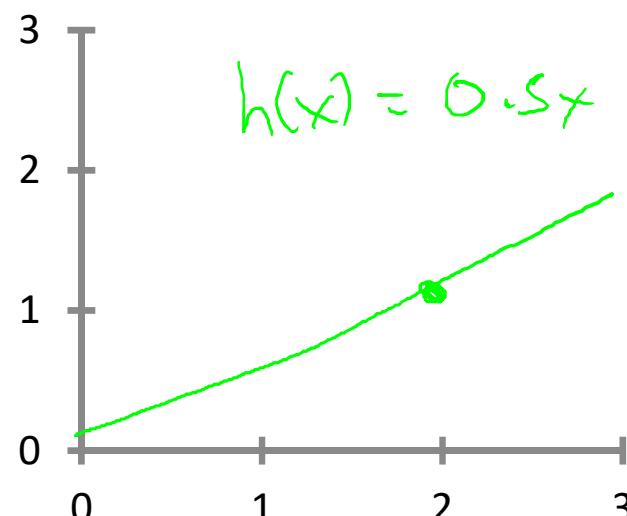
$\theta_i$ 's: Parameters

How to choose  $\theta_i$ 's ?

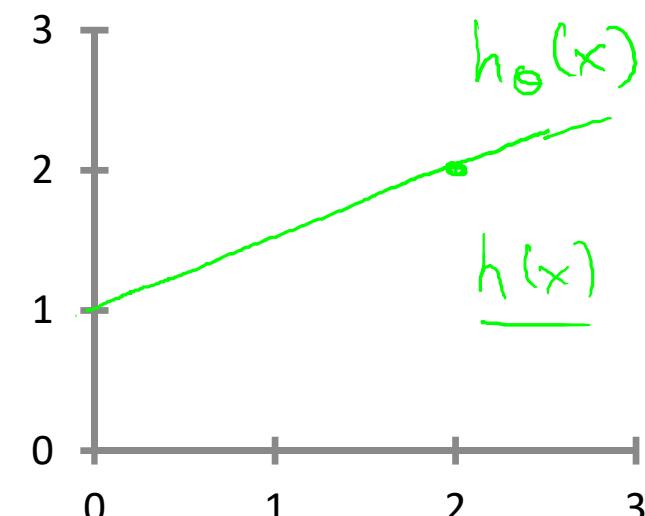
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$



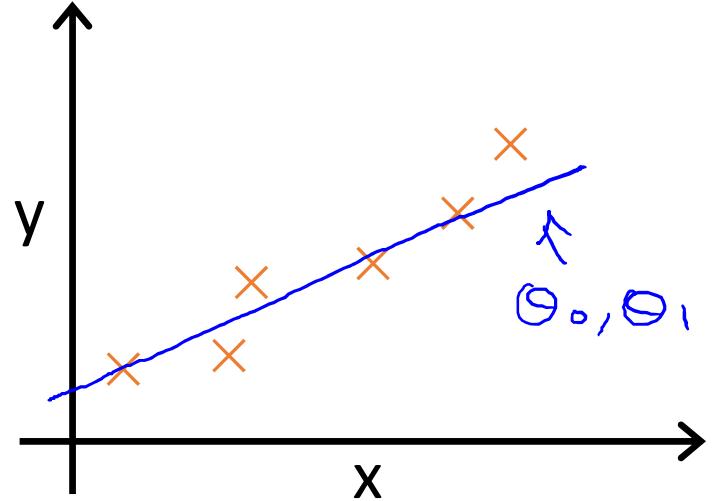
$$\begin{aligned}\theta_0 &= 1.5 \\ \theta_1 &= 0\end{aligned}$$



$$\begin{aligned}\theta_0 &= 0 \\ \theta_1 &= 0.5\end{aligned}$$



$$\begin{aligned}\theta_0 &= 1 \\ \theta_1 &= 0.5\end{aligned}$$



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

Simplified

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost Function:

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

Goal: minimize  $J(\theta_0, \theta_1)$

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

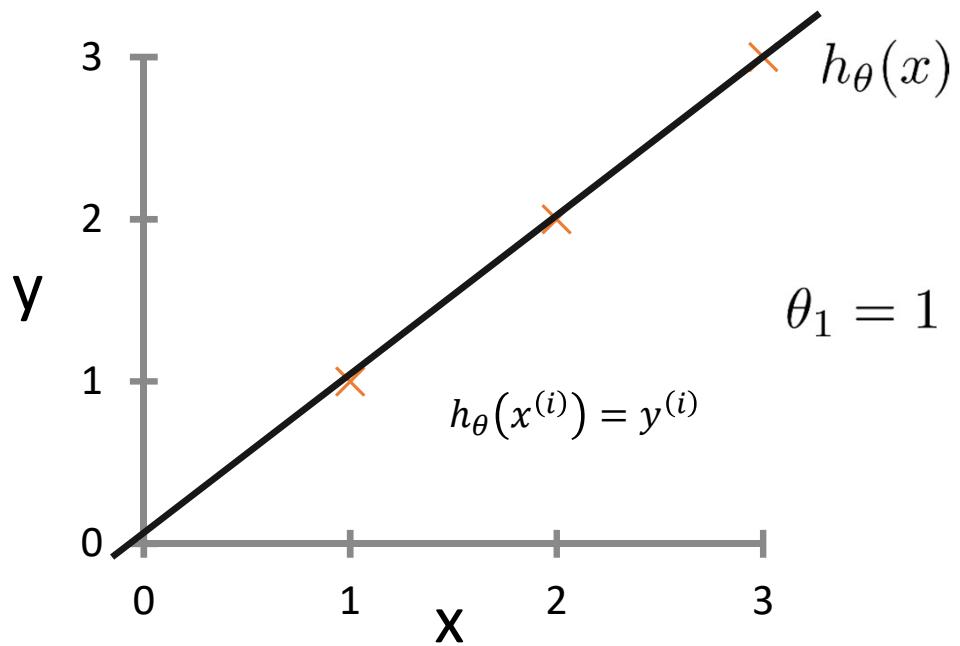
$$\theta_1$$

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

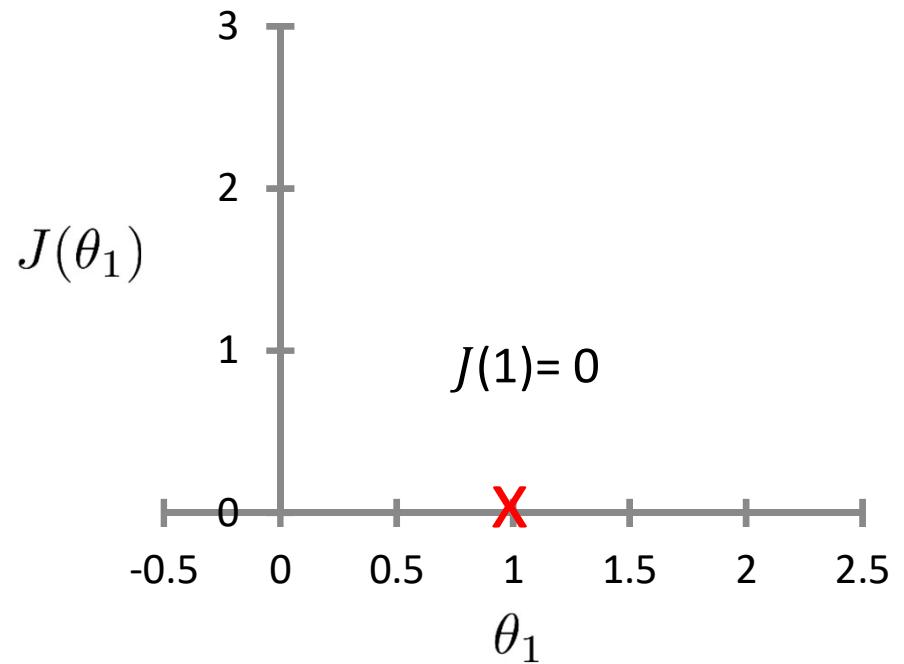
minimize  $J(\theta_1)$

Square error function

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



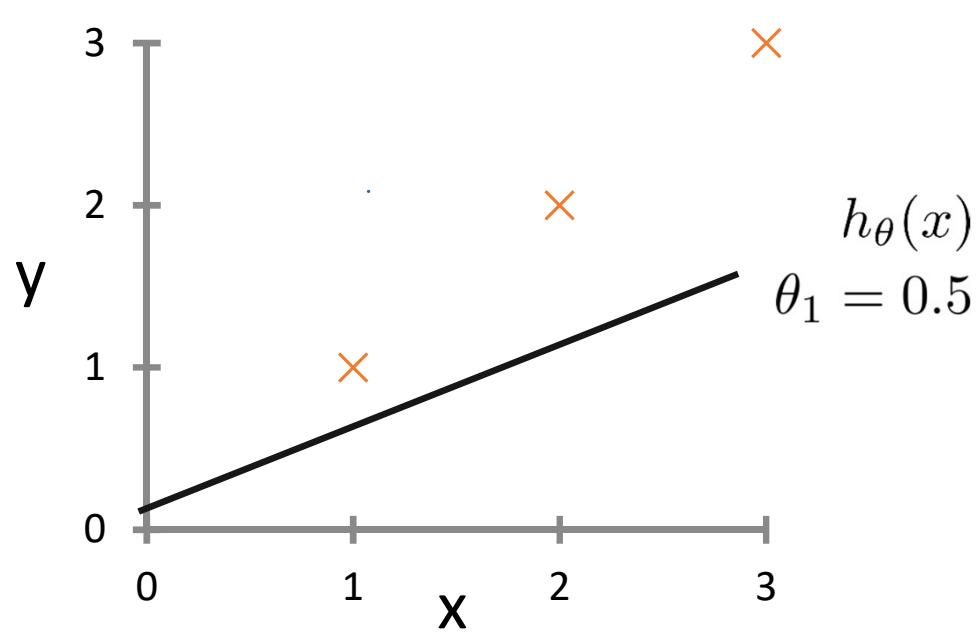
$J(\theta_1)$   
(function of the parameter  $\theta_1$ )



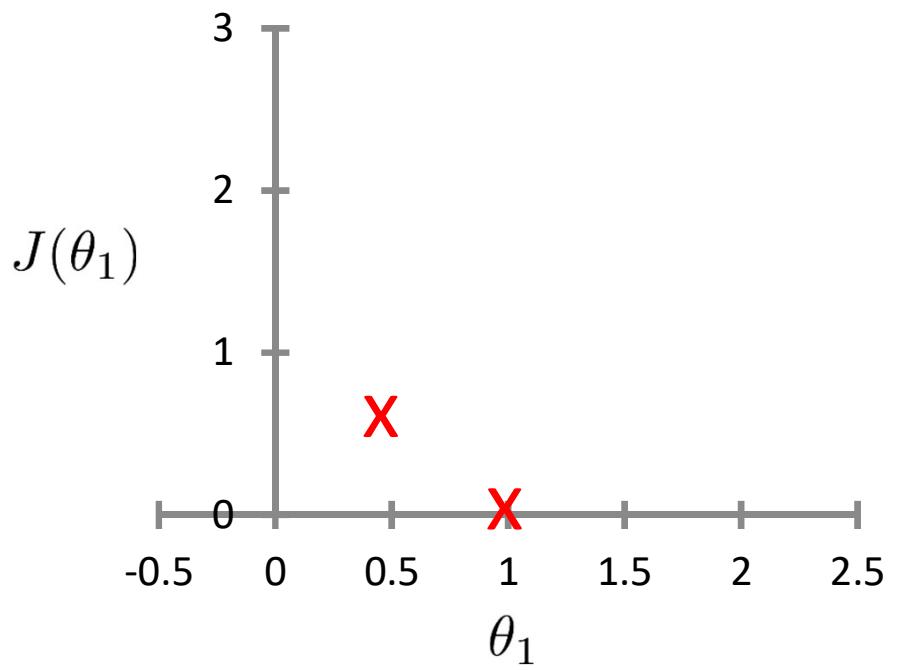
$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0 \end{aligned}$$

$\theta_1 = 0.5?$

$h_{\theta}(x)$   
(for fixed  $\theta_1$ , this is a function of  $x$ )



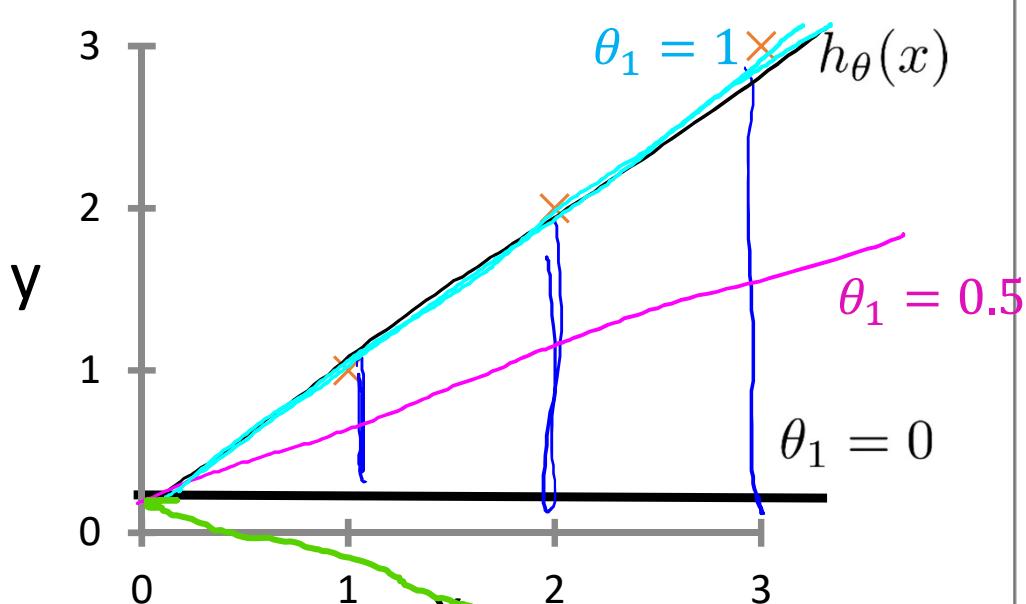
$J(\theta_1)$   
(function of the parameter  $\theta_1$ )



$$J(0.5) = \frac{1}{2m} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] = \frac{1}{2.3} (3.5) \approx 0.58$$

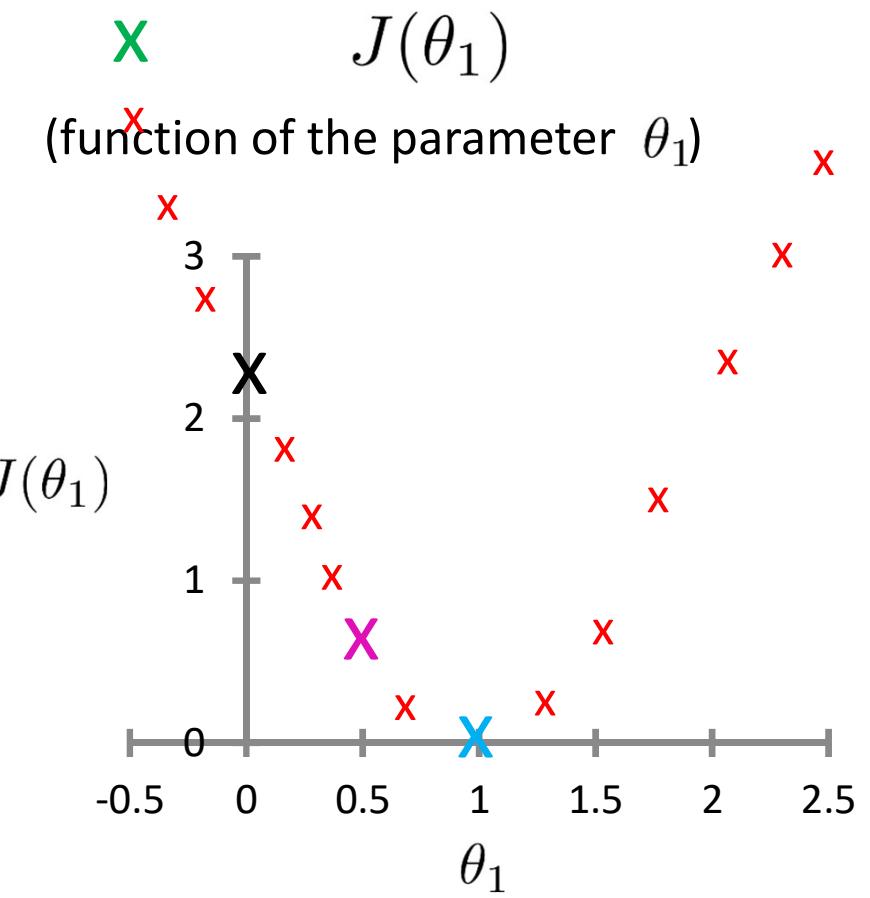
$\theta_1 = 0 ?$   
 $J(0) = ?$

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



$$J(0) = \frac{1}{2m} [1^2 + 2^2 + 3^2] = \frac{1}{2.3} (14) \approx 2.3$$

$$J(-0.5) = \frac{1}{2m} [1.5^2 + 3^2 + 4.5^2] = \frac{1}{2.3} (31.5) \approx 5.25$$



$$\min_{\theta_1} J(\theta_1)$$

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

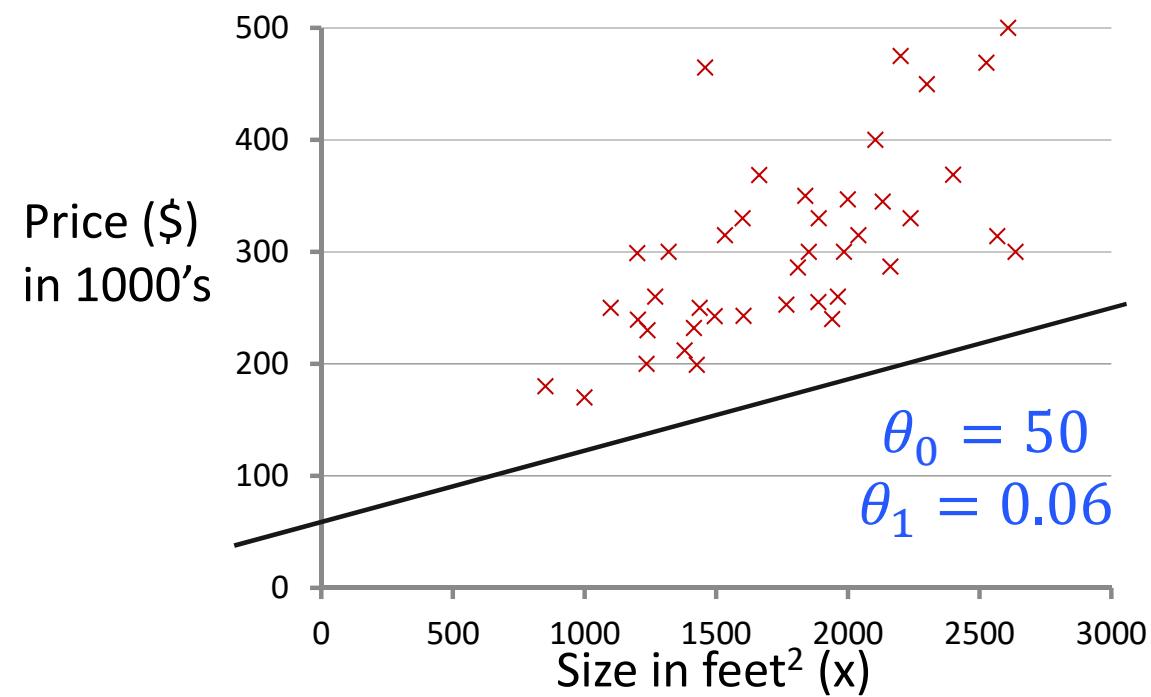
Parameters:  $\theta_0, \theta_1$

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

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

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

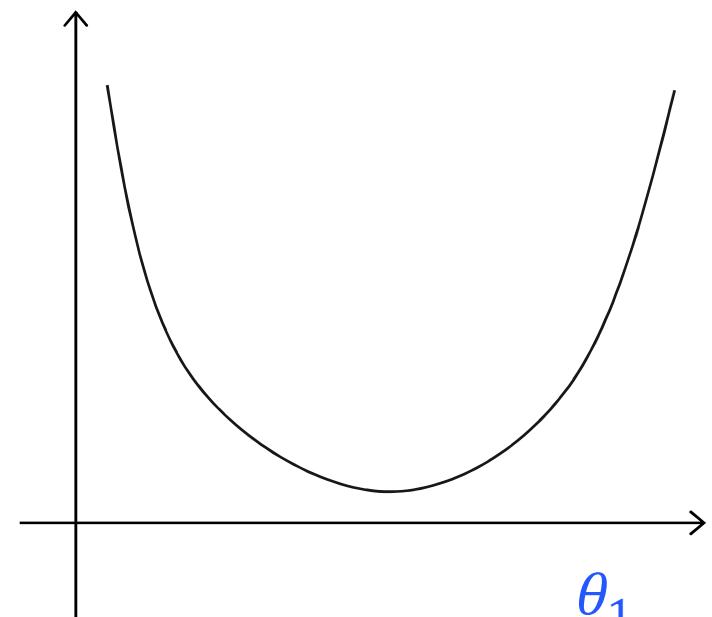
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



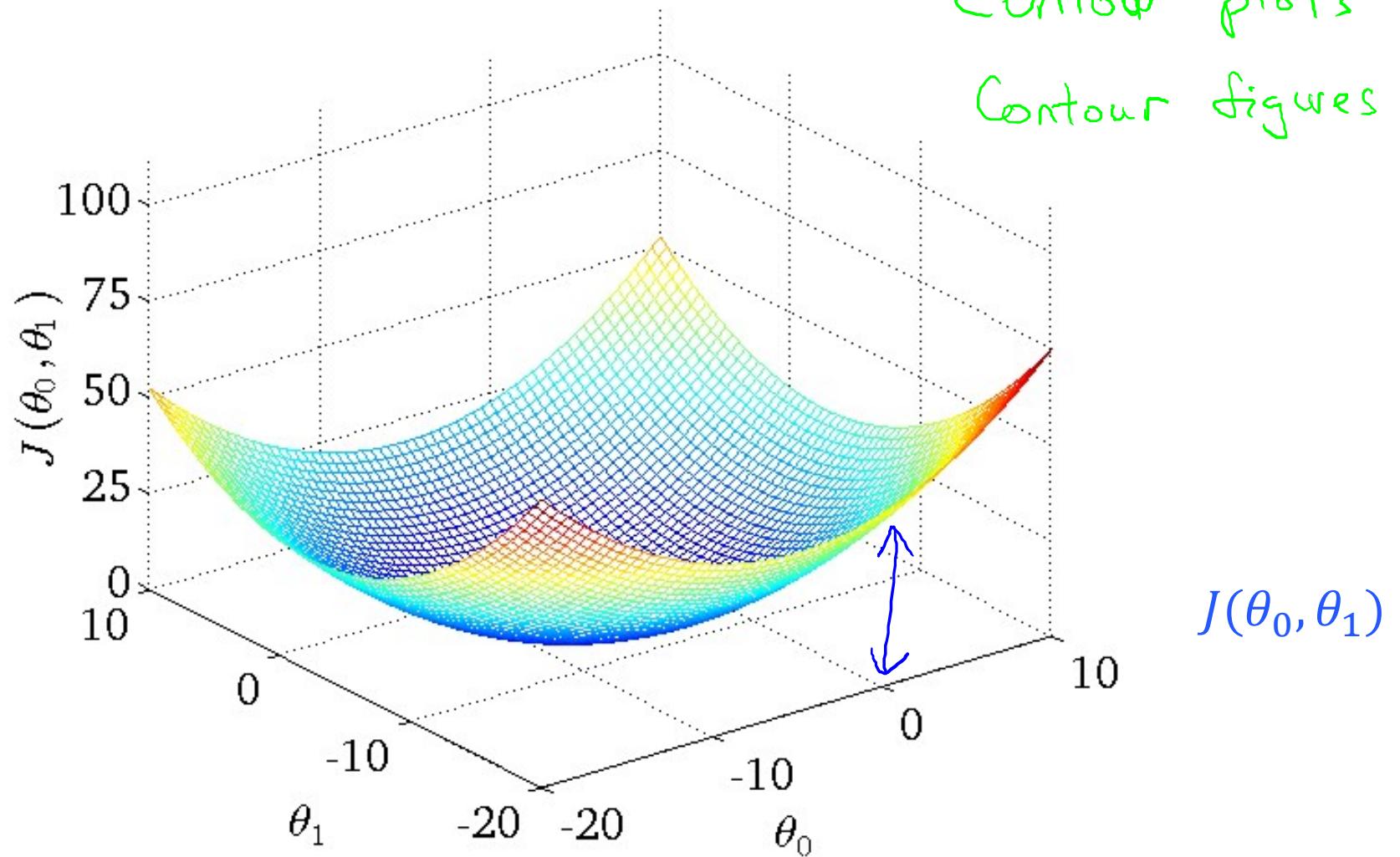
$$h_{\theta}(x) = 50 + 0.06x$$

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

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

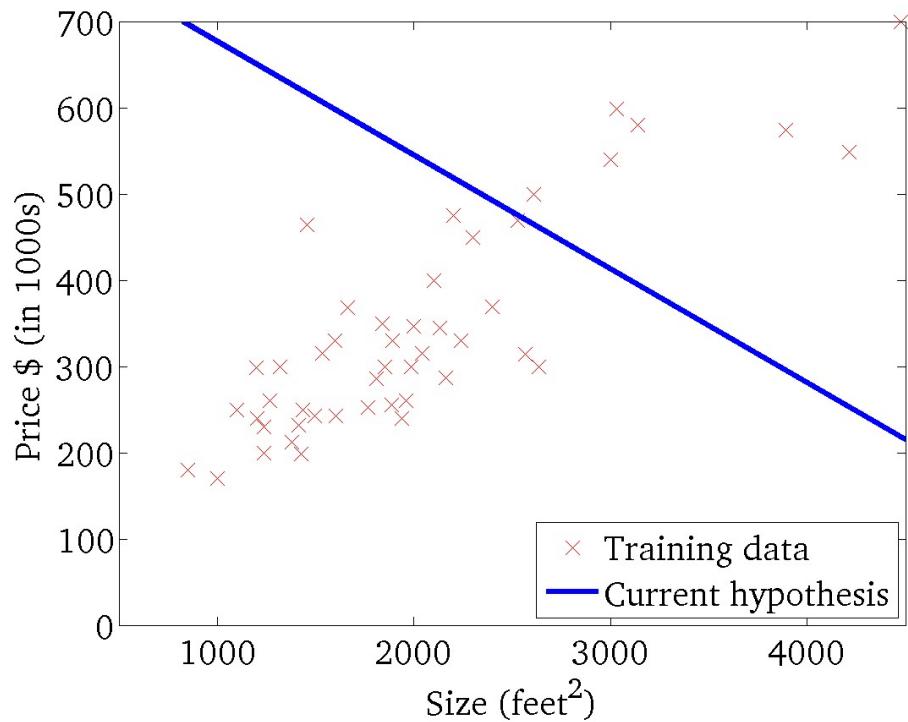


Contour plots  
Contour figures -



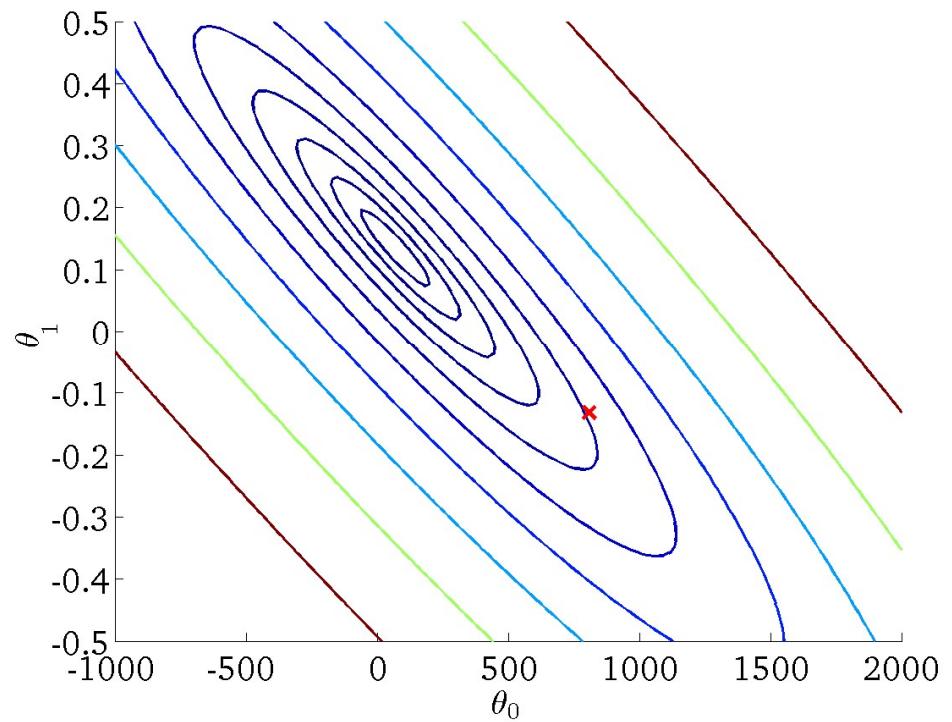
$$h_{\theta}(x)$$

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



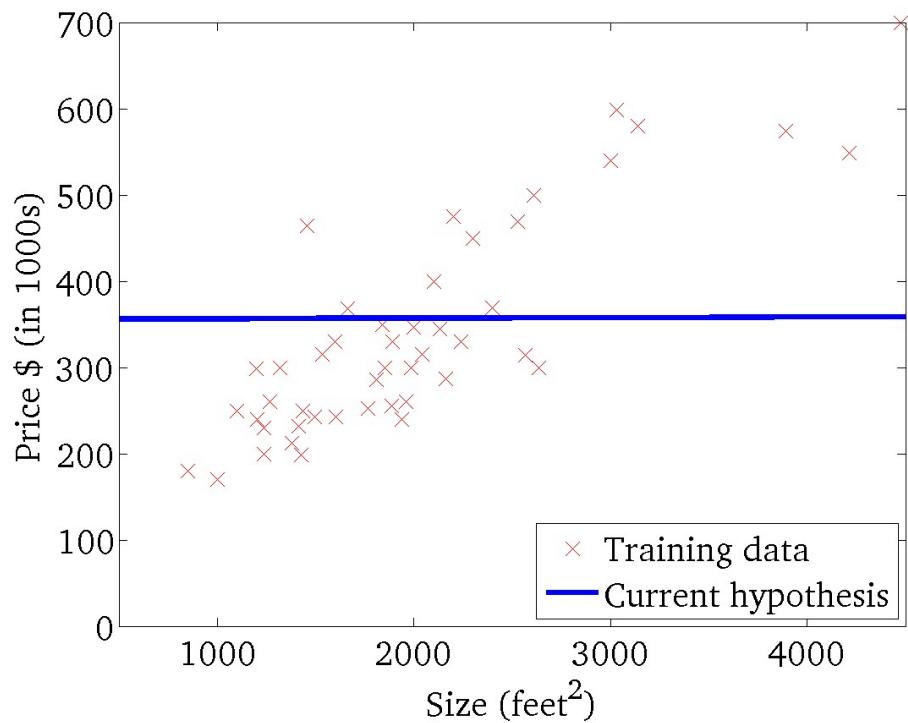
$$J(\theta_0, \theta_1)$$

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



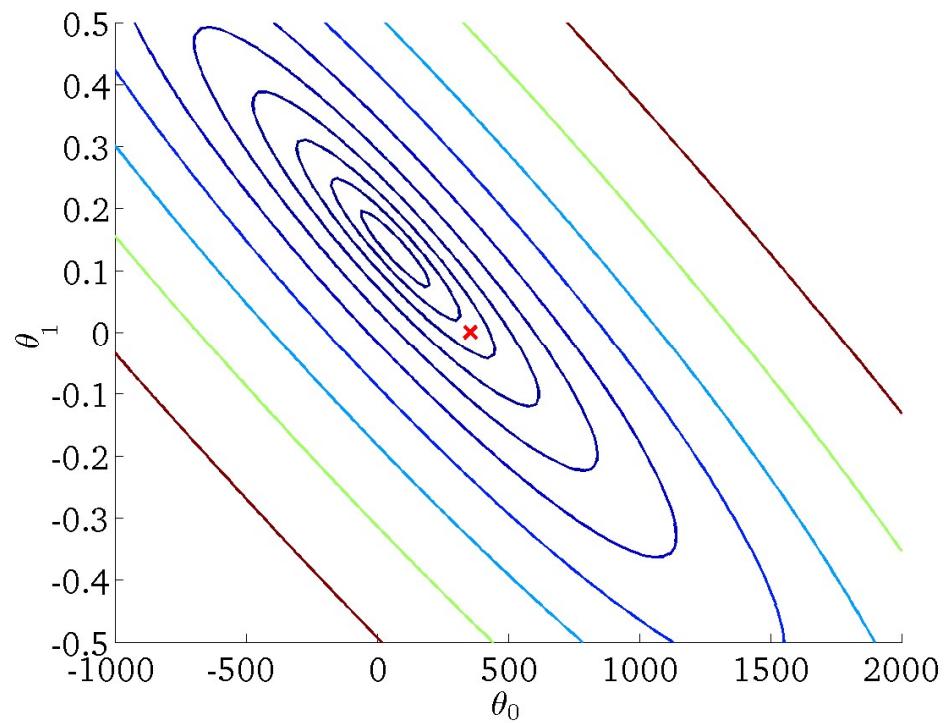
$$h_{\theta}(x)$$

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



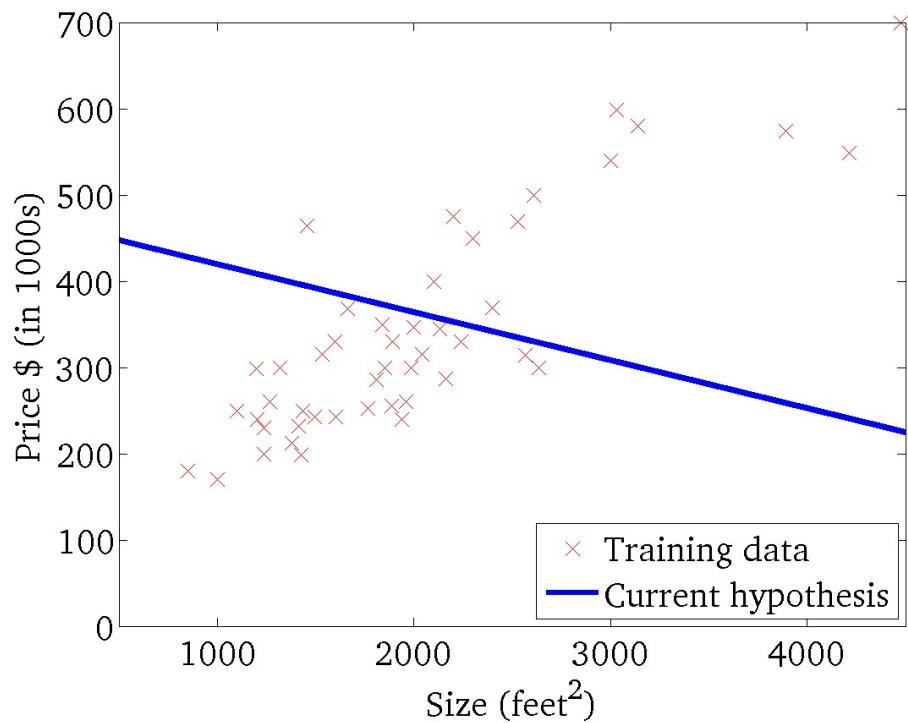
$$J(\theta_0, \theta_1)$$

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



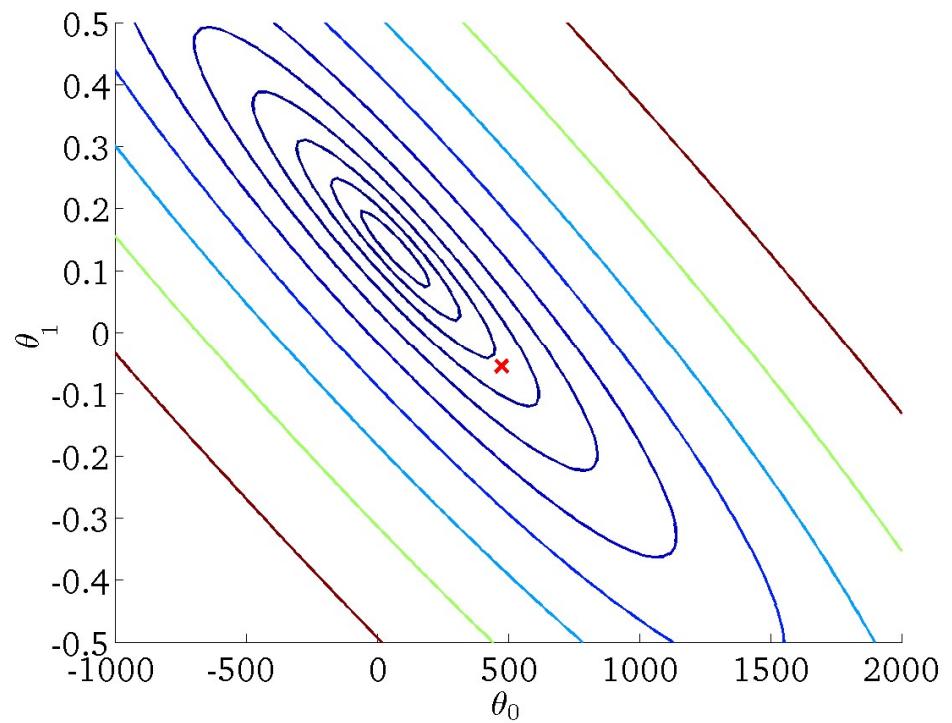
$$h_{\theta}(x)$$

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



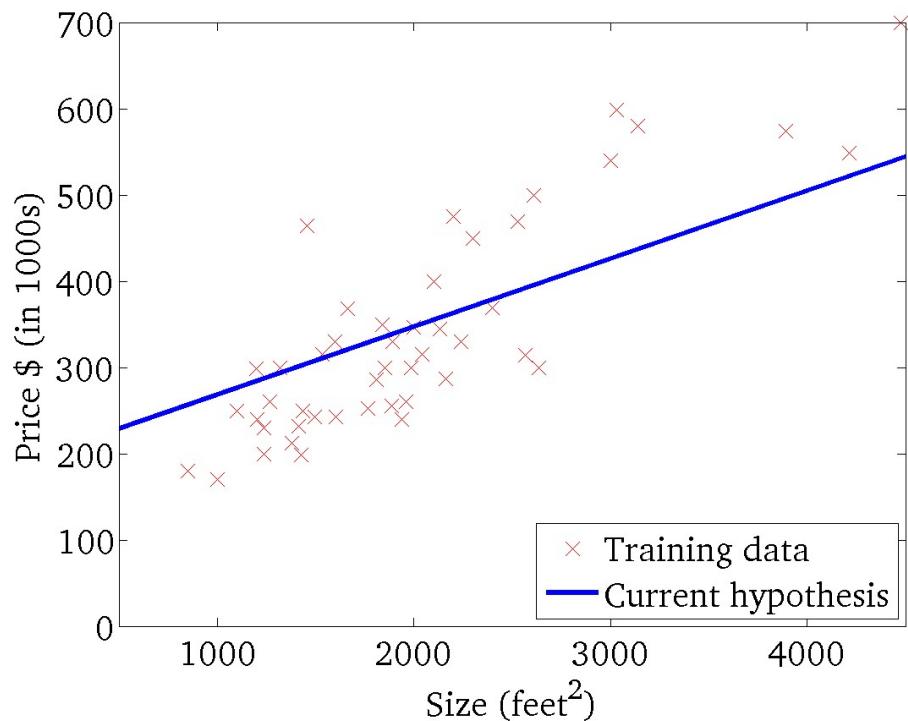
$$J(\theta_0, \theta_1)$$

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



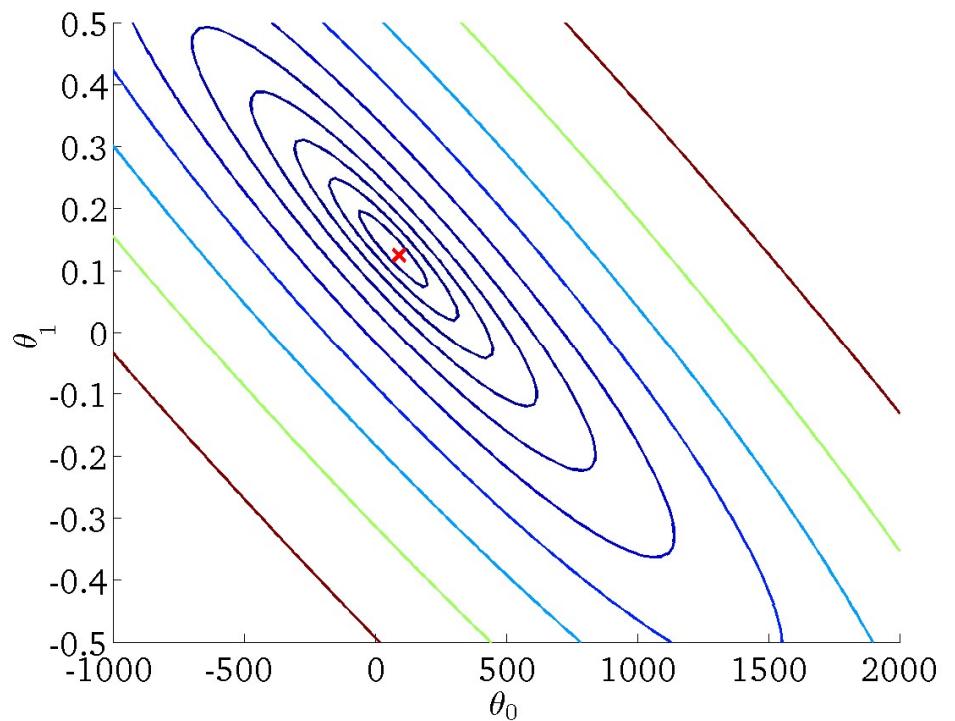
$$h_{\theta}(x)$$

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



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

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



# Lesson #1: Linear regression with one variable

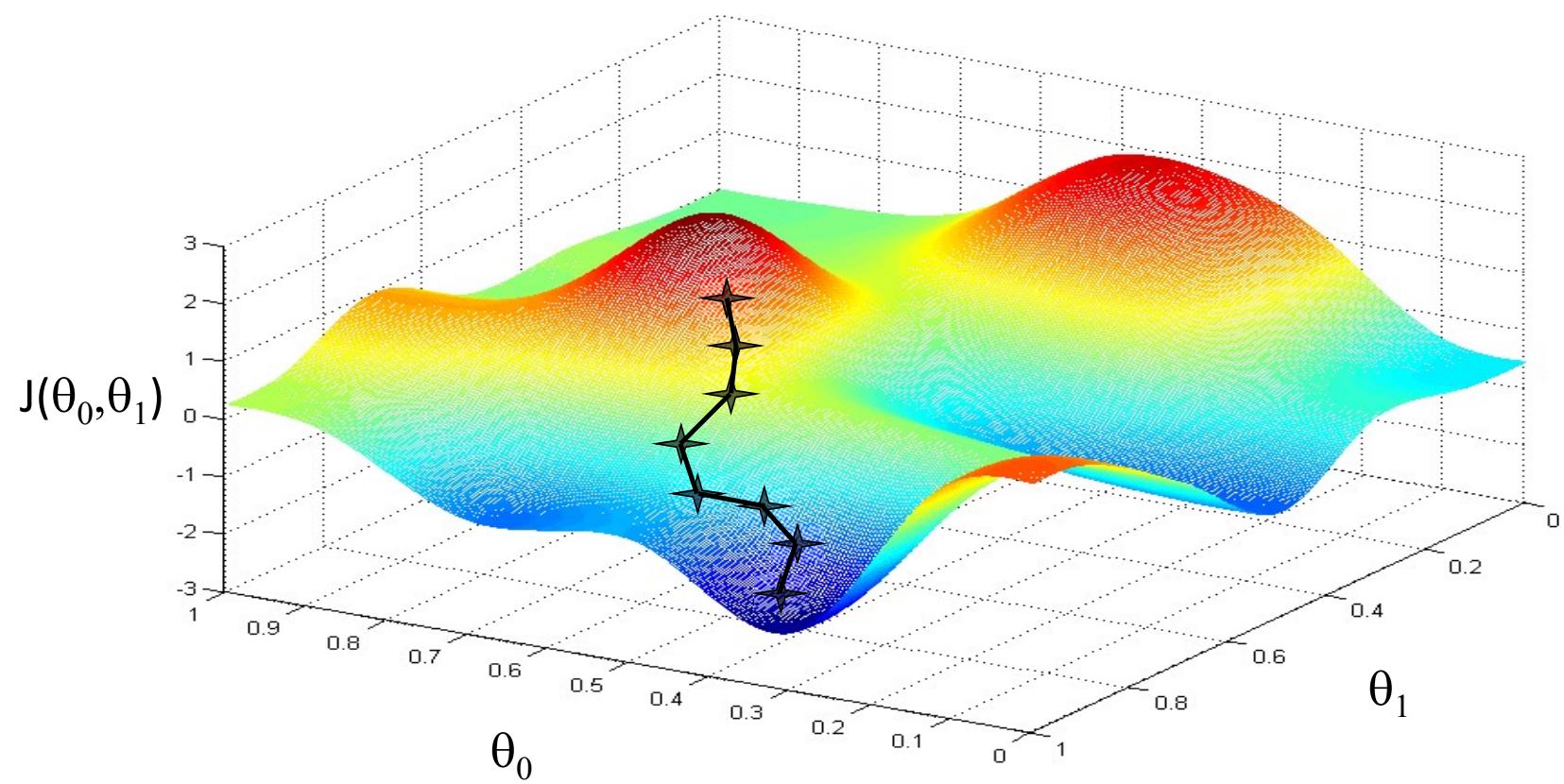
- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Cost function
  - Gradient descent
  - Gradient descent for linear regression

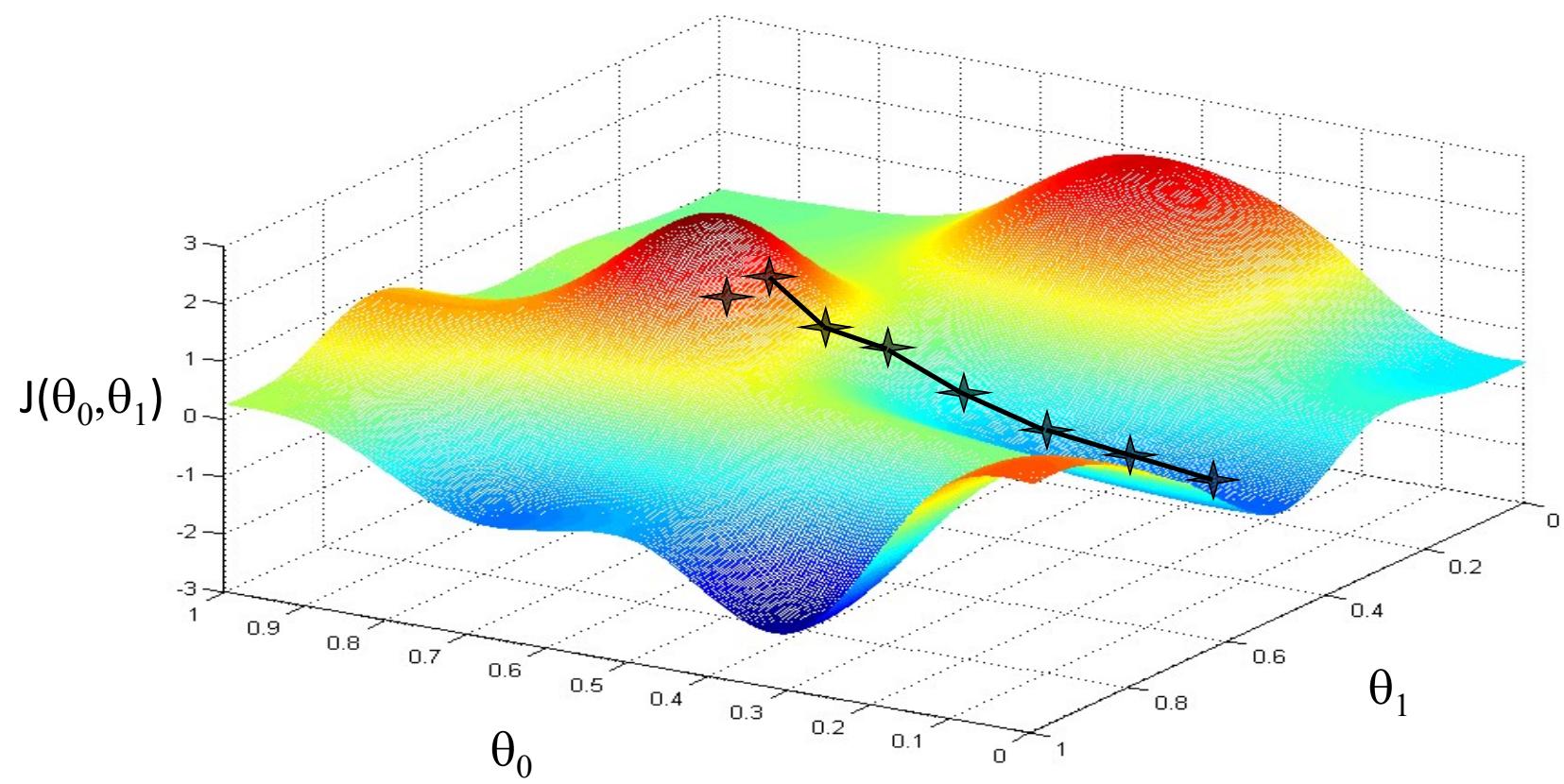
Have some function  $J(\theta_0, \theta_1)$      $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

Want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$      $\min_{\theta_0, \dots, \theta_n} J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

## Outline:

- Start with some  $\theta_0, \theta_1$  (say  $\theta_0 = 0, \theta_1 = 0$ )
- Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$   
until we hopefully end up at a minimum





# Gradient descent algorithm

repeat until convergence {

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

}

*$\alpha$ : learning rate*

---

Correct: Simultaneous update

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

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

$$\theta_0 := \text{temp0}$$

$$\theta_1 := \text{temp1}$$

Incorrect:

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

$$\theta_0 := \text{temp0}$$

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

$$\theta_1 := \text{temp1}$$

# Gradient descent algorithm

repeat until convergence {

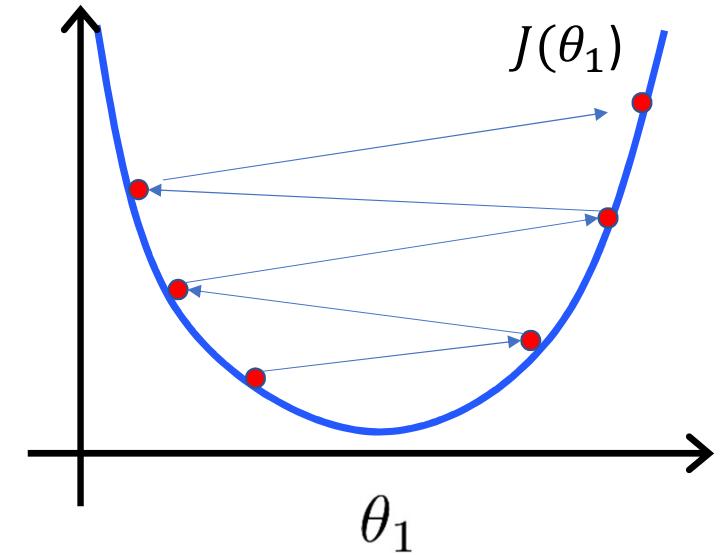
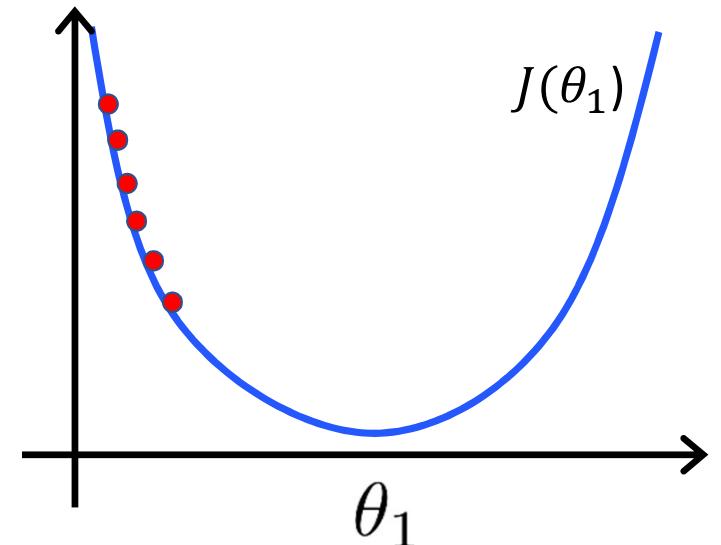
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad \begin{matrix} \text{(simultaneously update} \\ j = 0 \text{ and } j = 1 \end{matrix}$$

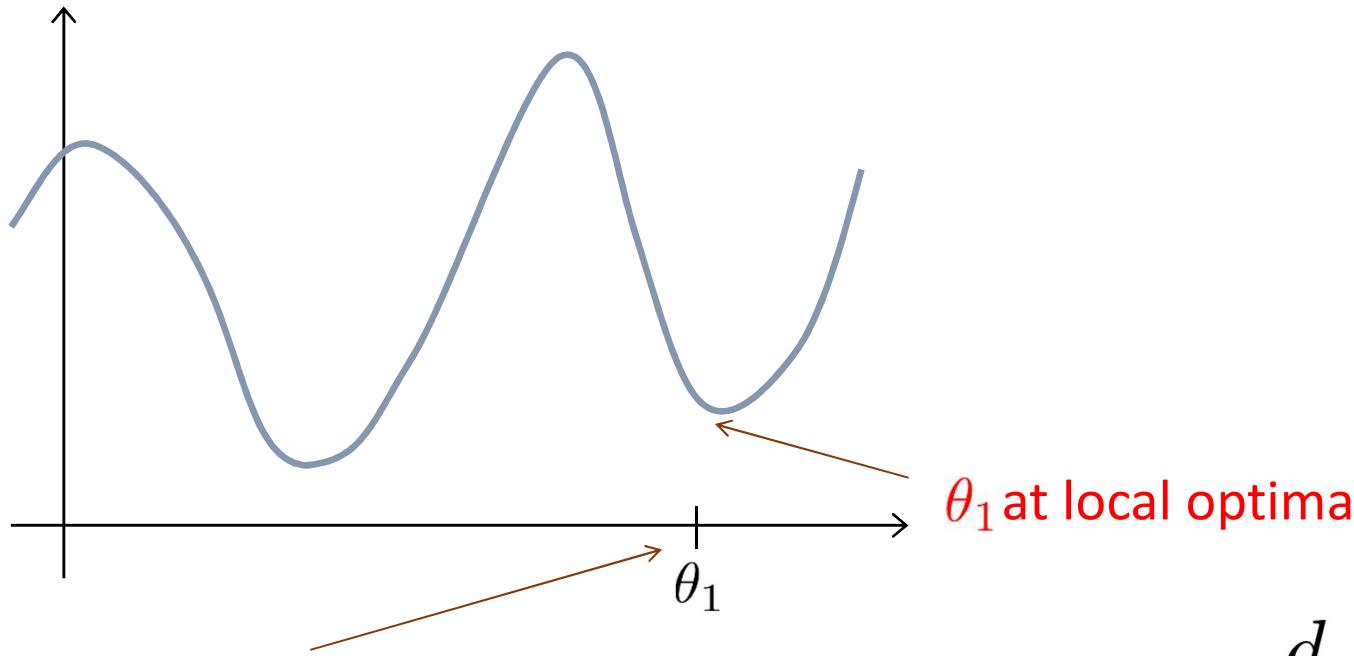
}

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

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.





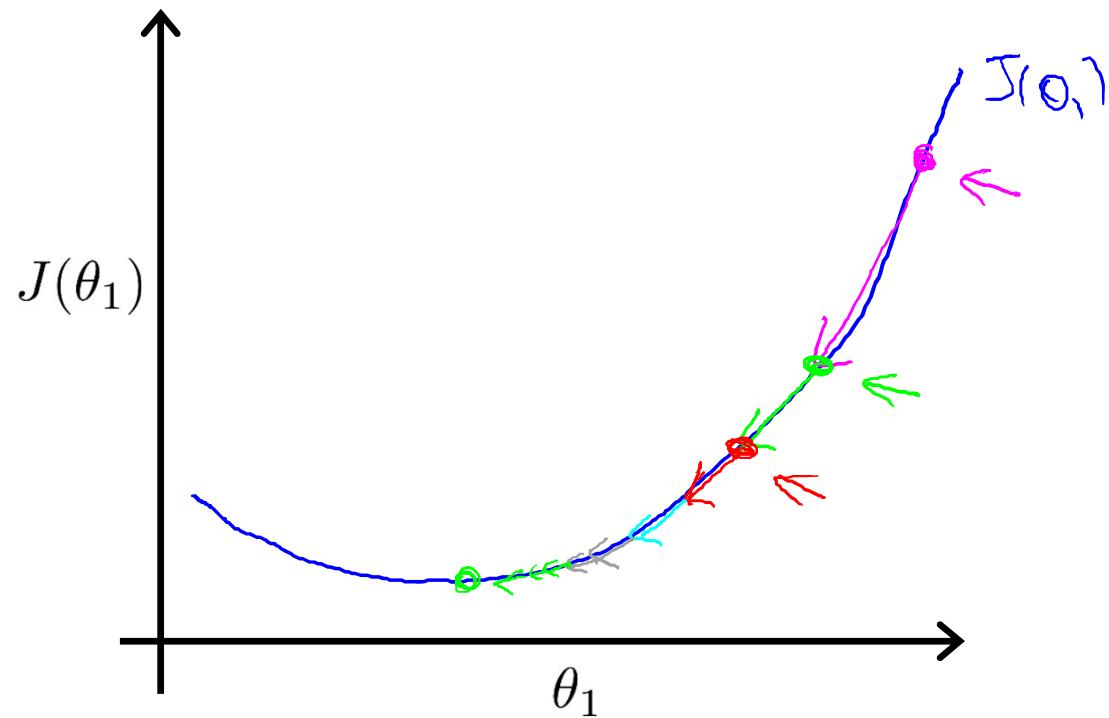
Current value of  $\theta_1$

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

Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.

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

As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.



# Lesson #1: Linear regression with one variable

- **Duration:** 2 hrs
- **Outline:**
  - Model representation
  - Cost function
  - Gradient descent
  - Gradient descent for linear regression

## Gradient descent algorithm

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

## Linear Regression Model

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

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

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

$$j = \mathbf{0}: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) =$$

$$j = \mathbf{1}: \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) =$$

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

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$j = 0: \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m (h_\theta(x^i) - y^{(i)})$$

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

# Gradient descent algorithm

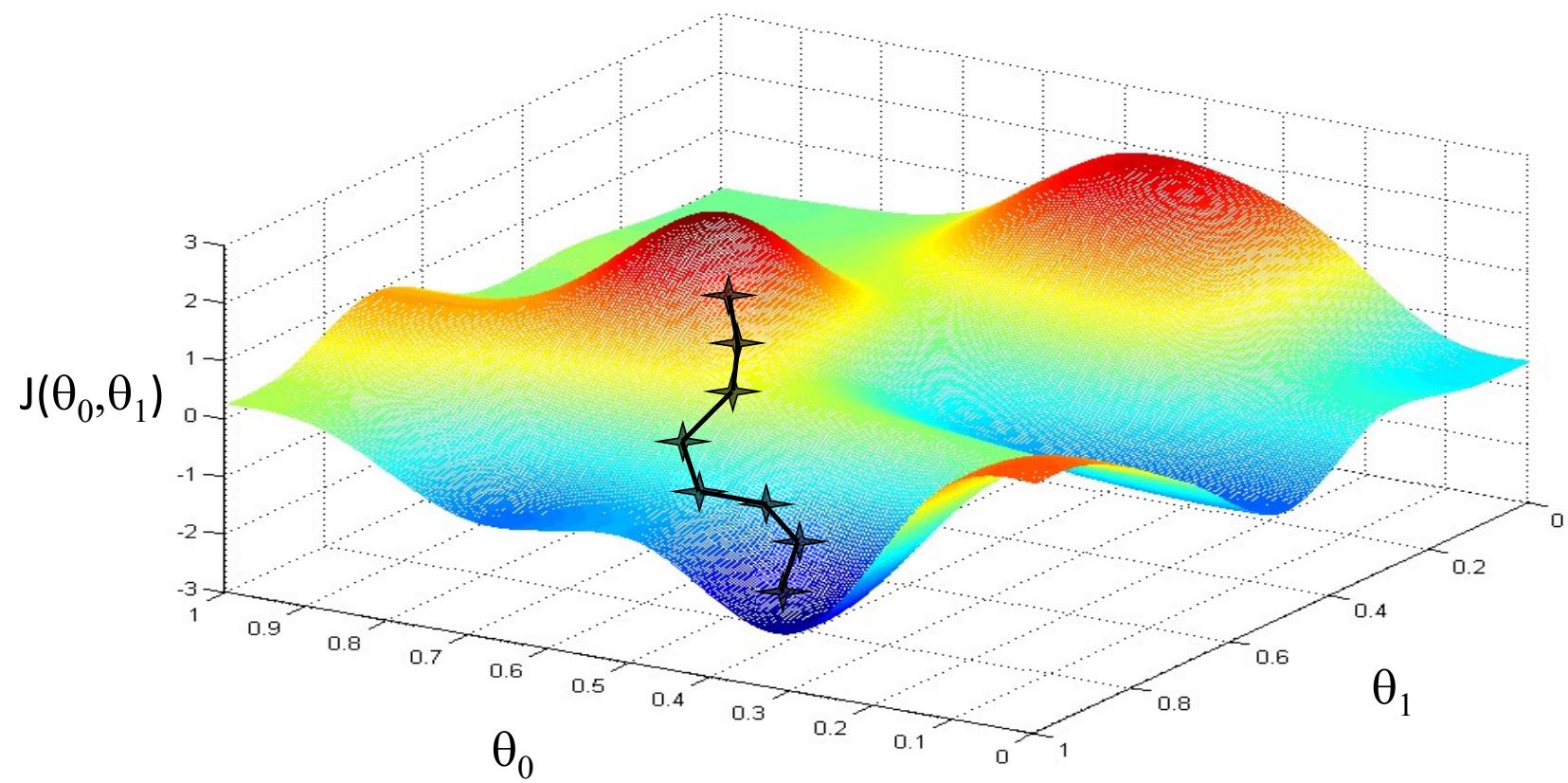
repeat until convergence {

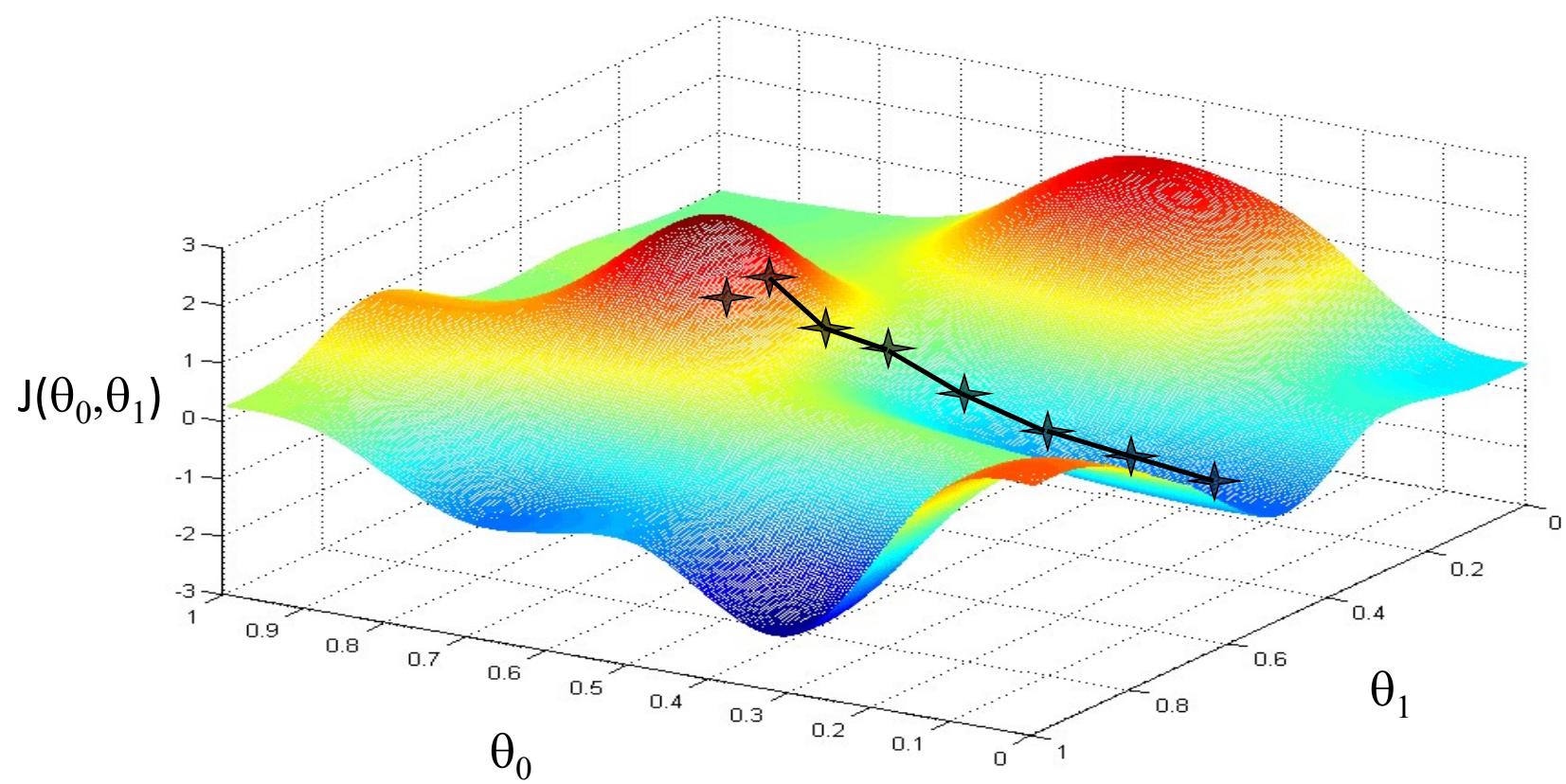
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$

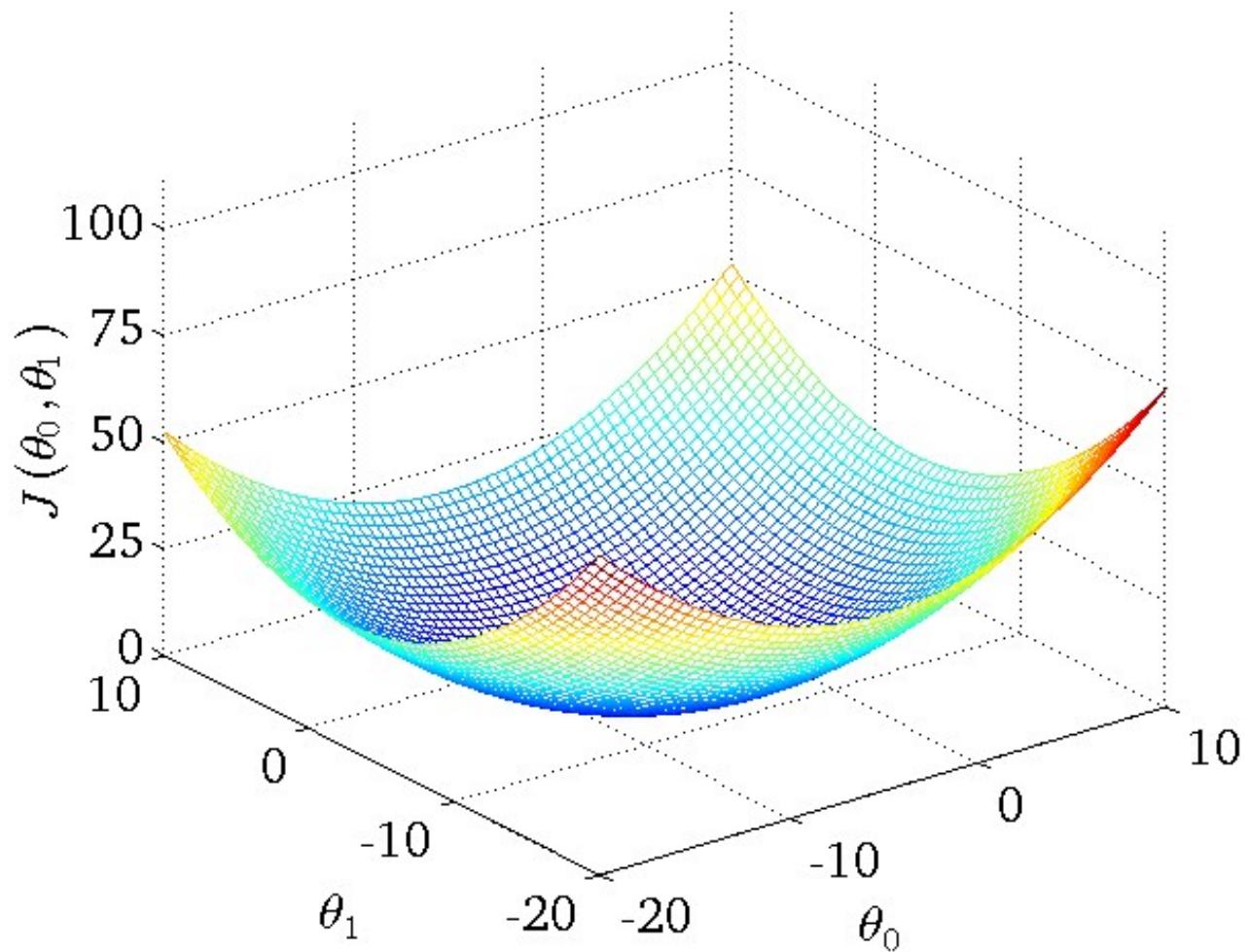
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

}

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

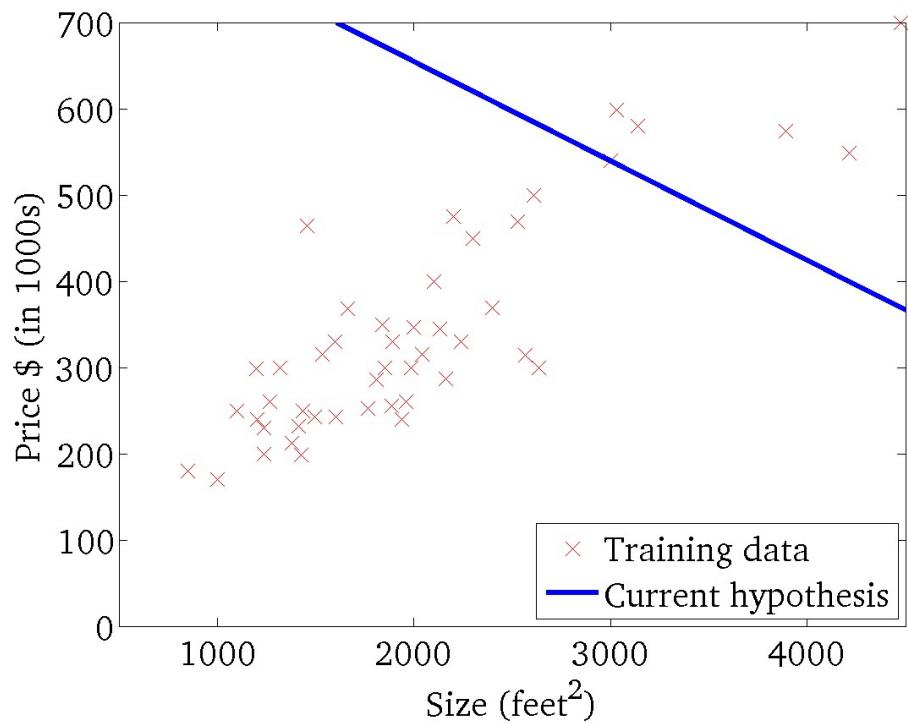






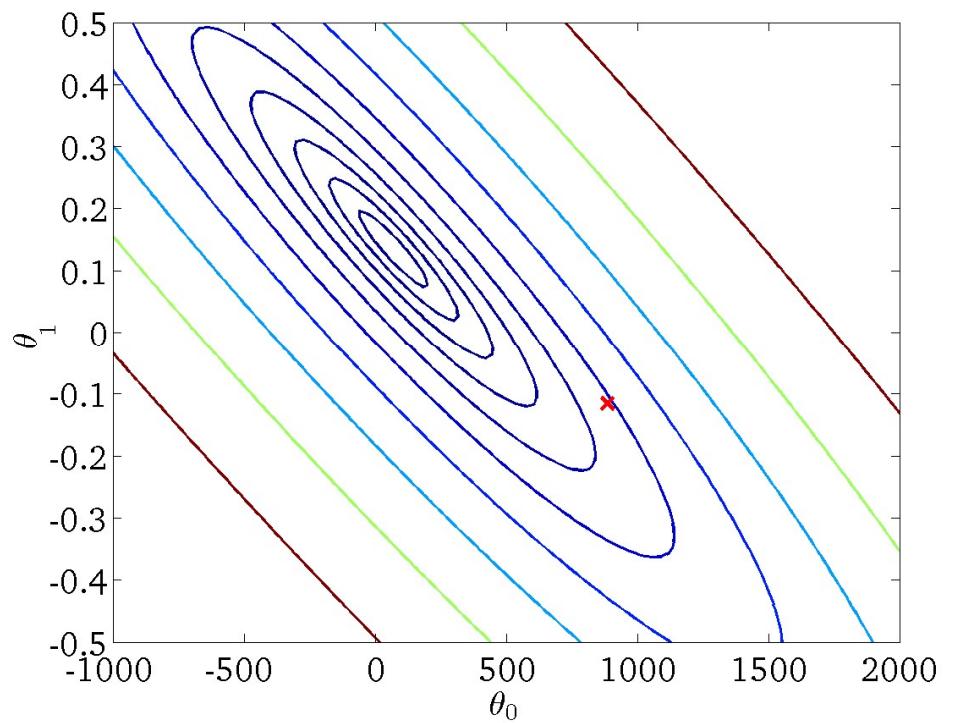
$$h_{\theta}(x)$$

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



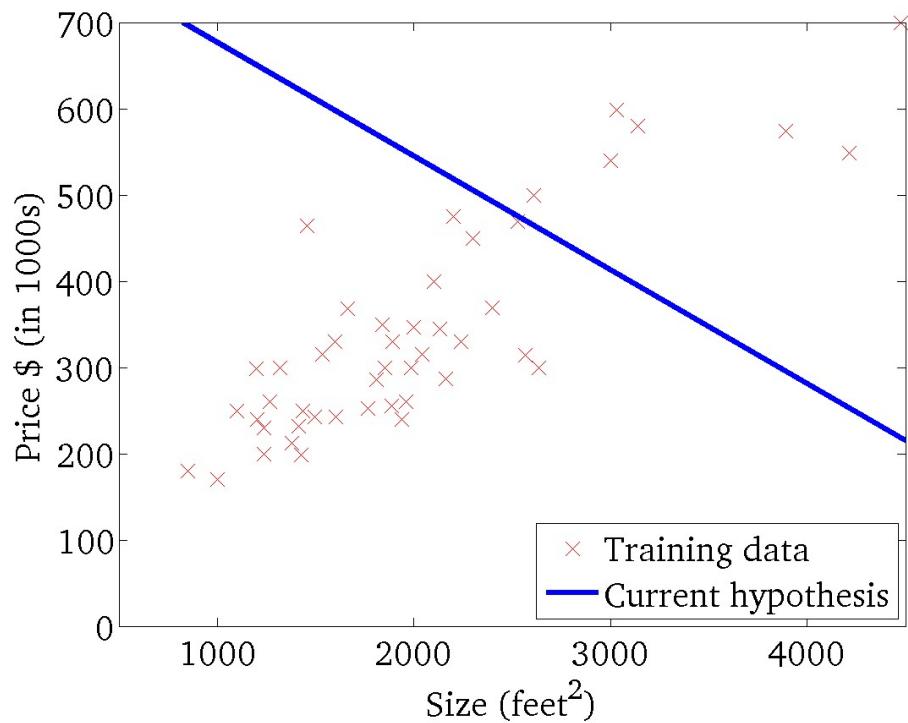
$$J(\theta_0, \theta_1)$$

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



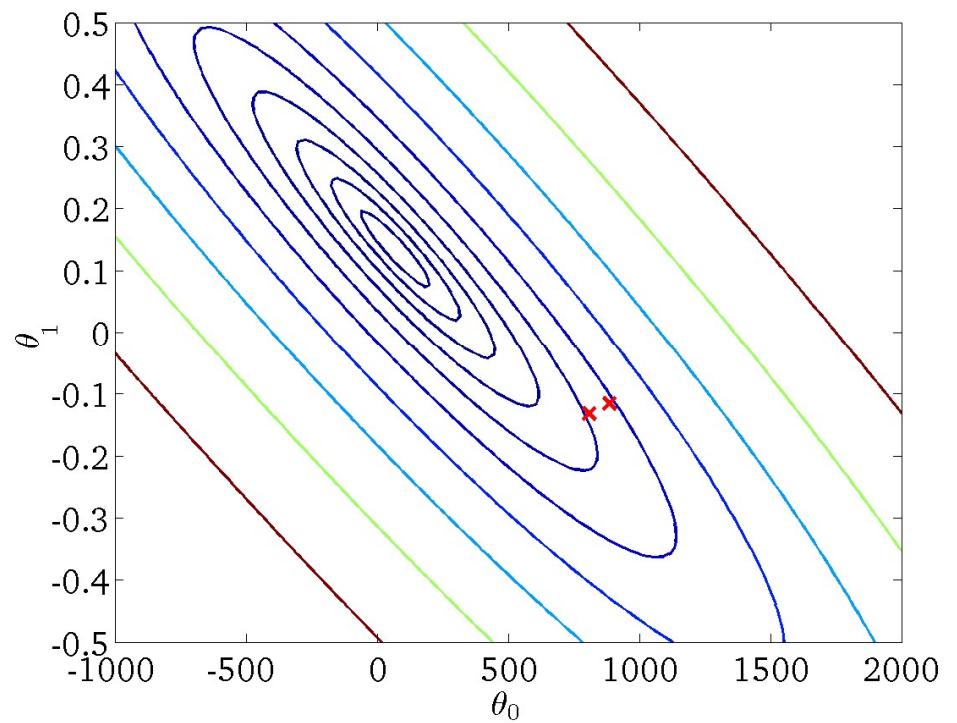
$$h_{\theta}(x)$$

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



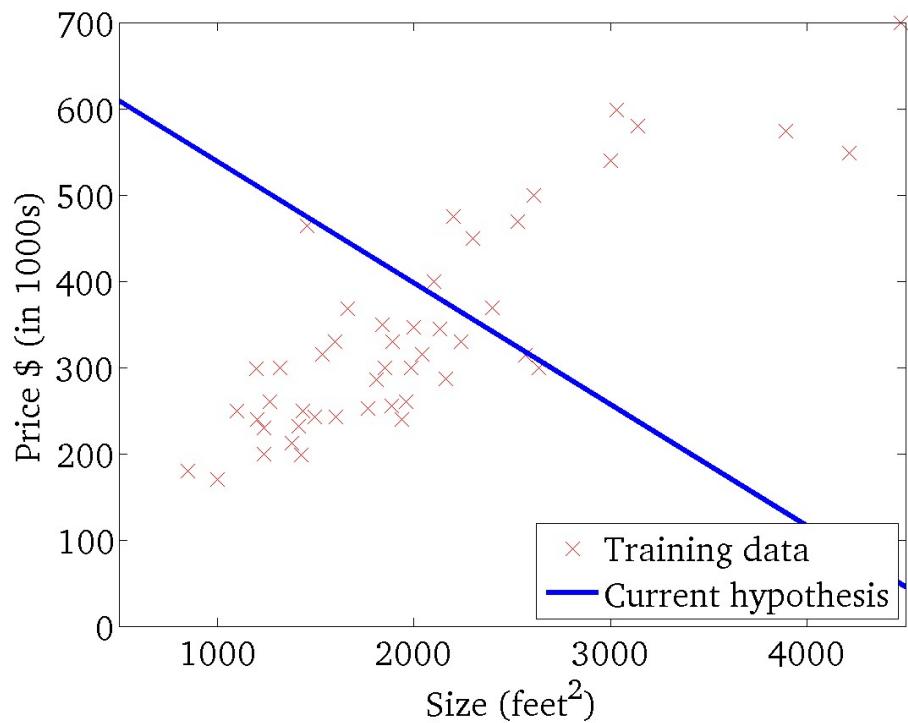
$$J(\theta_0, \theta_1)$$

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



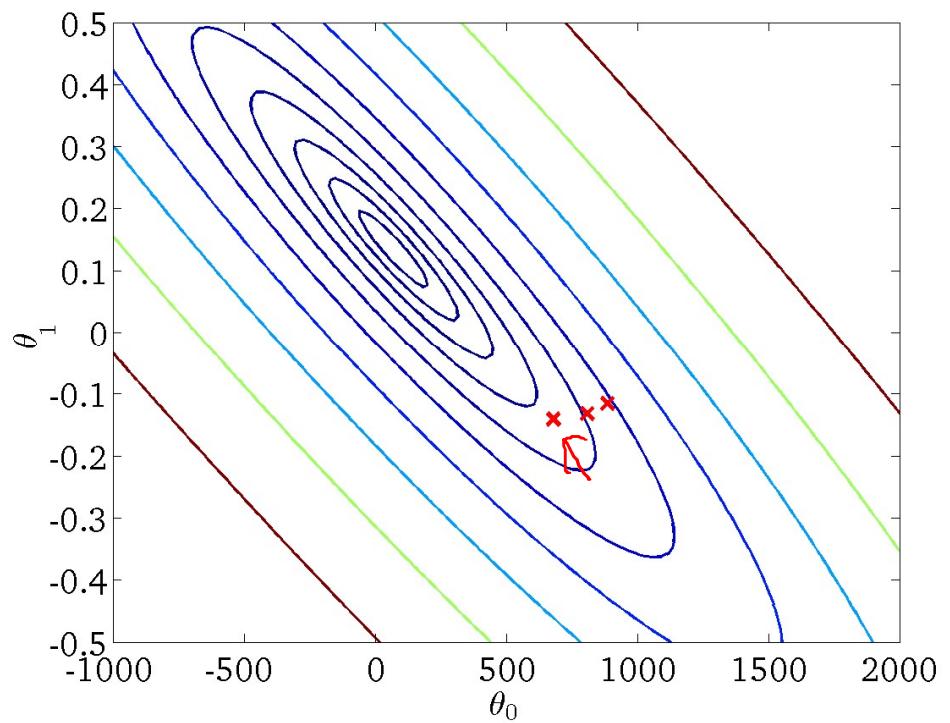
$$h_{\theta}(x)$$

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



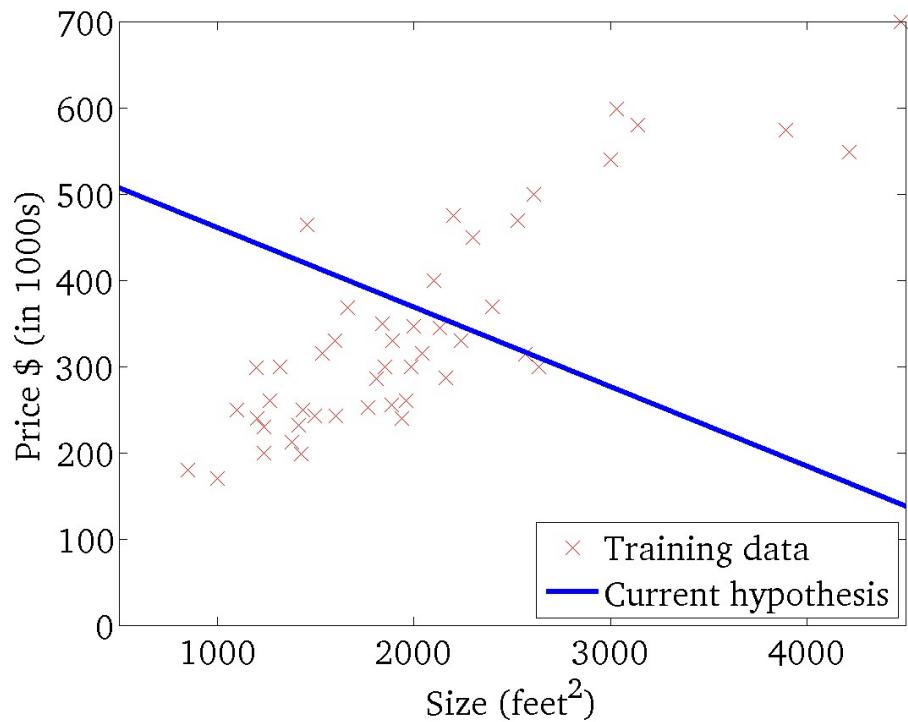
$$J(\theta_0, \theta_1)$$

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



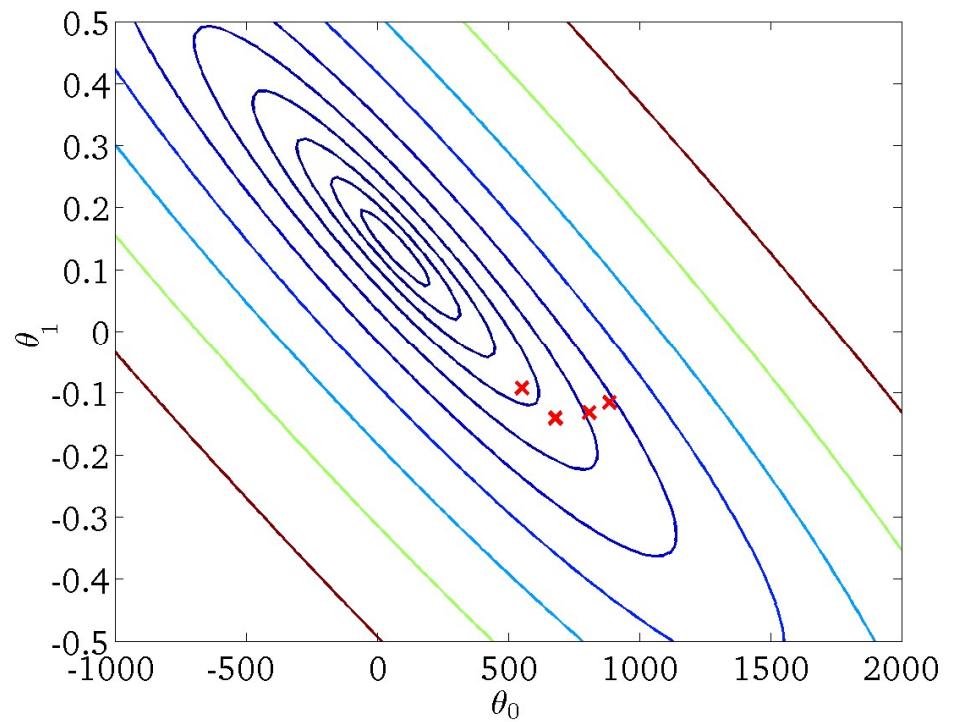
$$h_{\theta}(x)$$

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



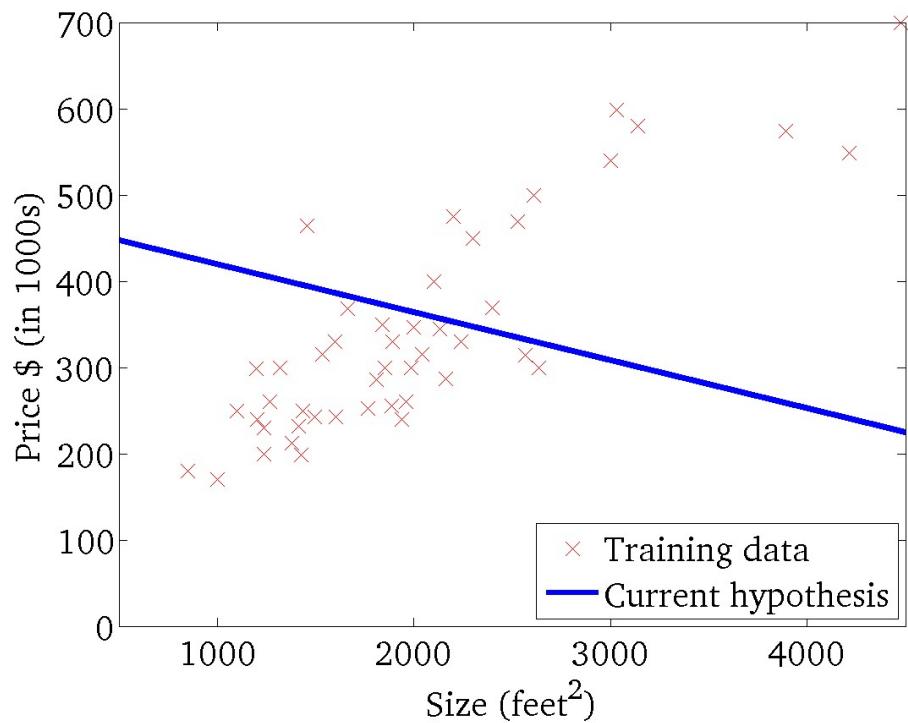
$$J(\theta_0, \theta_1)$$

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



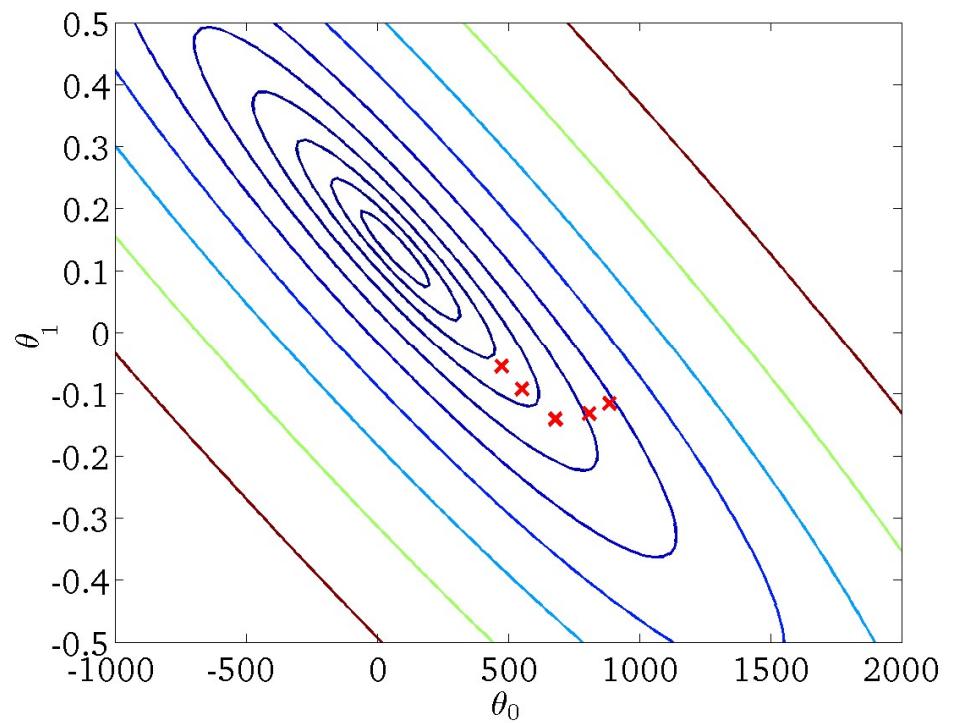
$$h_{\theta}(x)$$

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



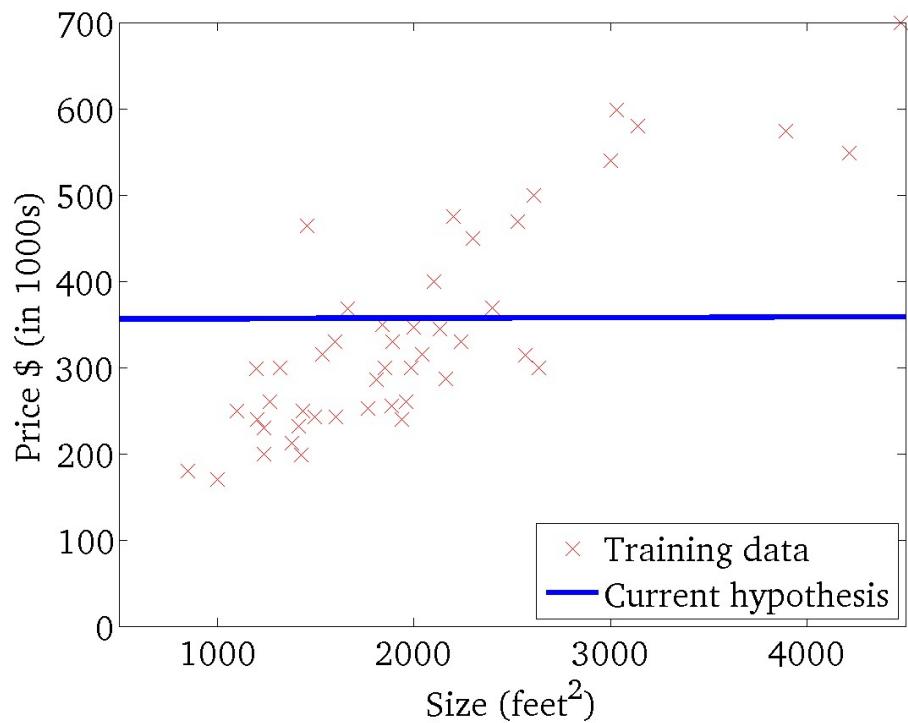
$$J(\theta_0, \theta_1)$$

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



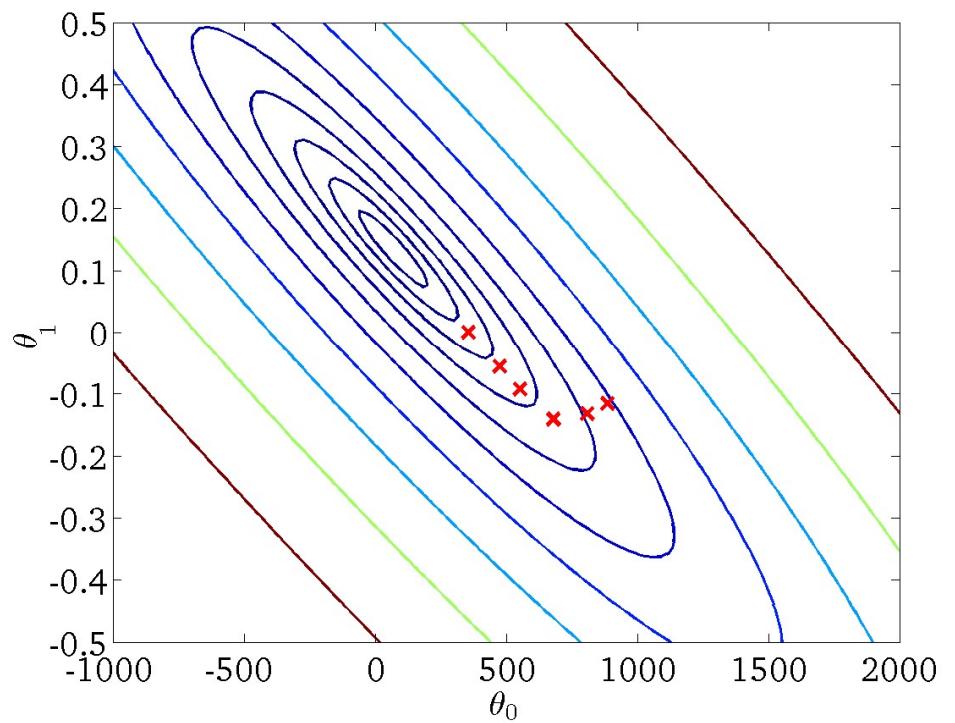
$$h_{\theta}(x)$$

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



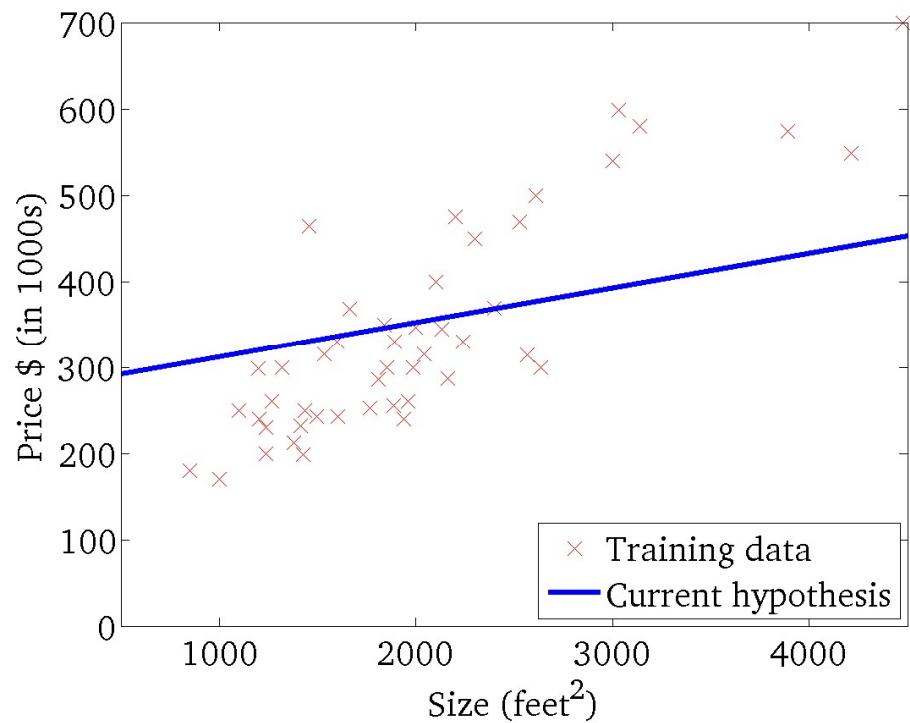
$$J(\theta_0, \theta_1)$$

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



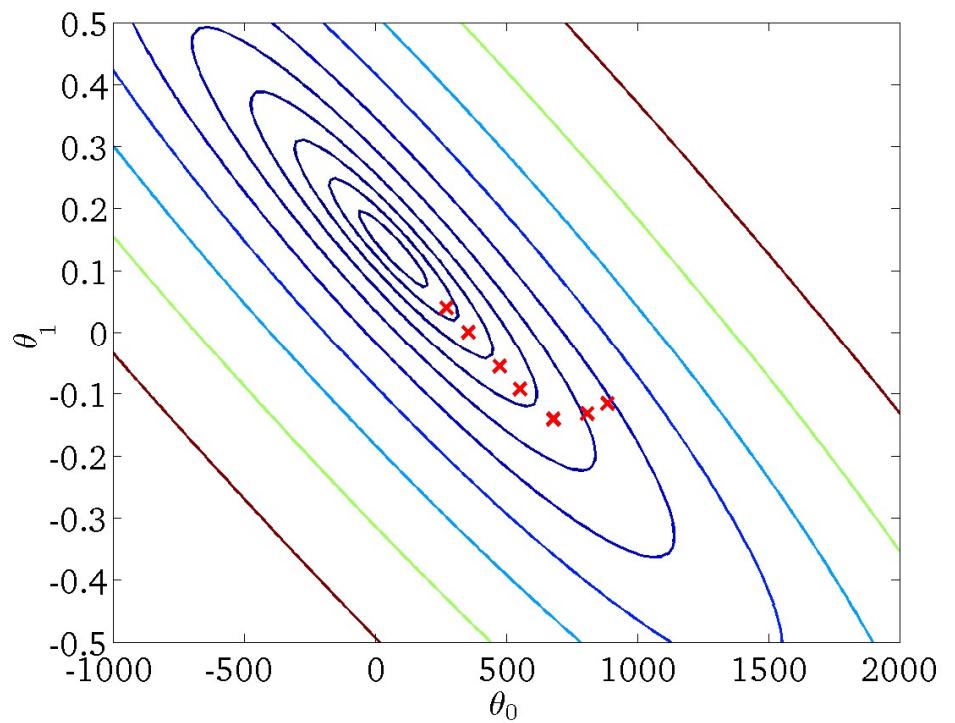
$$h_{\theta}(x)$$

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



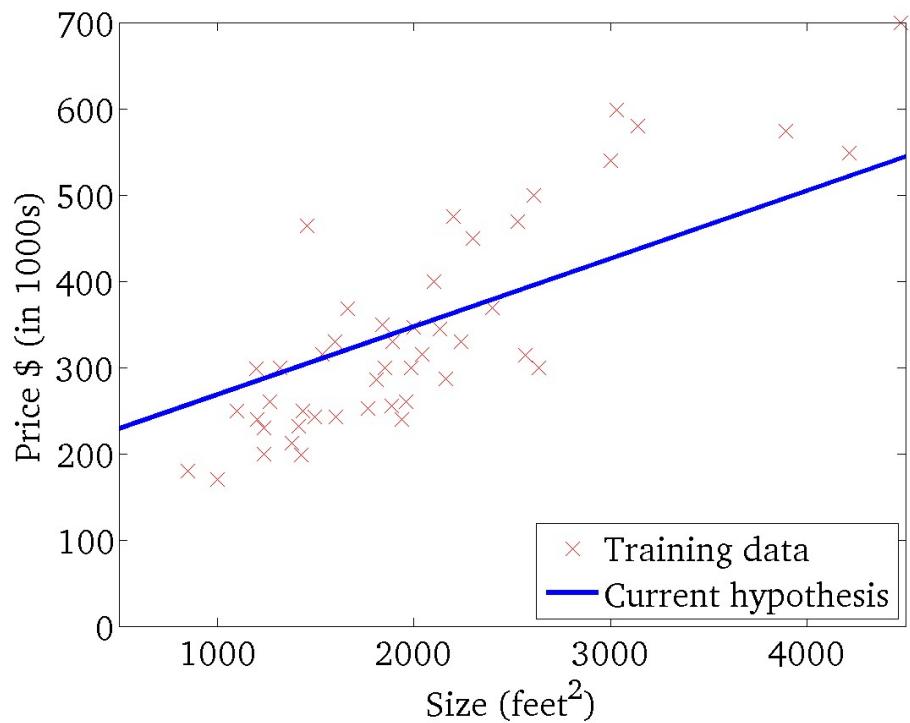
$$J(\theta_0, \theta_1)$$

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



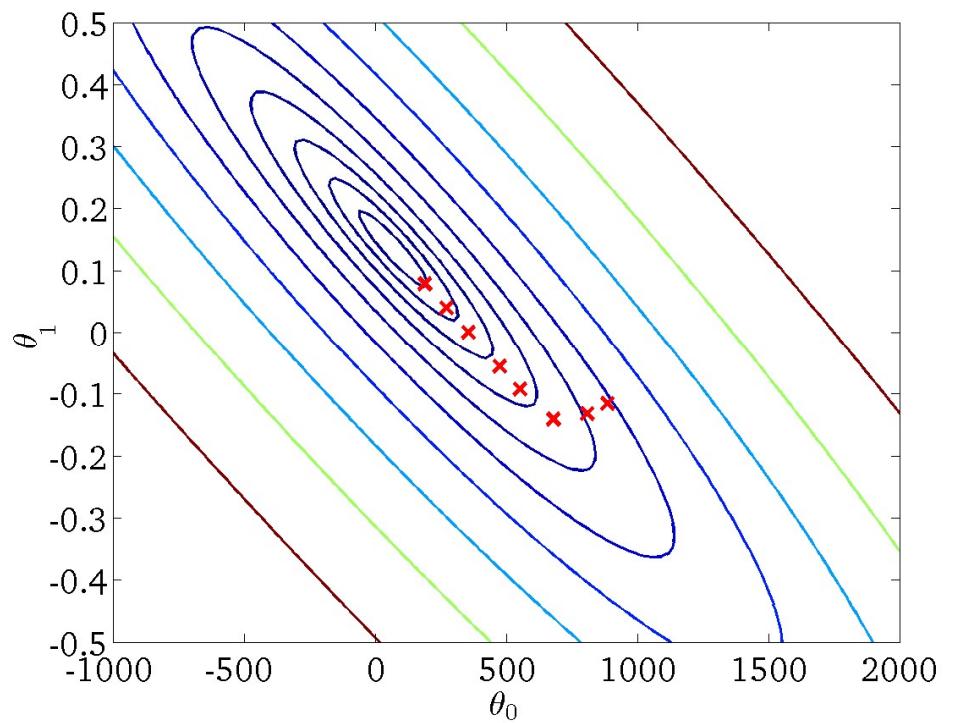
$$h_{\theta}(x)$$

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



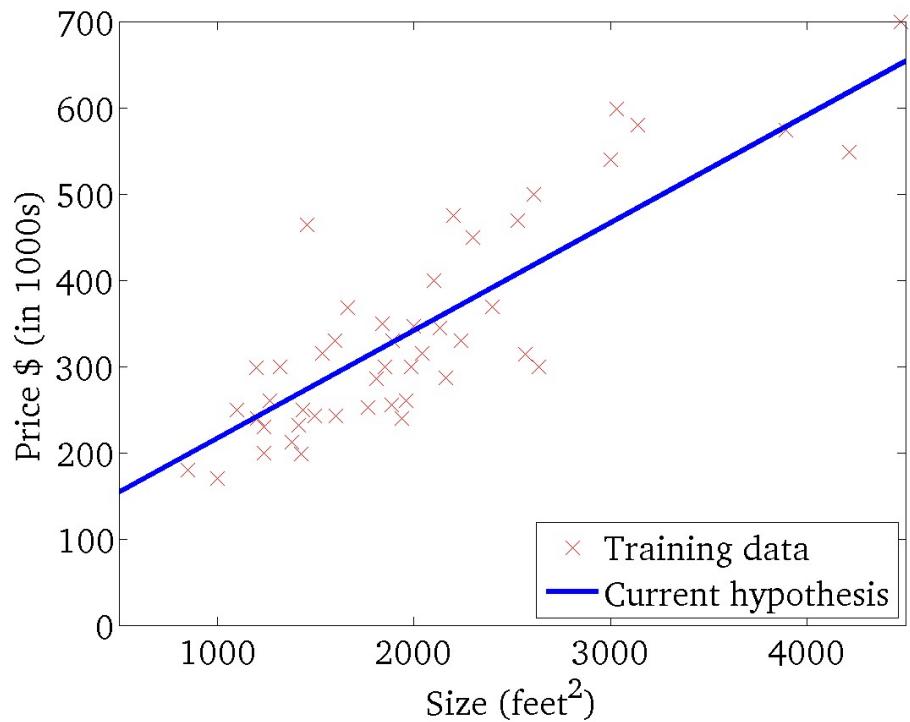
$$J(\theta_0, \theta_1)$$

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



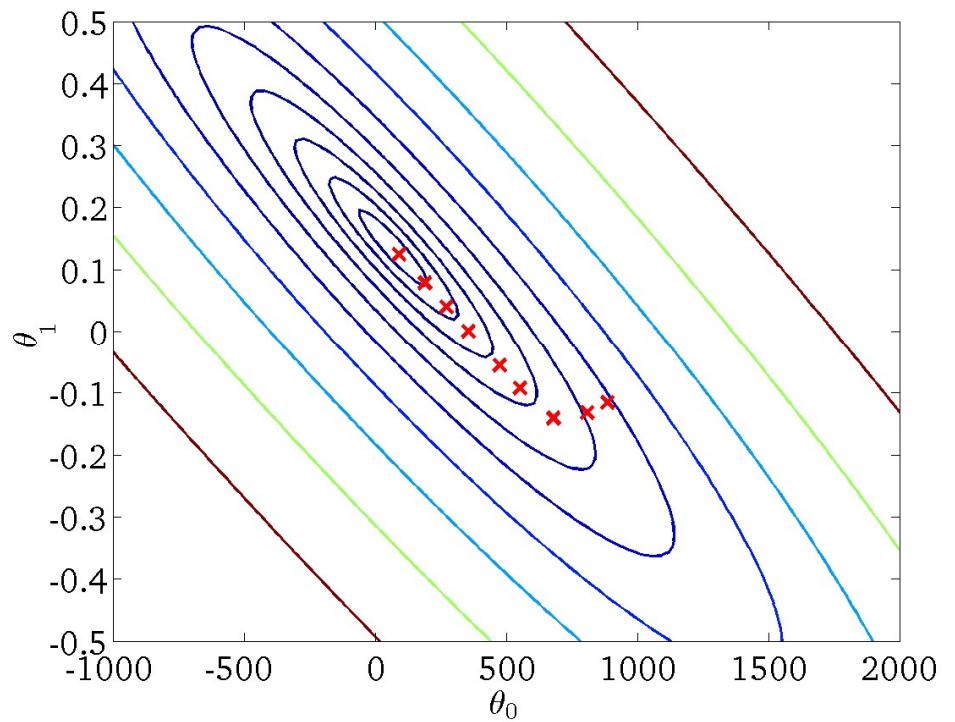
$$h_{\theta}(x)$$

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



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

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



## “Batch” Gradient Descent

“Batch”: Each step of gradient descent uses all the training examples.

## Stochastic Gradient Descent (SGD)

It updates the parameters for each training example, one by one.

## Stochastic Gradient Descent (SGD)

- Advantages:**
- faster than Batch GD in some problem.
  - the frequent updates allow us to have a pretty detailed rate of improvement.
- Disadvantages:**
- the frequent updates are more computationally expensive
  - The frequency of those updates can also result in noisy gradients

## Mini Batch Gradient Descent

A combination of the concepts of SGD and Batch Gradient Descent.

## Mini Batch Gradient Descent

- Splits the training dataset into small batches
  - Performs an update for each of these batches.
- Creates a balance between the robustness of stochastic gradient descent and the efficiency of batch gradient descent.

Mini-batch sizes range between 50 and 256.