

# Computer Vision (CS 419/619)

## Linear Regression and Gradient Descent

Anup Kumar Gupta

Under supervision of  
Dr. Puneet Gupta

# Linear Regression

## Supervised Learning

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

## Regression Problem

- Predict real-valued output

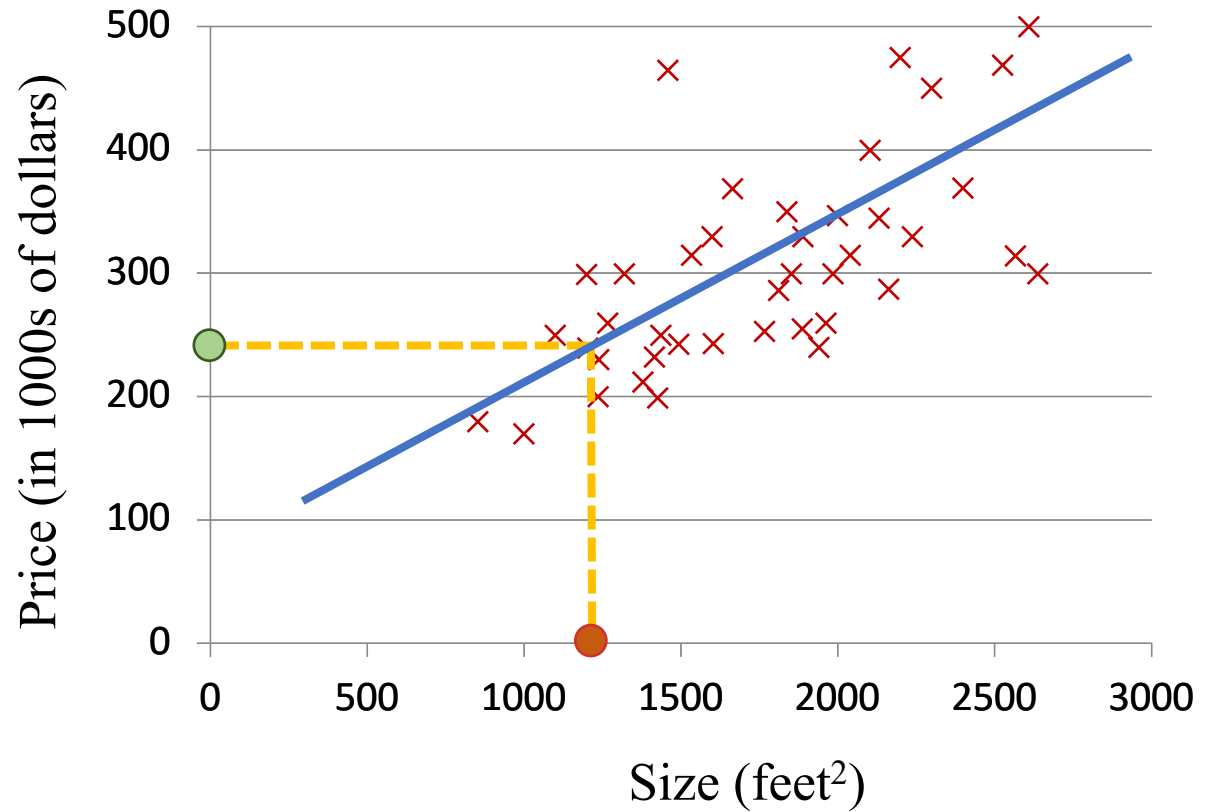


Fig. 1: Housing prices vs Size

# Training set

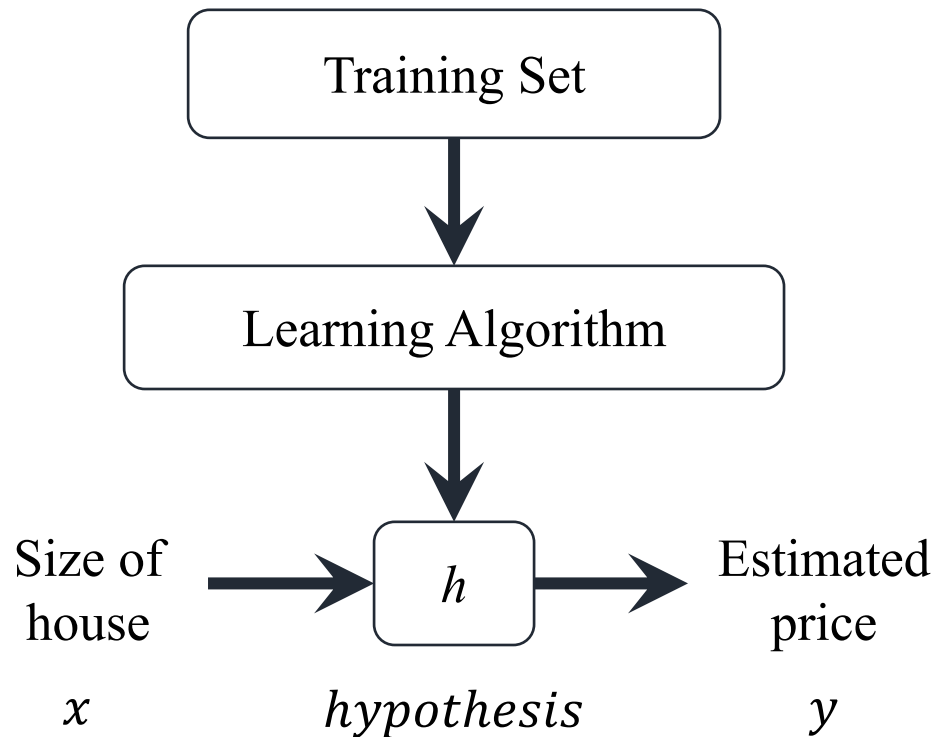
## Notations

- $m$  = number of training examples
- $x$  = “input” variable / features
- $y$  = “output” variable / “target” variable
- We will use  $(x, y)$  to denote one training example.
- We will use  $(x^{(i)}, y^{(i)})$  to denote the  $i^{\text{th}}$  training example.
- For example:  $x^{(1)} = 2104$  and  $y^{(1)} = 460$
- Note: superscript denotes index, and not power.

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

Tab. 1: Training Samples

# Training Procedure



$h$  maps the input  $x$  to the output  $y$

## How do we represent $h$ ?

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

As a shorthand we often remove the  $\theta$  notation:

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

$y$  is a linear function of  $x$ . It resembles:

$$y = c + m \cdot x$$

Equation of straight line.

- Linear Regression with one variable ( $x$ ).
- Univariate linear regression.

# Training procedure – How to compute the $\theta$ s

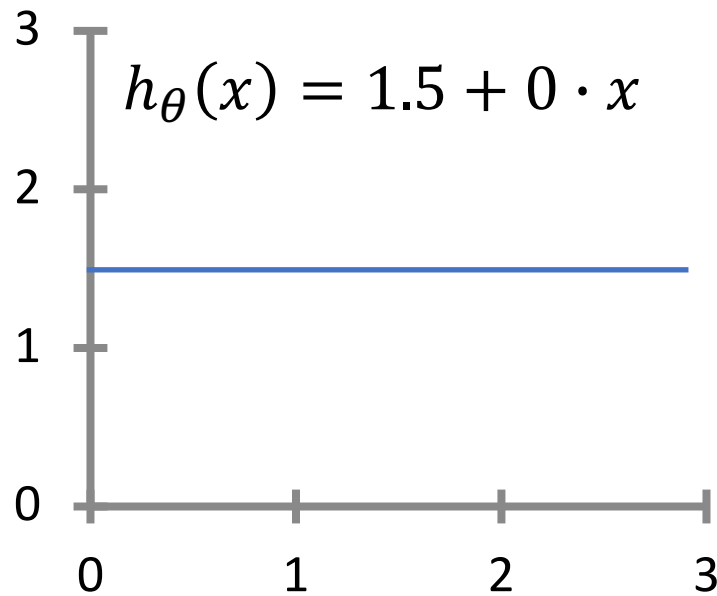
- $h_{\theta}(x) = \theta_0 + \theta_1 \cdot x$
- $\theta_1$  and  $\theta_2$  or in general  $\theta_i$  = parameters
- How to choose or find these  $\theta_i$ s ?
- If we have these  $\theta_i$ s and a new data point  $x$  is given to us,
- then we can calculate  $y$  using these parameters.

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

Tab. 1: Training Samples

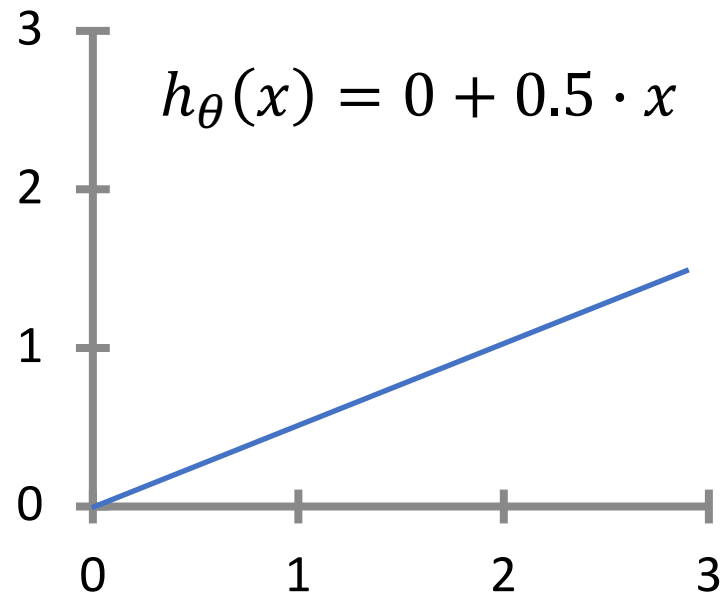
# Variations of $h_\theta$ with variations in $\theta_i$

$$h_\theta(x) = \theta_0 + \theta_1 \cdot x$$



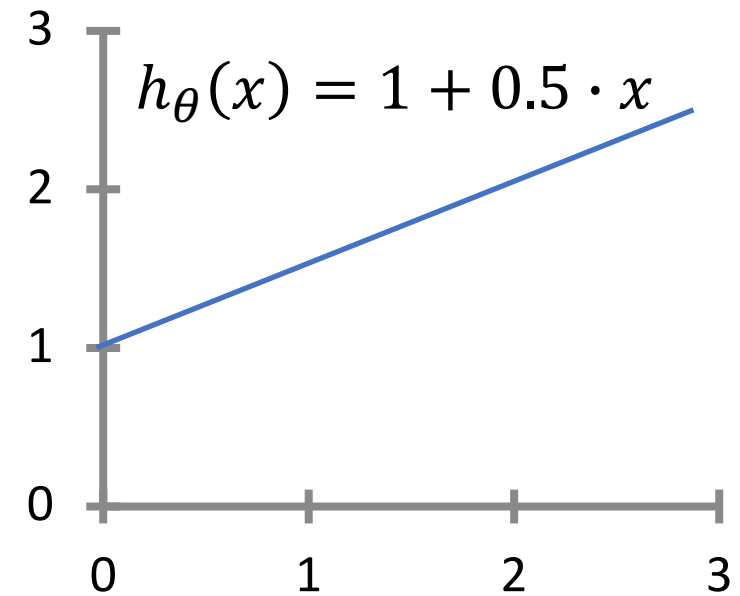
$$\theta_0 = 1.5$$

$$\theta_1 = 0$$



$$\theta_0 = 0$$

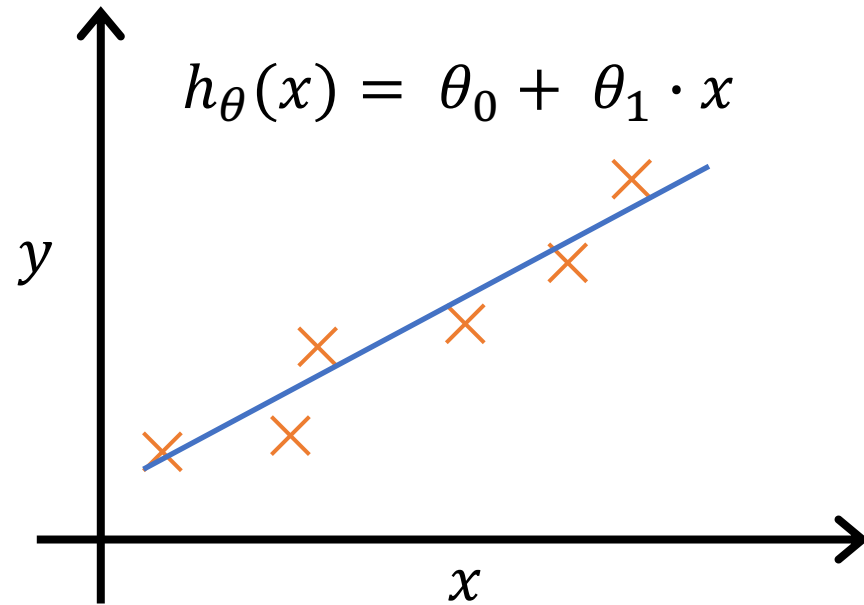
$$\theta_1 = 0.5$$



$$\theta_0 = 1$$

$$\theta_1 = 0.5$$

# Deciding and Choosing the parameters $\theta_i$ s



The idea is to choose  $\theta_0$  and  $\theta_1$  such that  $h_{\theta}(x)$  is close to  $y$  for given training sample pairs  $(x, y)$

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

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

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

$J(\theta_0, \theta_1)$  is called as **cost function**.

It is also called as **squared error function**.

# Hypothesis function (Simplified)

- Hypothesis

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

- 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

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

- Simplified Hypothesis ( $\theta_0 = 0$ )

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

- Parameter

$$\theta_1$$

- Cost function

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

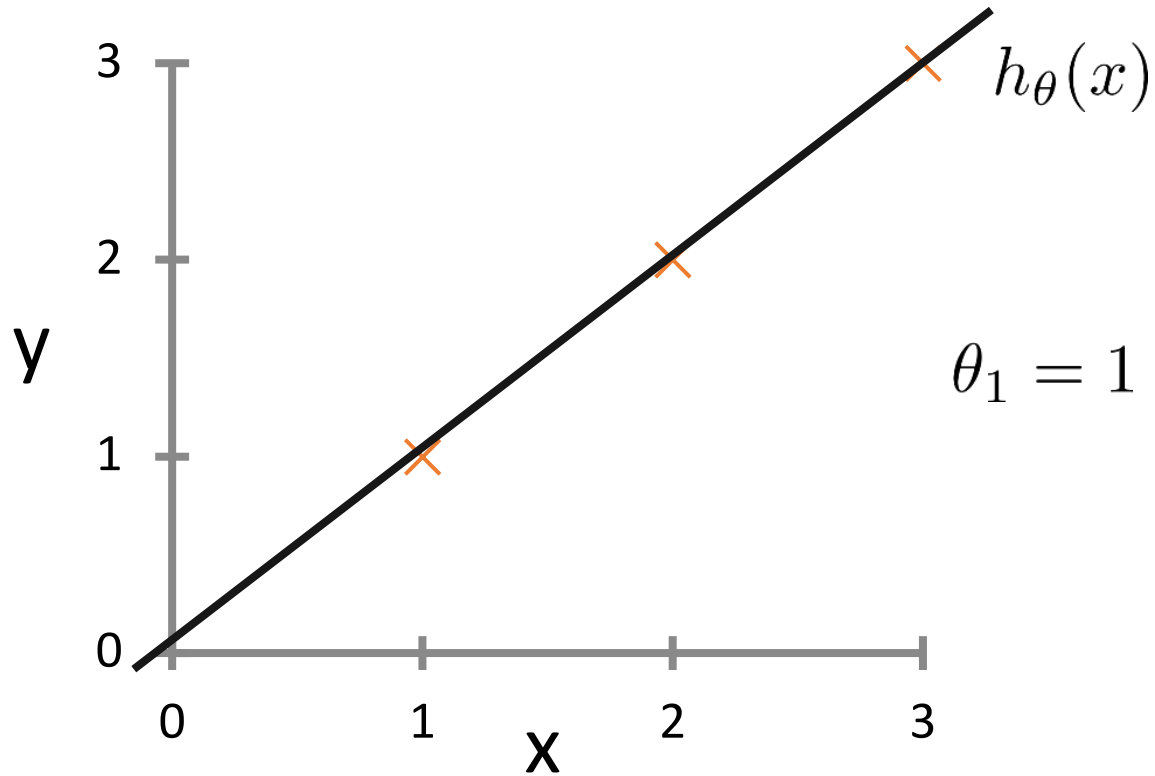
- Goal

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



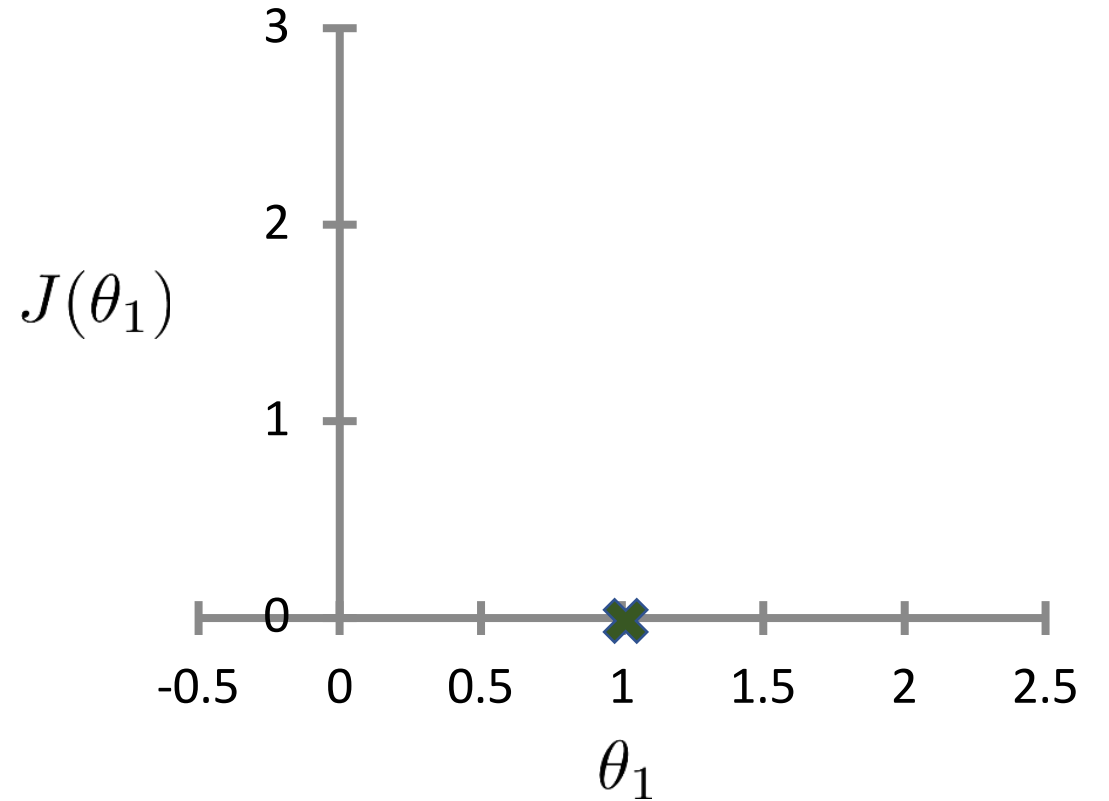
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )



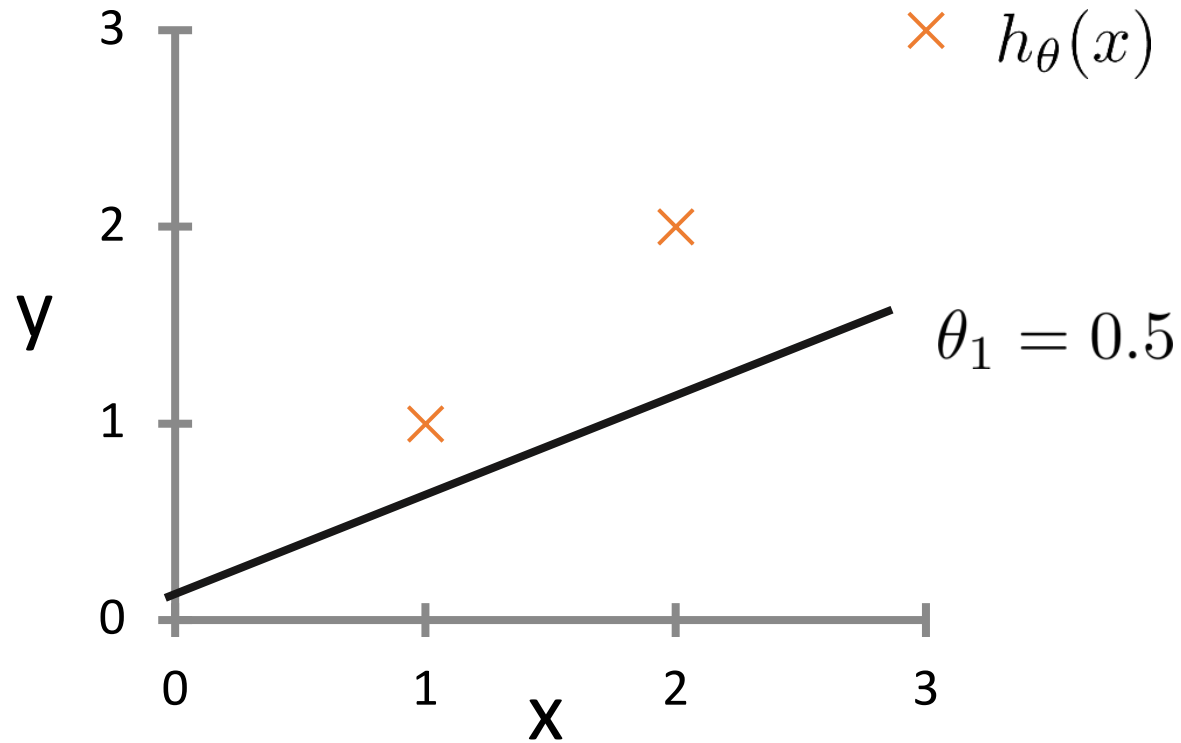
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 \cdot x^{(i)} - y^{(i)})^2$$

$$J(\theta_1 = 1) = \frac{1}{2 * 3} \{(1 \times 1 - 1)^2 + (1 \times 2 - 2)^2 + (1 \times 3 - 3)^2\}$$

$$J(\theta_1 = 1) = 0$$

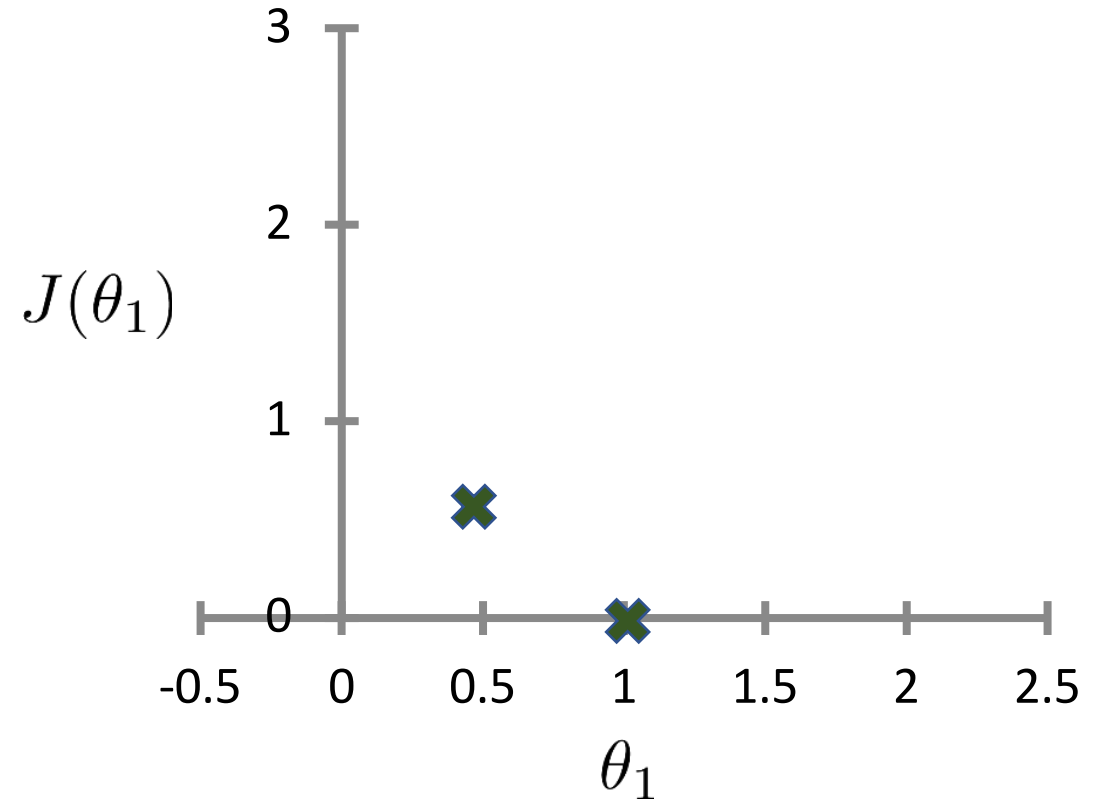
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )

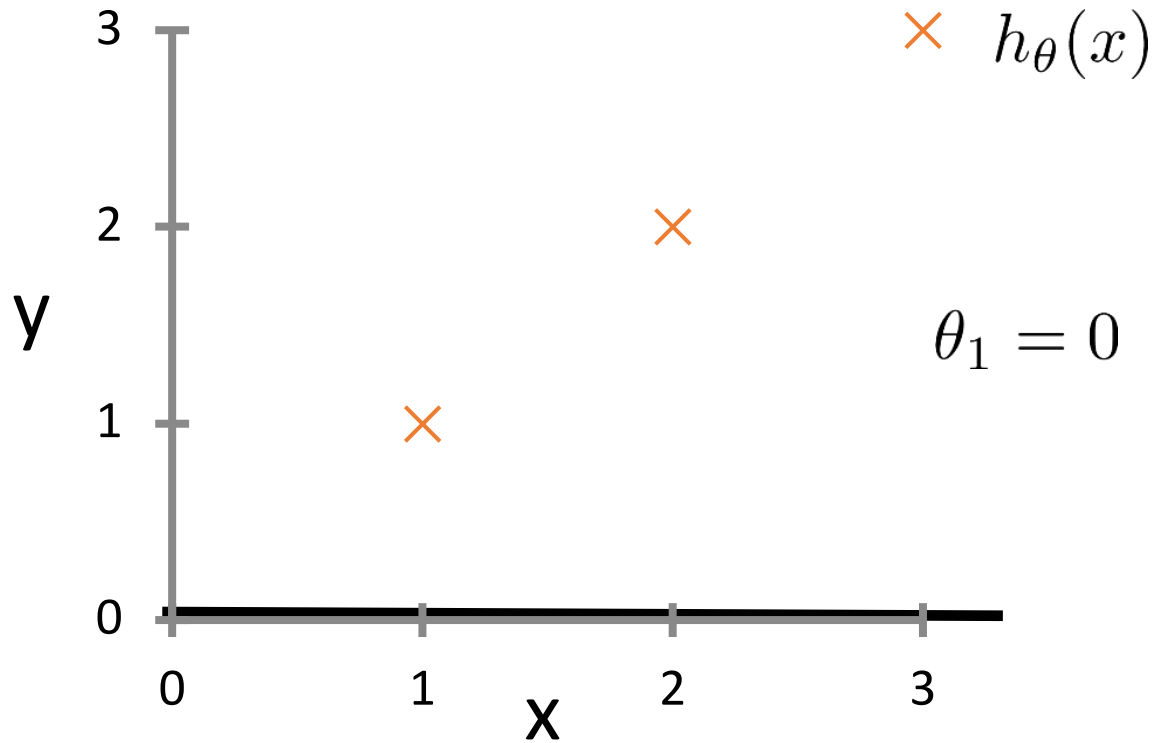


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 \cdot x^{(i)} - y^{(i)})^2$$

$$J(\theta_1 = 0.5) = \frac{1}{2 * 3} \{ (0.5 \times 1 - 1)^2 + (0.5 \times 2 - 2)^2 + (0.5 \times 3 - 3)^2 \} \quad J(\theta_1 = 0.5) = 3.5/6 = 0.58$$

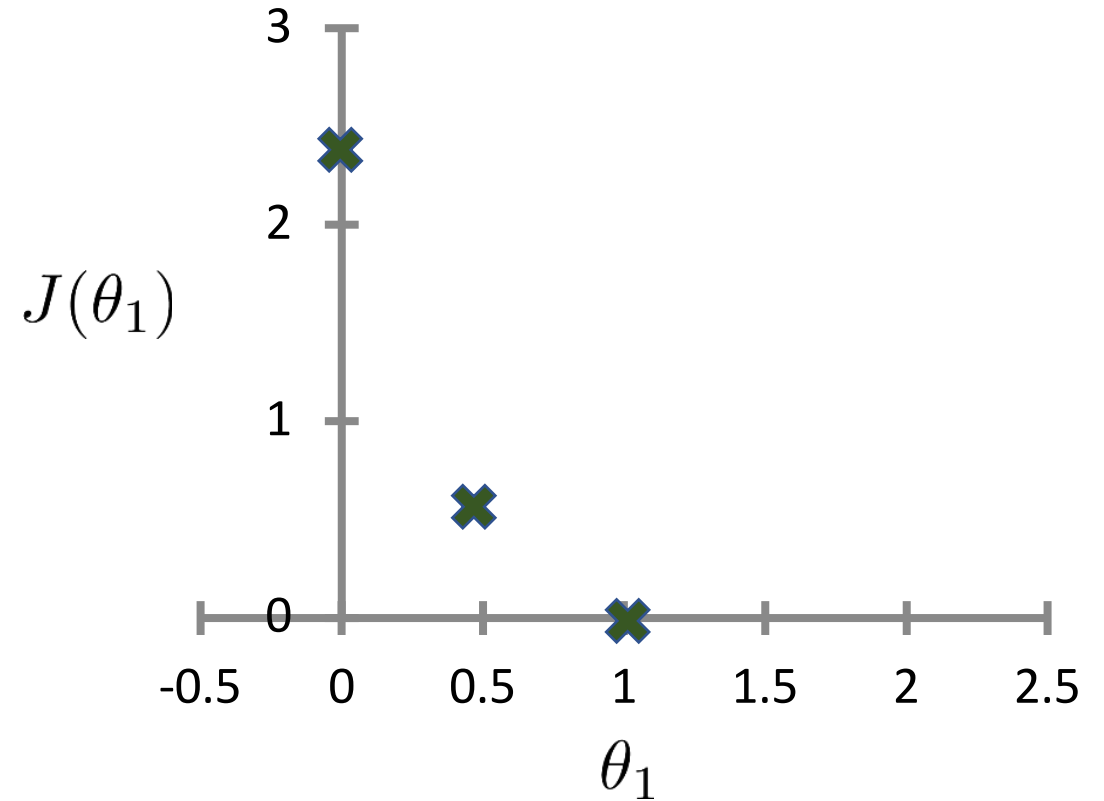
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )

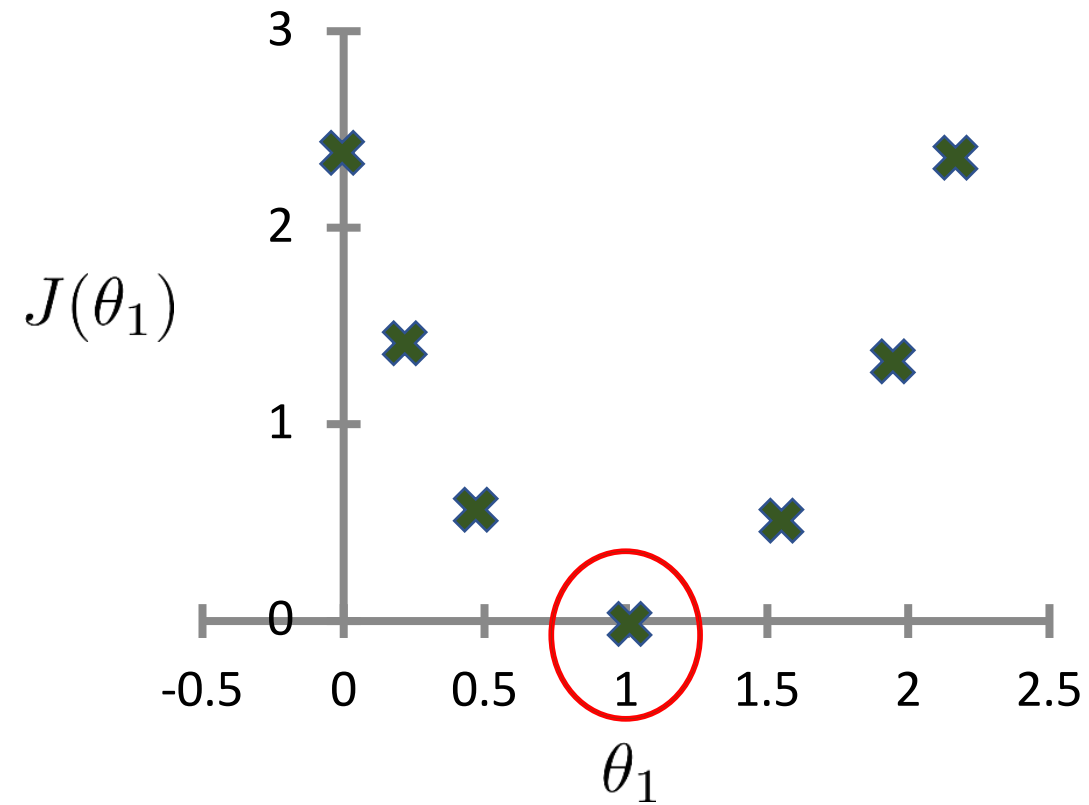


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (\theta_1 \cdot x^{(i)} - y^{(i)})^2$$

$$J(\theta_1 = 0) = \frac{1}{2 * 3} \{(0 - 1)^2 + (0 - 2)^2 + (0 - 3)^2\}$$

$$J(\theta_1 = 0) = 14/6 = 2.3$$

Which one to choose?



The value of  $\theta_1$  which gives minimum  $J(\theta_1)$

# Hypothesis function (Summary)

- Hypothesis

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

- Parameters

$$\theta_0, \theta_1$$

- Cost function

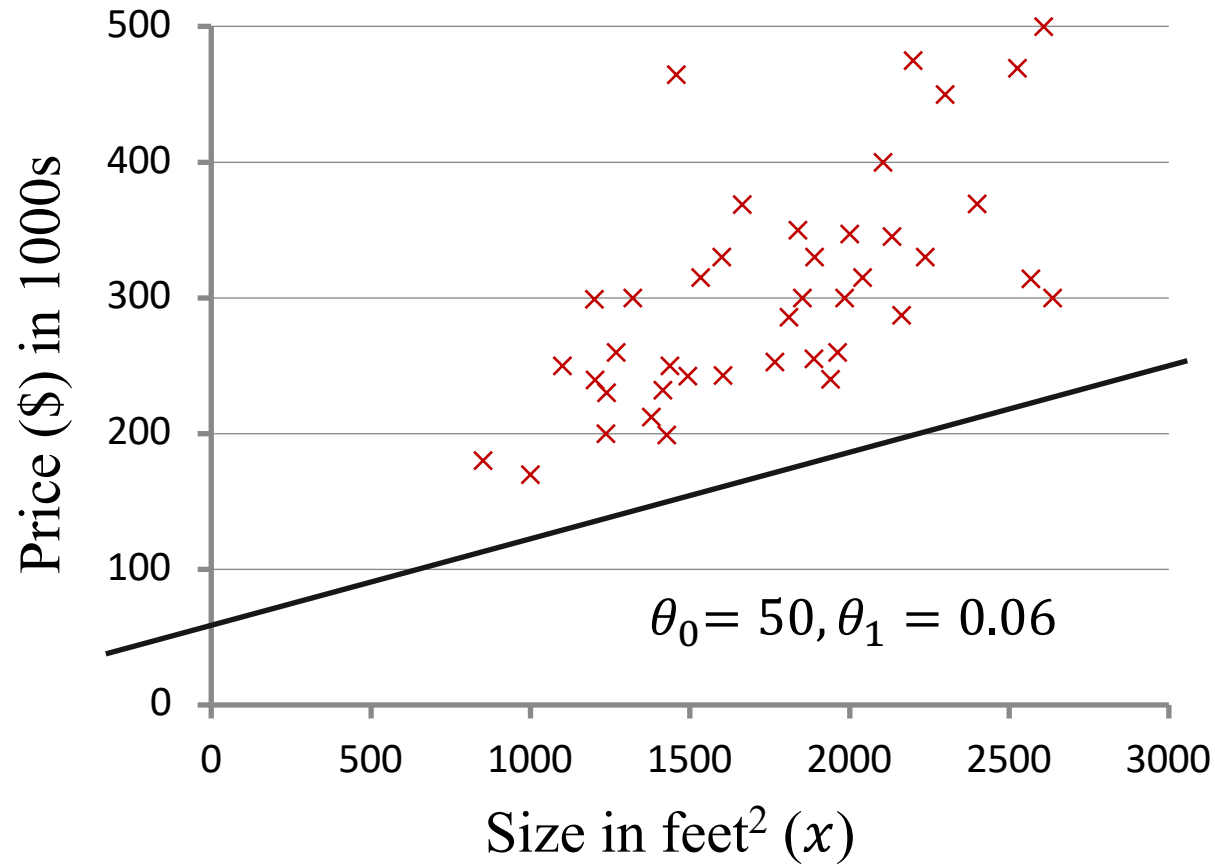
$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

- Goal

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

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

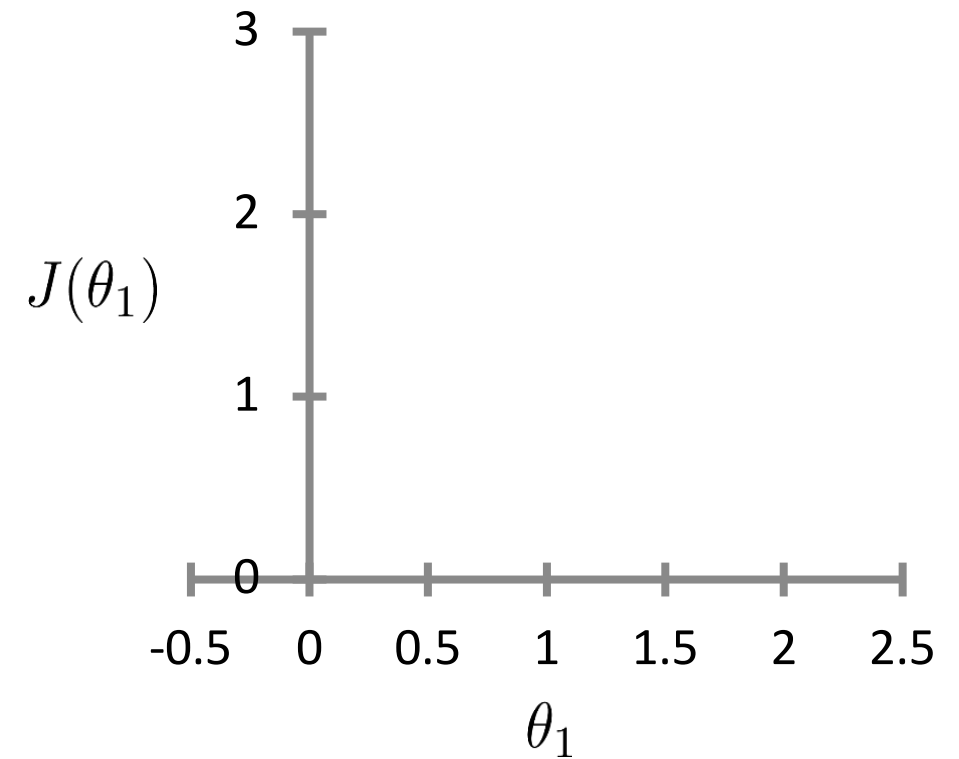
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

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

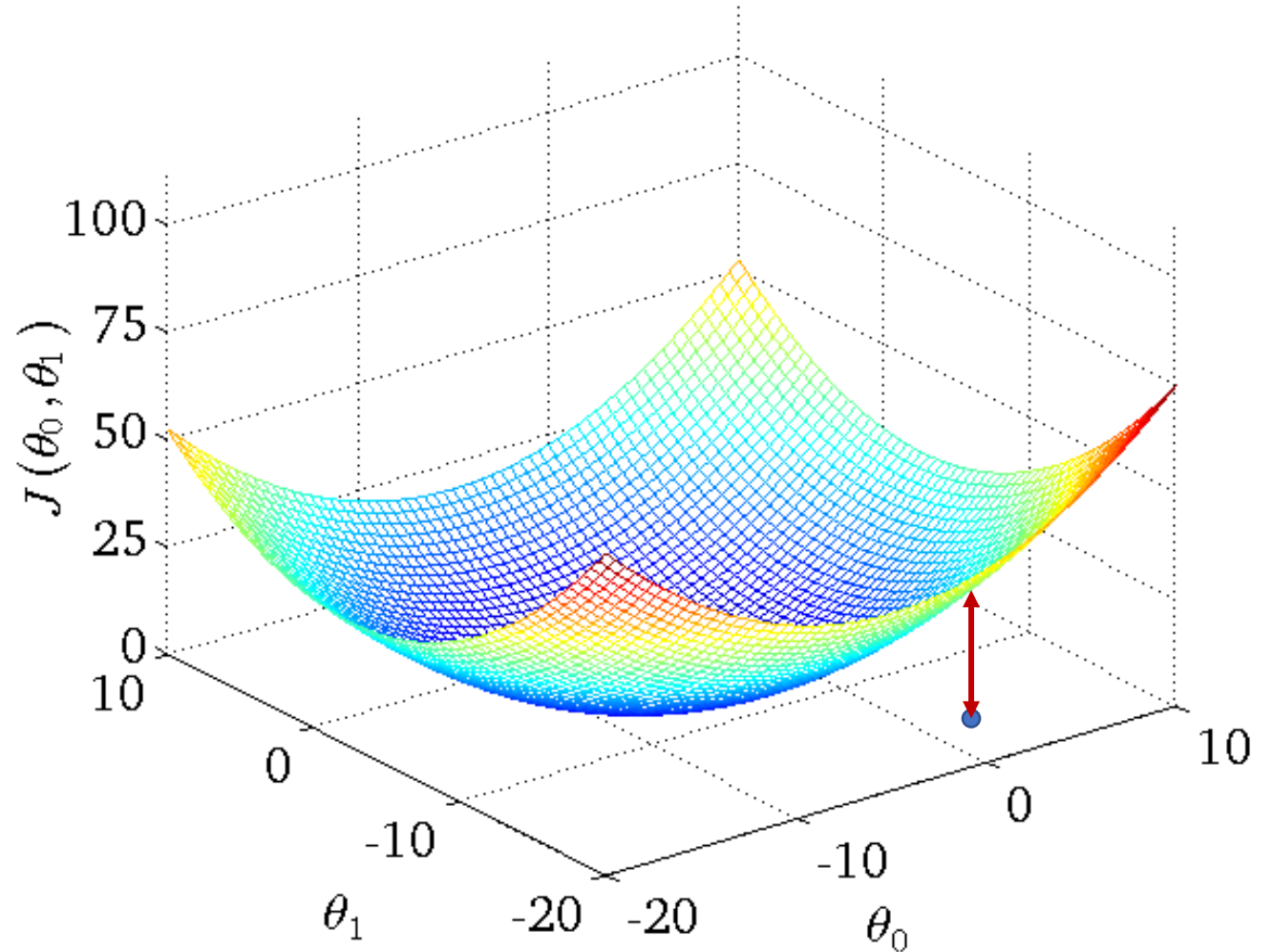
(function of the parameters  $\theta_0, \theta_1$ )



- But now we have two parameters  $\theta_0, \theta_1$ .
- Hence, we need to add an extra axis to the plot.

# Loss Function in 3-D (Example)

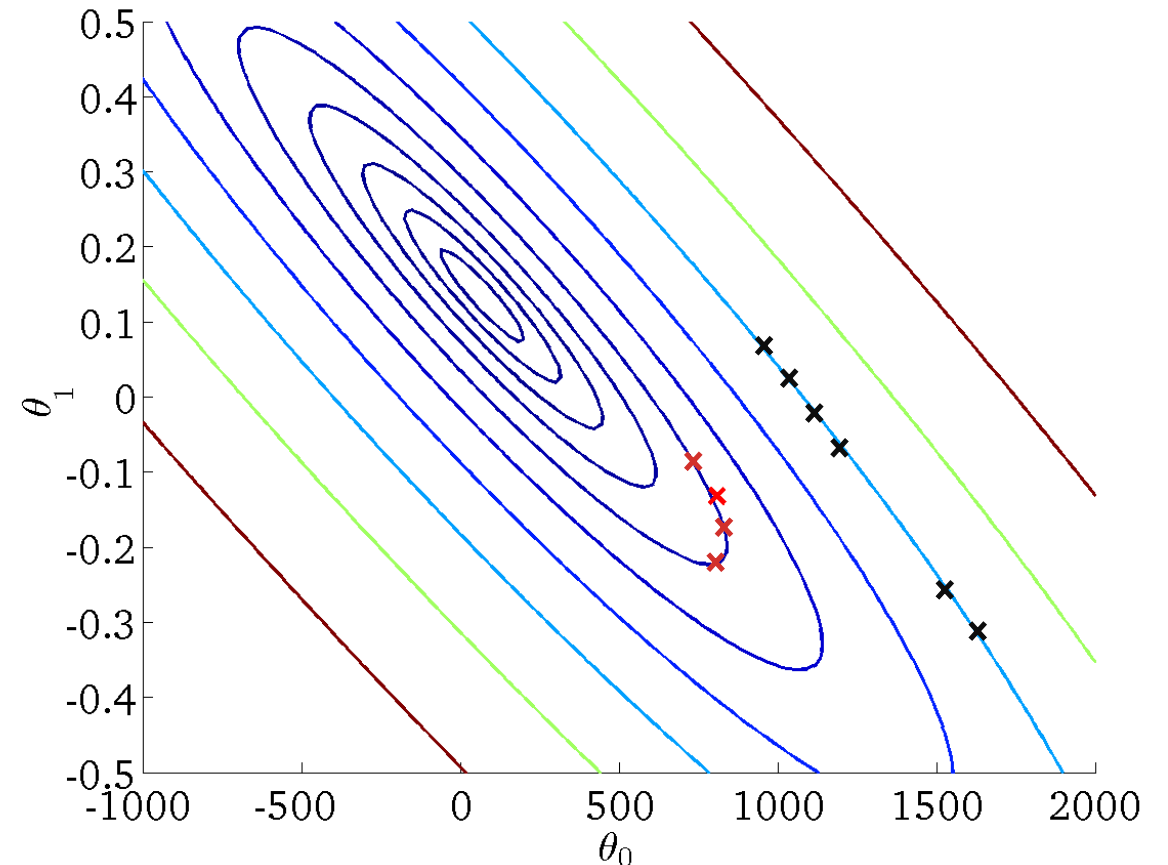
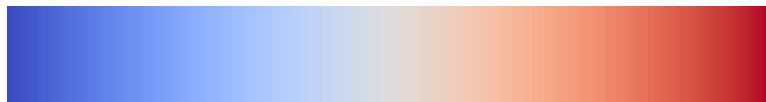
- This is a 3D surface plot
- The axes denotes:
  - both the parameters ( $\theta_0, \theta_1$ )
  - the loss function  $J(\theta_0, \theta_1)$
- As we vary the values of  $\theta_0$  and  $\theta_1$  we get different values of the loss function  $J(\theta_0, \theta_1)$
- *[Informal]* The height of the surface from the  $(\theta_0, \theta_1)$  denotes the value of loss at that point.



Demo: <https://academo.org>

# Contour Plots

- This is a contour plot of a 3D surface.
- The axes denotes the parameters  $\theta_0$  and  $\theta_1$ .
- An ellipse (or oval) is a set of points that takes on the same value.
- *[Generally]* A color towards the shade of blue (cold color) denotes a lower value, whereas a color towards the red shade (warm color) denotes a higher value.
- For example:

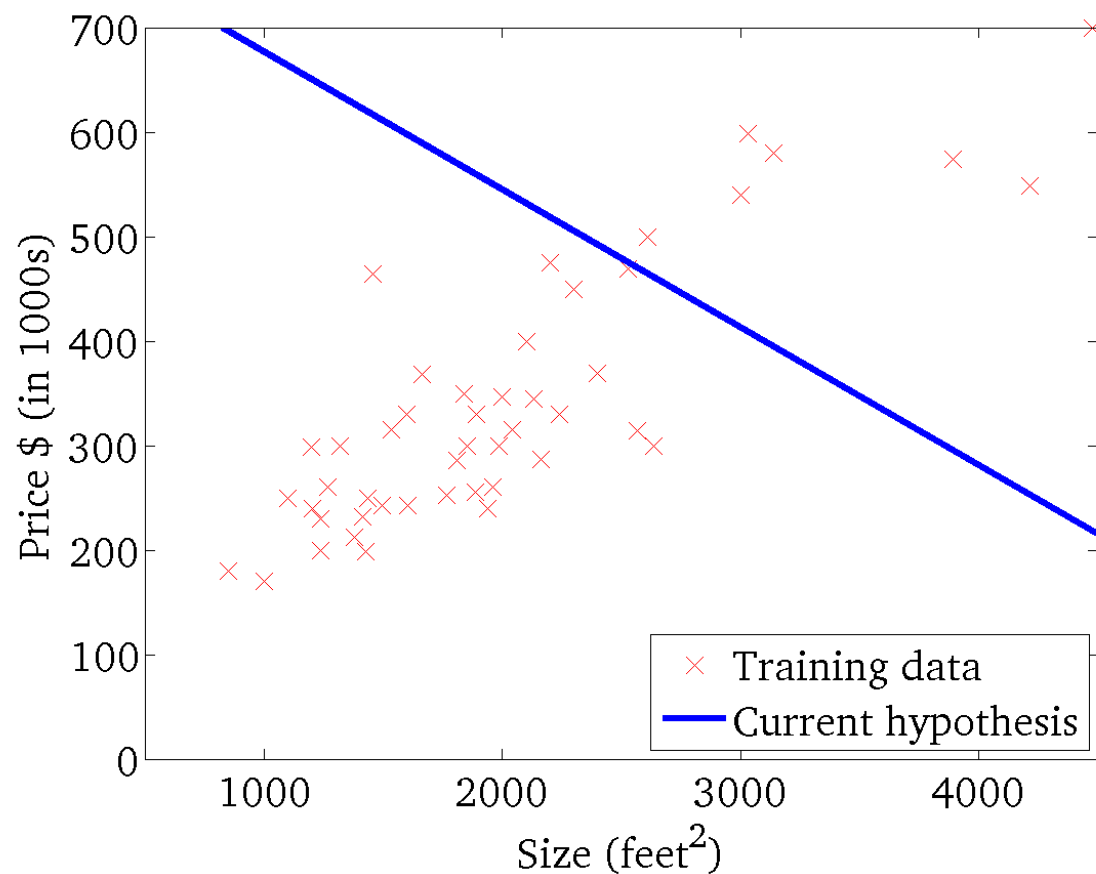


- These all points will have same value as they lie on the same ellipse



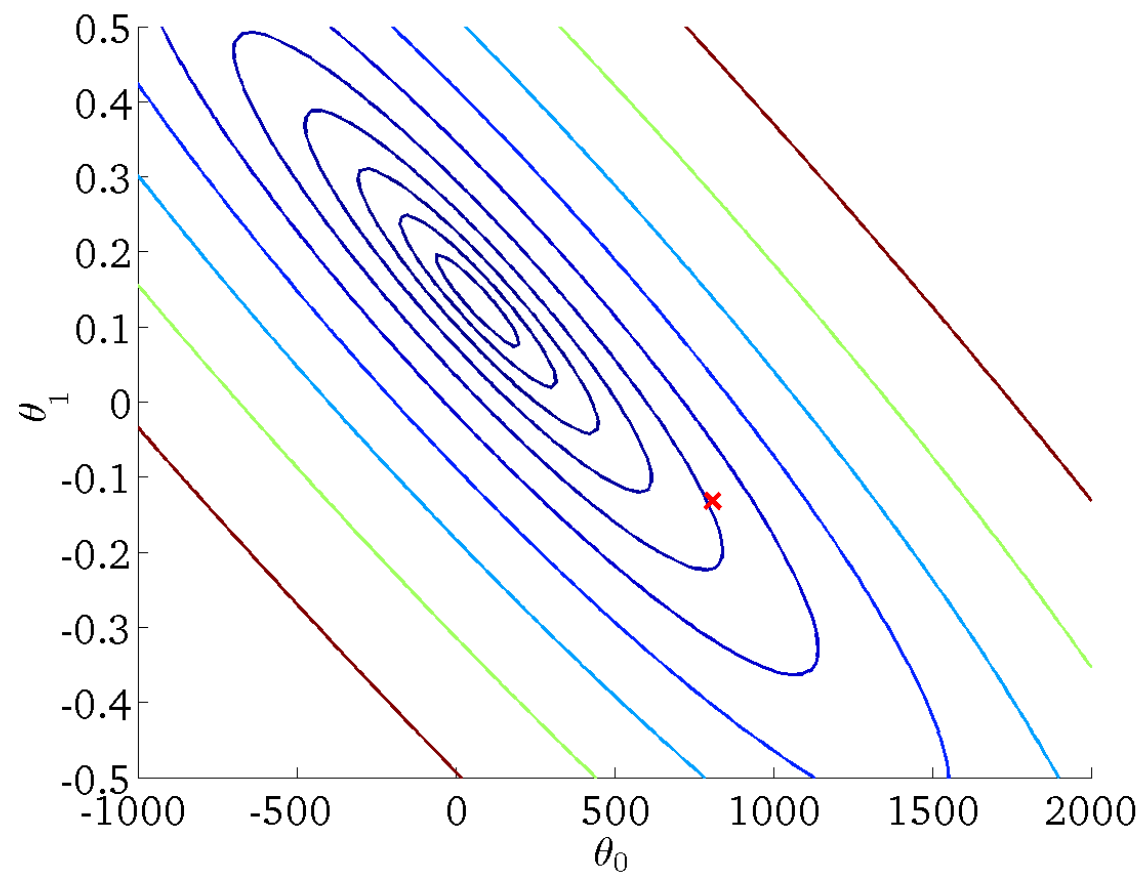
$$h_{\theta}(x)$$

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



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

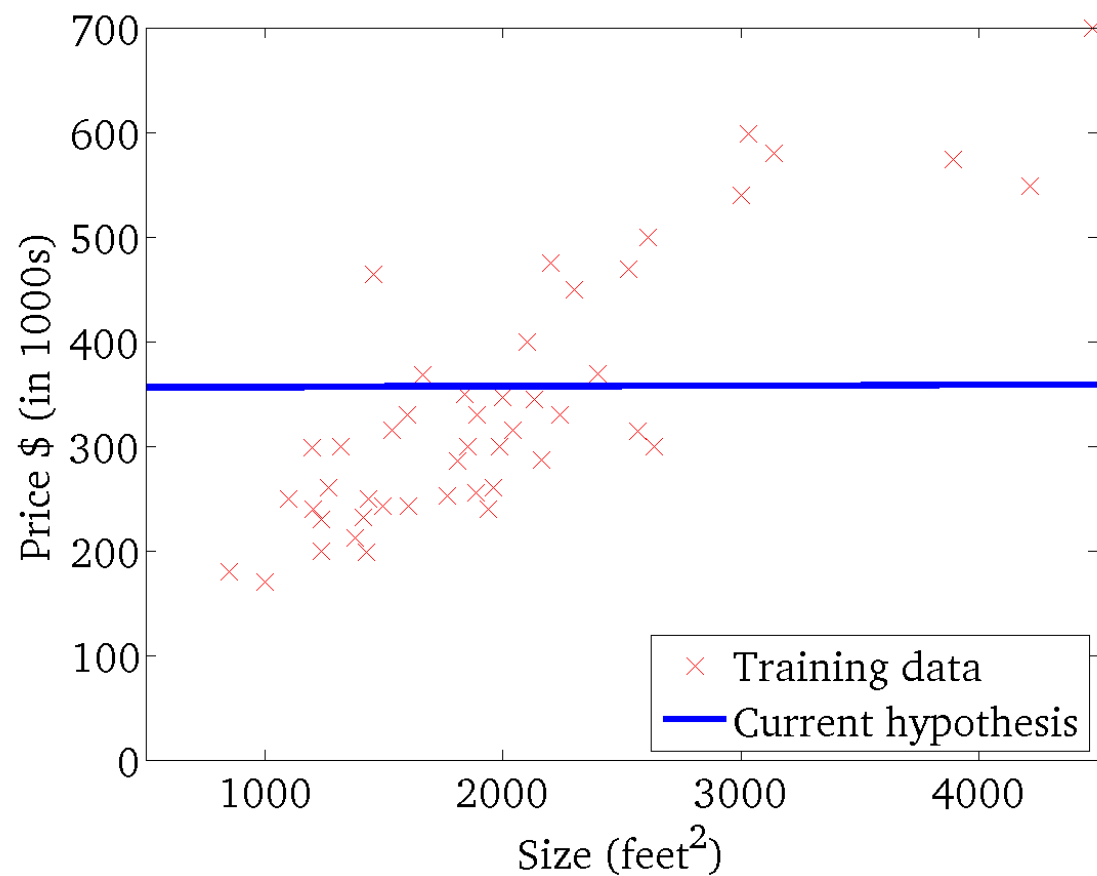
(function of the parameters  $\theta_0, \theta_1$ )



$$\theta_0 \approx 800 \text{ and } \theta_1 \approx -0.15$$

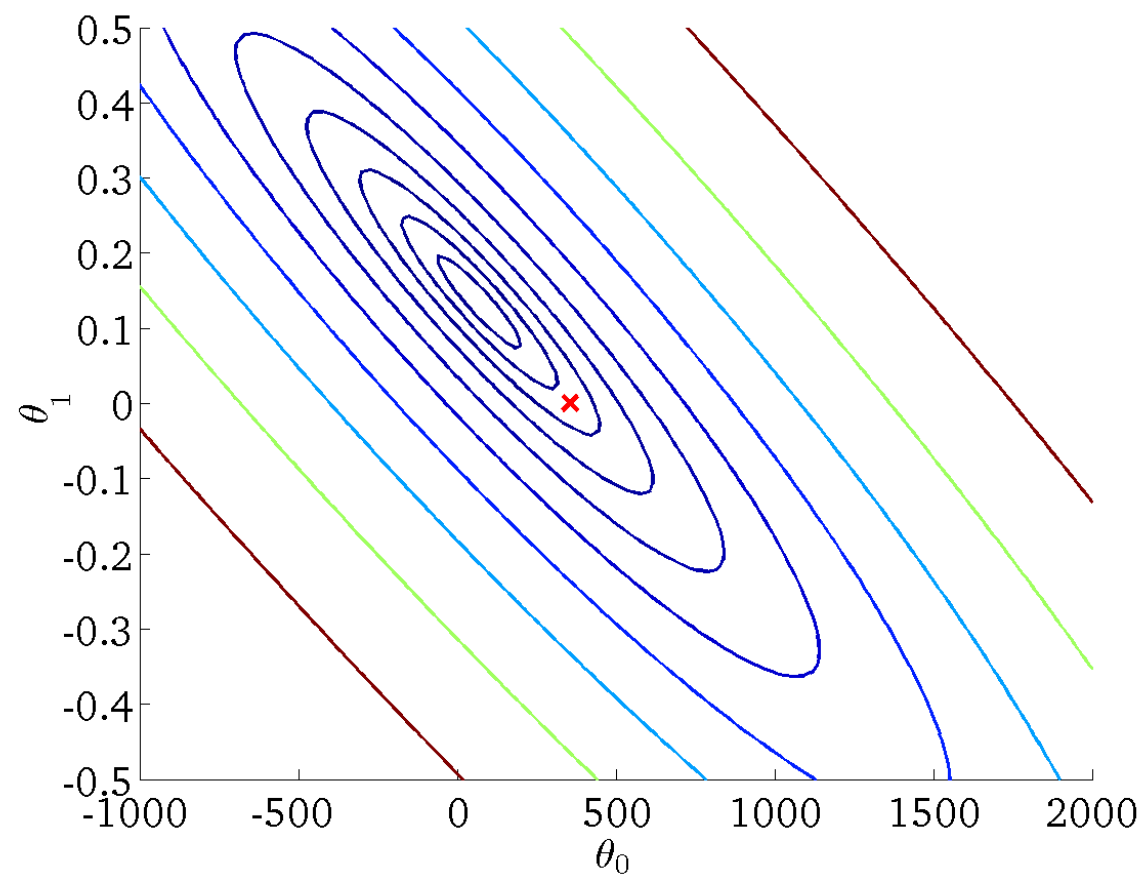
$$h_{\theta}(x)$$

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



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

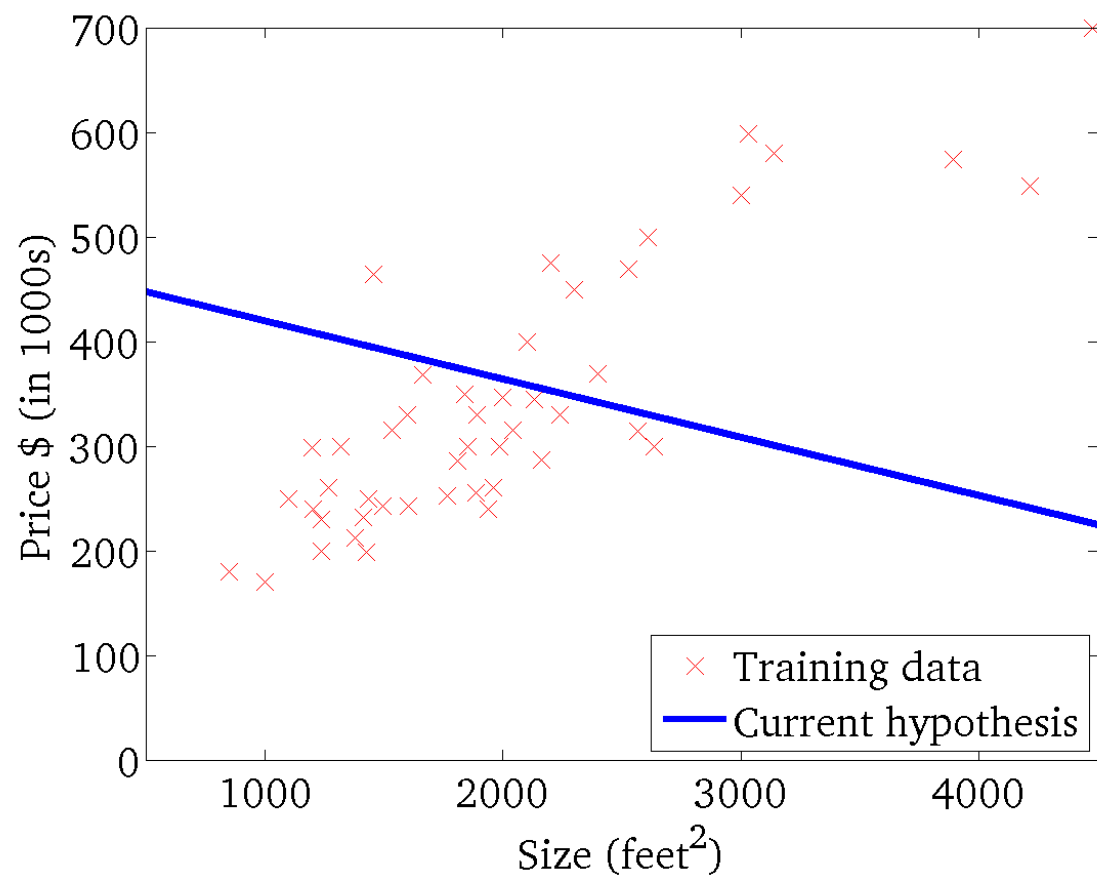
(function of the parameters  $\theta_0, \theta_1$ )



$$\theta_0 \approx 360 \text{ and } \theta_1 \approx 0$$

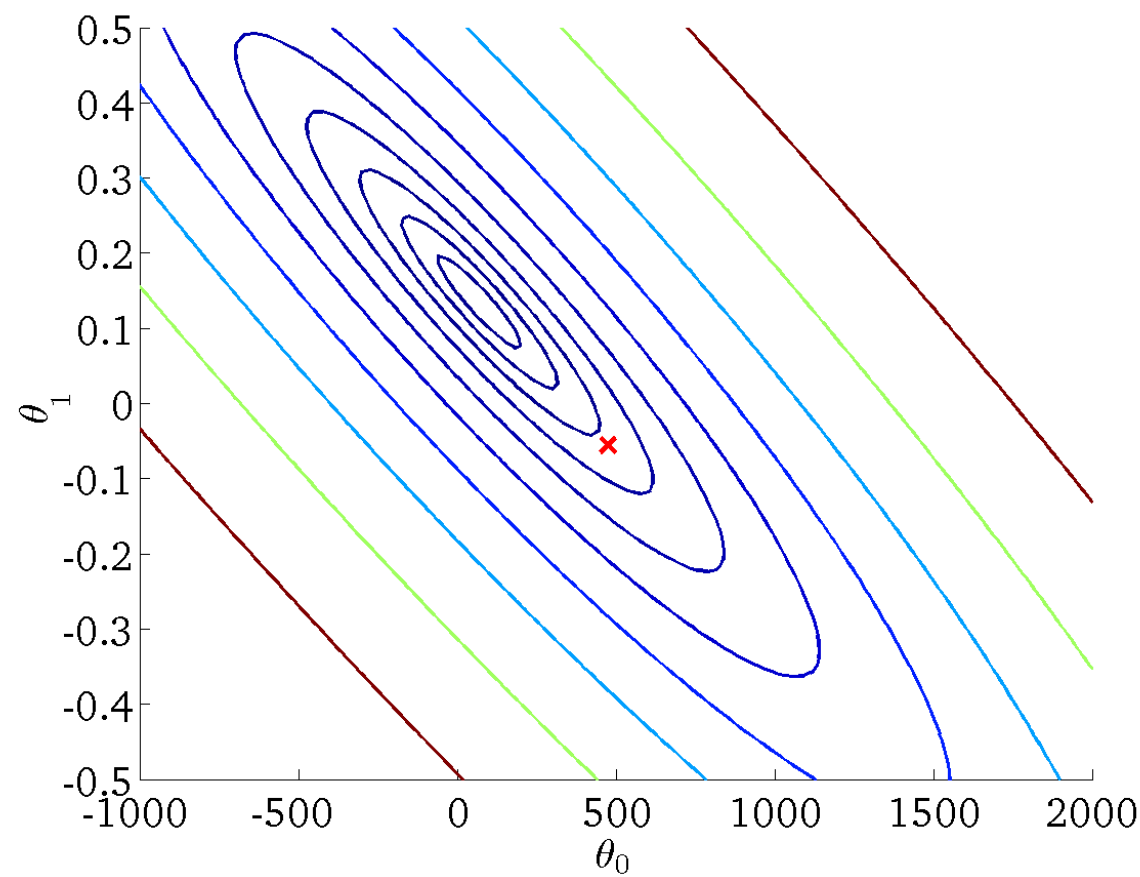
$$h_{\theta}(x)$$

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



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

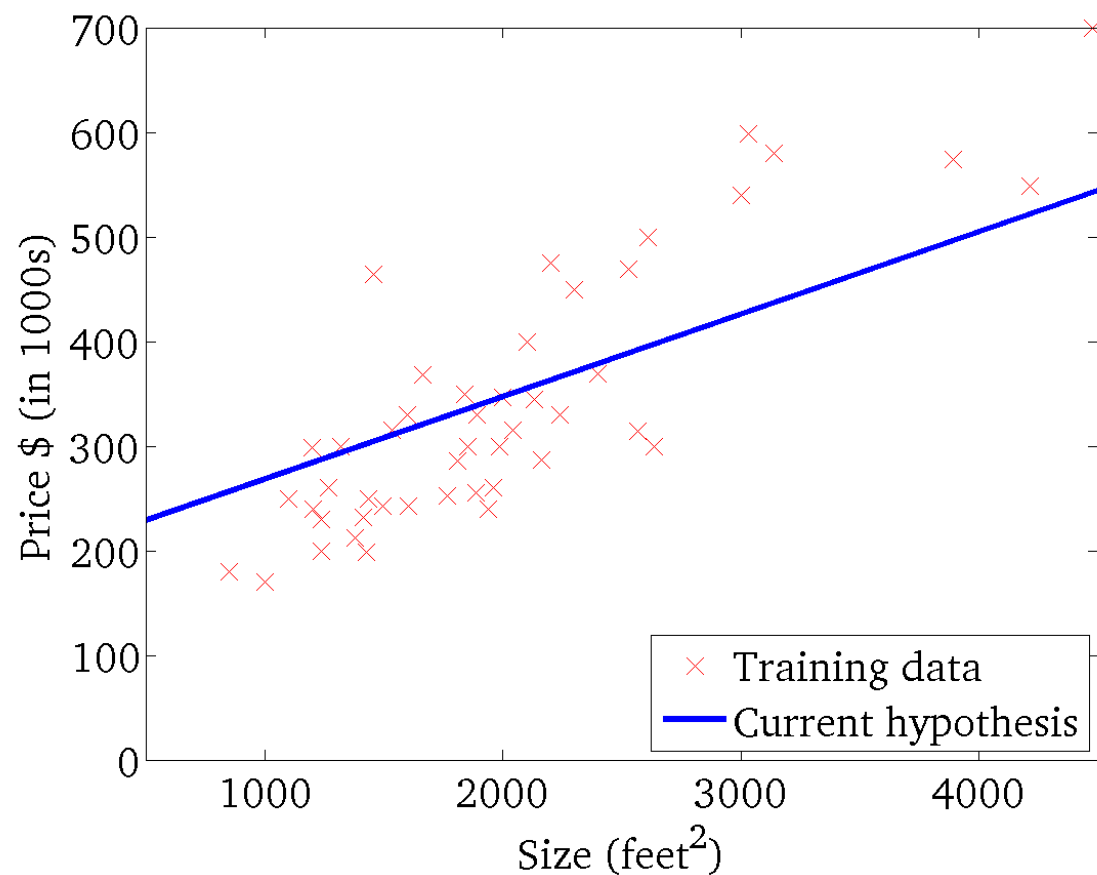
(function of the parameters  $\theta_0, \theta_1$ )



$$\theta_0 \approx 510 \text{ and } \theta_1 \approx -0.02$$

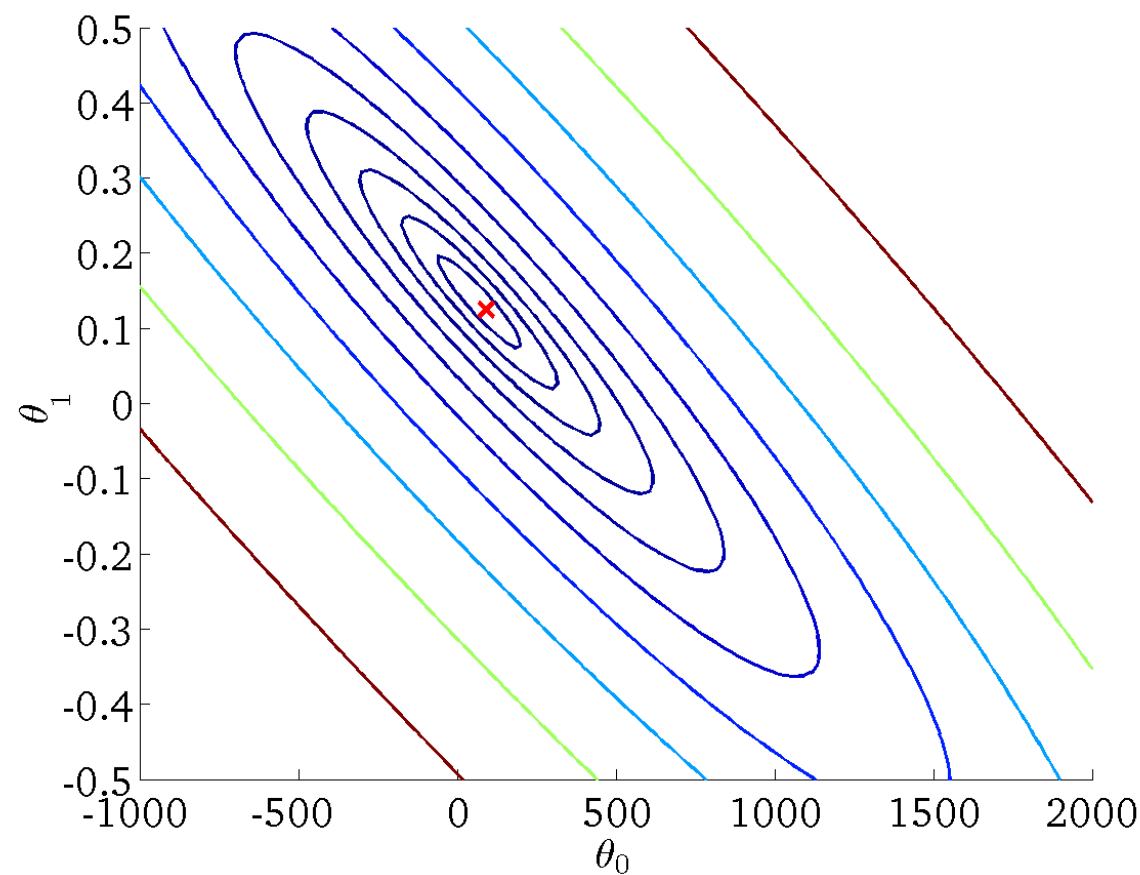
$$h_{\theta}(x)$$

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



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

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



$$\theta_0 \approx 480 \text{ and } \theta_1 \approx 0.14$$

Not minimum, but close to minimum.

# Gradient Descent Algorithm

- Cost function

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

- Goal

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

## Algorithm

1. Start with some  $\theta_0, \theta_1$ .
2. Keep changing  $\theta_0, \theta_1$  to reduce  $J(\theta_0, \theta_1)$  until we hopefully end up at a minimum.

- Generalized Cost function

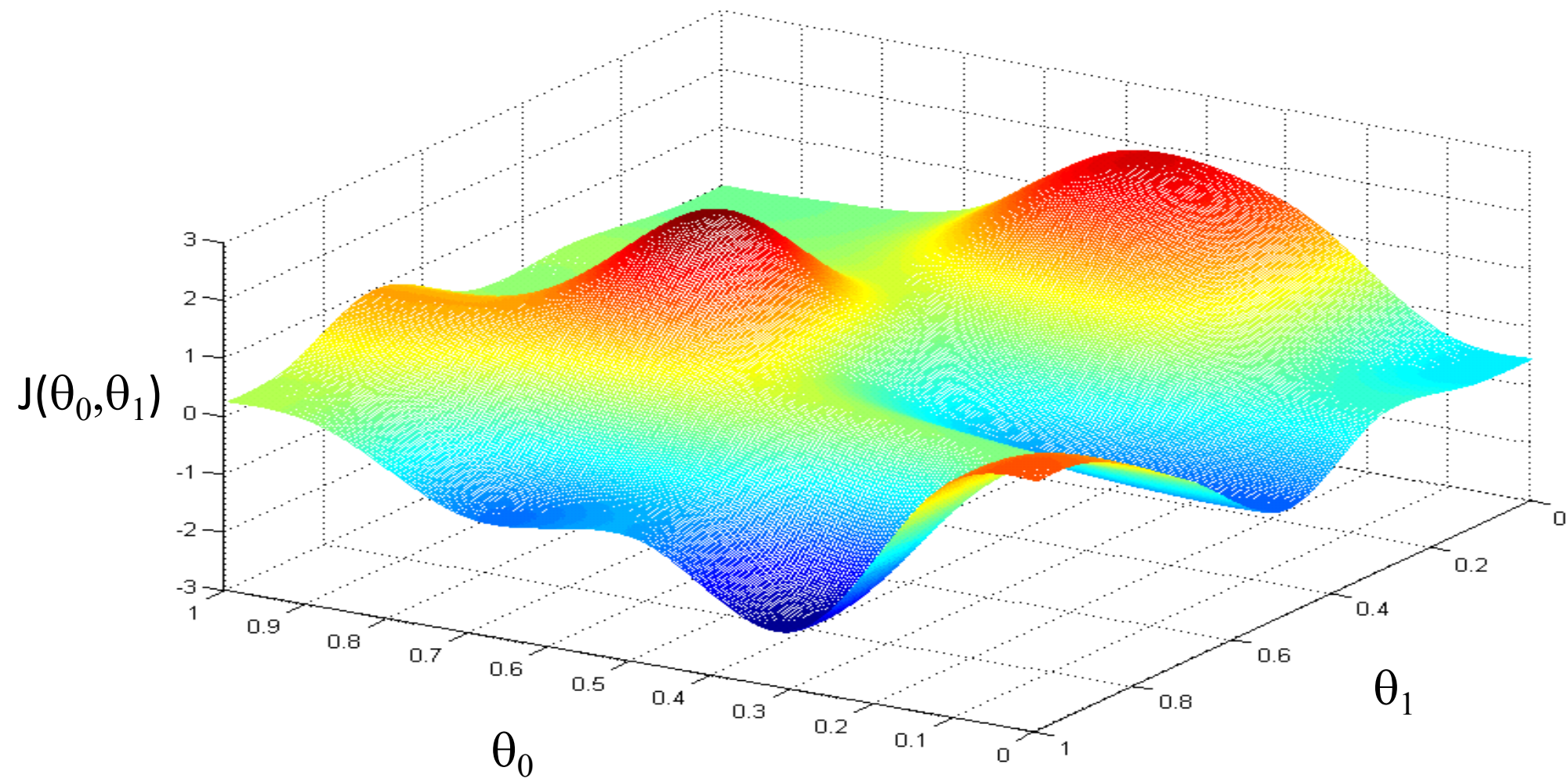
$$J(\theta_0, \theta_1, \theta_2 \dots \theta_n)$$

- Goal

$$\underset{\theta_0 \dots \theta_n}{\text{minimize}} J(\theta_0, \theta_1, \theta_2 \dots \theta_n)$$

## Algorithm

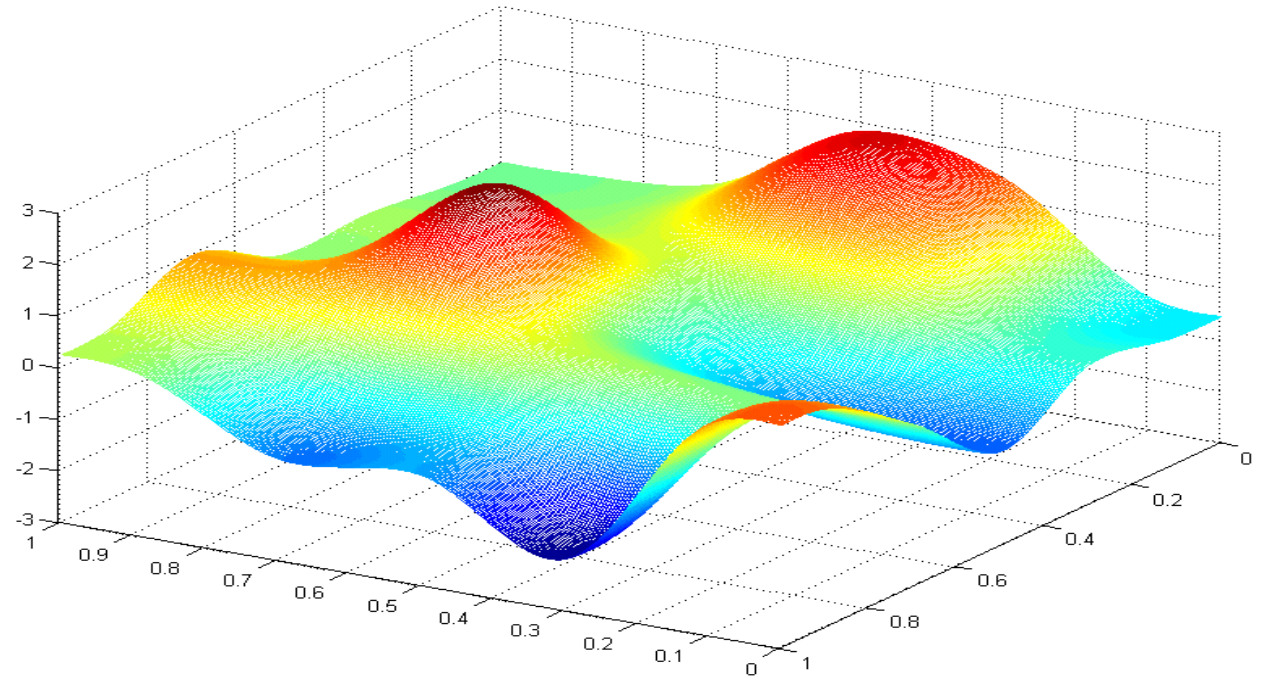
1. Start with some  $\theta_0, \theta_1, \theta_2 \dots \theta_n$ .
2. Keep changing  $\theta_0, \theta_1, \theta_2 \dots \theta_n$  to reduce  $J(\theta_0, \theta_1, \theta_2 \dots \theta_n)$  until we hopefully end up at a minimum.



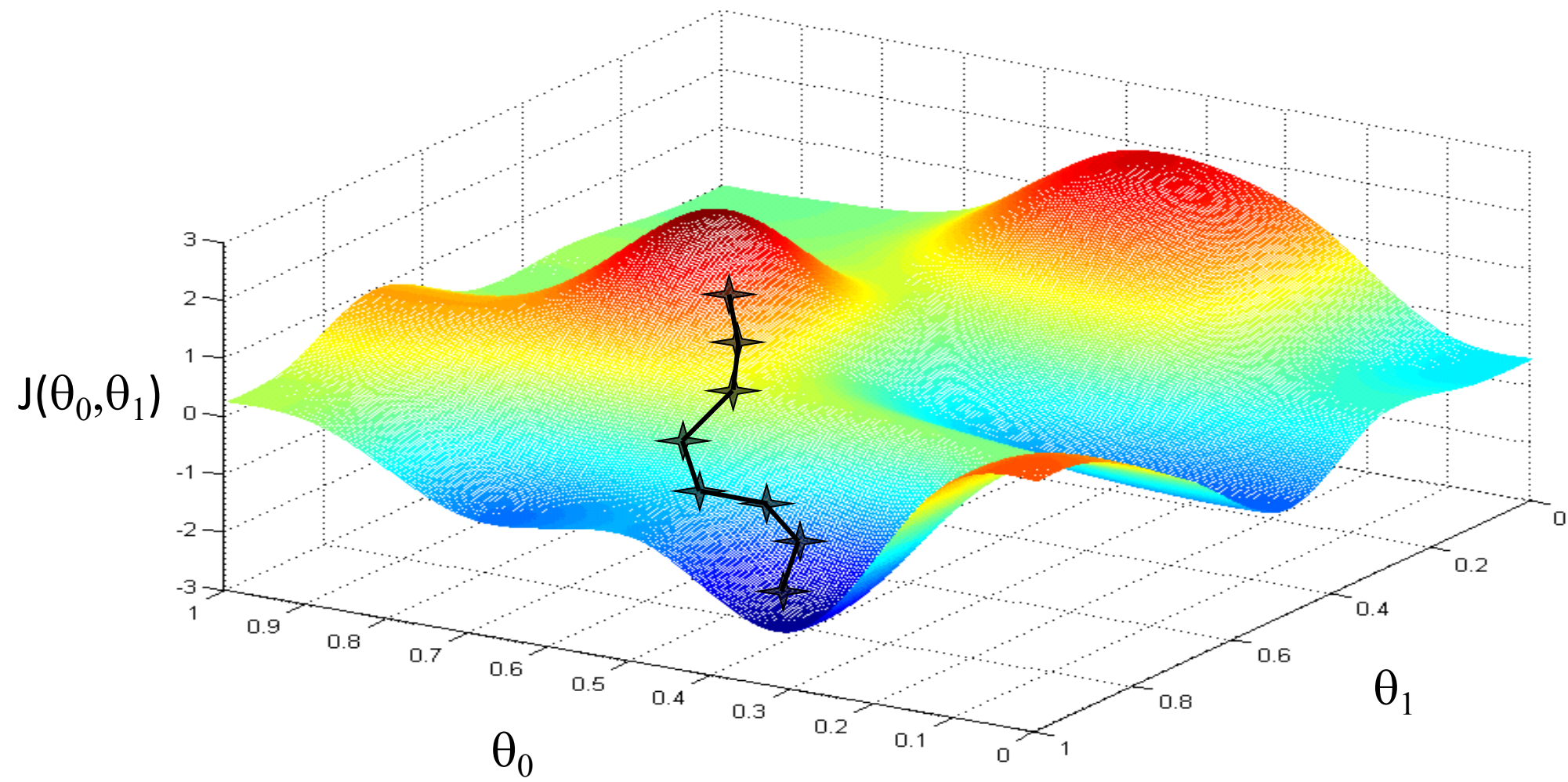
# Gradient Descent Algorithm – Intuitive Approach

- Assume we are in some grassy park, with two hill like structures.
- Assume that we are physically standing at that point on the hill, on this little red hill in the park.
- In gradient descent, we take a 360 degrees spin and ponder:

“If we were to take a little steps in some direction, so as we want to go downhill as quickly as possible, What direction should we choose” ?

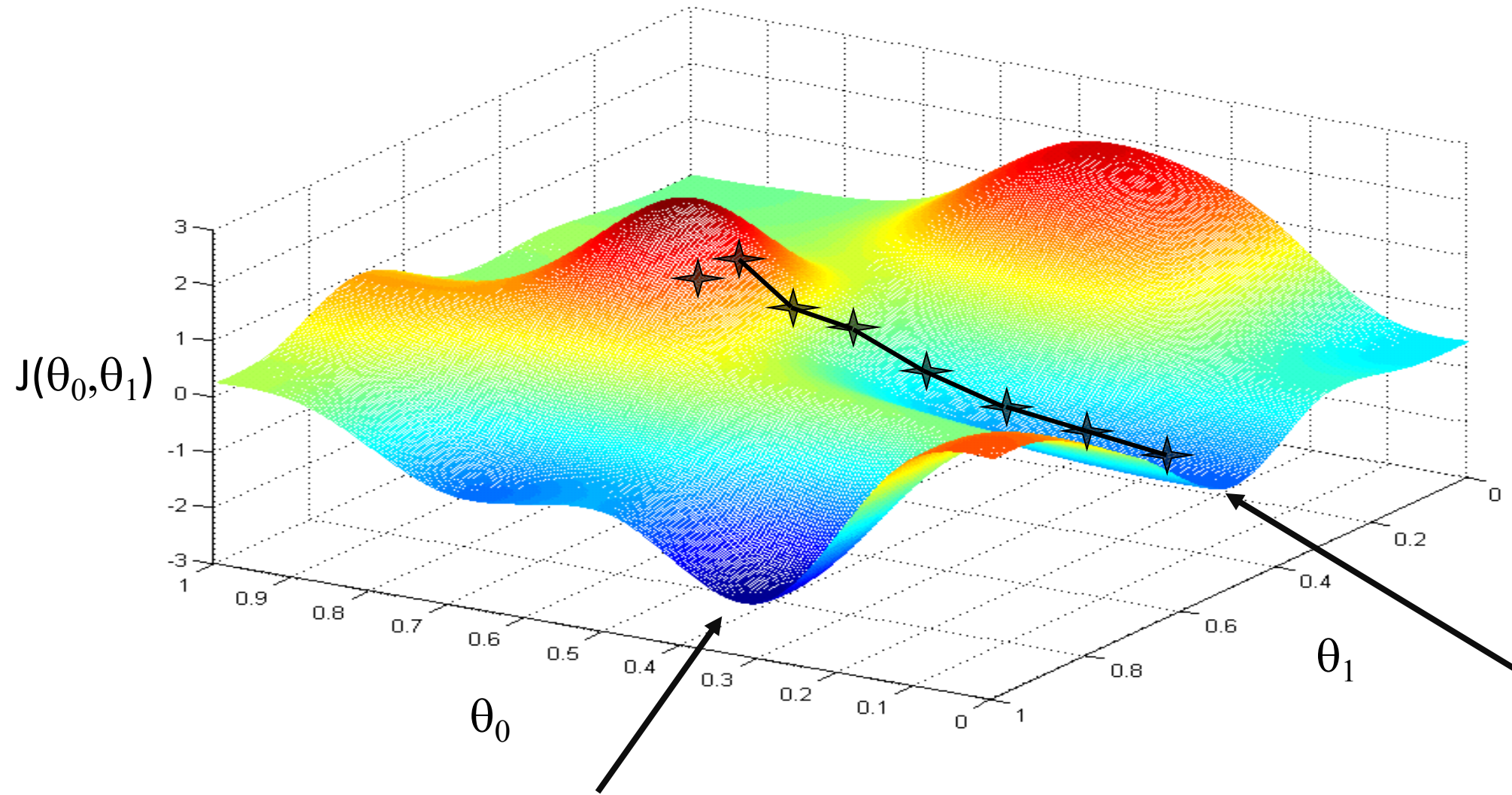








# Gradient Descent Algorithm – Initial point matters



# Gradient Descent Algorithm

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

- $\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$  – gradient or derivative
- $\alpha$  – learning rate (positive constant)

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 – 1-Dimensional

repeat until convergence {

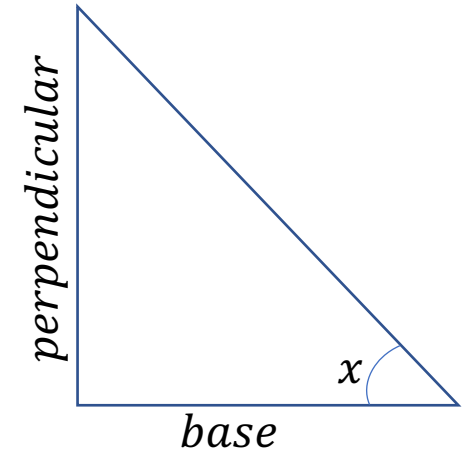
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_j} J(\theta_1)$$

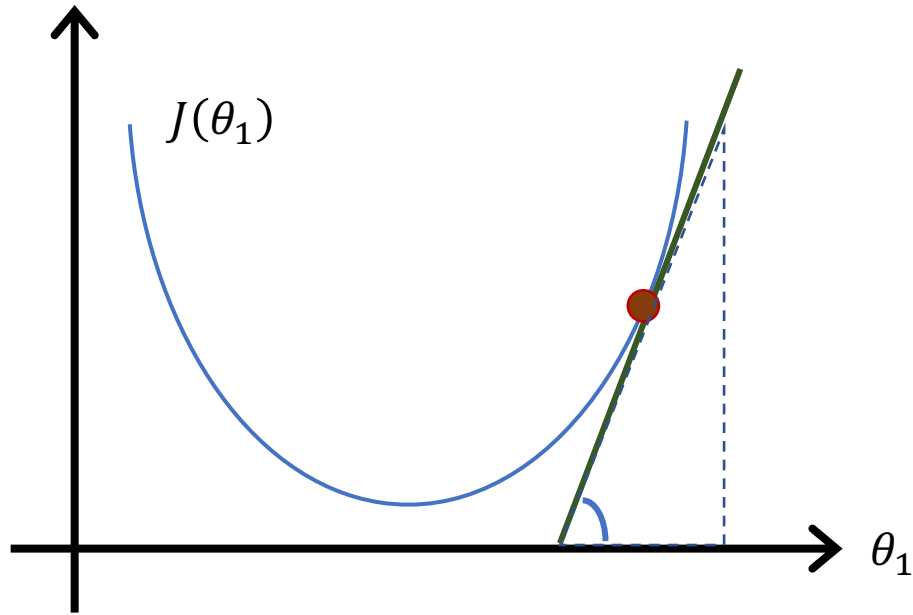
}

- $\frac{\partial}{\partial \theta_j}$  – partial derivative
- $\frac{d}{d\theta_j}$  – derivative

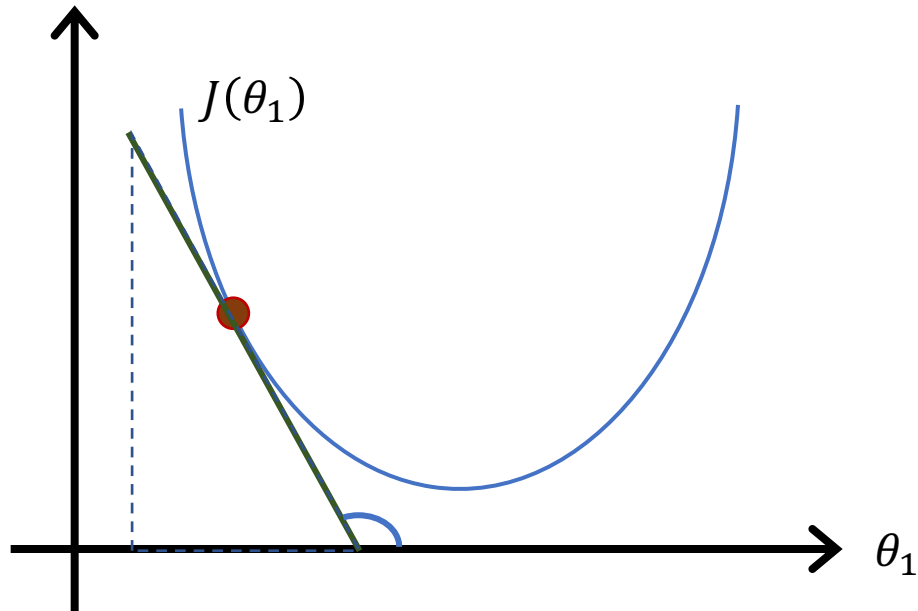
Few more pointers:

- $\tan x = \frac{\text{perpendicular}}{\text{base}}$
- $\tan x$  is positive in 1<sup>st</sup> and 3<sup>rd</sup> quadrant.
- $\tan x$  is positive when  $0^\circ < x < 90^\circ$  and  $180^\circ < x < 270^\circ$ .





- $\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_j} J(\theta_1)$
- $\frac{d}{d\theta_j} J(\theta_1)$  is positive
- $\theta_1 := \theta_1 - \alpha$  (positive)
- $\theta_1$  decreases

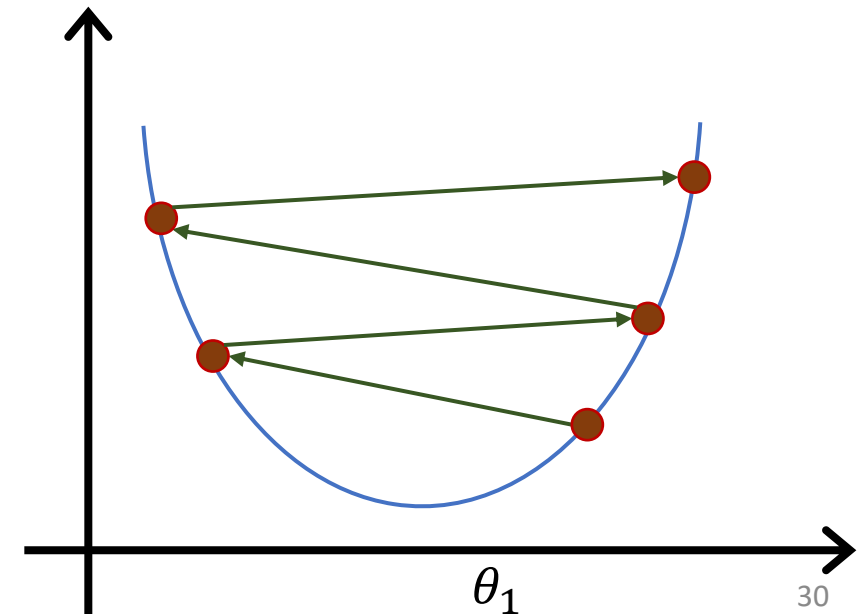
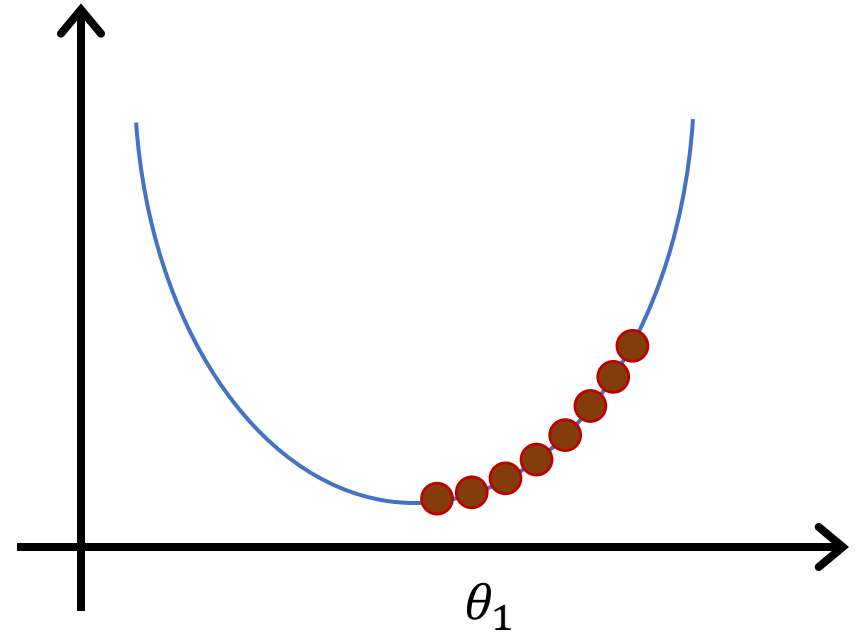


- $\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_j} J(\theta_1)$
- $\frac{d}{d\theta_j} J(\theta_1)$  is negative
- $\theta_1 := \theta_1 - \alpha$  (negative)
- $\theta_1$  increases

# Effect of Learning Rate ( $\alpha$ )

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

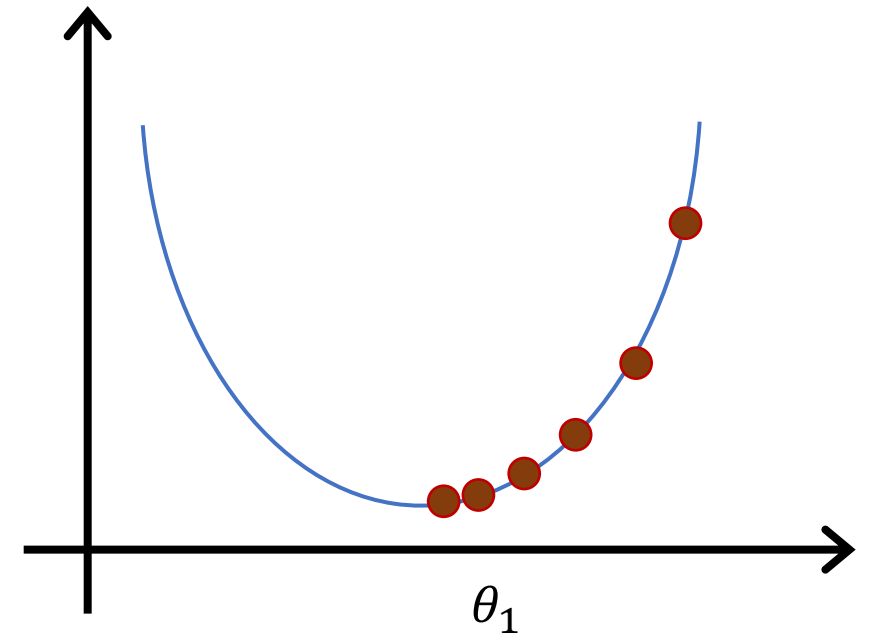
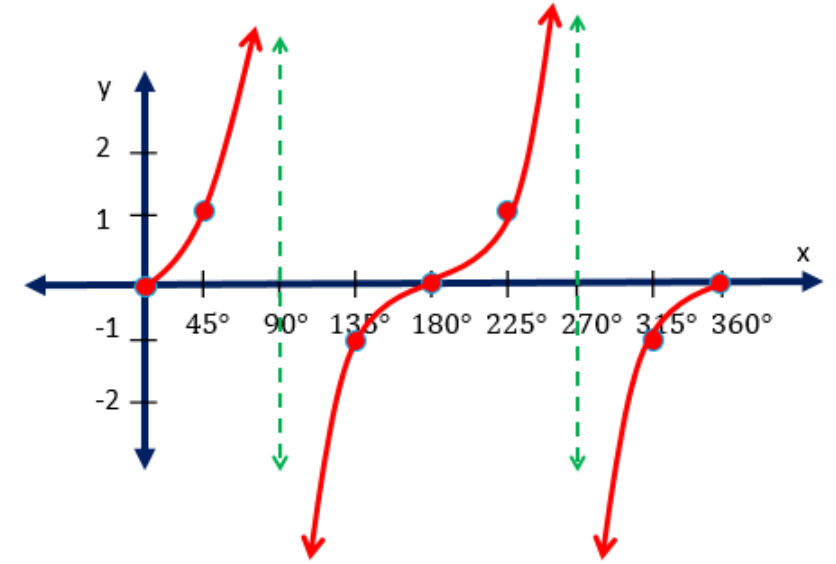
- If  $\alpha$  is too small, gradient descent can be slow.
- If  $\alpha$  is too large, gradient descent can overshoot the minimum. It may fail to converge, or even diverge.



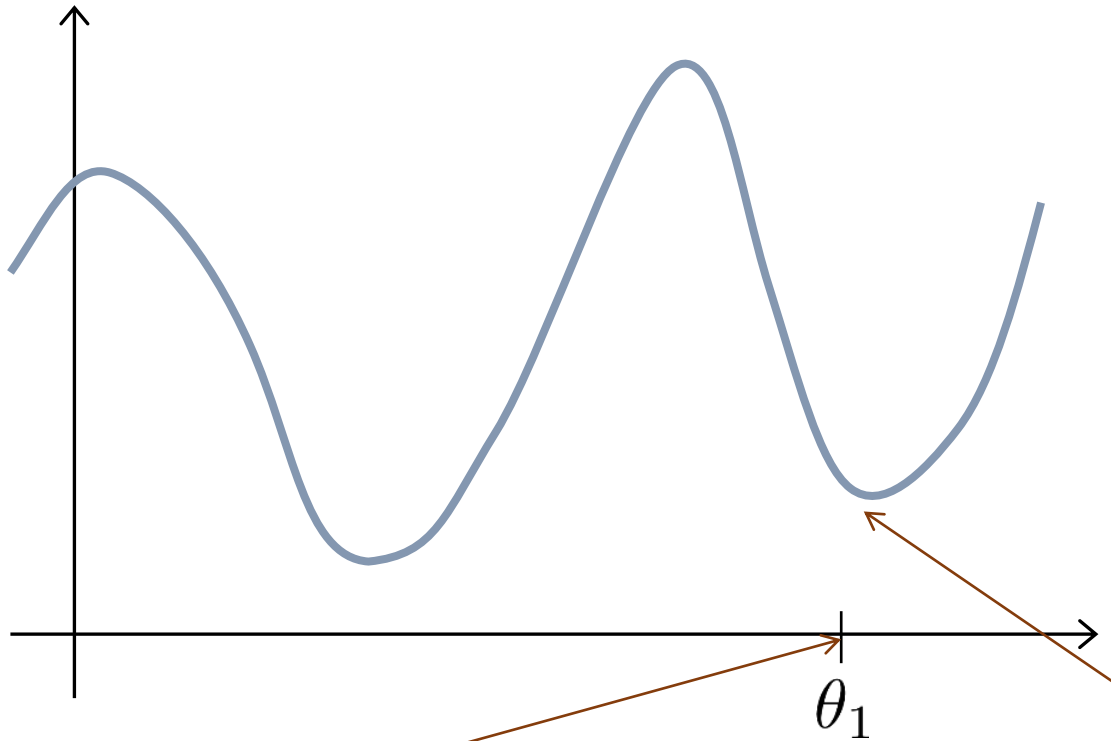
# Effect of Learning Rate ( $\alpha$ )

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

- Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.
- As we approach a local minimum, gradient descent will automatically take smaller steps. So, no need to decrease  $\alpha$  over time.



# Behavior at Local Minima



Current value of  $\theta_1$

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

At local optima:  $\frac{d}{d\theta_j} J(\theta_1) = 0$

No update will occur

$\theta_1$  at local optima

# Linear Regression and Gradient Descent Algorithm

## Gradient Descent Algorithm

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

## Linear Regression

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

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

We will use the Gradient Descent Algorithm to

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



# Calculating the Gradients

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) &= \frac{\partial}{\partial \theta_j} \left( \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right) \\ &= \frac{\partial}{\partial \theta_j} \left( \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 \cdot x^{(i)} - y^{(i)})^2 \right)\end{aligned}$$

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

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

# Gradient descent algorithm on Linear Regression

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

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

repeat until convergence {

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

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

}

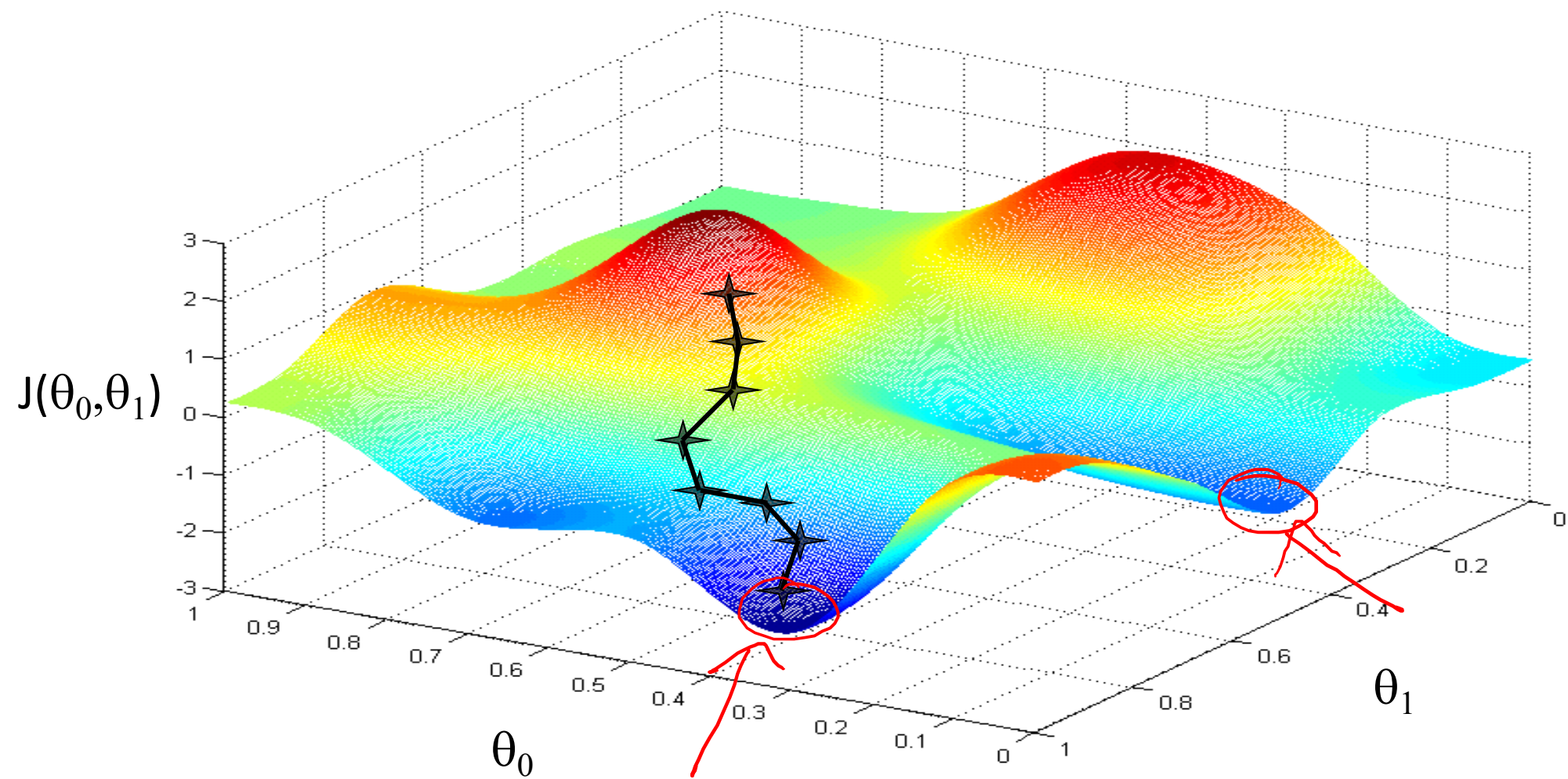
# update  $\theta_0, \theta_1$  simultaneously

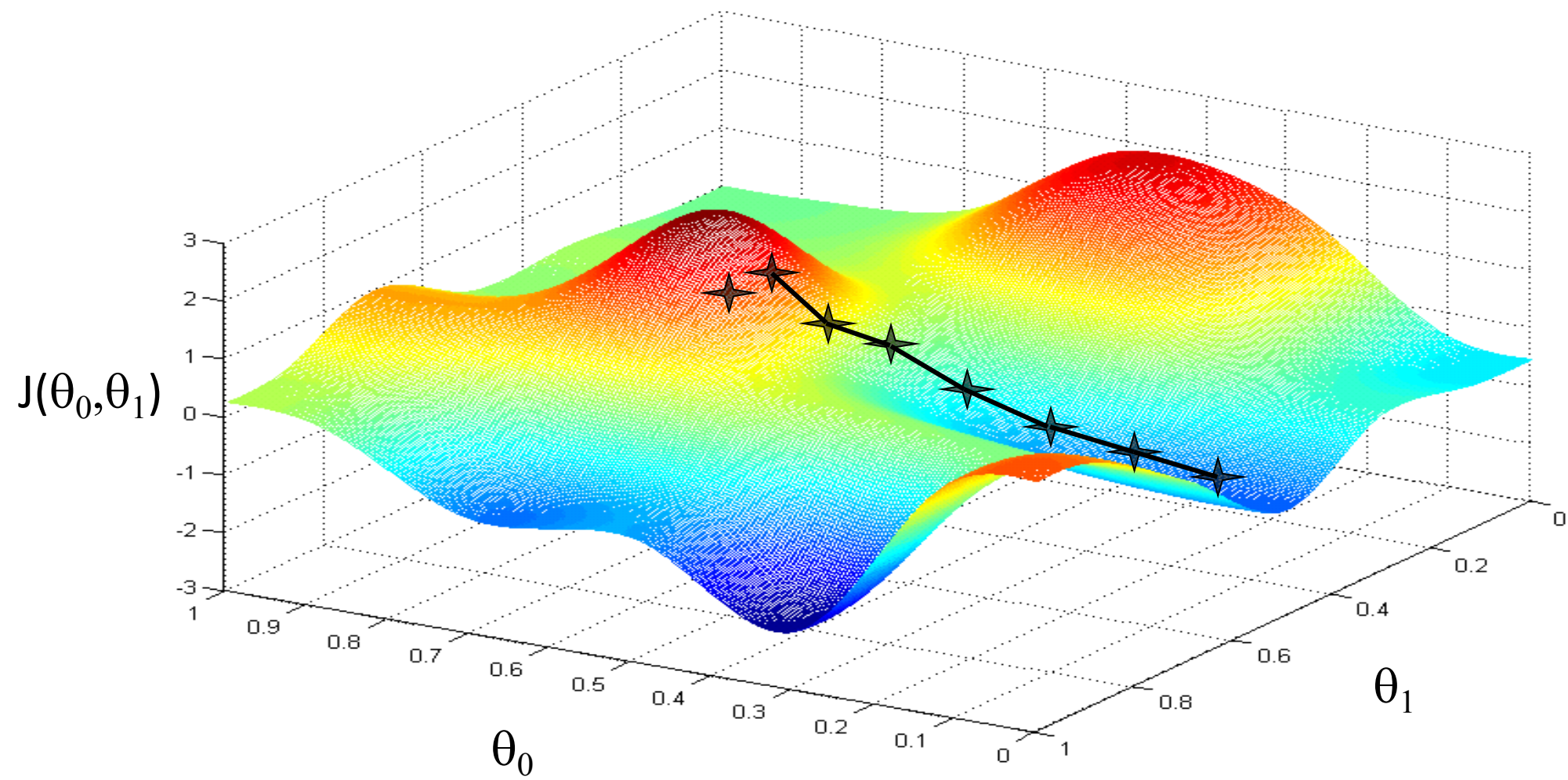
repeat until convergence {

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

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

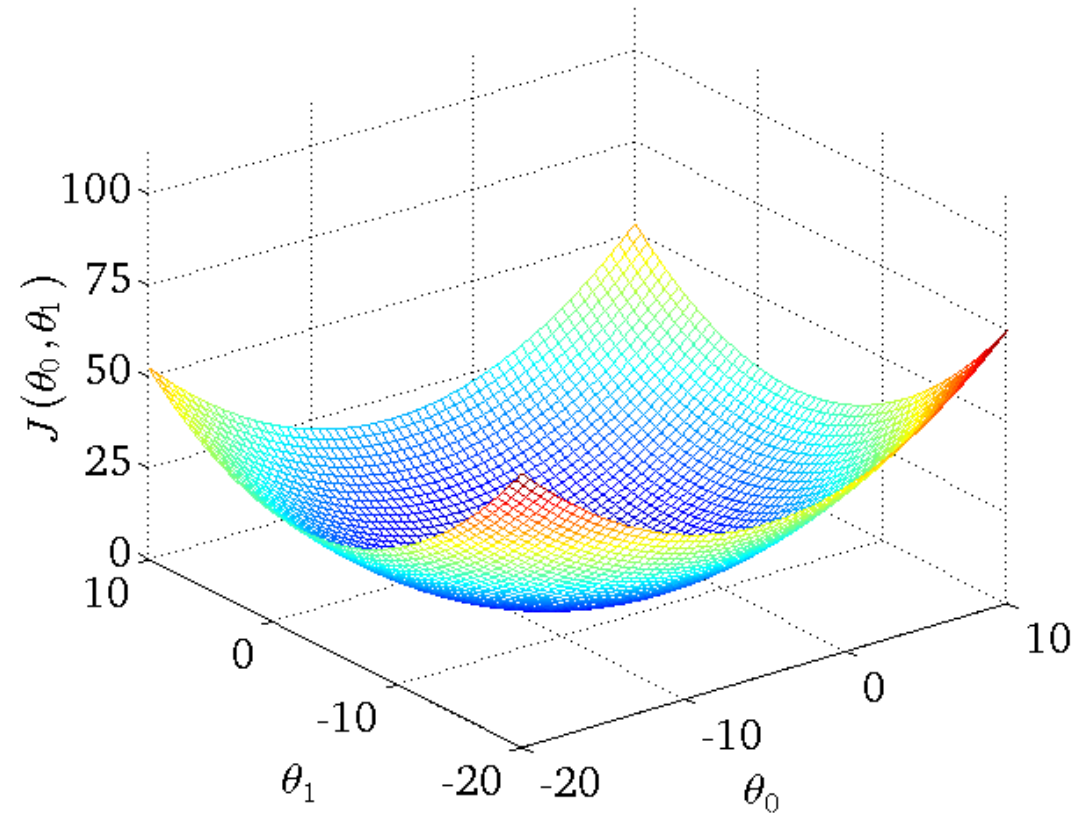
}





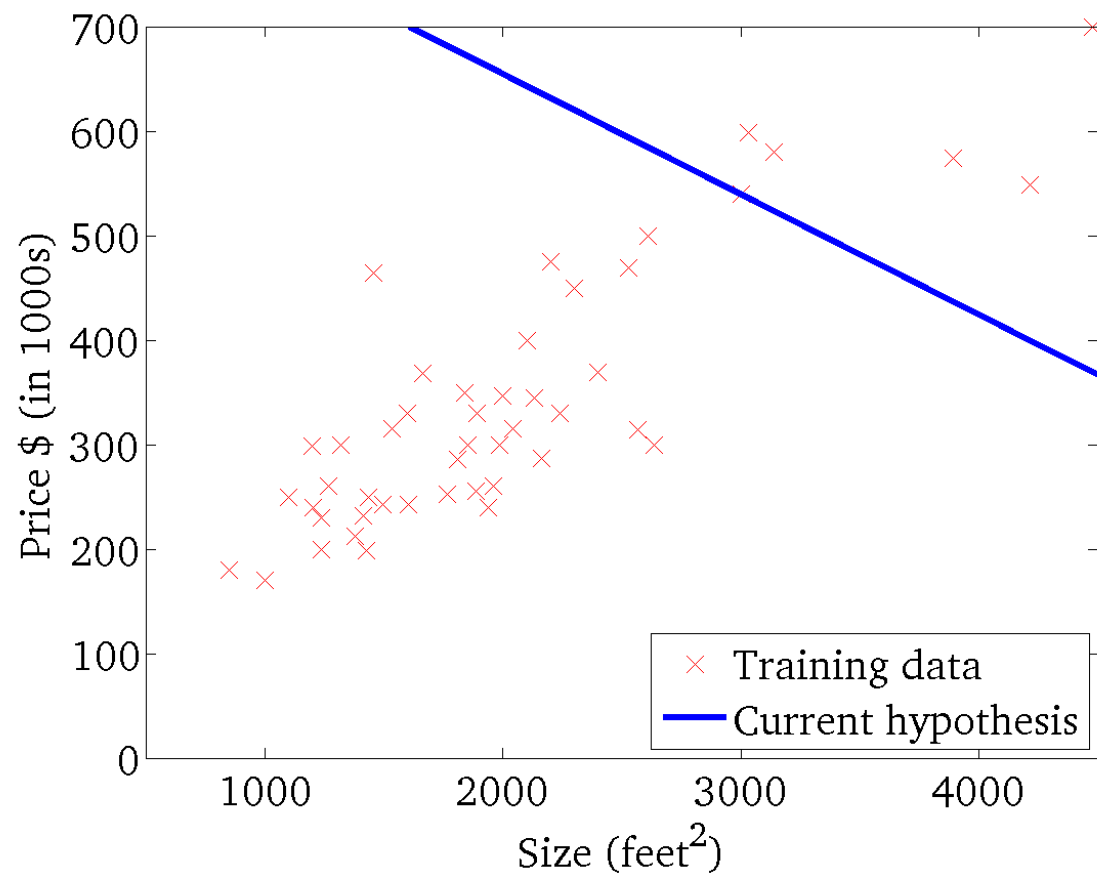
# Gradient descent algorithm on Linear Regression

- Cost function for linear regression will always be a **bowl shaped**.
- Formally, these kind of functions are called as **convex function**.
- A convex function will always have only one local optima, which is also the global optima.



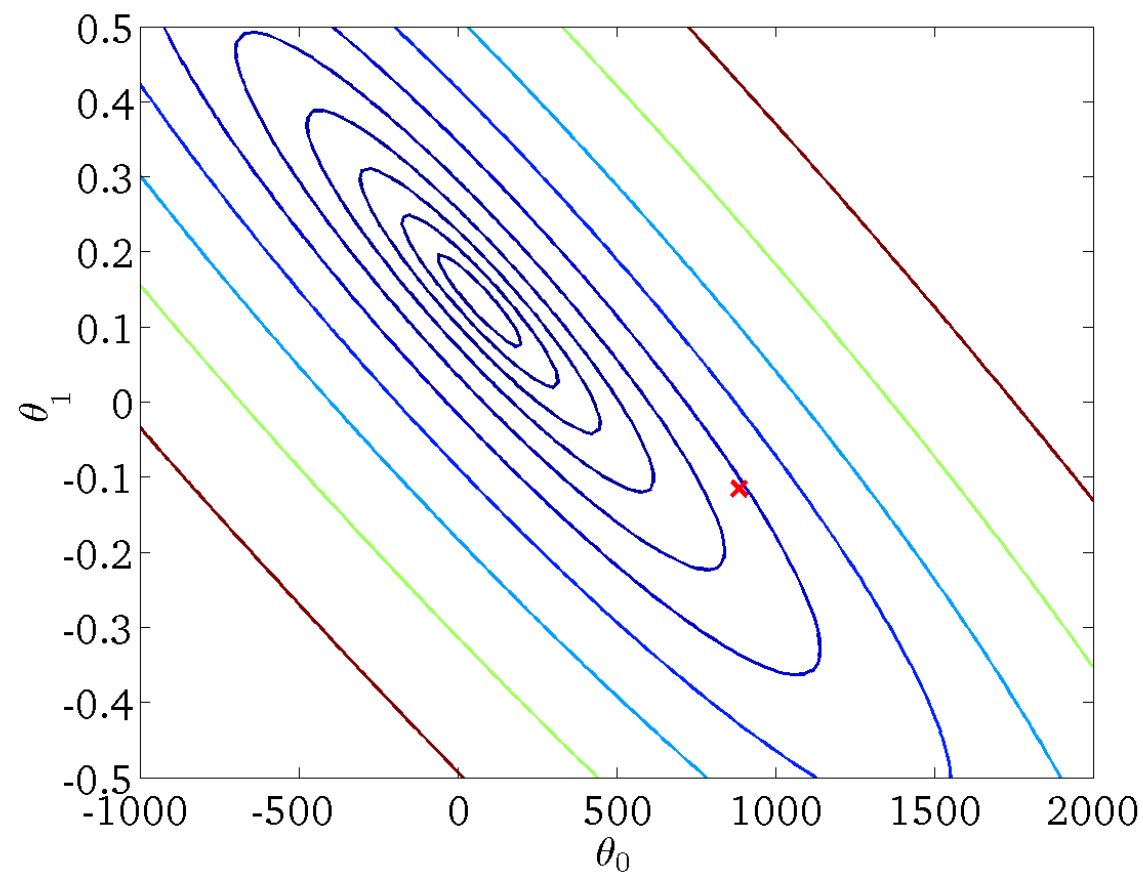
$$h_{\theta}(x)$$

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



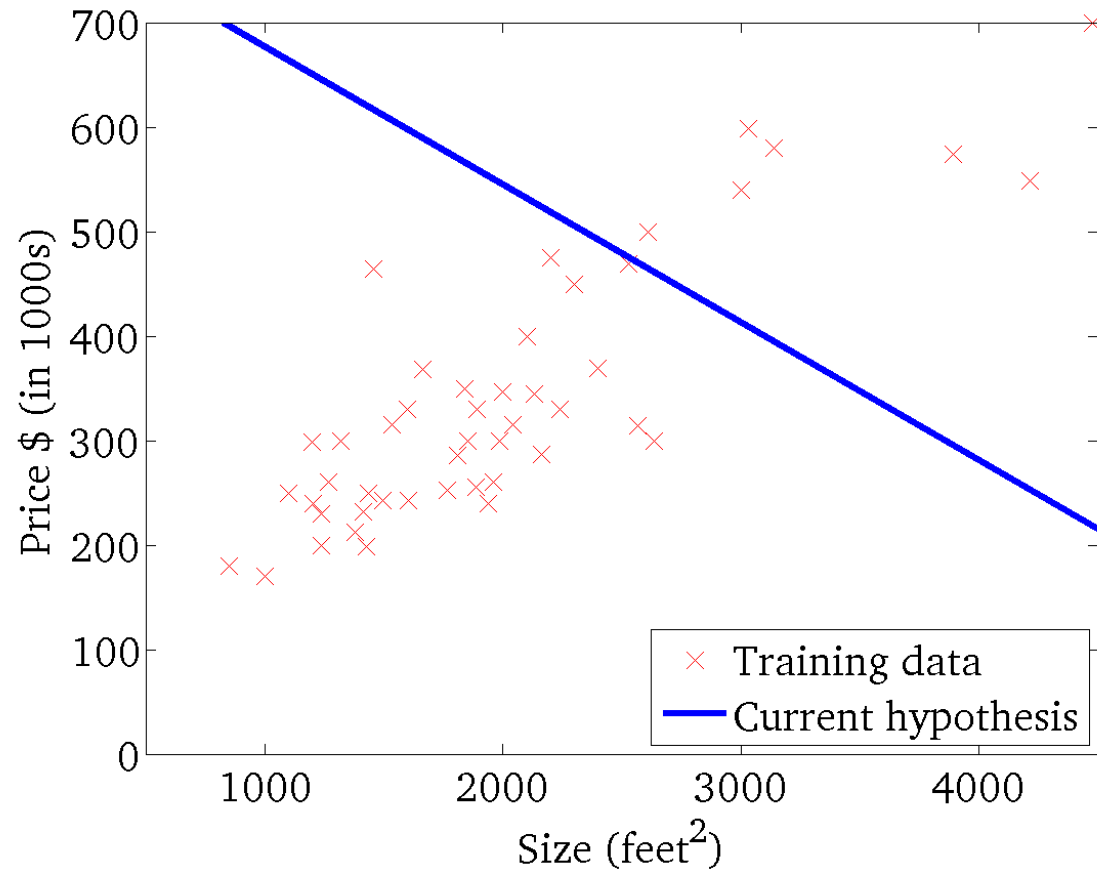
$$J(\theta_0, \theta_1)$$

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



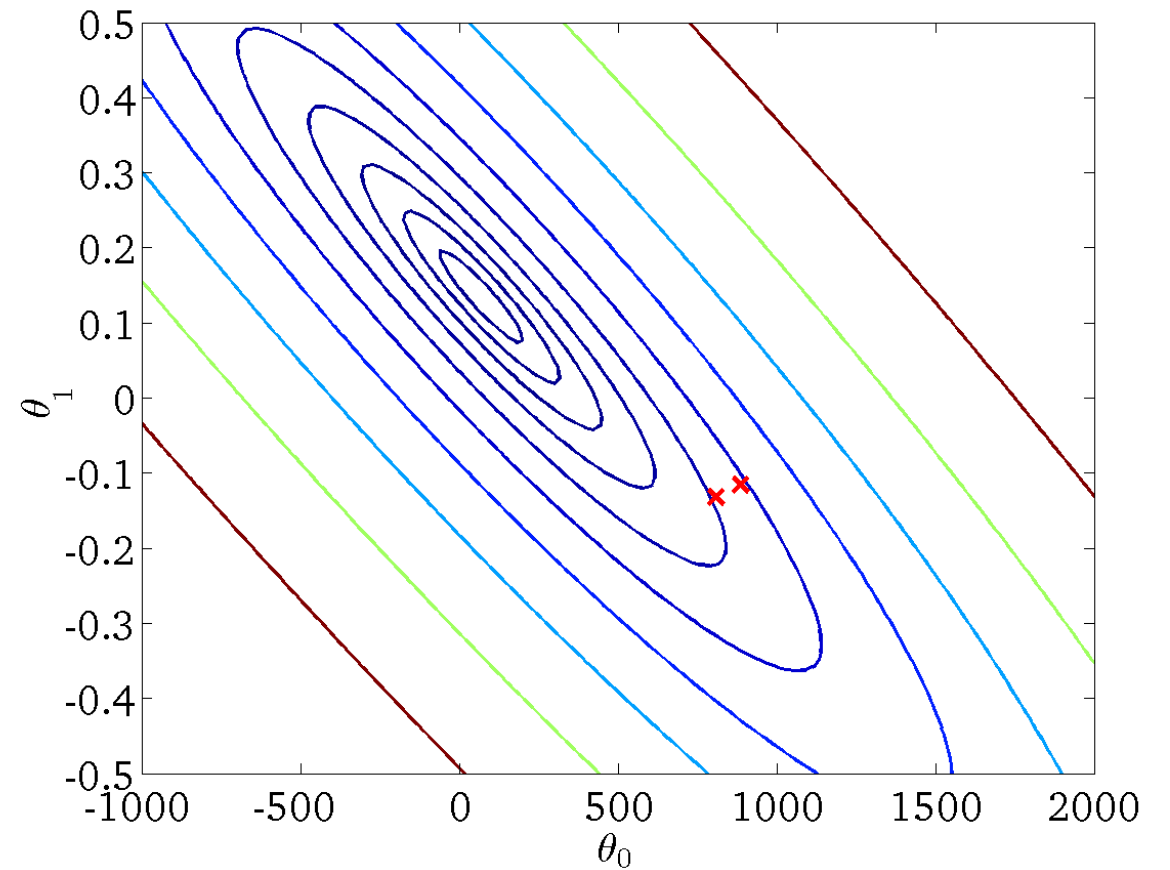
$$h_{\theta}(x)$$

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



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

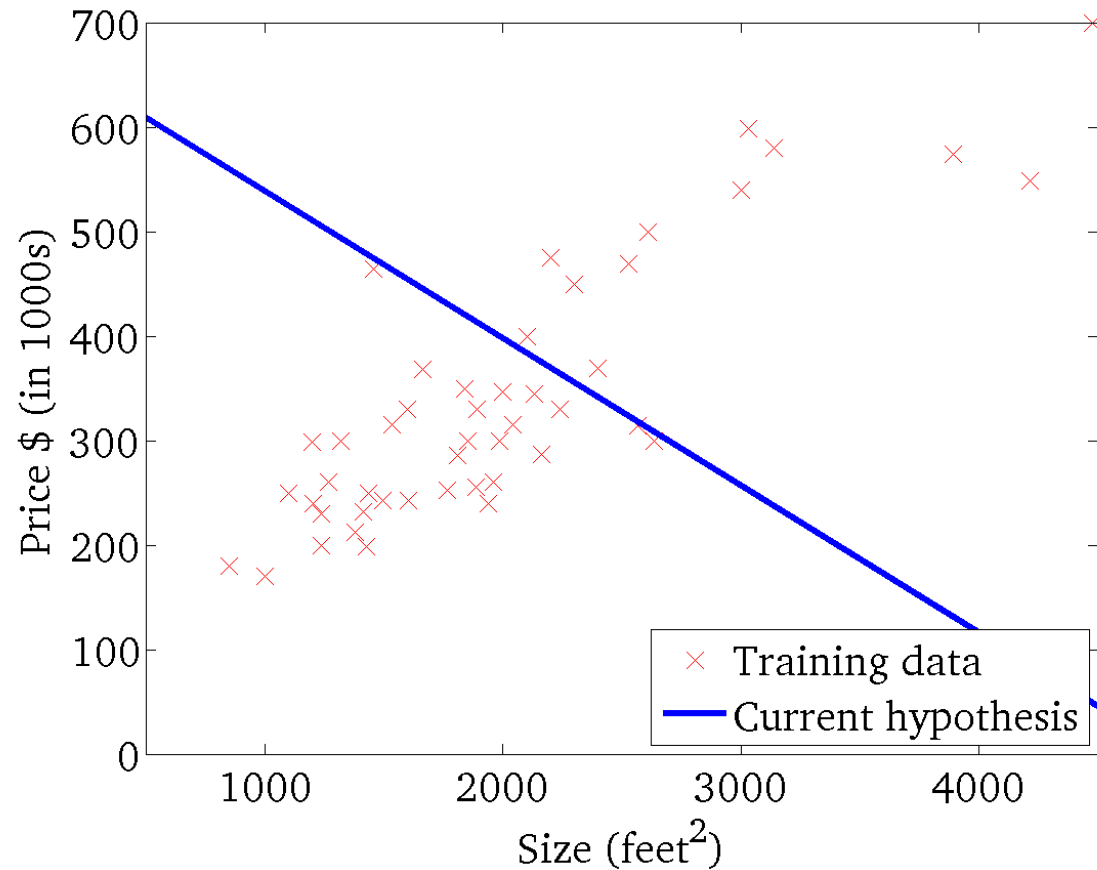
(function of the parameters  $\theta_0, \theta_1$ )





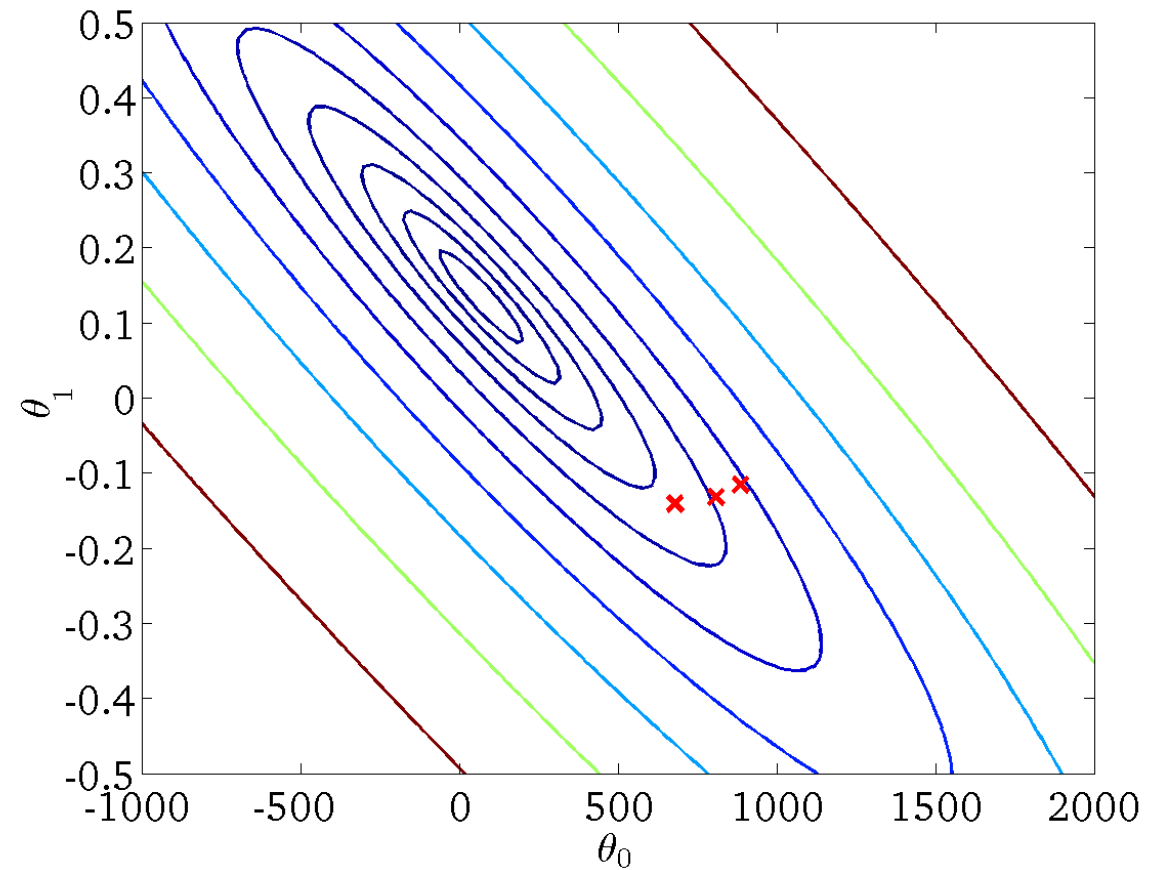
$$h_{\theta}(x)$$

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



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

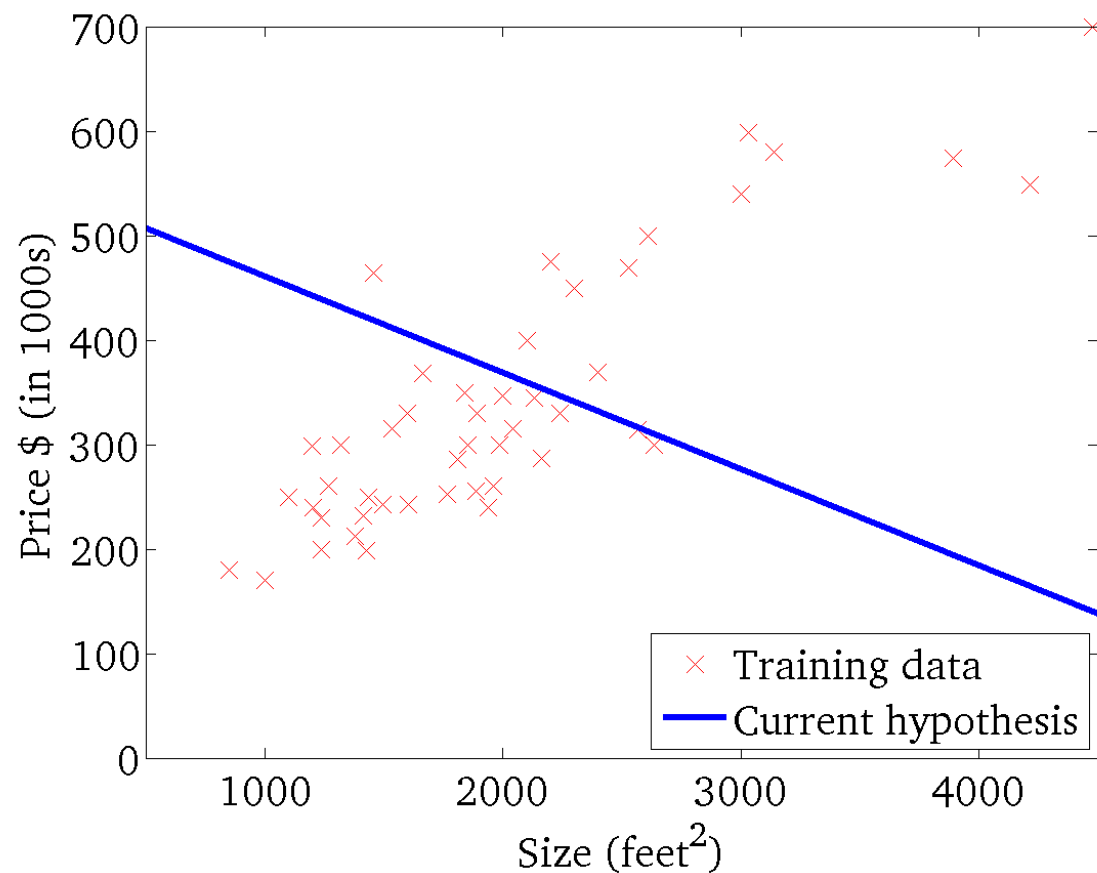
(function of the parameters  $\theta_0, \theta_1$ )





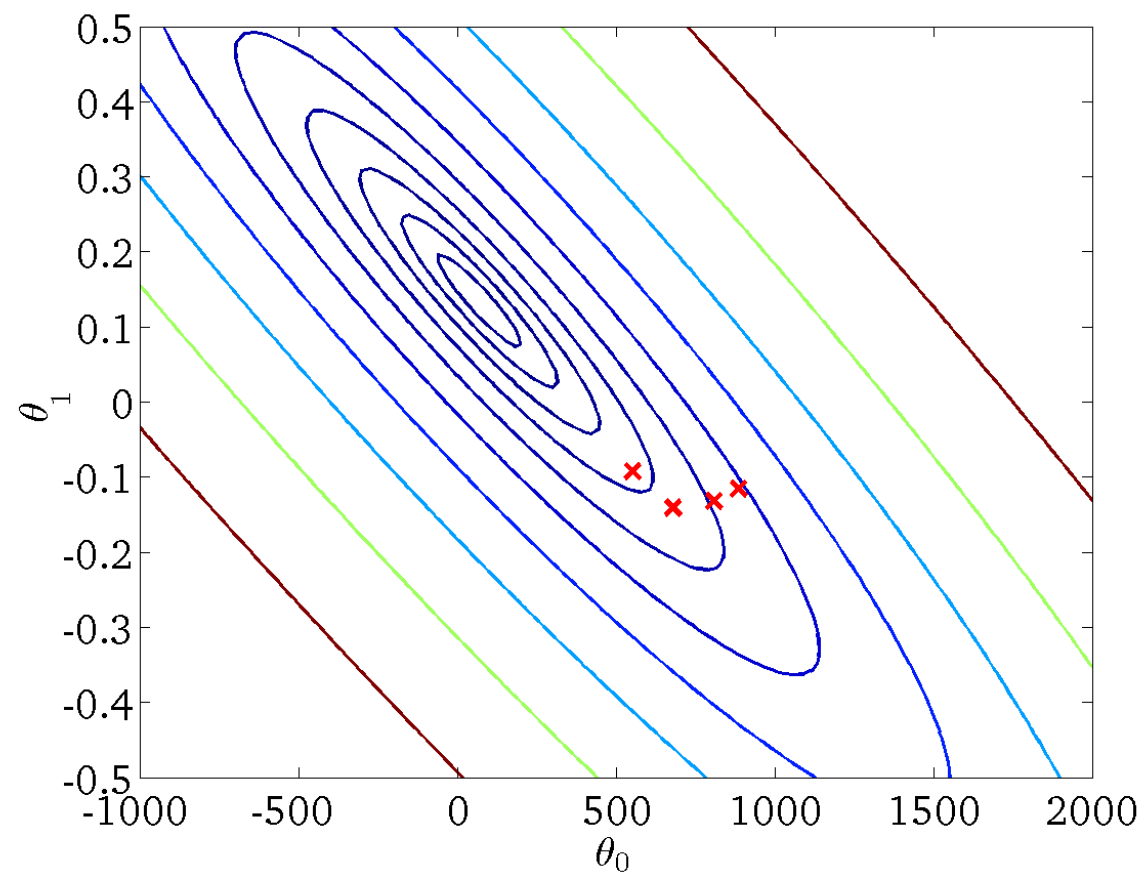
$$h_{\theta}(x)$$

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



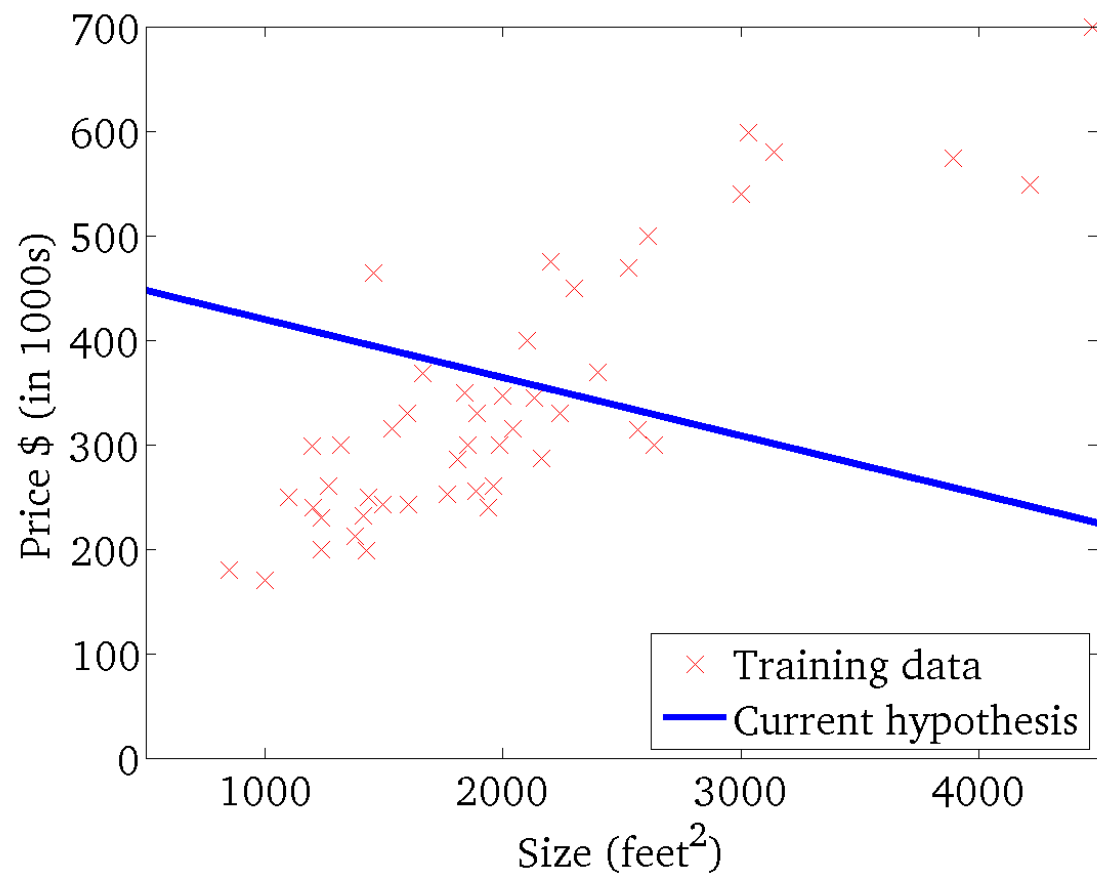
$$J(\theta_0, \theta_1)$$

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



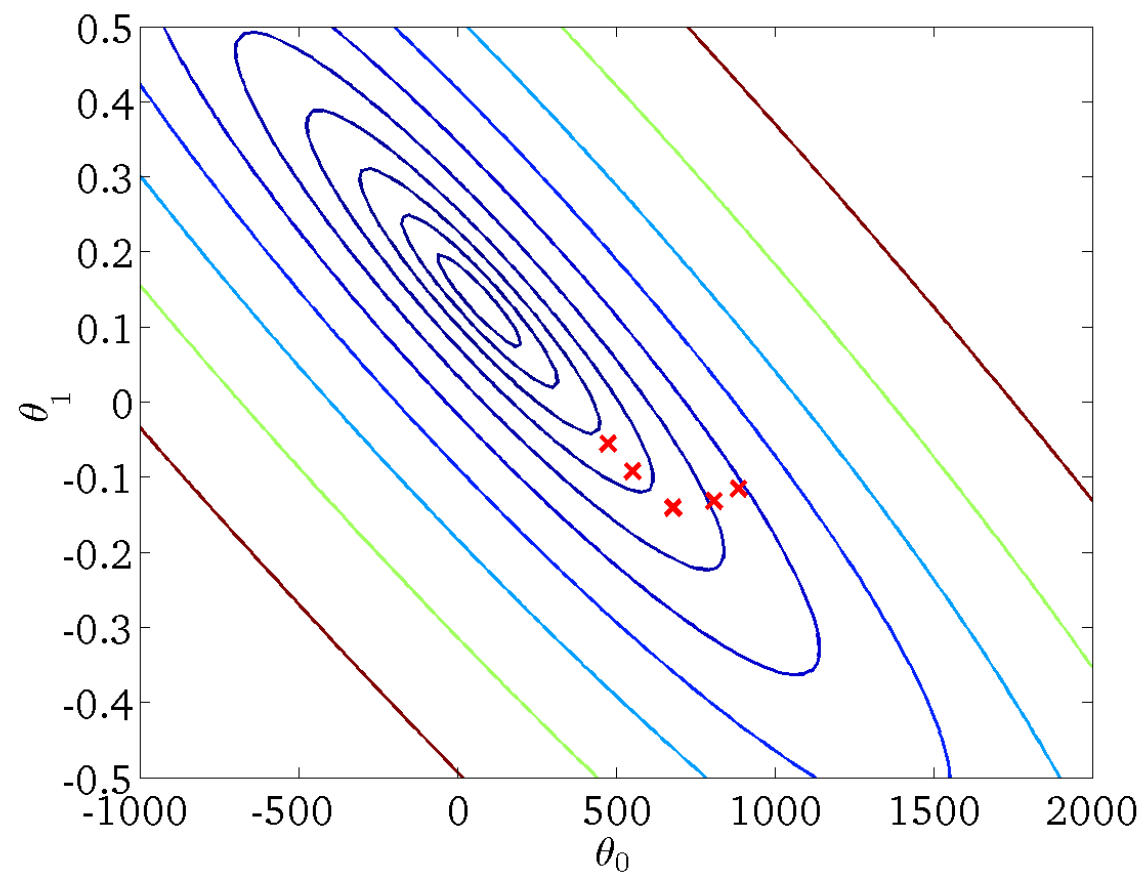
$$h_{\theta}(x)$$

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



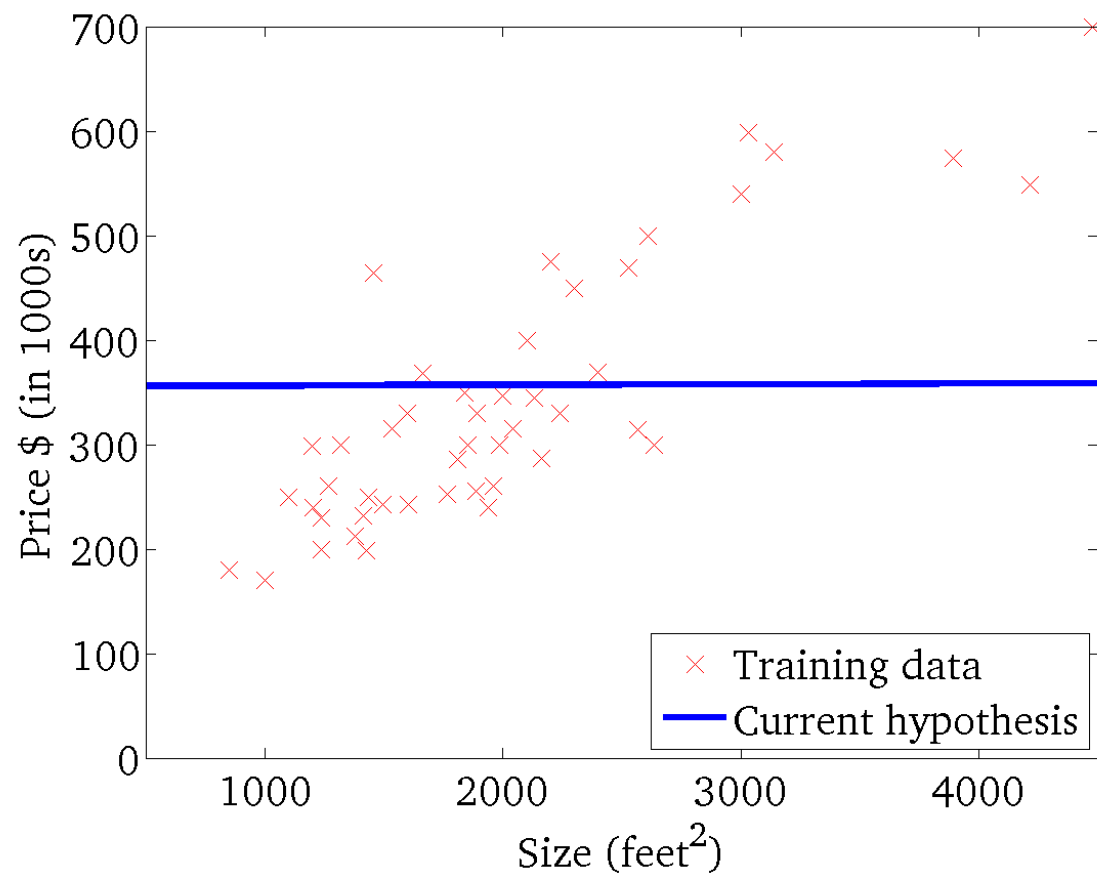
$$J(\theta_0, \theta_1)$$

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



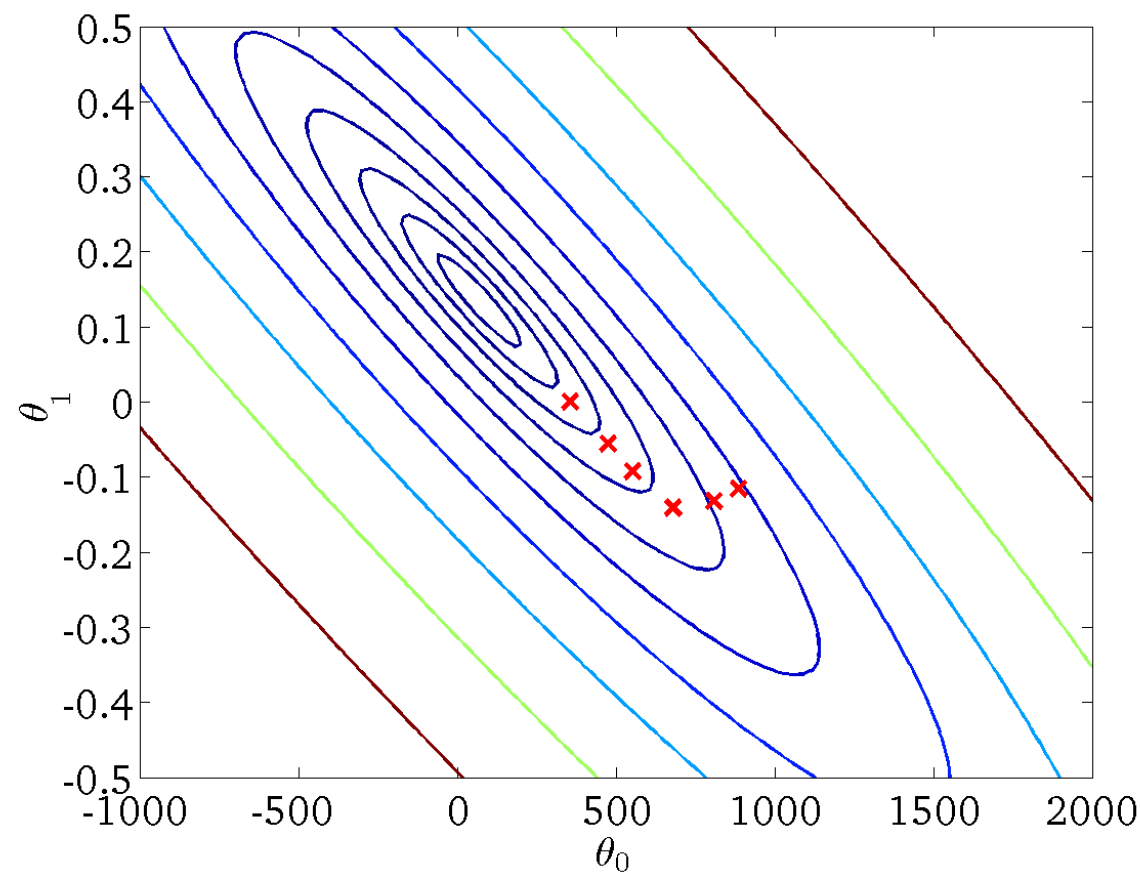
$$h_{\theta}(x)$$

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



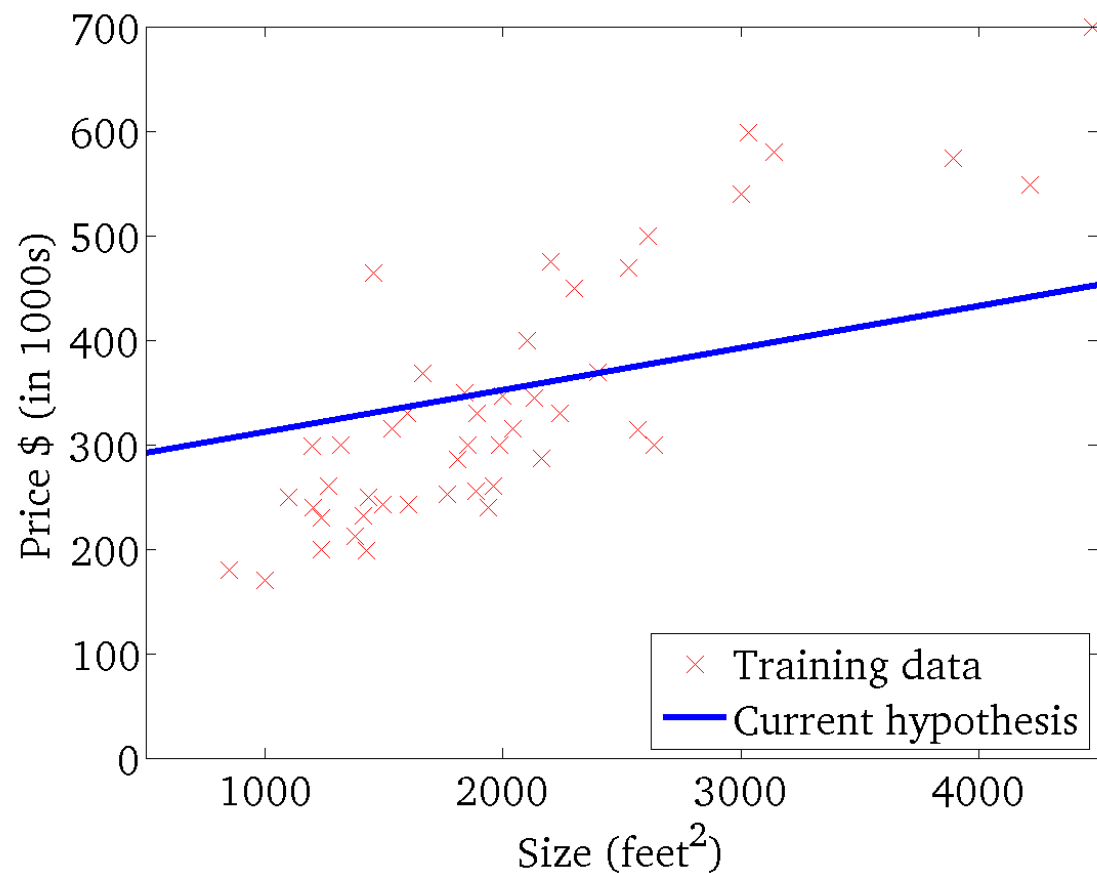
$$J(\theta_0, \theta_1)$$

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



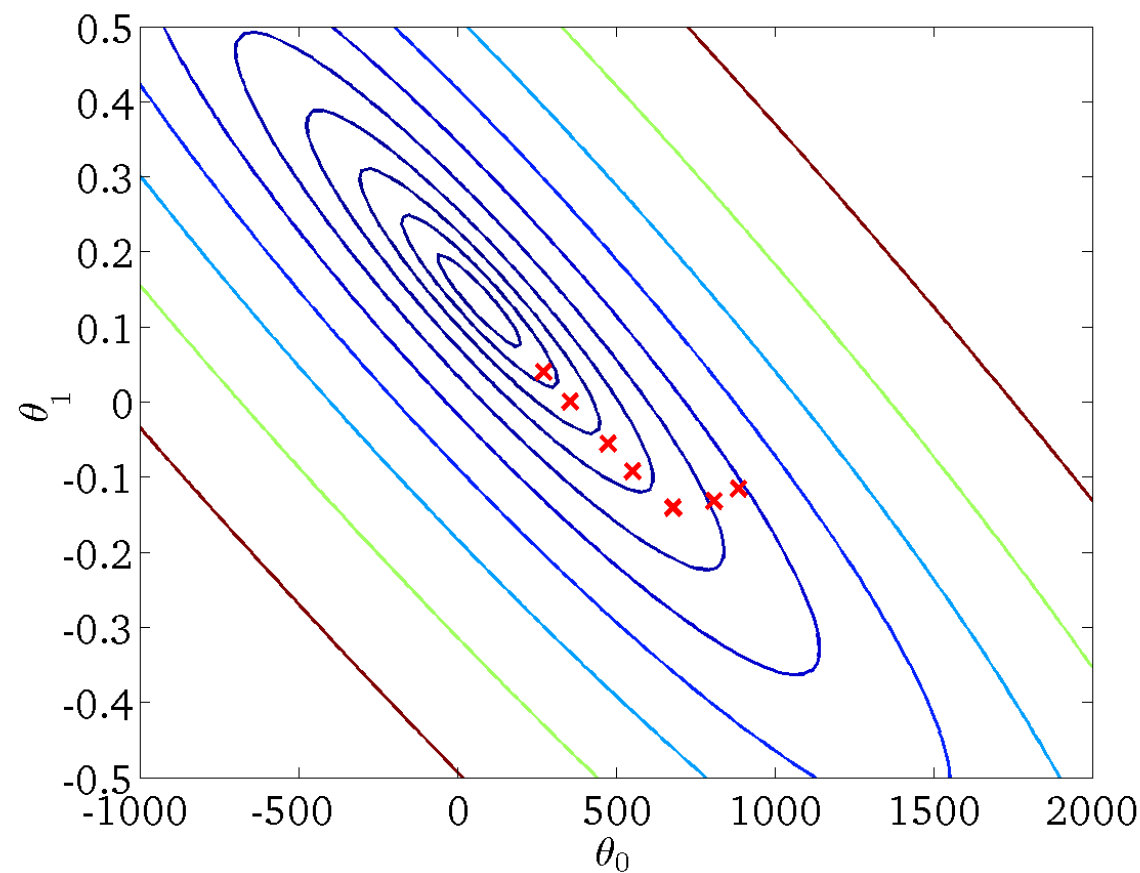
$$h_{\theta}(x)$$

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



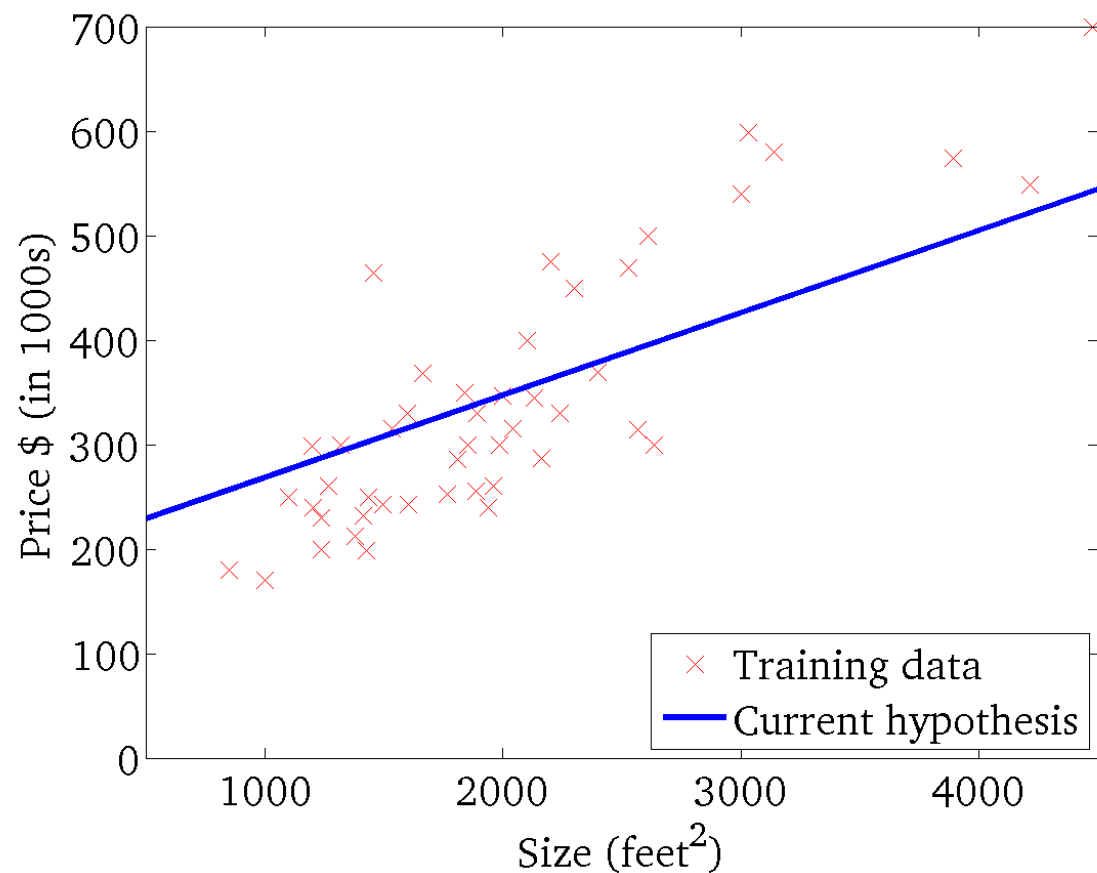
$$J(\theta_0, \theta_1)$$

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



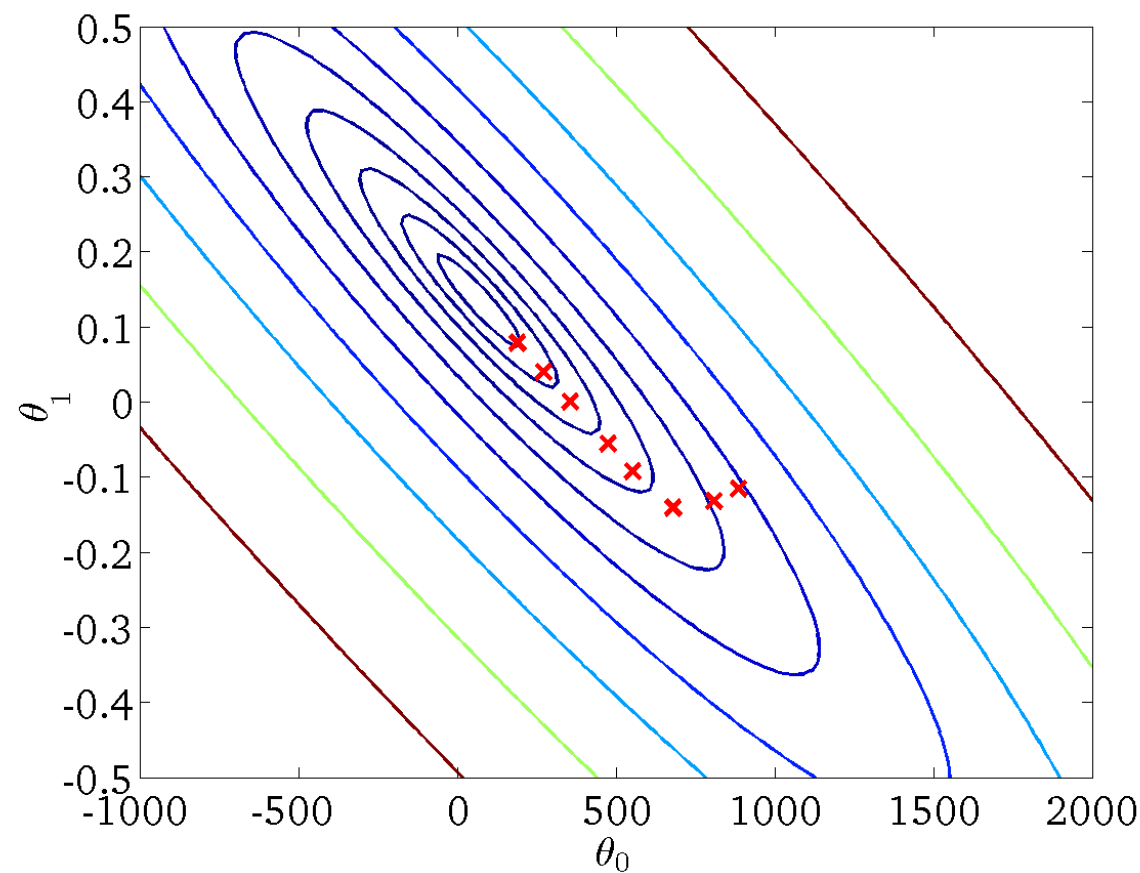
$$h_{\theta}(x)$$

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



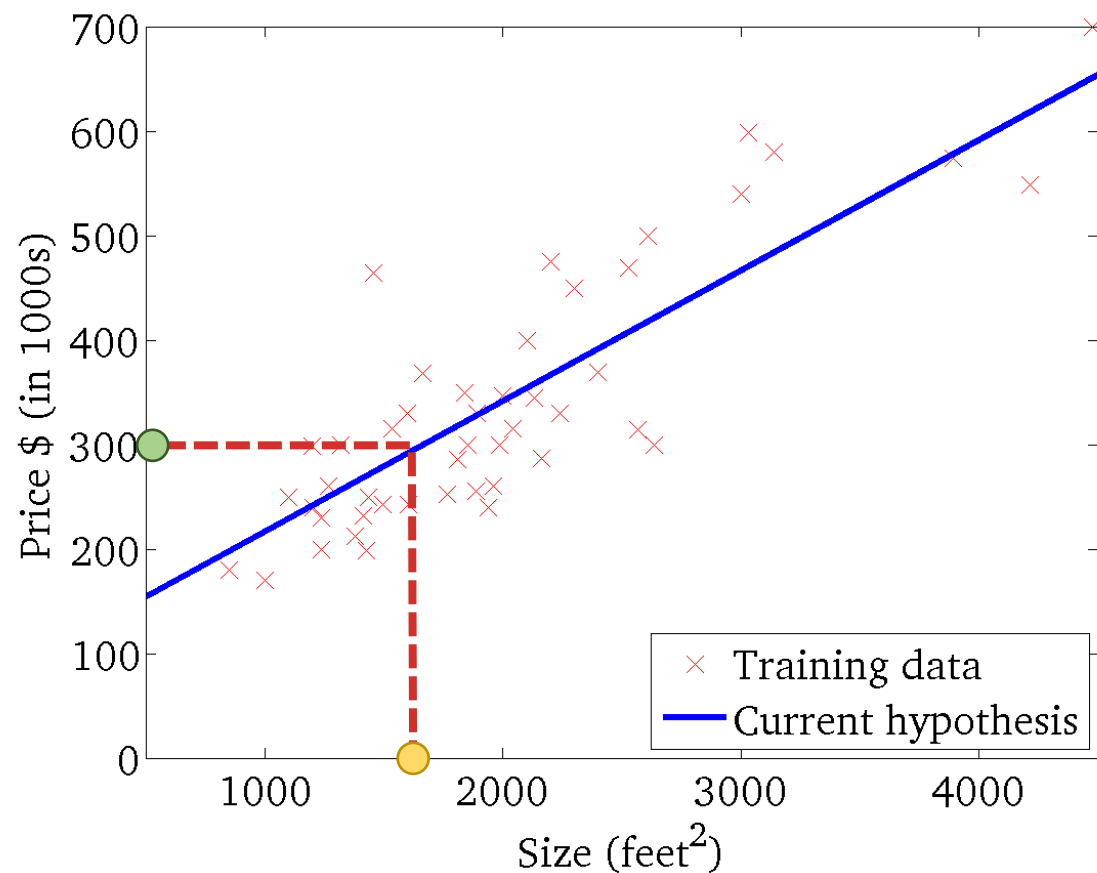
$$J(\theta_0, \theta_1)$$

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



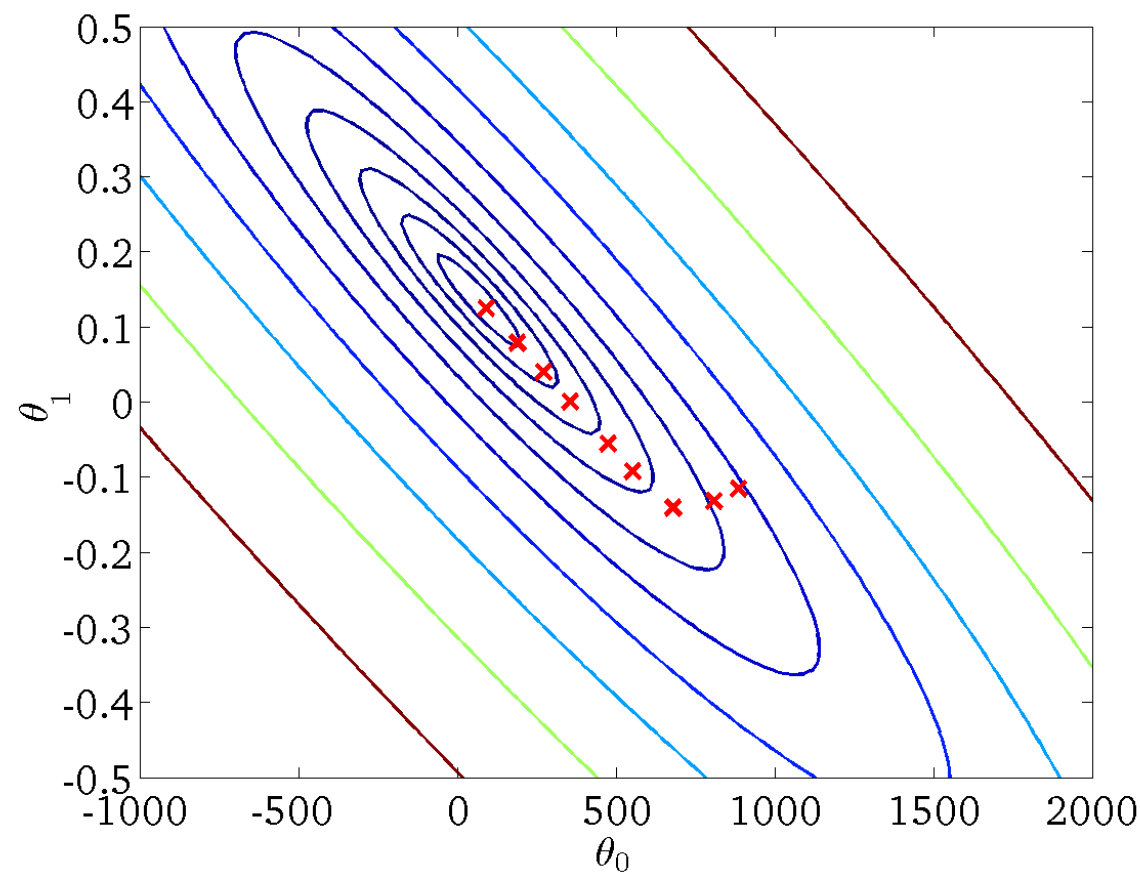
$$h_{\theta}(x)$$

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



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

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



Why  
“Linear Regression”  
and not  
“Affine Regression”?

## Reference

- CS229: Machine Learning, <https://cs229.stanford.edu/>
- Machine Learning, <https://www.coursera.org/learn/machine-learning>