



# Linear regression with one variable

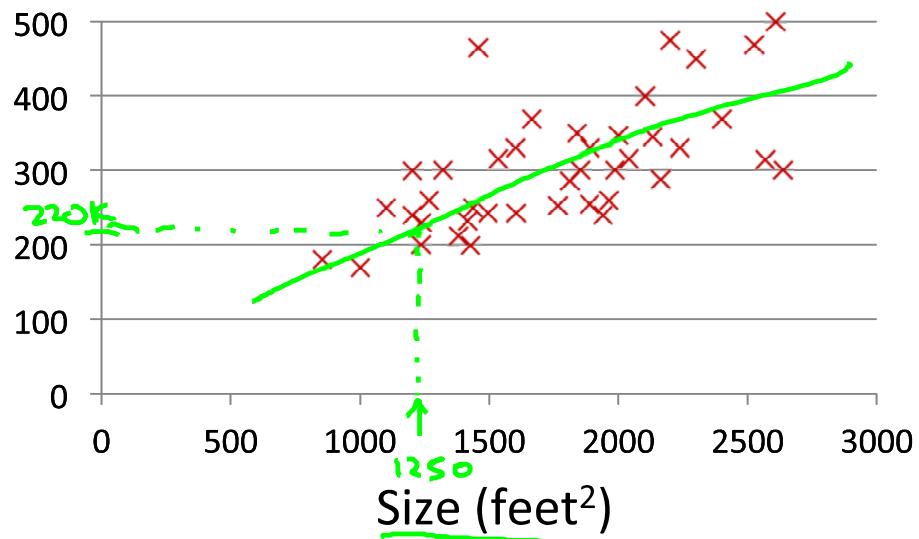
---

## Model representation

Machine Learning

# 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
...	...
i	↑

Notation:

- m = Number of training examples
- x's = "input" variable / features
- y's = "output" variable / "target" variable

$(x, y)$  - one training example  
 $(x^{(i)}, y^{(i)})$  - i<sup>th</sup> training example

$$\left| \begin{array}{l} x^{(1)} = 2104 \\ x^{(2)} = 1416 \\ y^{(1)} = 460 \\ \uparrow \end{array} \right.$$

Andrew |

Training Set

Learning Algorithm

Size of house

x

hypothesis

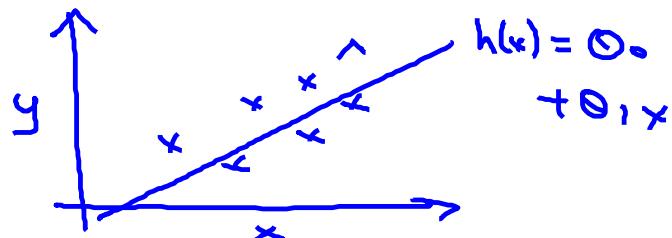
Estimated price  
(estimated value of y)

h maps from x's to y's.

How do we represent  $h$  ?

$$h_{\Theta}(x) = \underline{\underline{\Theta_0 + \Theta_1 x}}$$

Shorthand:  $\hat{h}(x)$

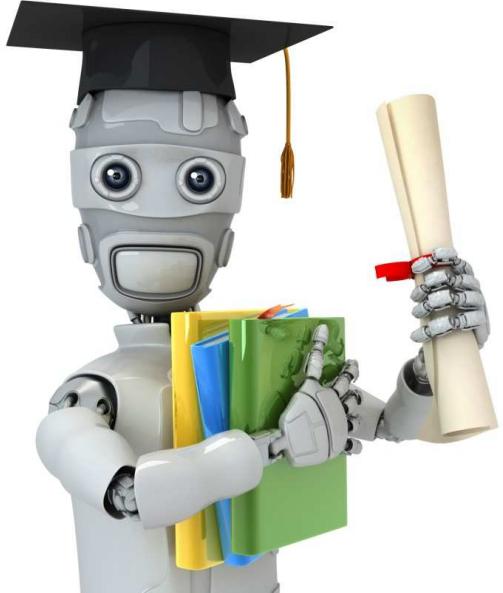


Linear regression with one variable.  
Univariate linear regression.

one variable

b7

Andrew |



# Linear regression with one variable

---

## Cost function

Machine Learning

# Training Set

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

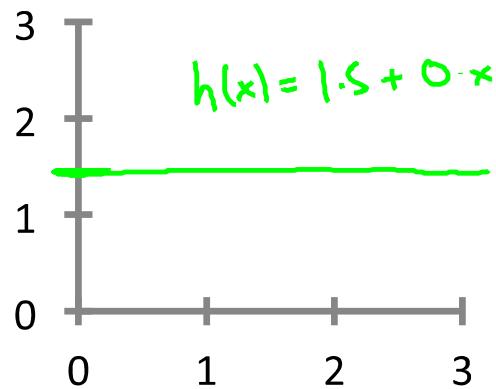
$m = 47$

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

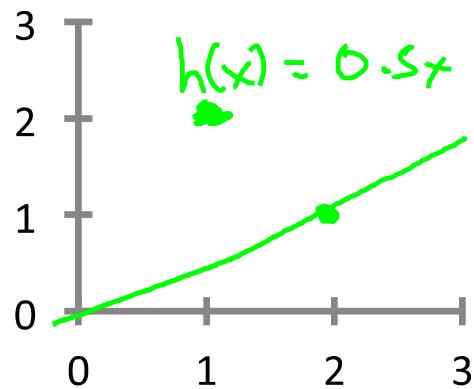
$\theta_i$ 's: Parameters

How to choose  $\theta_i$ 's ?

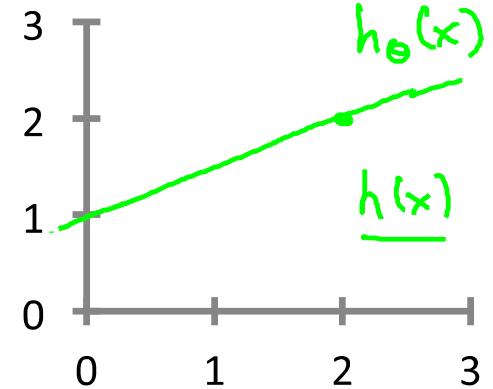
$$\underline{h_{\theta}(x) = \theta_0 + \theta_1 x}$$



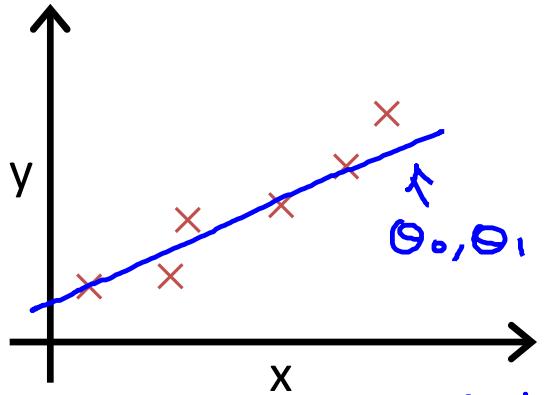
$$\begin{aligned}\rightarrow \theta_0 &= 1.5 \\ \rightarrow \theta_1 &= 0\end{aligned}$$



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



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



↑

minimize  
 $\theta_0, \theta_1$

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

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

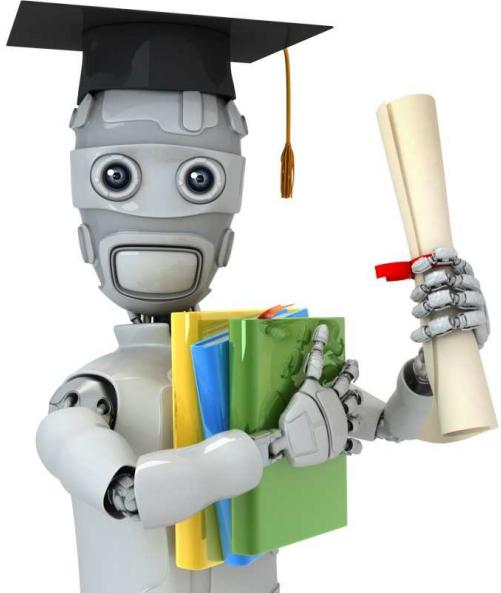
$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

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

$x, y$

minimize  $\underline{\theta_0, \theta_1}$   $\underbrace{J(\theta_0, \theta_1)}$   
 Cost function

Squared error function



# Linear regression with one variable

---

## Cost function intuition I

Machine Learning

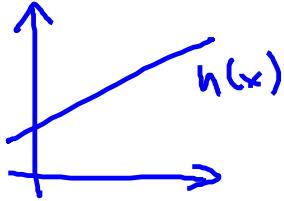
## Simplified

Hypothesis:

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

Parameters:

$$\underline{\theta_0, \theta_1}$$



Cost Function:

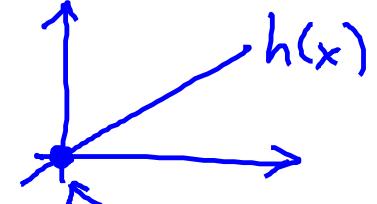
$$\rightarrow J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

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

$$\underline{\theta_0 = 0}$$

$$\underline{\theta_1}$$

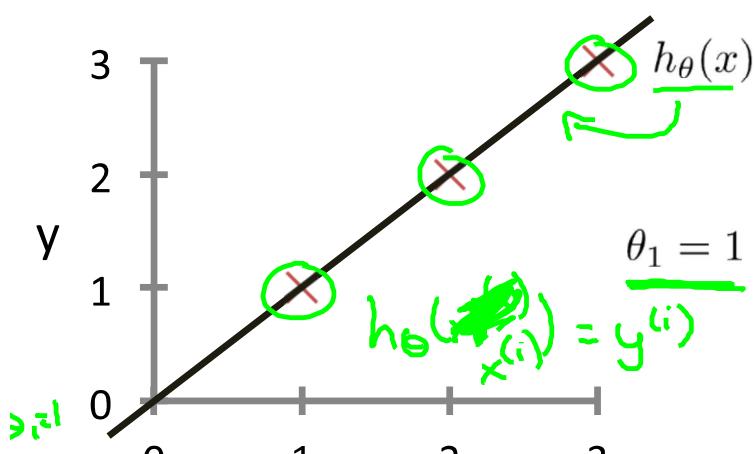


$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\underset{\theta_1}{\text{minimize}} \underline{J(\theta_1)} \quad \underline{\theta_0, x^{(i)}}$$

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

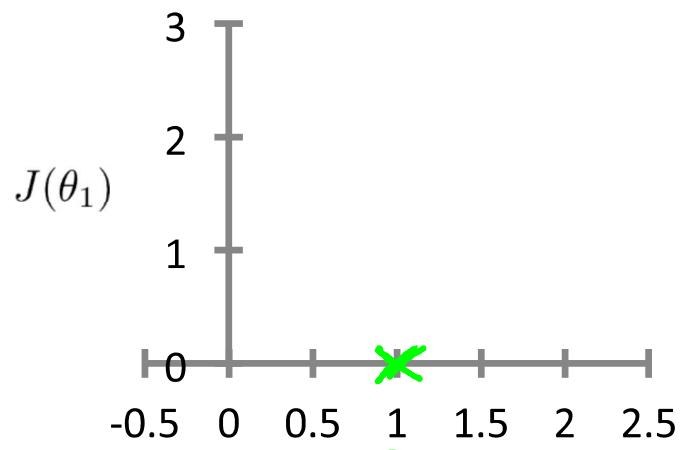
(for fixed  $\theta_1$ , this is a function of x)



$$\begin{aligned} J(\theta_1) &= \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2 \approx \frac{1}{3m} (0^2 + 0^2 + 0^2) = 0^2 \end{aligned}$$

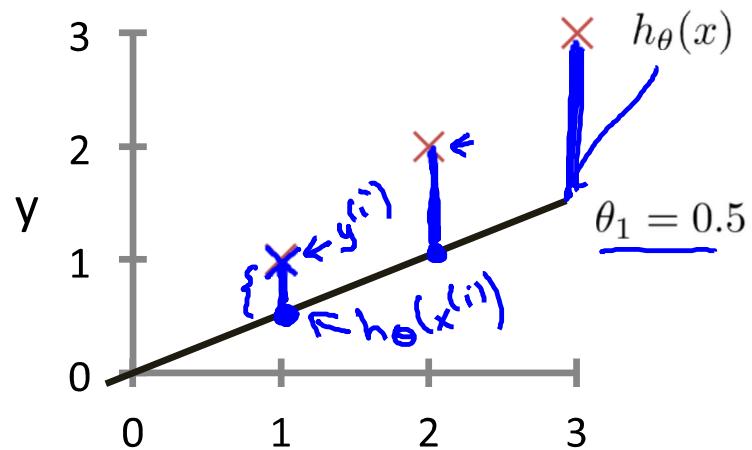
$\rightarrow \underline{J(\theta_1)}$

(function of the parameter  $\theta_1$ )



$h_{\theta}(x)$

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

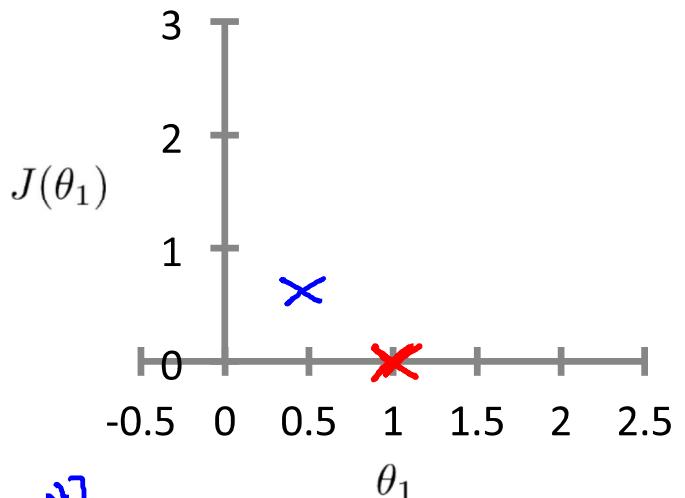


$$J(0.5) = \frac{1}{2m} \left[ (0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right]$$

$$= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} \approx \underline{0.58}$$

$J(\theta_1)$

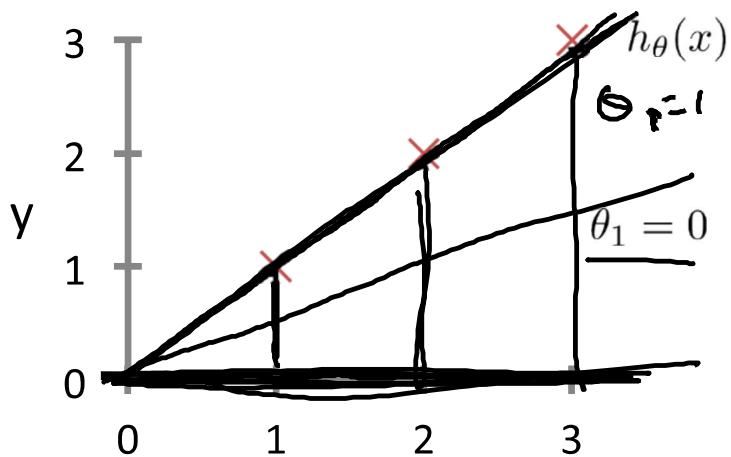
(function of the parameter  $\theta_1$ )



$$\Theta_1 = 0^2$$

$$J(0) = ?$$

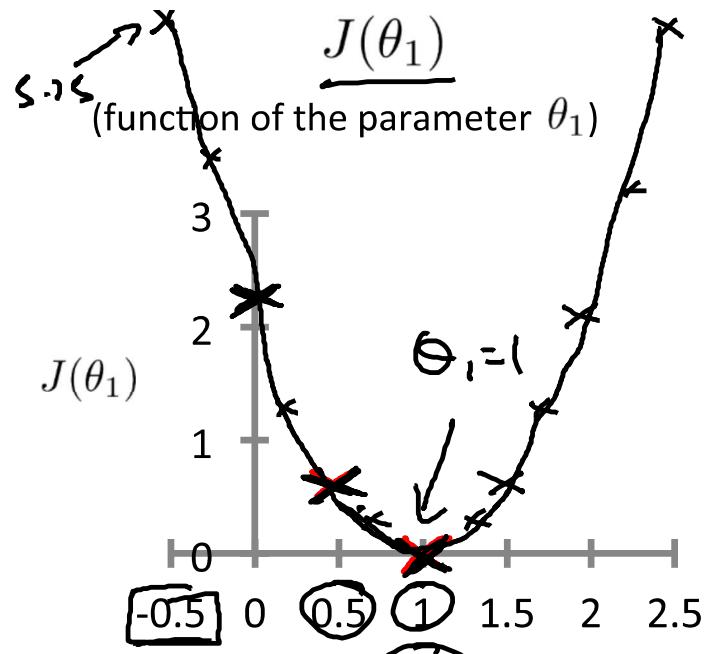
$h_{\theta}(x)$   
(for fixed  $\theta_1$ , this is a function of  $x$ )



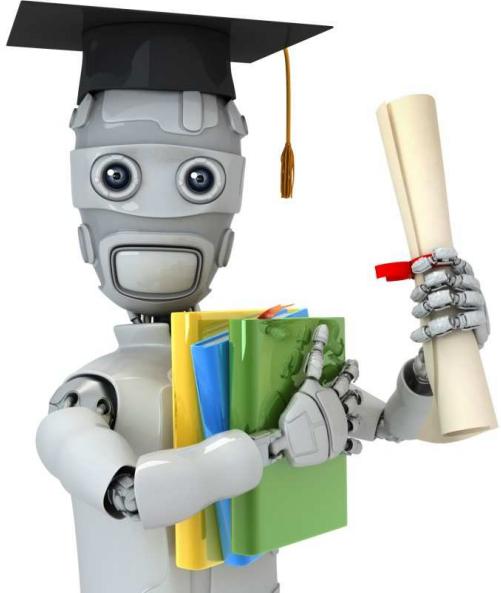
$$\mathcal{J}(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2) \\ = \frac{1}{6} \cdot 14 \approx 2.3$$

$$h(x) = -0.5x$$

minimize  $\mathcal{J}(\theta_1)$



Andrew |



# Linear regression with one variable

---

## Cost function intuition II

Machine Learning

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

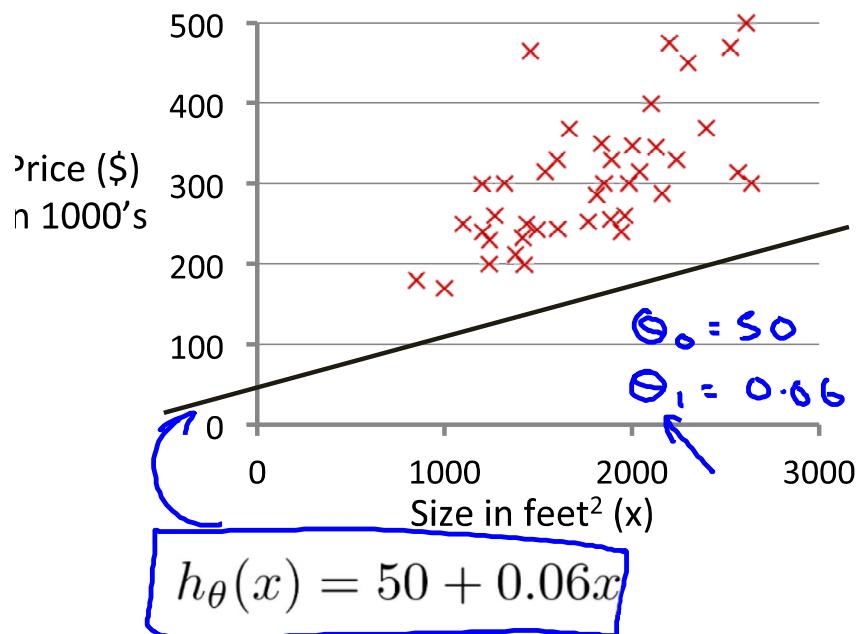
Parameters:  $\underline{\theta_0, \theta_1}$

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

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

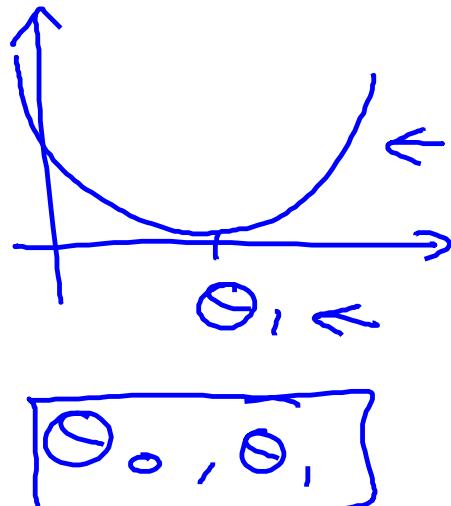
$$\underline{h_{\theta}(x)}$$

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



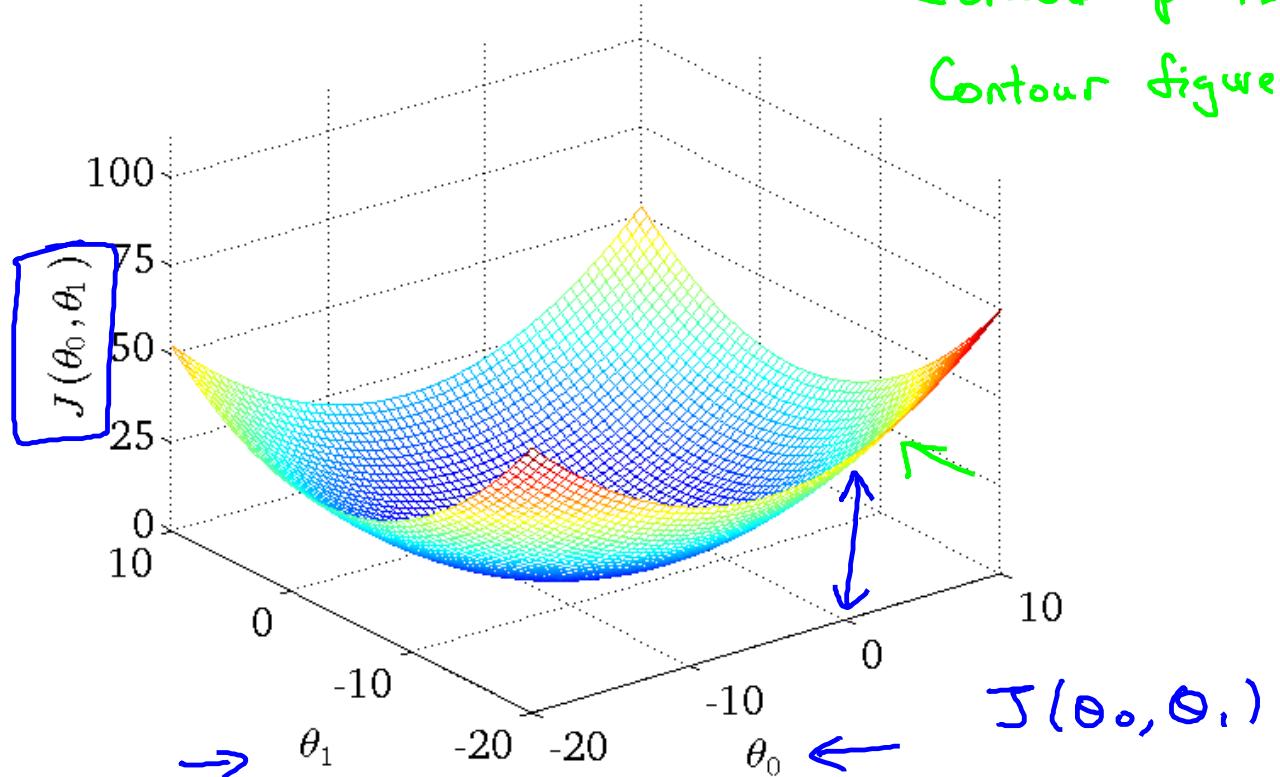
$$\underline{J(\theta_0, \theta_1)}$$

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



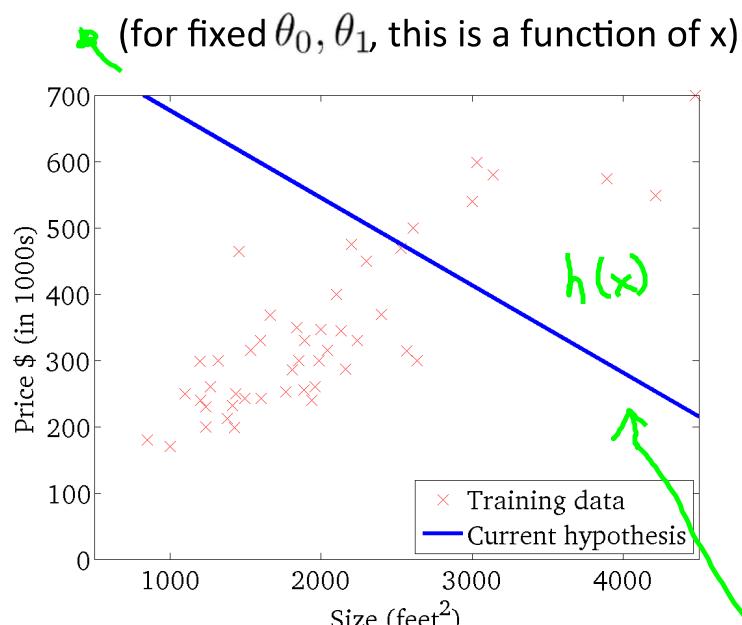
Contour plots

Contour figures -

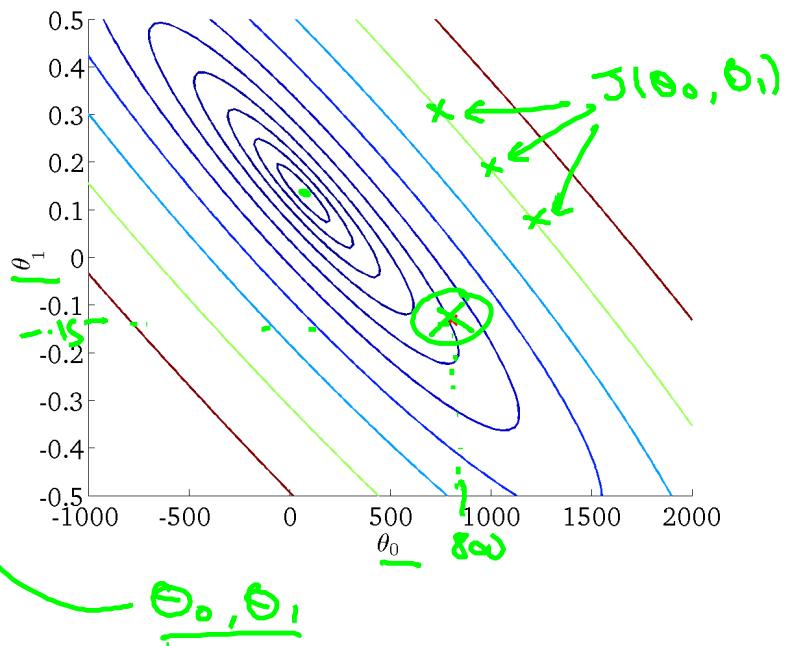


$$h_{\theta}(x)$$

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

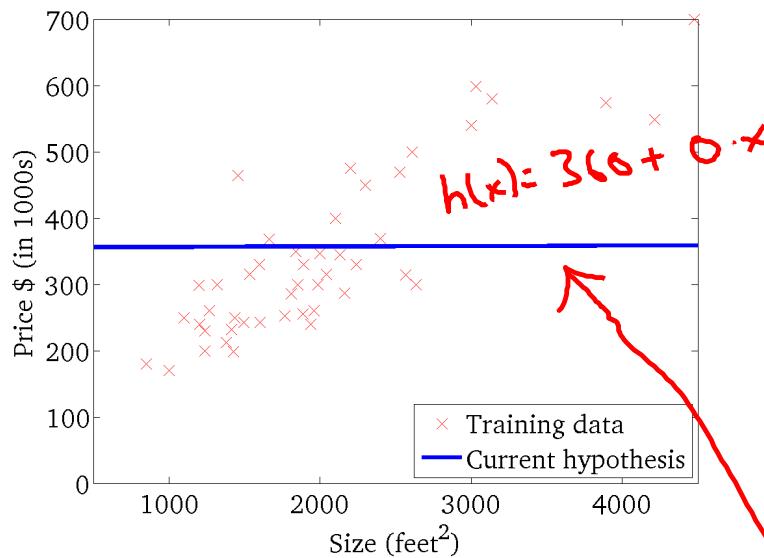


(function of the parameters  $\theta_0, \theta_1$ )



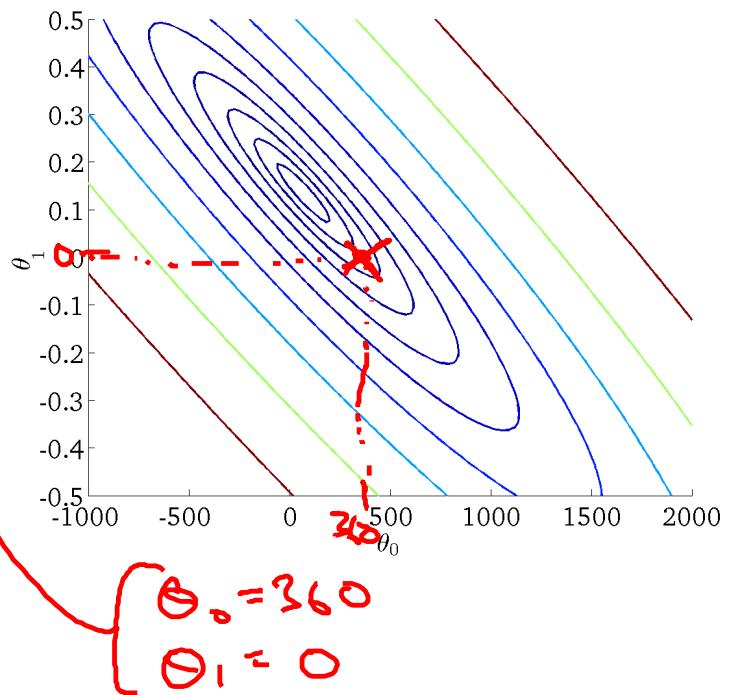
$$h_{\theta}(x)$$

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



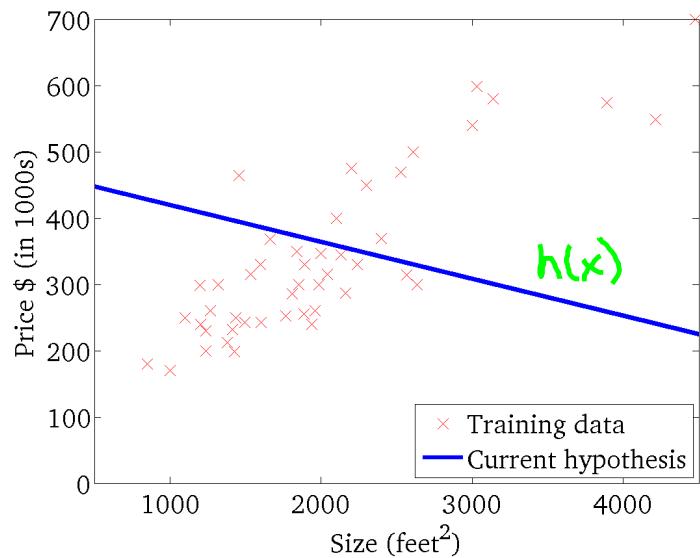
$$J(\theta_0, \theta_1)$$

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



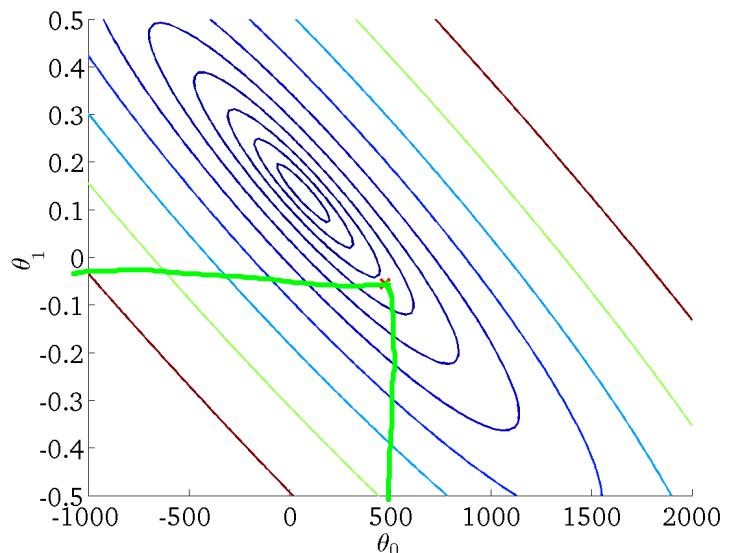
$$h_{\theta}(x)$$

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



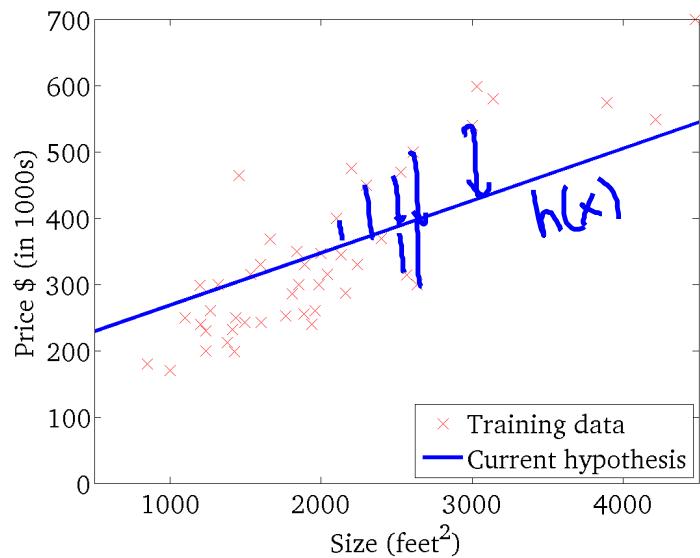
$$J(\theta_0, \theta_1)$$

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



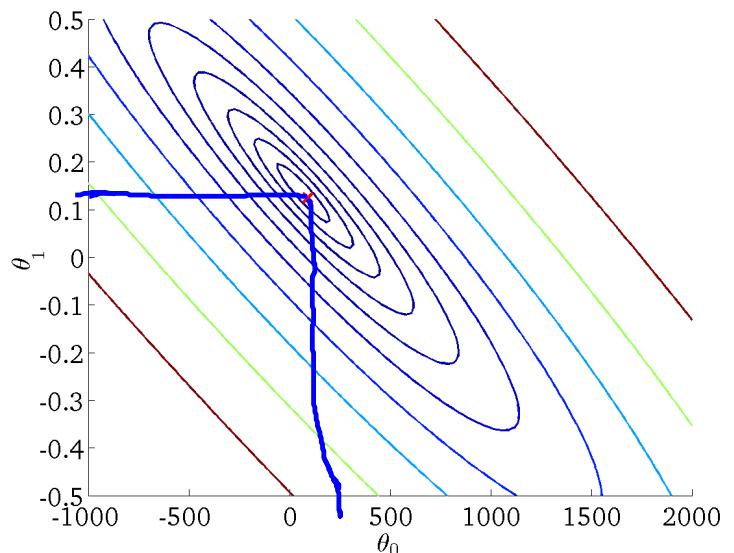
$$h_{\theta}(x)$$

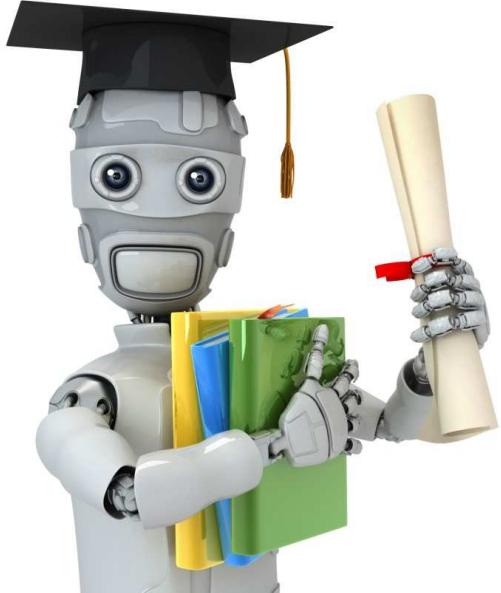
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

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





Linear regression  
with one variable

---

# Gradient descent

Machine Learning

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

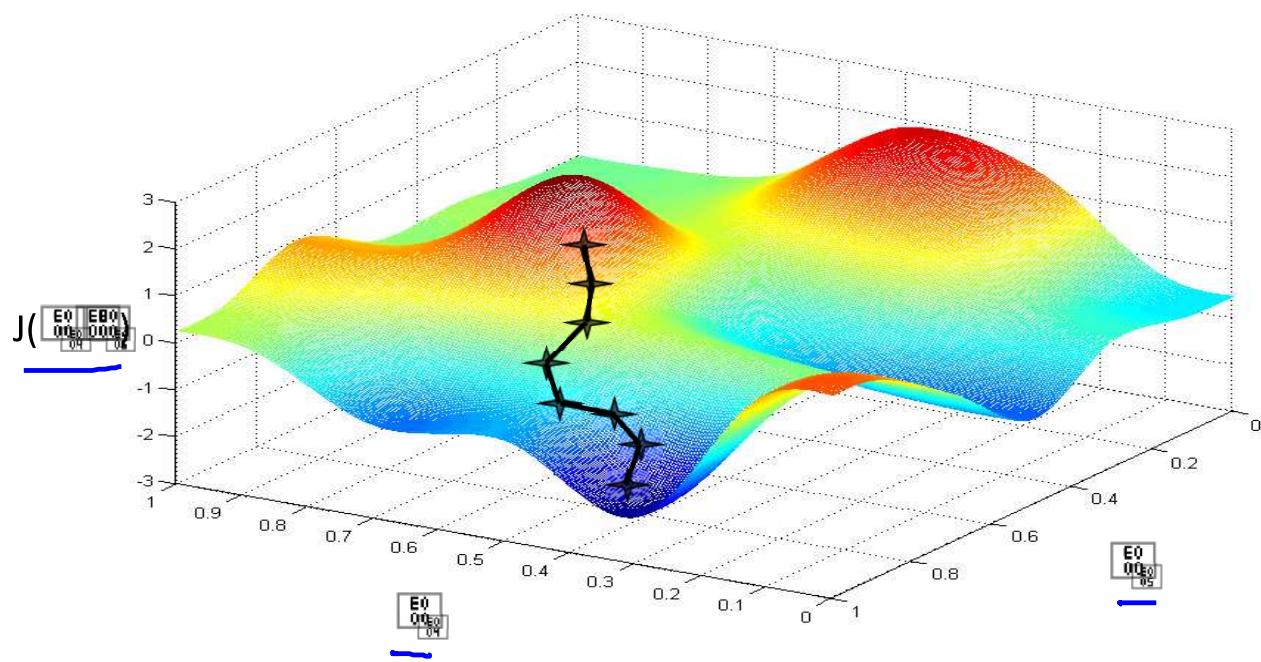
Want  $\min_{\theta_0, \theta_1} \underline{J(\theta_0, \theta_1)}$   $\min_{\theta_0, \dots, \theta_n} \underline{J(\theta_0, \dots, \theta_n)}$

## Outline:

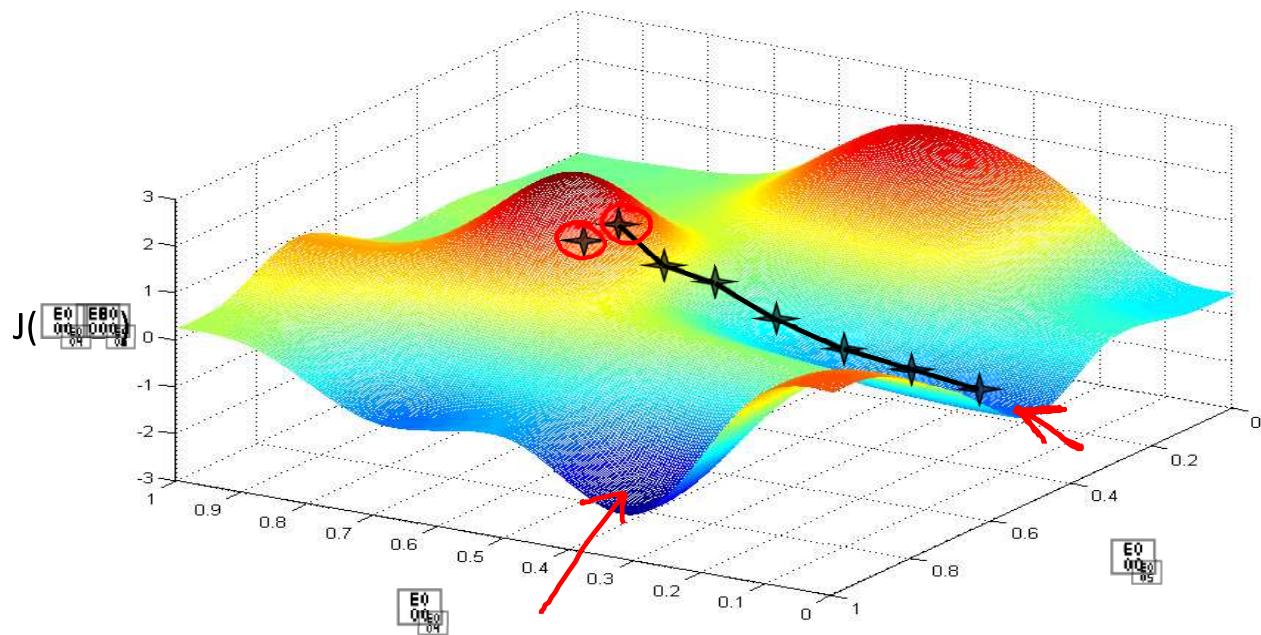
Start with some  $\underline{\theta_0, \theta_1}$  (say  $\theta_0 = 0, \theta_1 = 0$ )

Keep changing  $\underline{\theta_0, \theta_1}$  to reduce  $\underline{J(\theta_0, \theta_1)}$

until we hopefully end up at a minimum



Andrew |



Andrew |

# Gradient descent algorithm

repeat until convergence {  
 } →  $\theta_0, \theta_1$   
 } →  $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$   
 } → learning rate

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

Simultaneously update  
 $\theta_0$  and  $\theta_1$

→  $a := b$   
 $a := a + 1$

→  $a = b$  ←  
 $a = a + 1$  X

## Correct: Simultaneous update

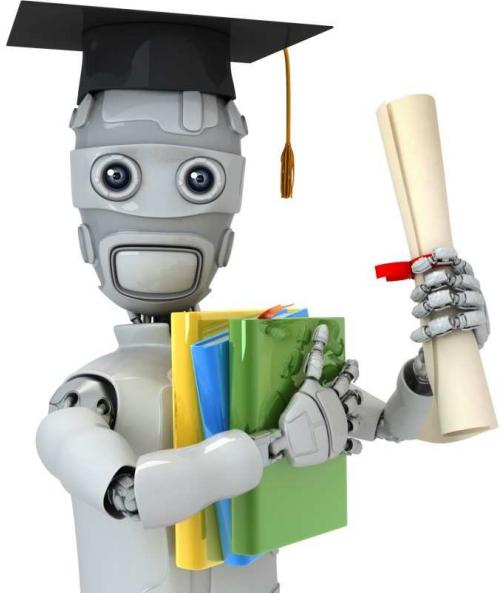
→  $\text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
 →  $\text{temp1} := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 →  $\theta_0 := \text{temp0}$   
 →  $\theta_1 := \text{temp1}$



## Incorrect:

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





# Linear regression with one variable

---

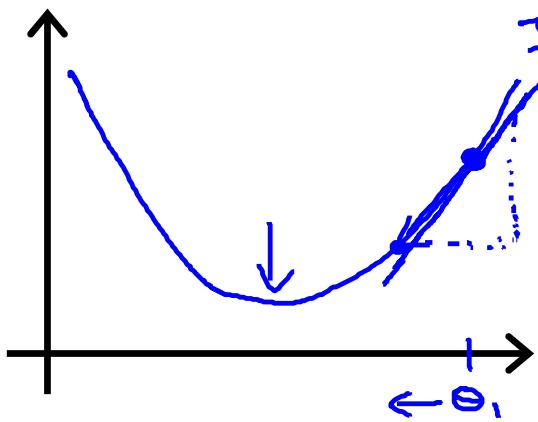
## Gradient descent intuition

Machine Learning

# Gradient descent algorithm

repeat until convergence {  
     $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$       (simultaneously update  
    }  
         $\theta_0$  and  $\theta_1$ )  
    }      learning rate  
                    derivative

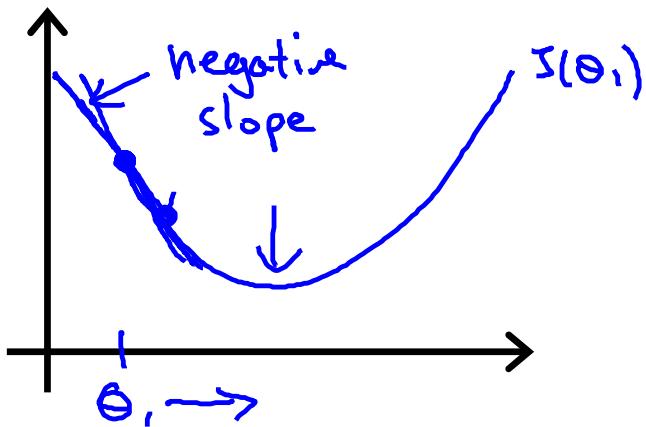
$$\min_{\theta_1} J(\theta_1) \quad \theta_1 \in \mathbb{R}$$



$(\theta_1 \in \mathbb{R})$

$$\theta_1 := \theta_1 - \frac{\alpha}{\frac{d}{d\theta_1} J(\theta_1)} \geq 0$$

$$\theta_1 := \theta_1 - \underline{\alpha} \cdot (\text{positive number})$$

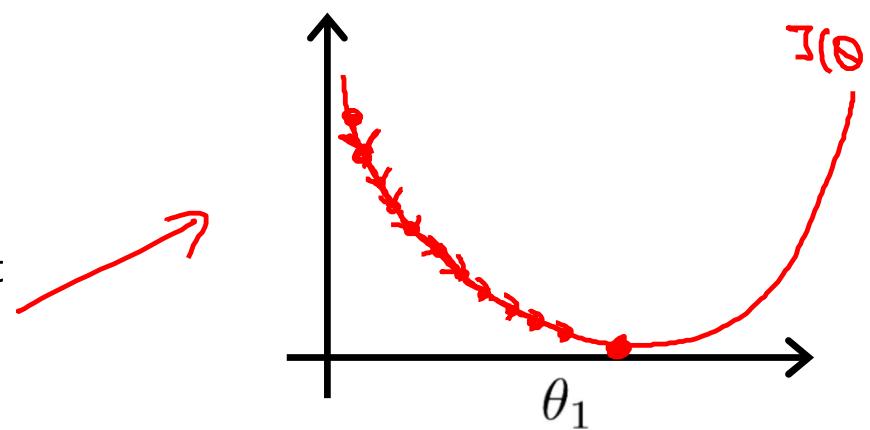


$$\frac{\frac{\partial}{\partial \theta_1} J(\theta_1)}{\leq 0}$$

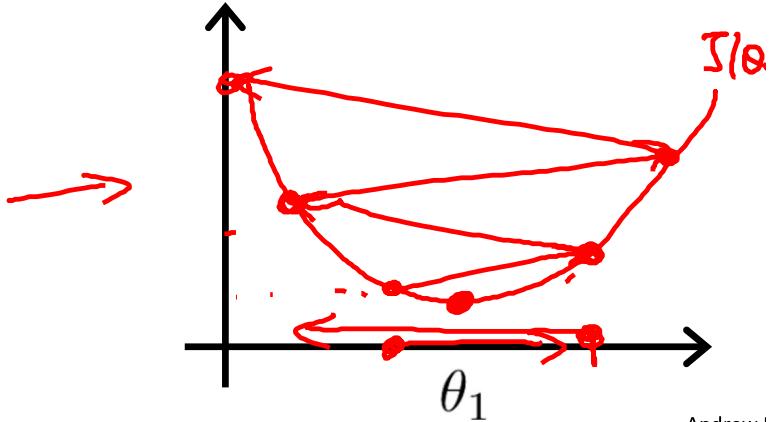
$$\theta_1 := \theta_1 - \underline{\alpha} \uparrow (\text{negative number})$$

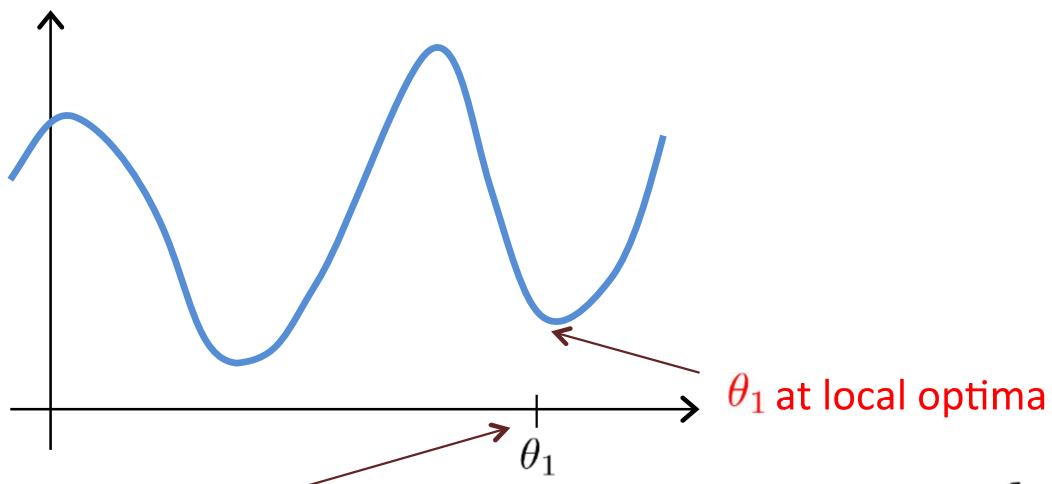
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

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.



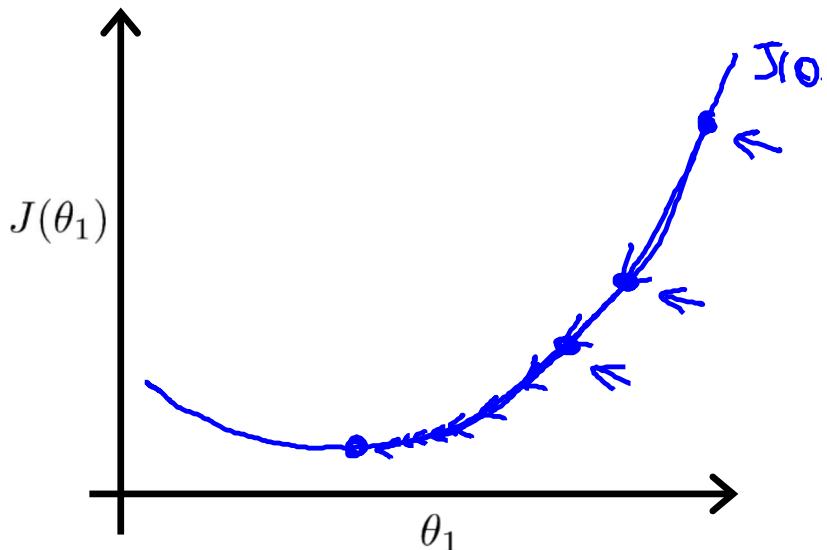


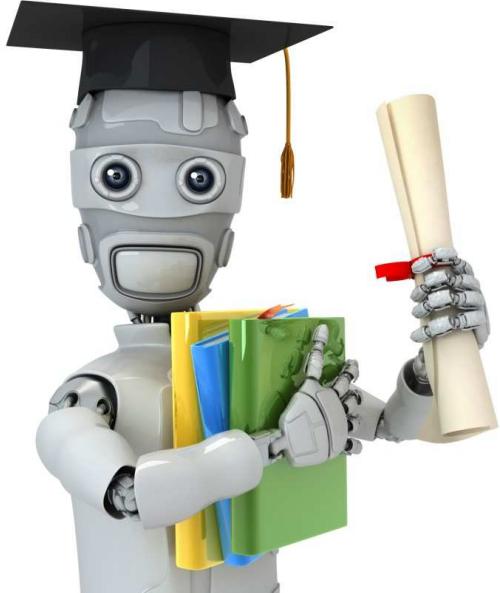
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

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.





# Linear regression with one variable

---

## Gradient descent for linear regression

Machine Learning

## 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

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

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

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

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

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

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

## Gradient descent algorithm

repeat until convergence {

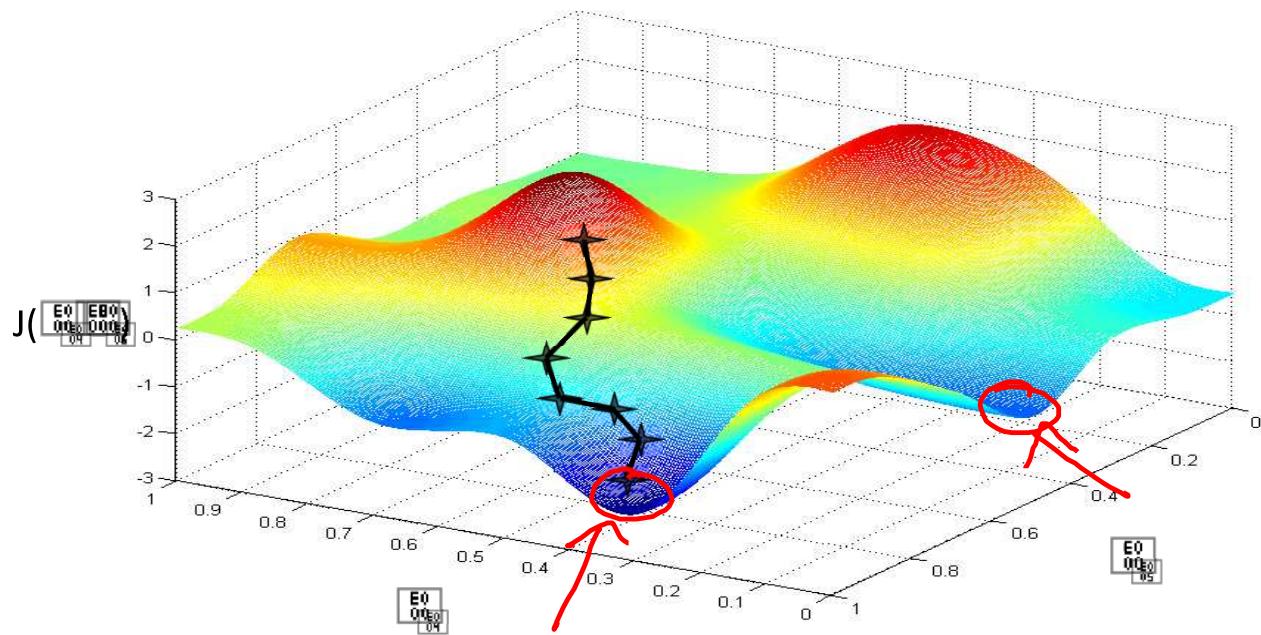
$$\theta_0 := \theta_0 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \right]$$
$$\theta_1 := \theta_1 - \alpha \left[ \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right]$$

}

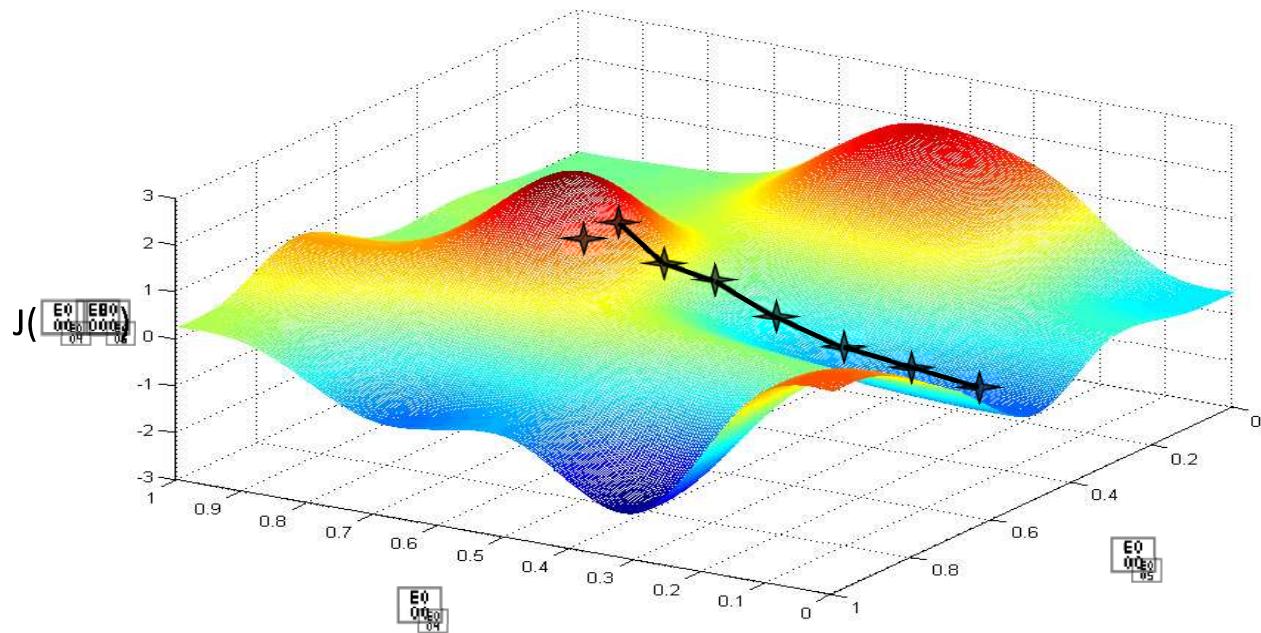
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

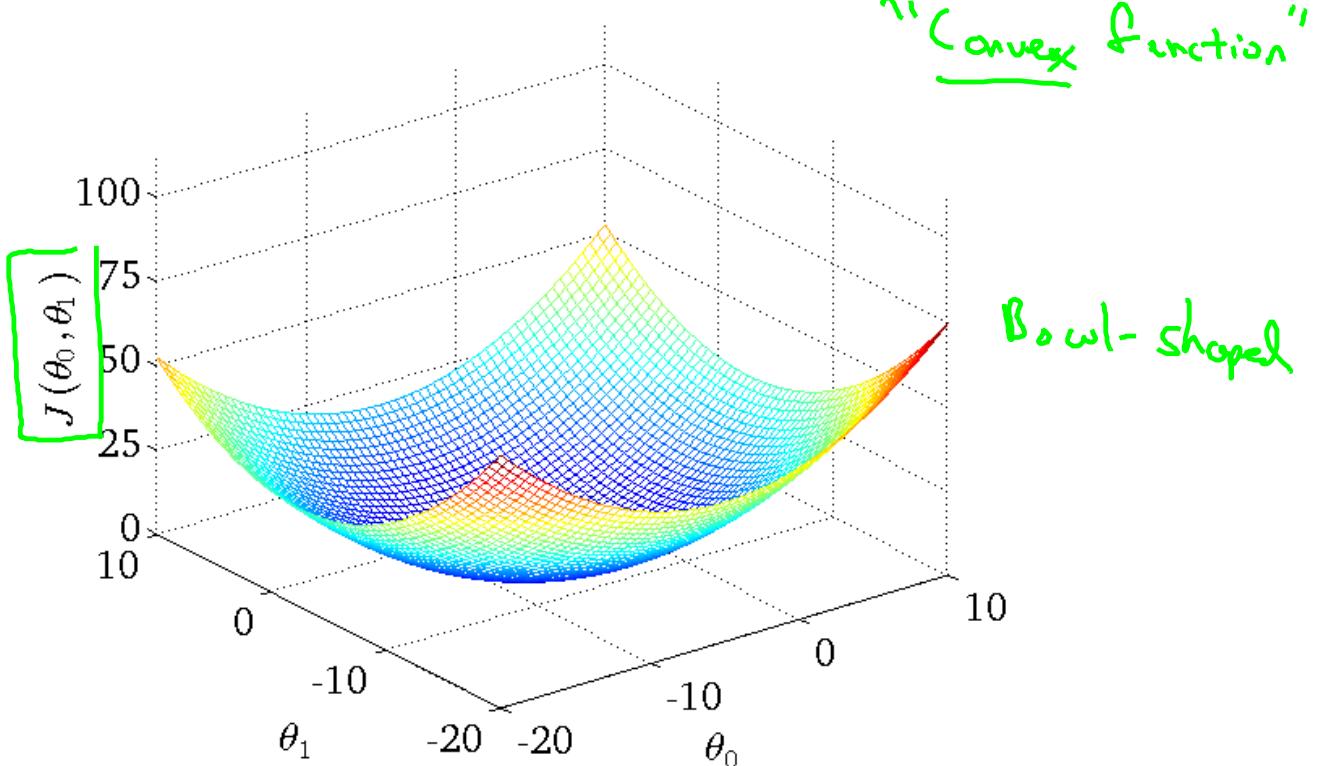
$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$



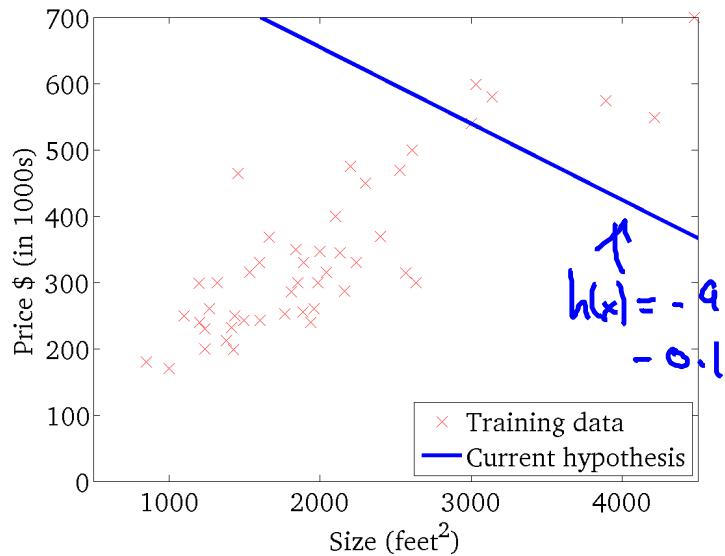
Andrew |



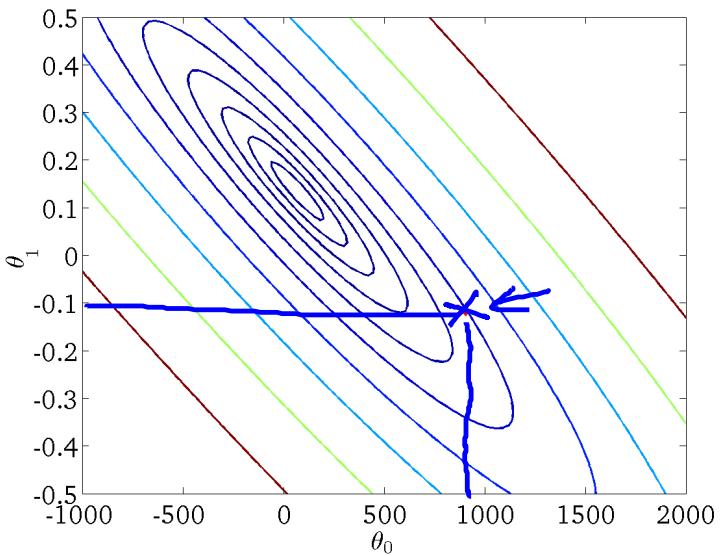
Andrew |



$h_{\theta}(x)$   
 (for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )

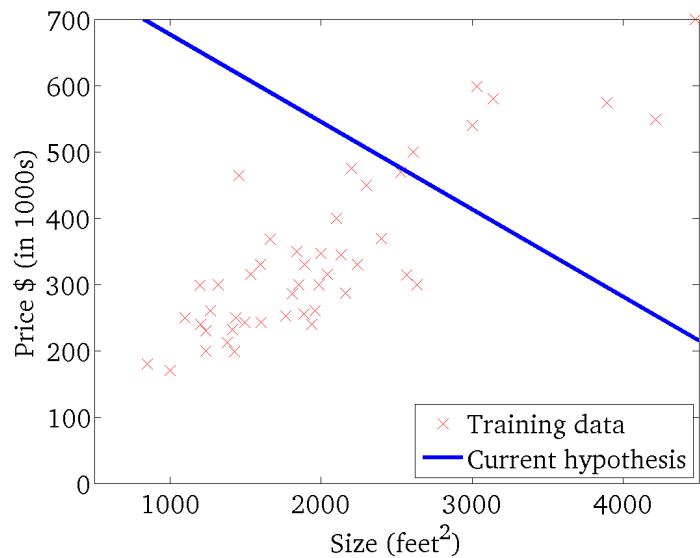


$J(\theta_0, \theta_1)$   
 (function of the parameters  $\theta_0, \theta_1$ )



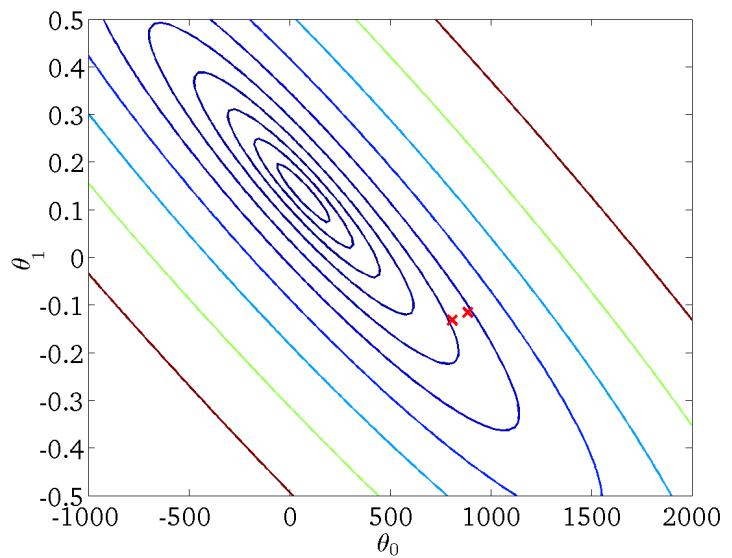
$$h_{\theta}(x)$$

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



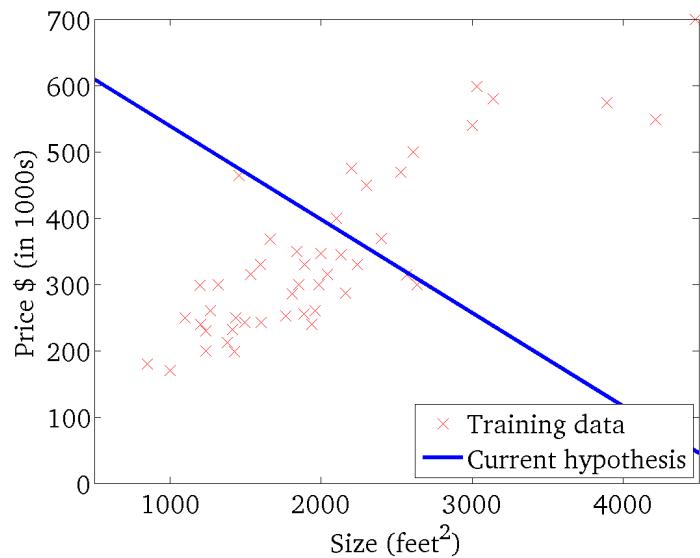
$$J(\theta_0, \theta_1)$$

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



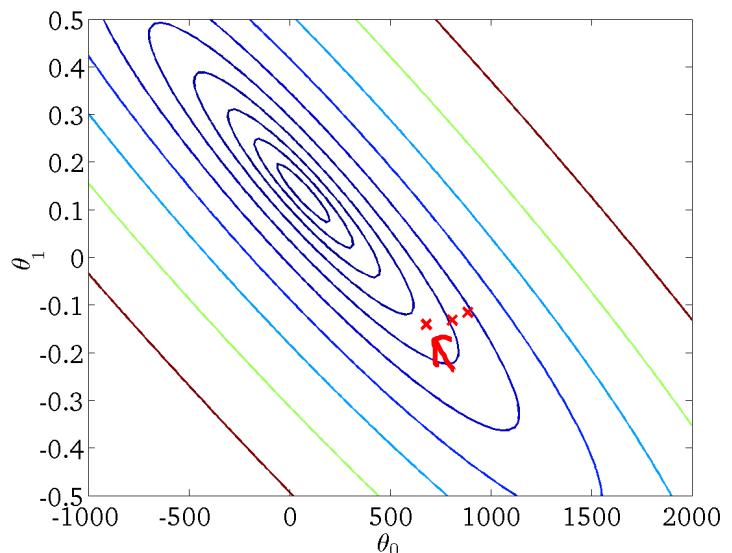
$$h_{\theta}(x)$$

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



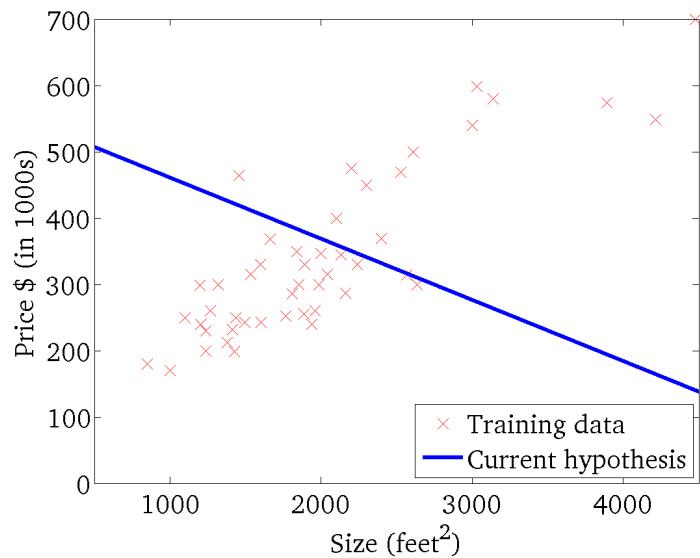
$$J(\theta_0, \theta_1)$$

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



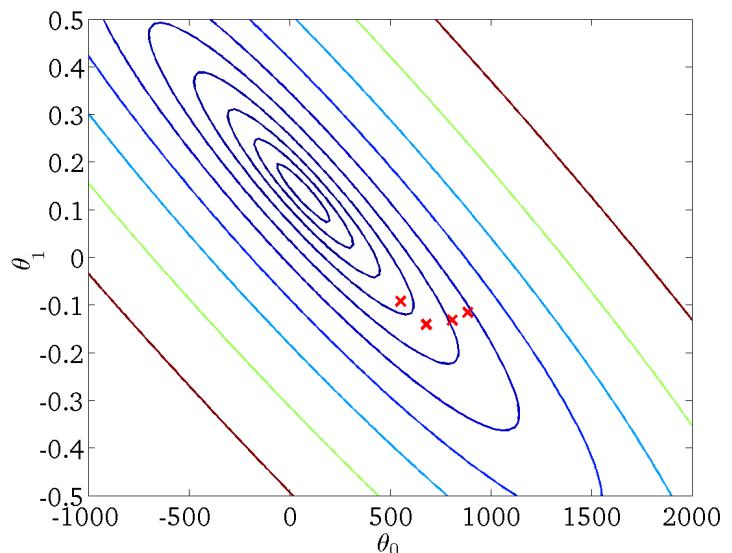
$$h_{\theta}(x)$$

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



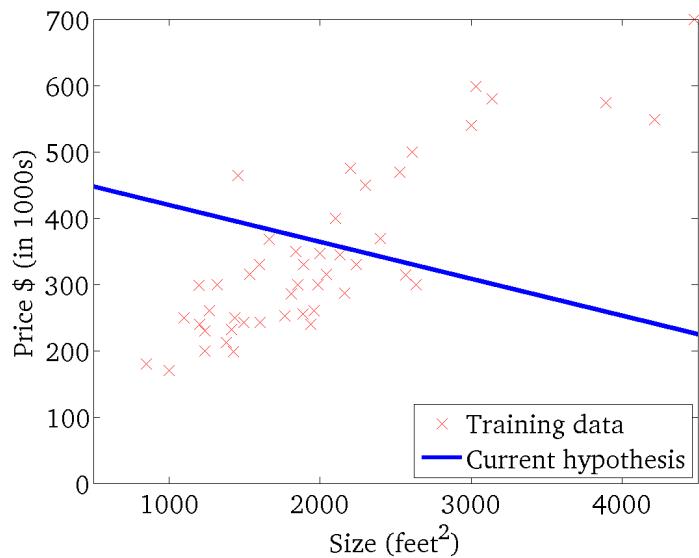
$$J(\theta_0, \theta_1)$$

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



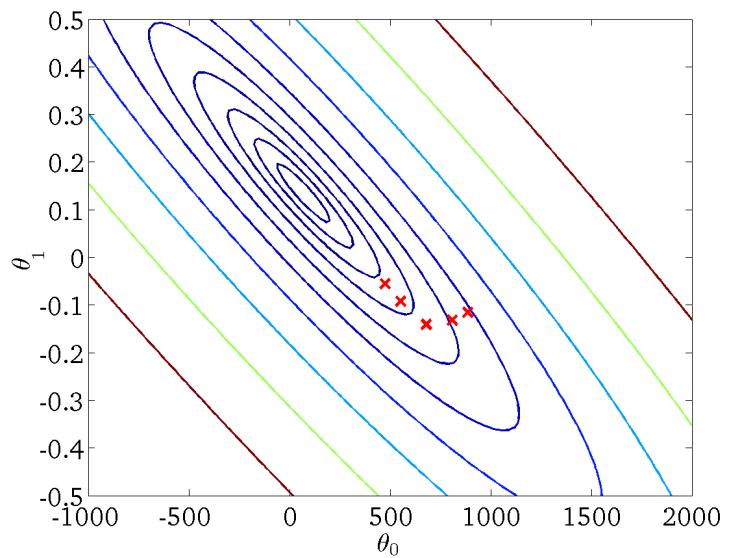
$$h_{\theta}(x)$$

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



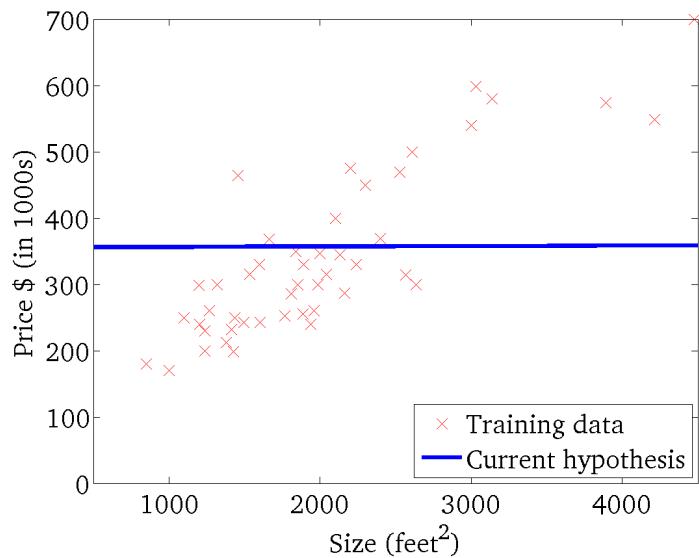
$$J(\theta_0, \theta_1)$$

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



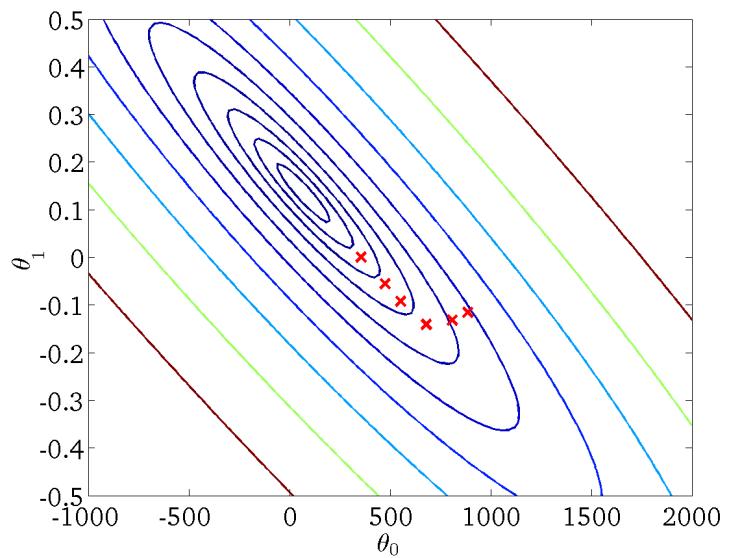
$$h_{\theta}(x)$$

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



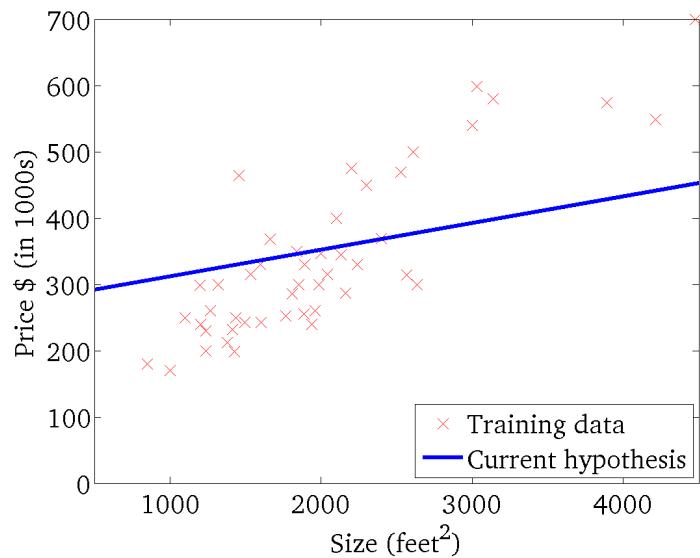
$$J(\theta_0, \theta_1)$$

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



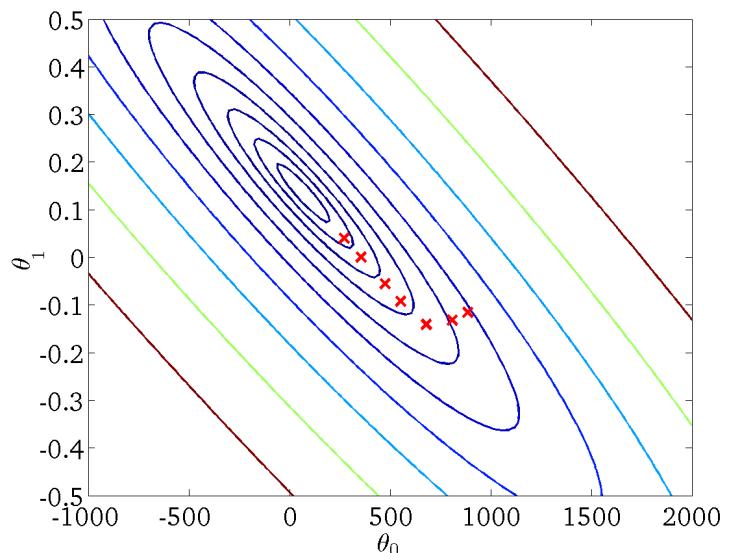
$$h_{\theta}(x)$$

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



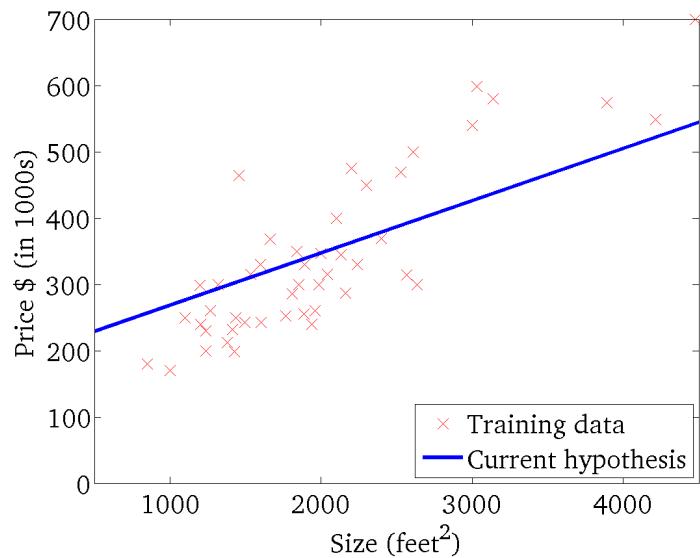
$$J(\theta_0, \theta_1)$$

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



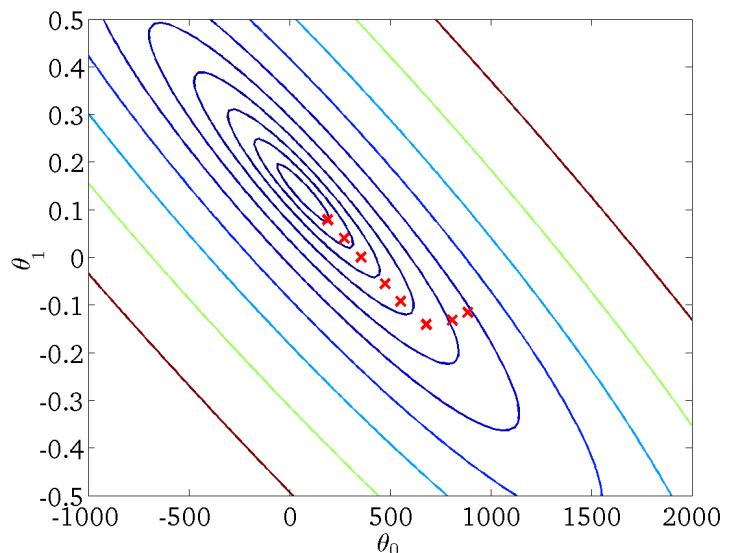
$$h_{\theta}(x)$$

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



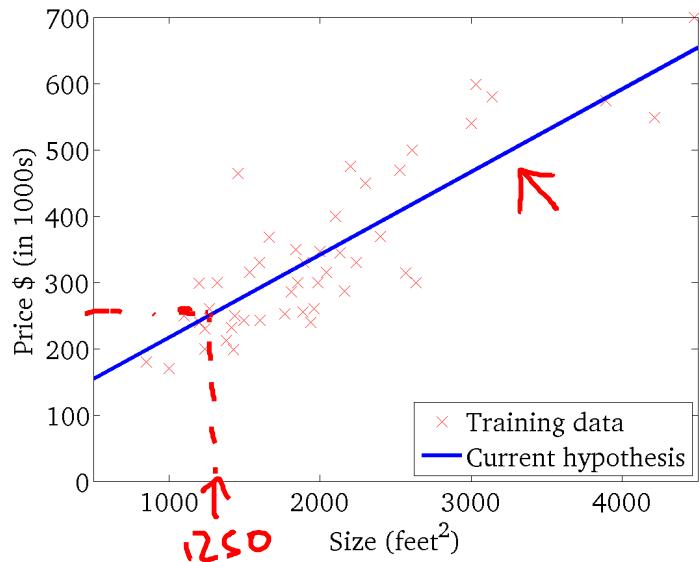
$$J(\theta_0, \theta_1)$$

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



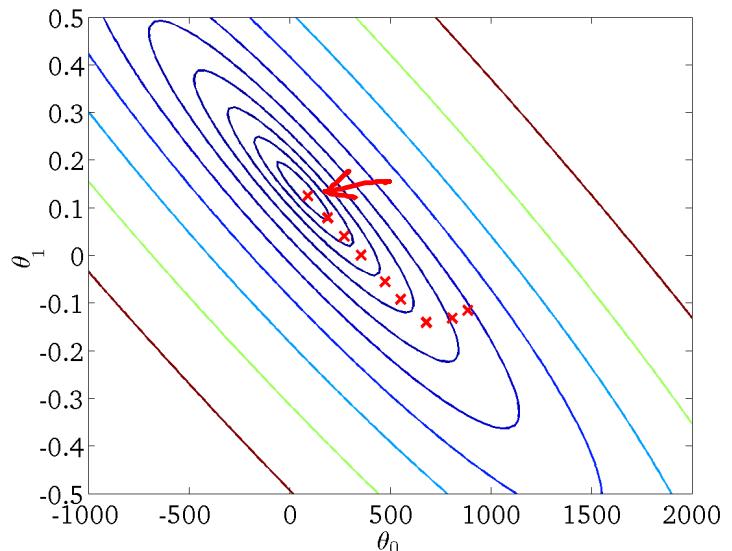
$$h_{\theta}(x)$$

(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.

$$\rightarrow \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})$$