

CS1138

# Machine Learning

Lecture : Gradient Descent

(Slide Credits: Andrew Ng)

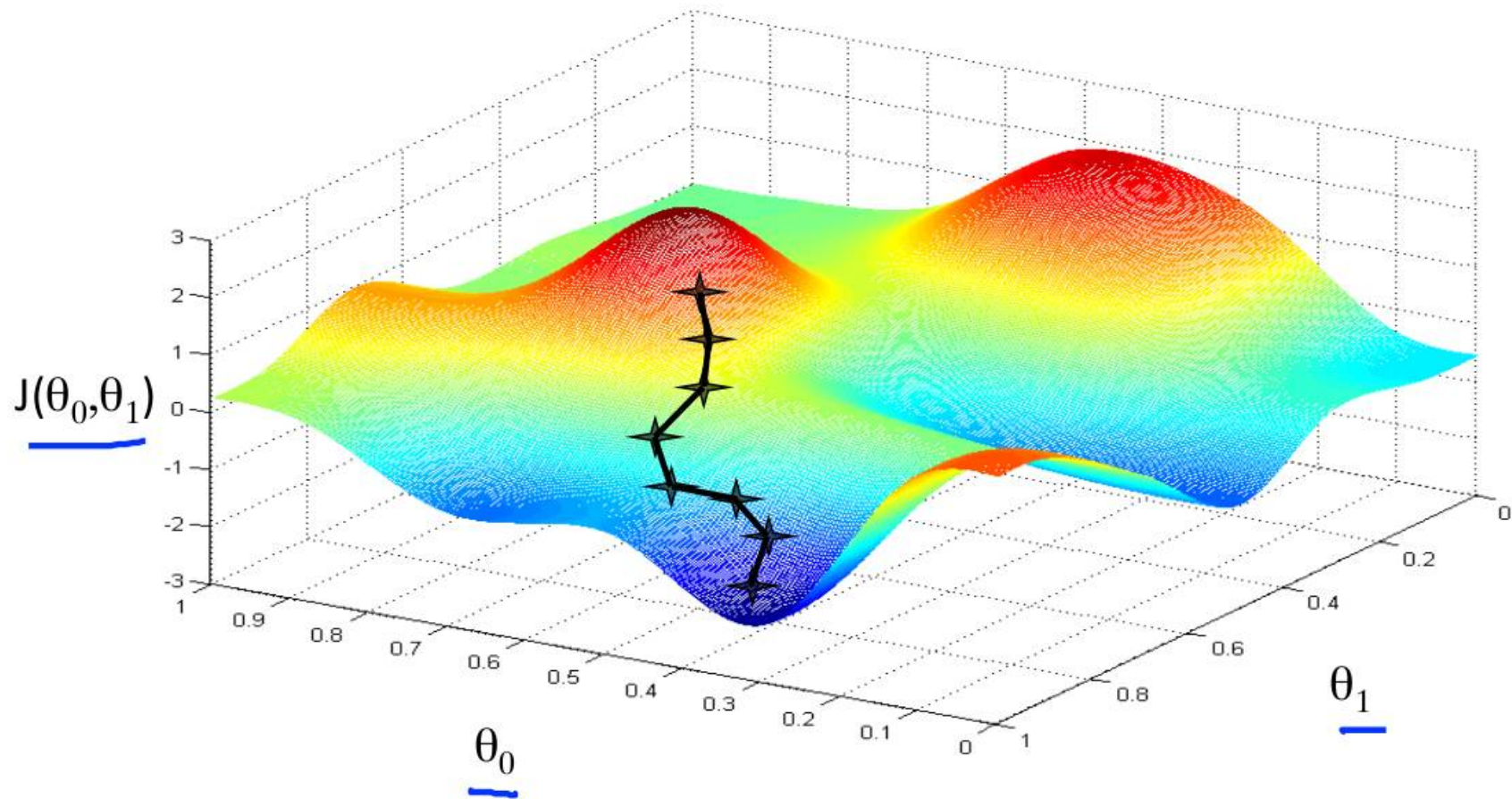
Arpan Gupta

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

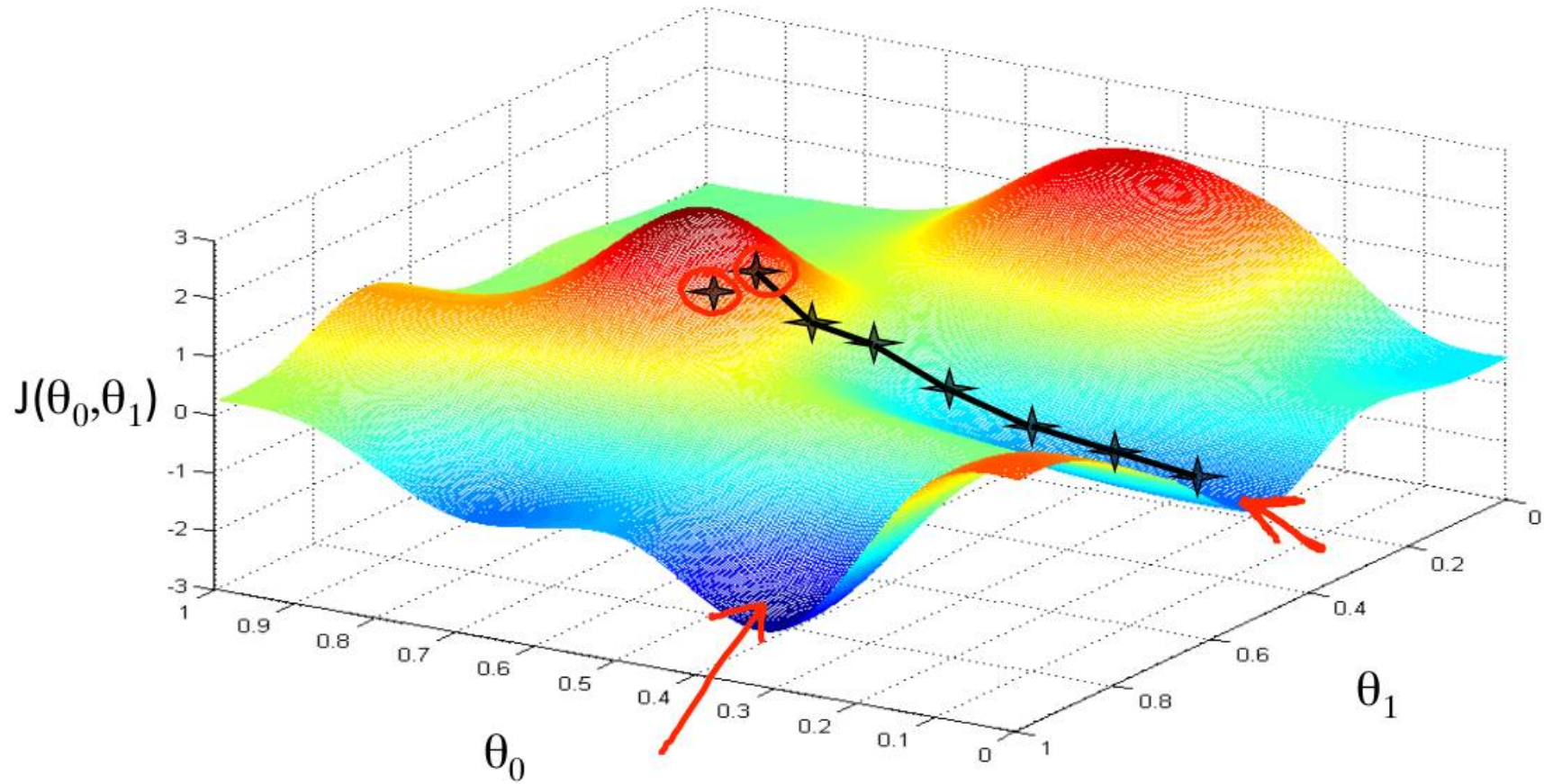
Want  $\min_{\theta_0, \theta_1} J(\theta_0, \theta_1)$   $\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$

## 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



Case: When gradient descent reaches global minima



Case: When gradient descent does not reach global minima

# Gradient descent algorithm

repeat until convergence {

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

}

---

Correct: Simultaneous update

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

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

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

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



# Gradient descent algorithm

$\theta_0, \theta_1$

Assignment

$a := b$   
 $a := a + 1$

Truth assertion

$a = b$

$a = a + 1$  ✗

repeat until convergence {

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

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

Simultaneously update  
 $\theta_0$  and  $\theta_1$

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 {

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

(simultaneously update  
 $j = 0$  and  $j = 1$ )

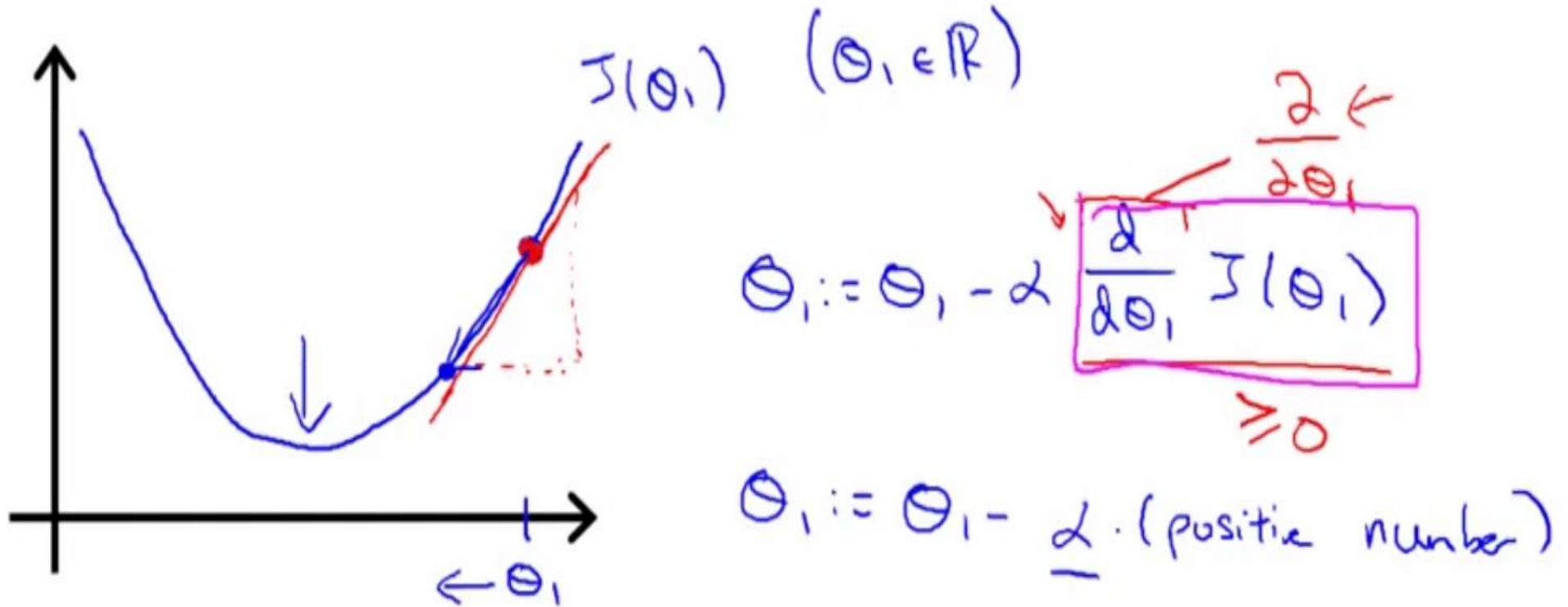
learning  
rate

derivative

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

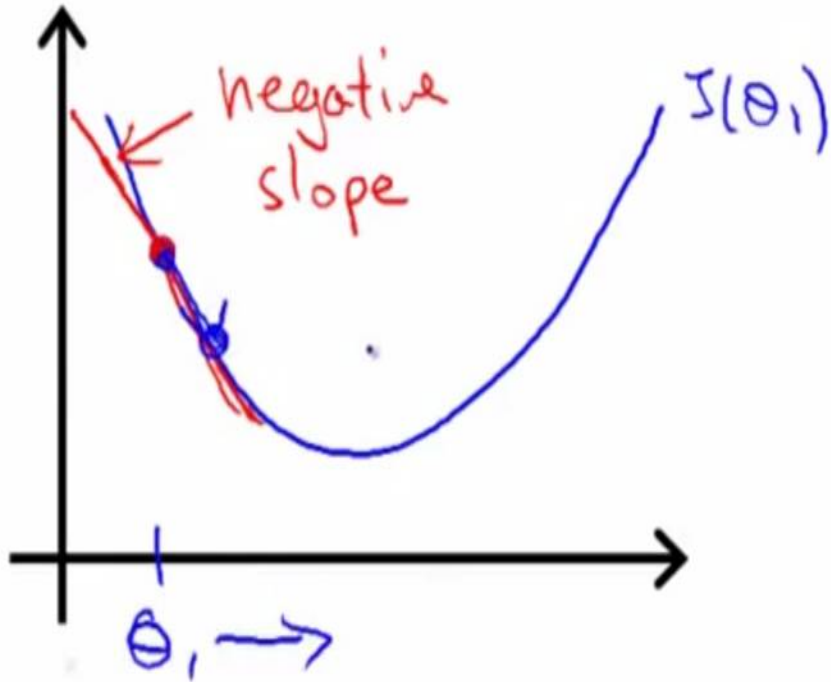
$$\theta_1 \in \mathbb{R}.$$

Assuming that  $J$  is a function of single parameter  $\Theta_1$   
 When slope is positive (move left).





Assuming that  $J$  is a function of single parameter  $\Theta_1$   
When slope is negative (move right).



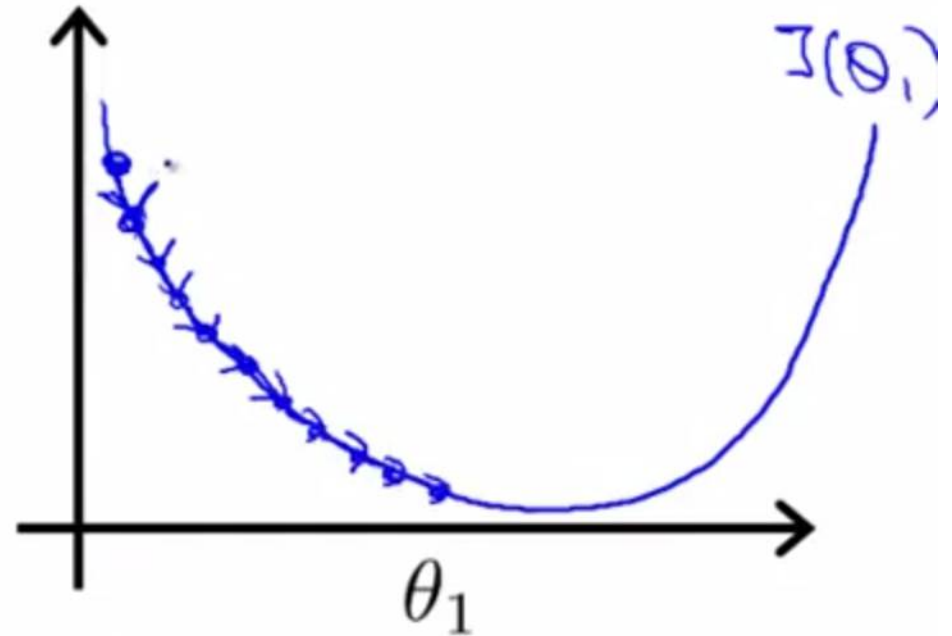
$$\frac{\frac{\partial}{\partial \theta_1} J(\theta_1)}{\leq 0}$$
$$\theta_1 := \theta_1 - \alpha (\text{negative number})$$

Andrew Ng

## Effect of Learning Rate ( $\alpha$ )

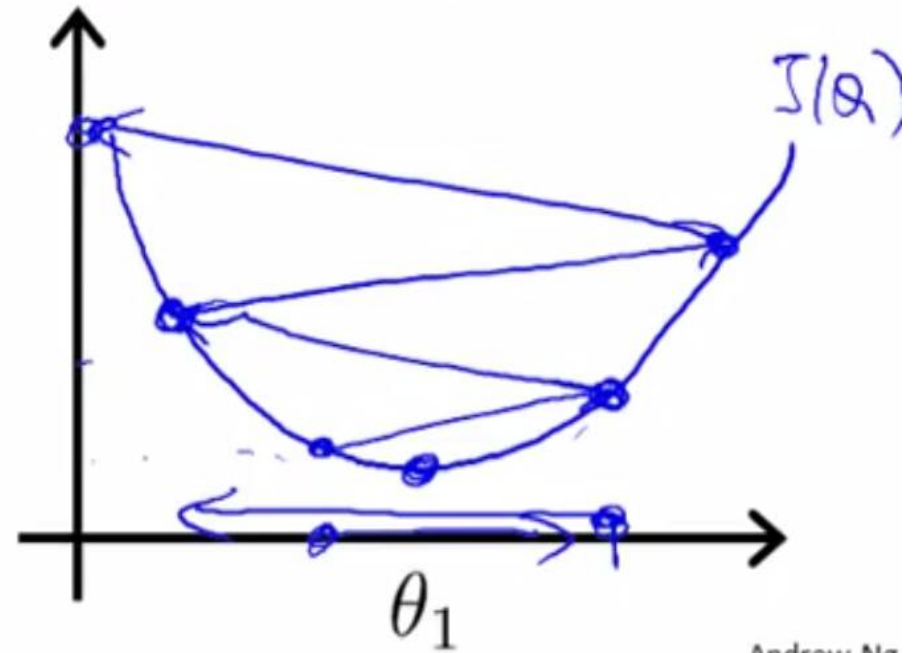
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

If  $\alpha$  is too small, gradient descent can be slow.



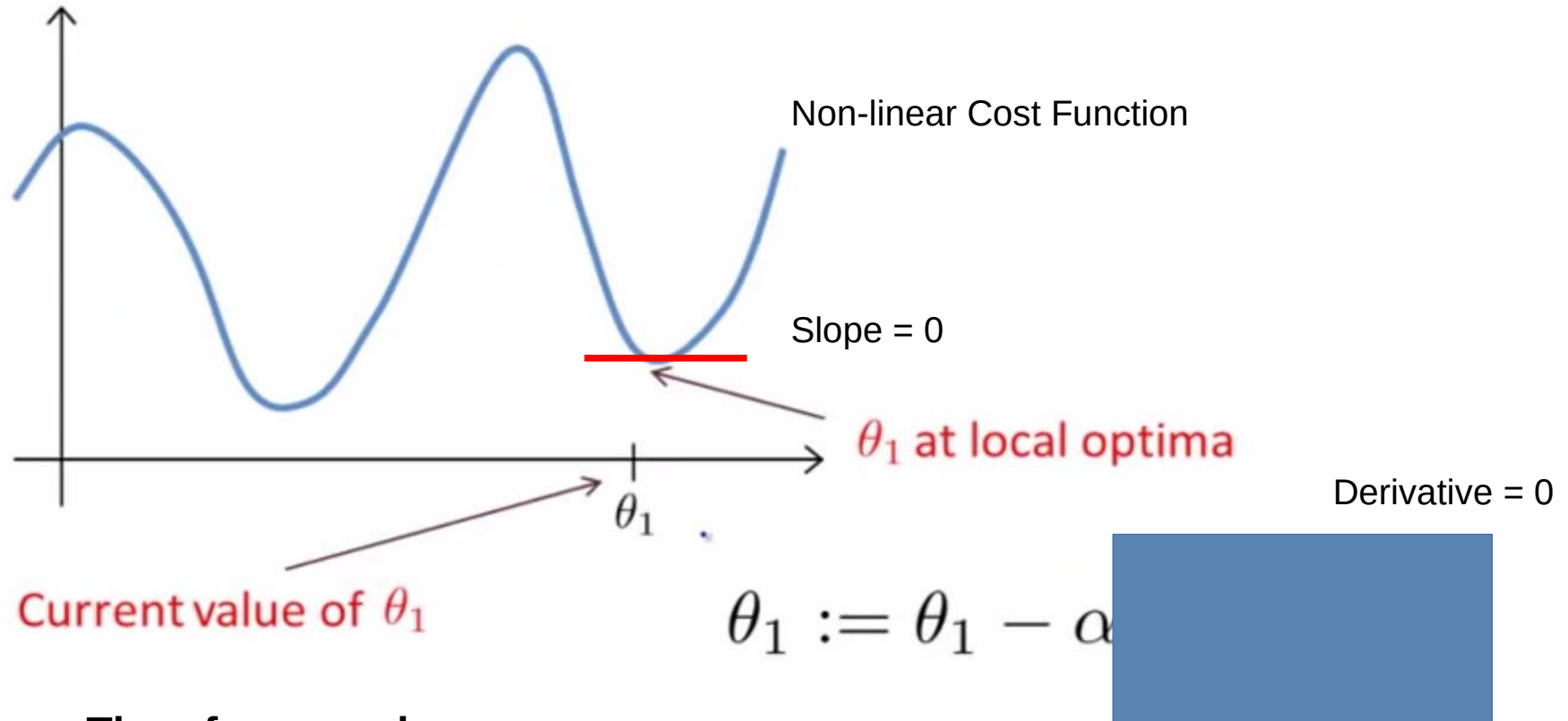
## Effect of Learning Rate ( $\alpha$ )

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



Andrew Ng

What will one step of gradient descent do if  $\Theta_1$  is already at the local minima?

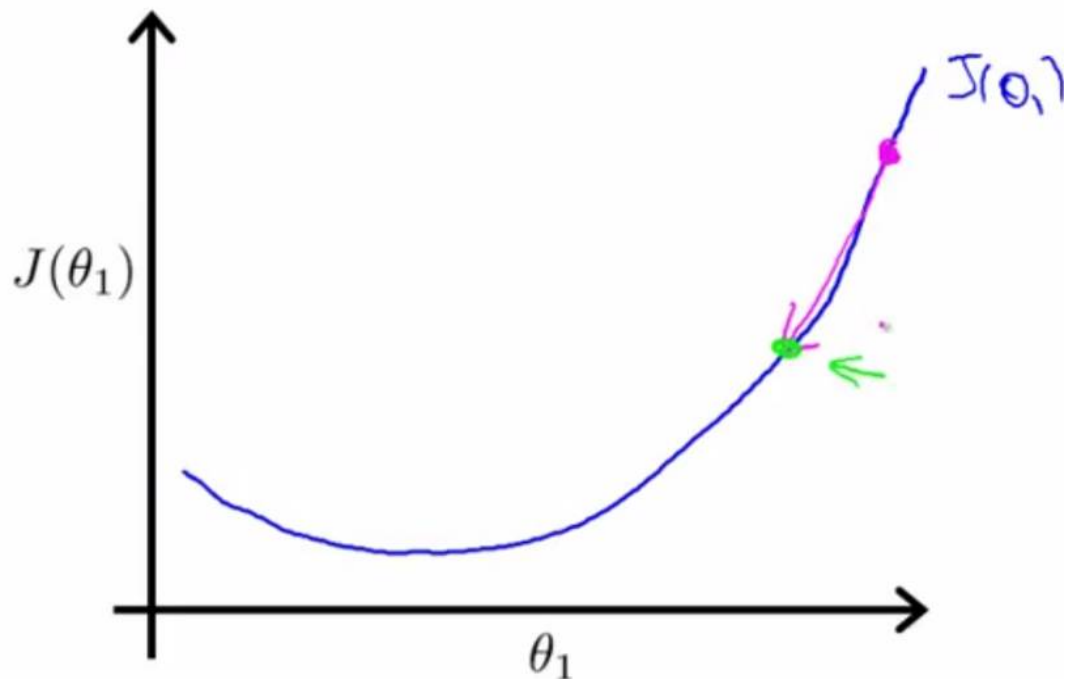


Therefore, no change.

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.

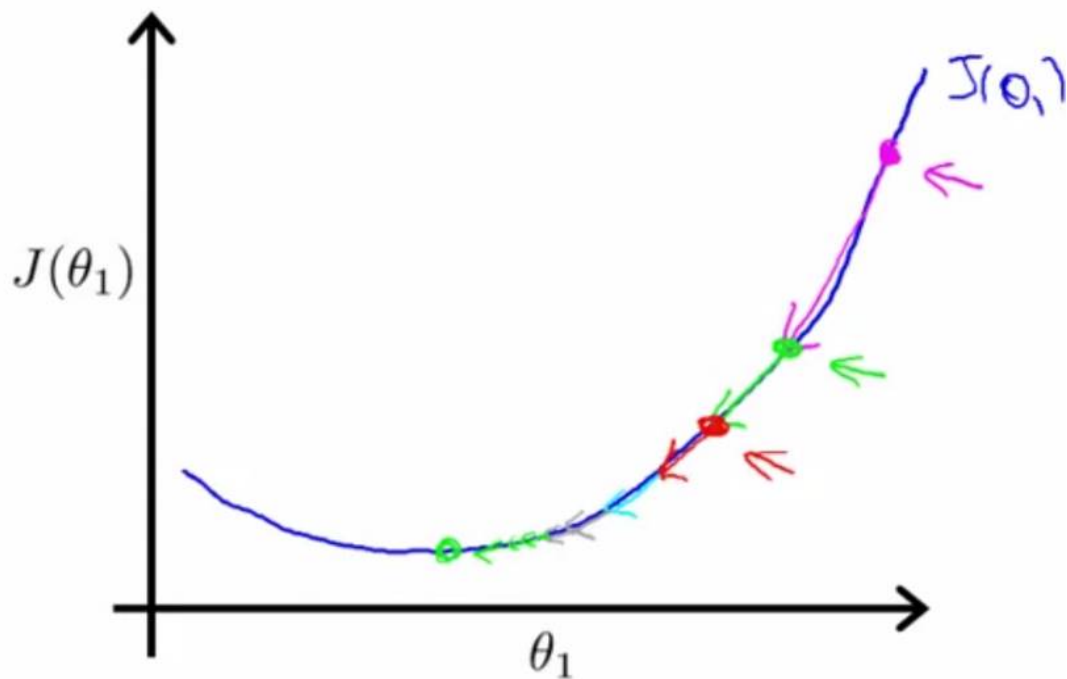




Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.

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

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



# Gradient Descent for linear regression

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

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

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

# Gradient descent algorithm

repeat until convergence {

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

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

}

## Gradient descent algorithm

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

repeat until convergence {

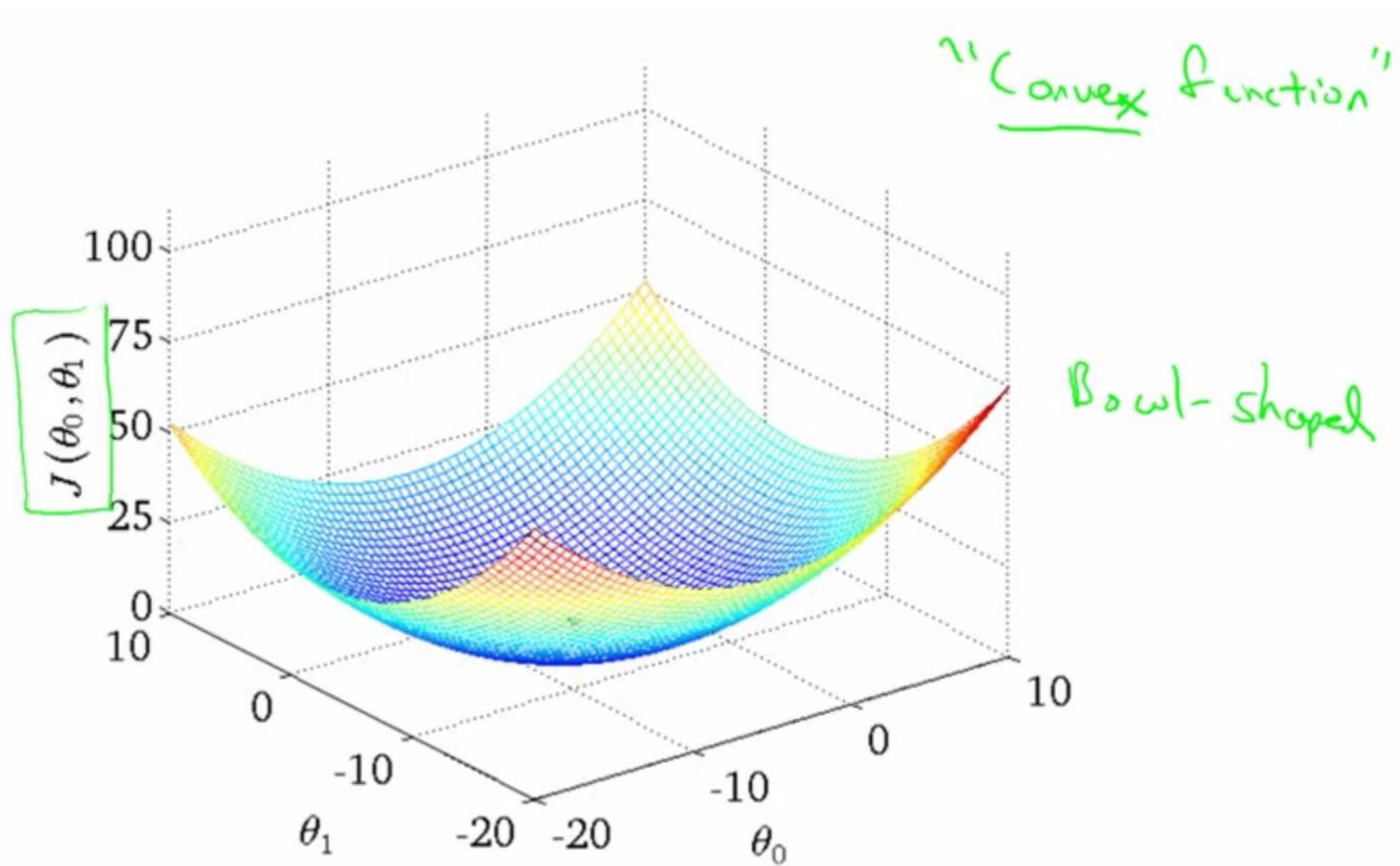
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

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

}

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

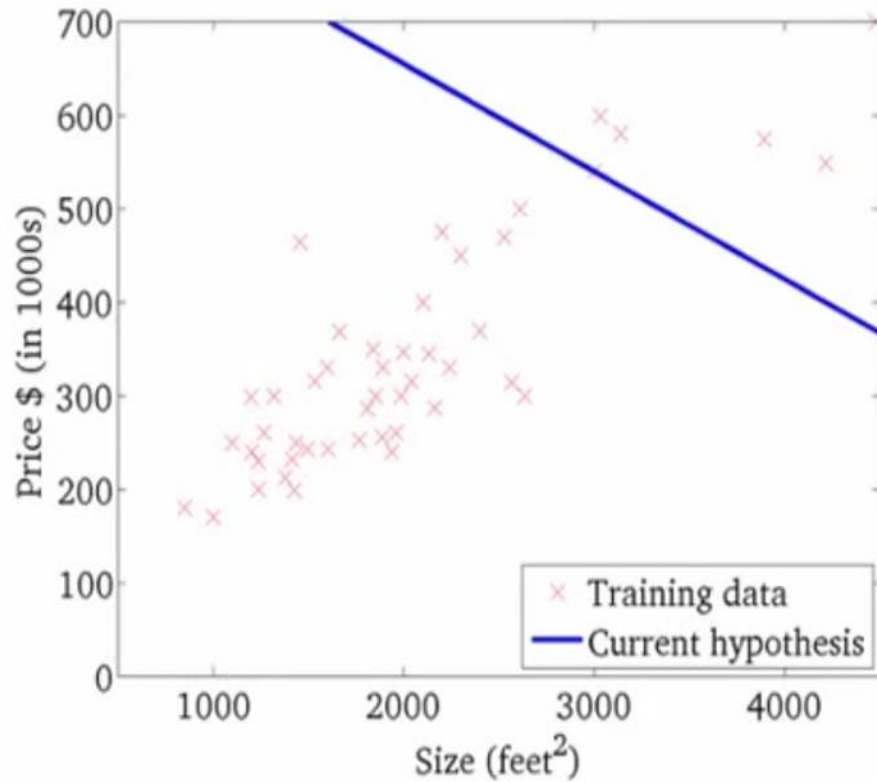
$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$$





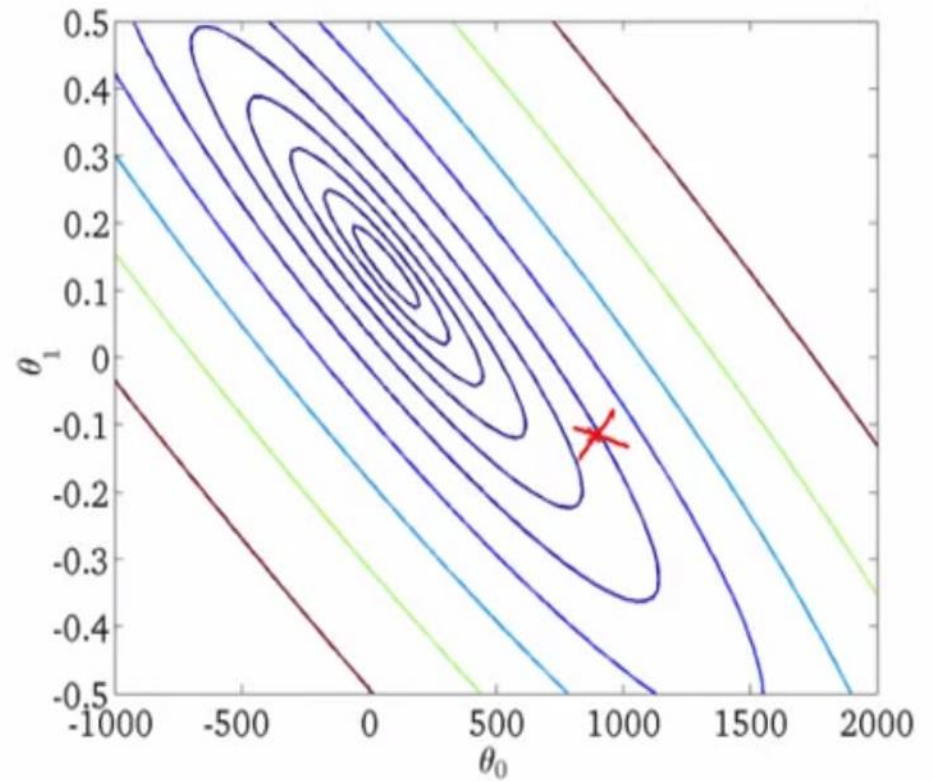
$$\underline{h_{\theta}(x)}$$

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



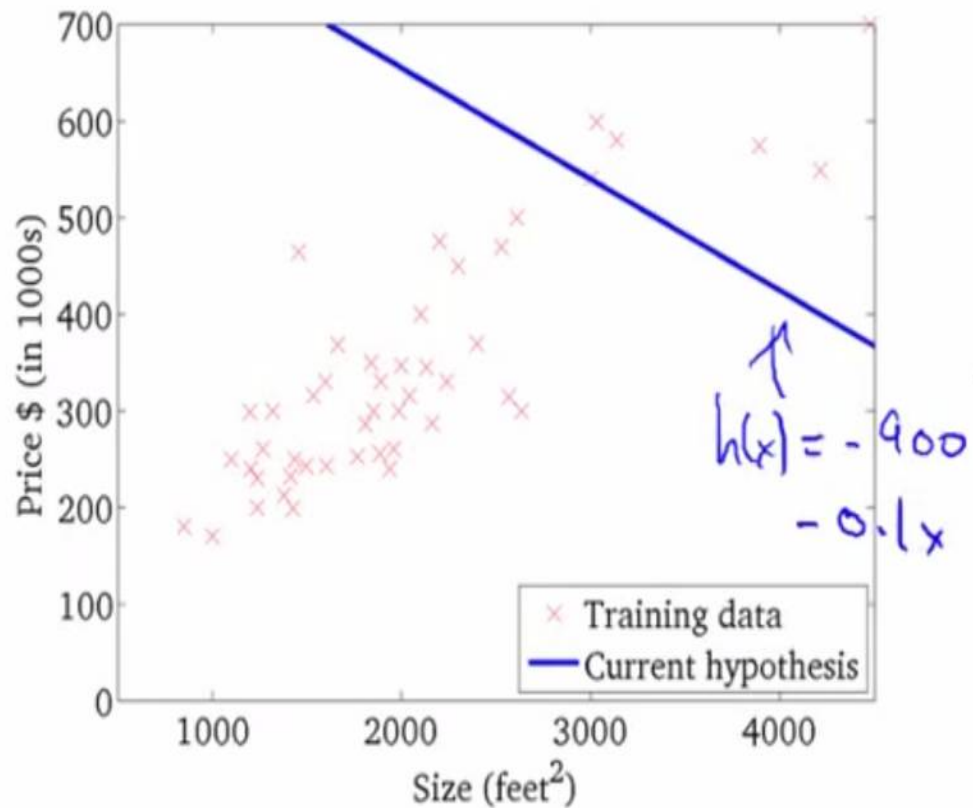
$$\underline{J(\theta_0, \theta_1)}$$

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



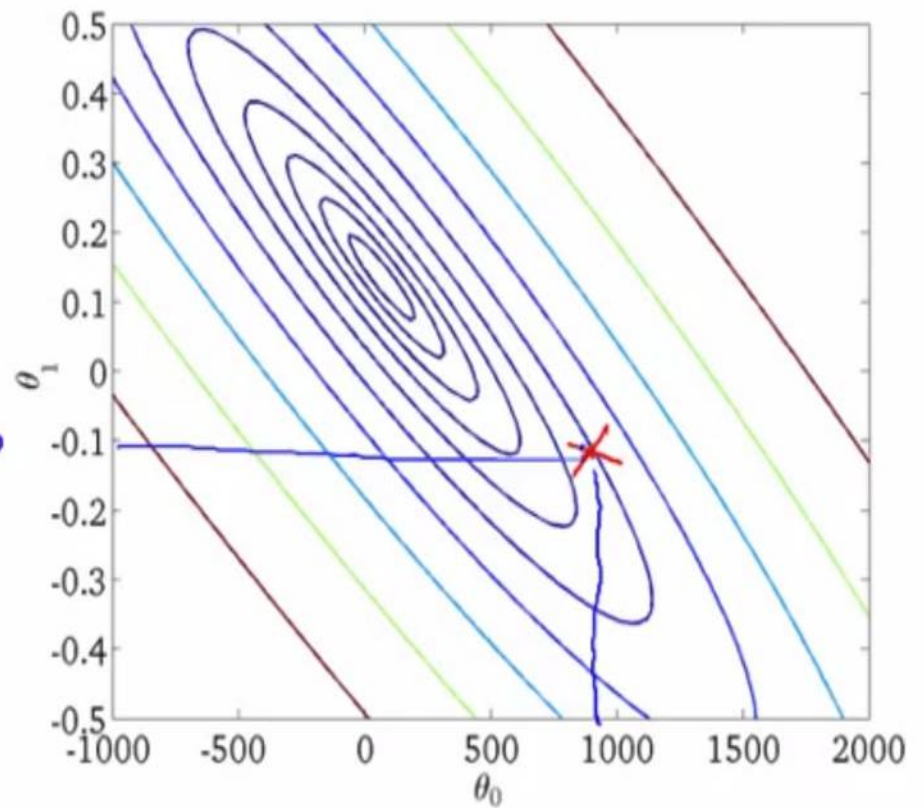
$$h_{\theta}(x)$$

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



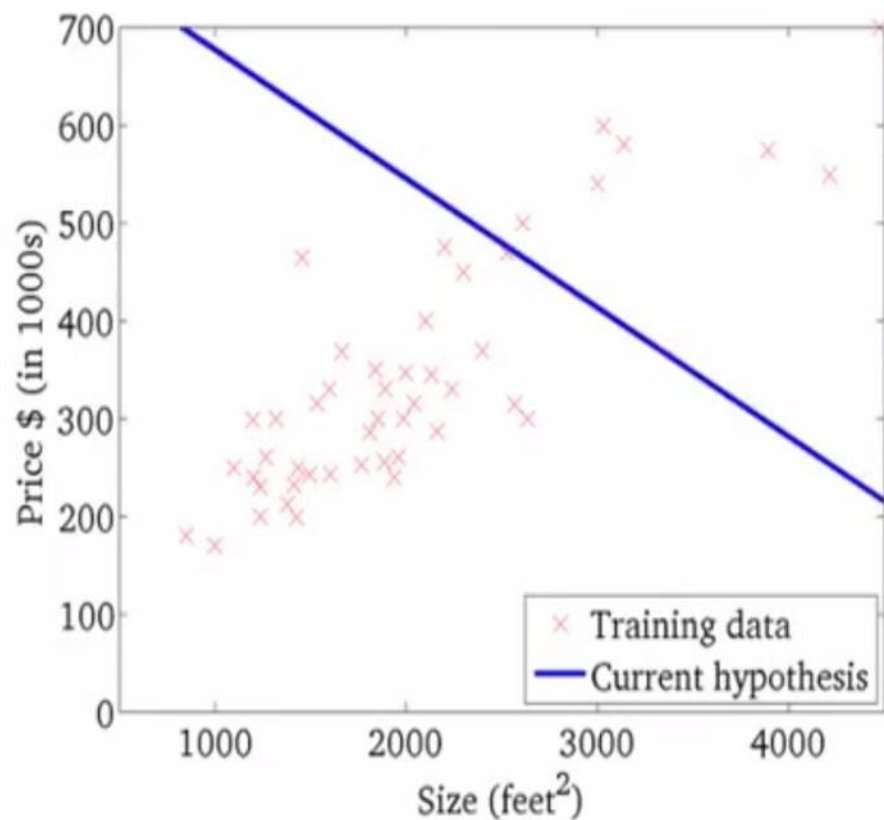
$$J(\theta_0, \theta_1)$$

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



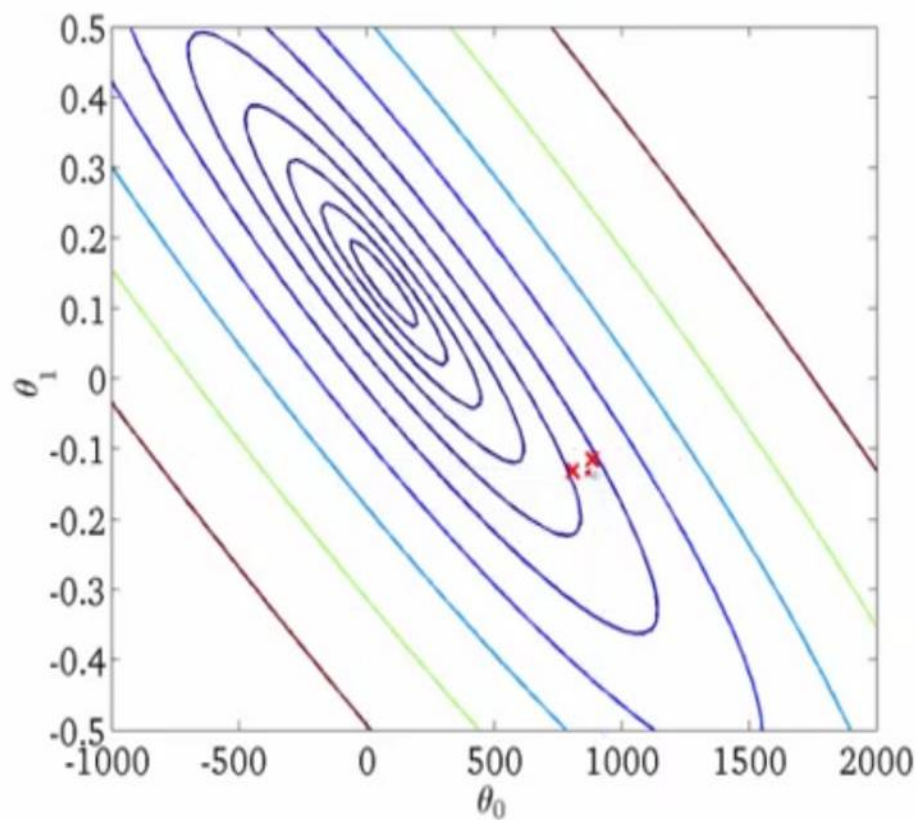
$$h_{\theta}(x)$$

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



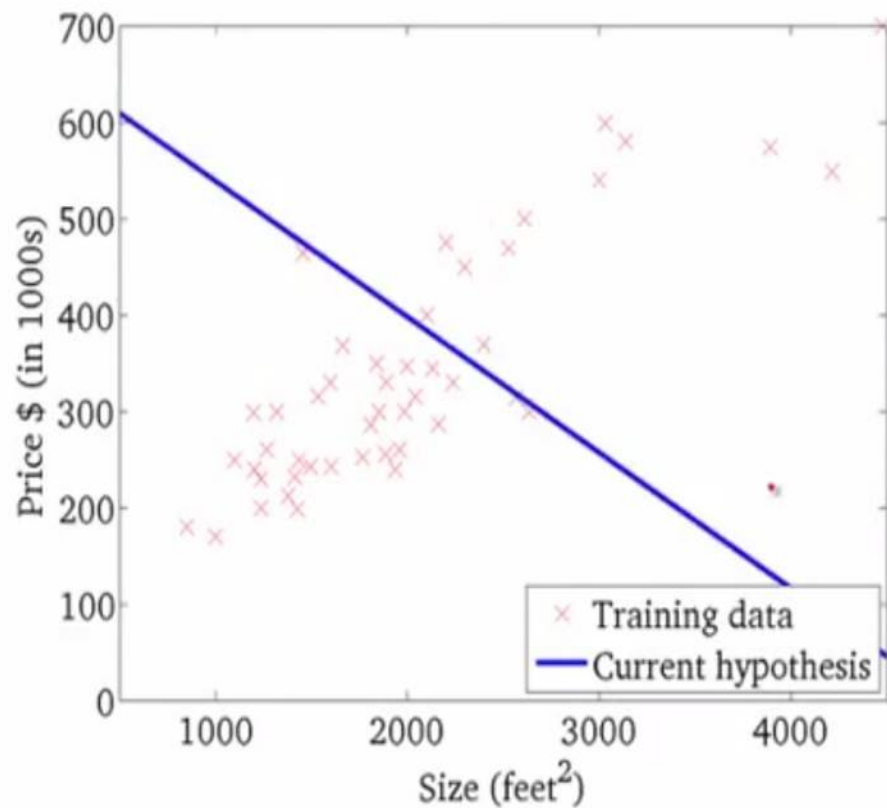
$$J(\theta_0, \theta_1)$$

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



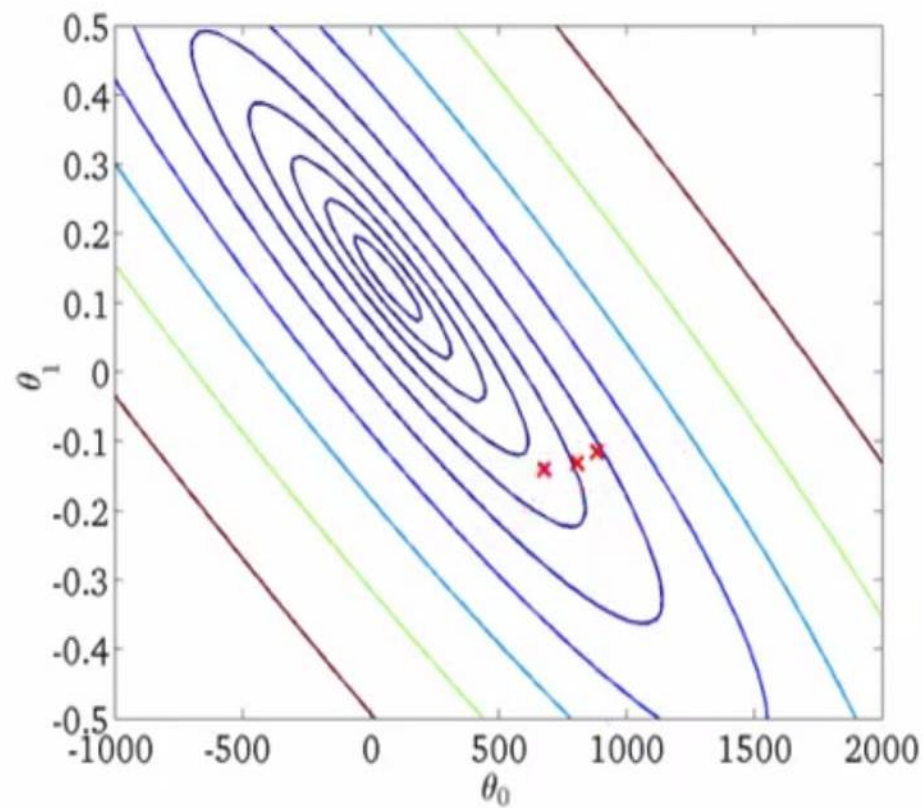
$$h_{\theta}(x)$$

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



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

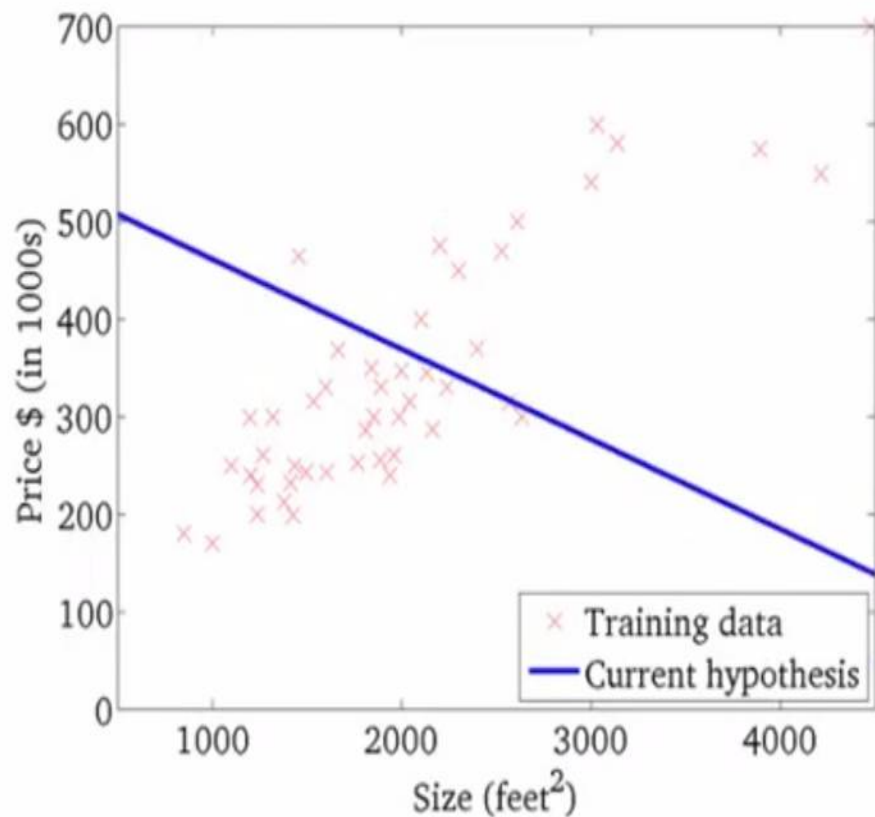
(function of the parameters  $\theta_0, \theta_1$ )





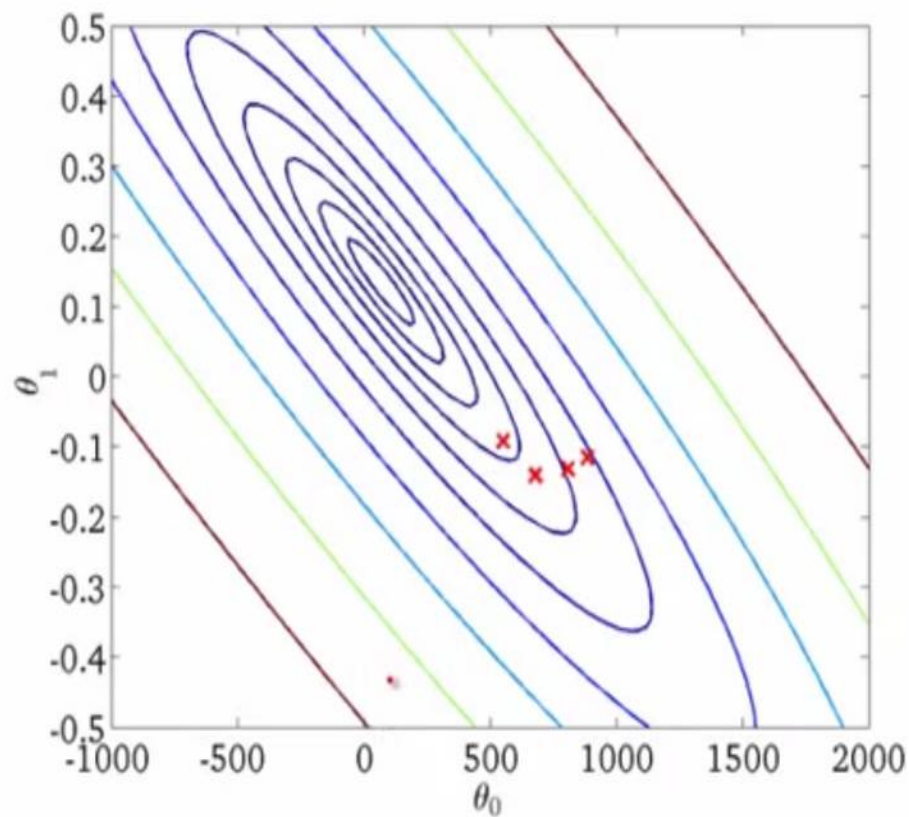
$$h_{\theta}(x)$$

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



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

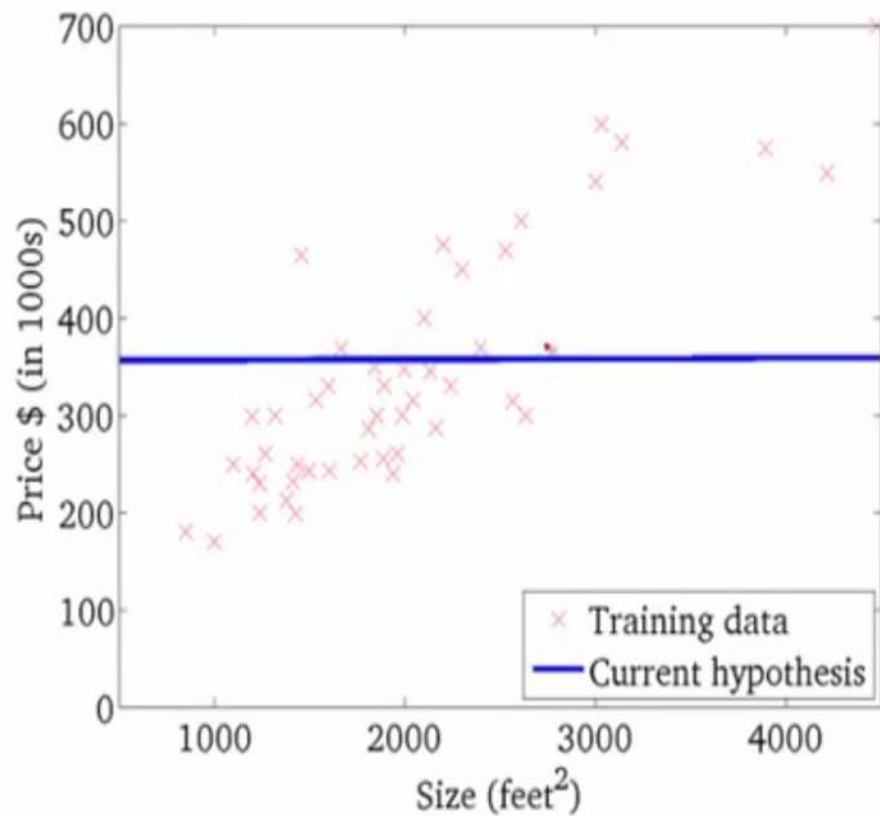
(function of the parameters  $\theta_0, \theta_1$ )





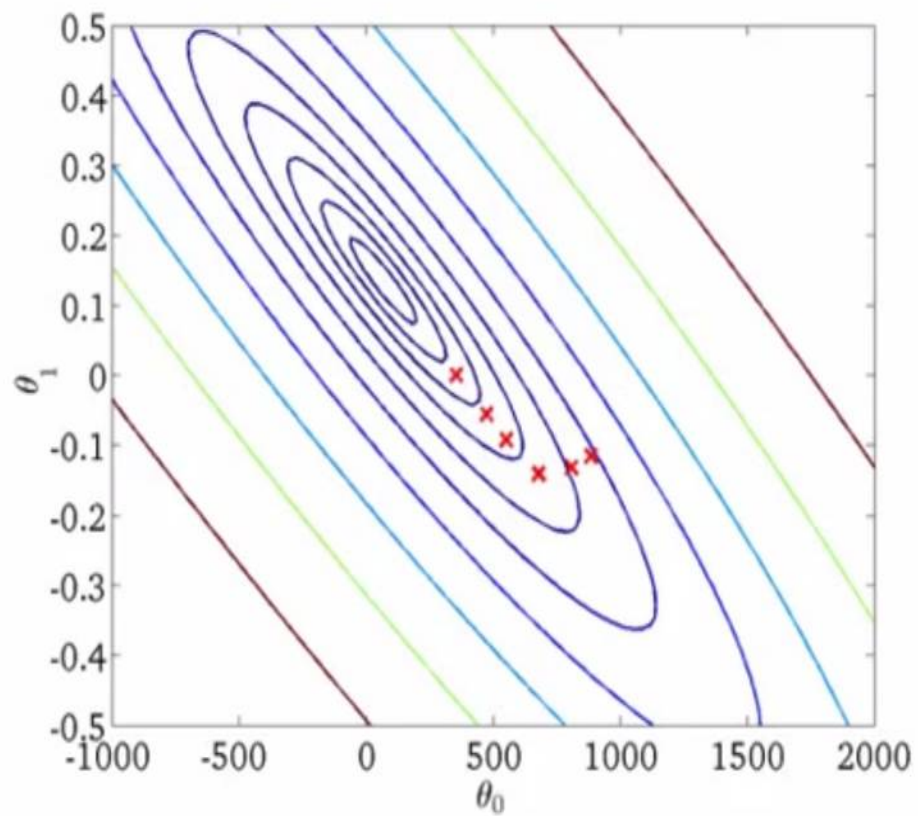
$$h_{\theta}(x)$$

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



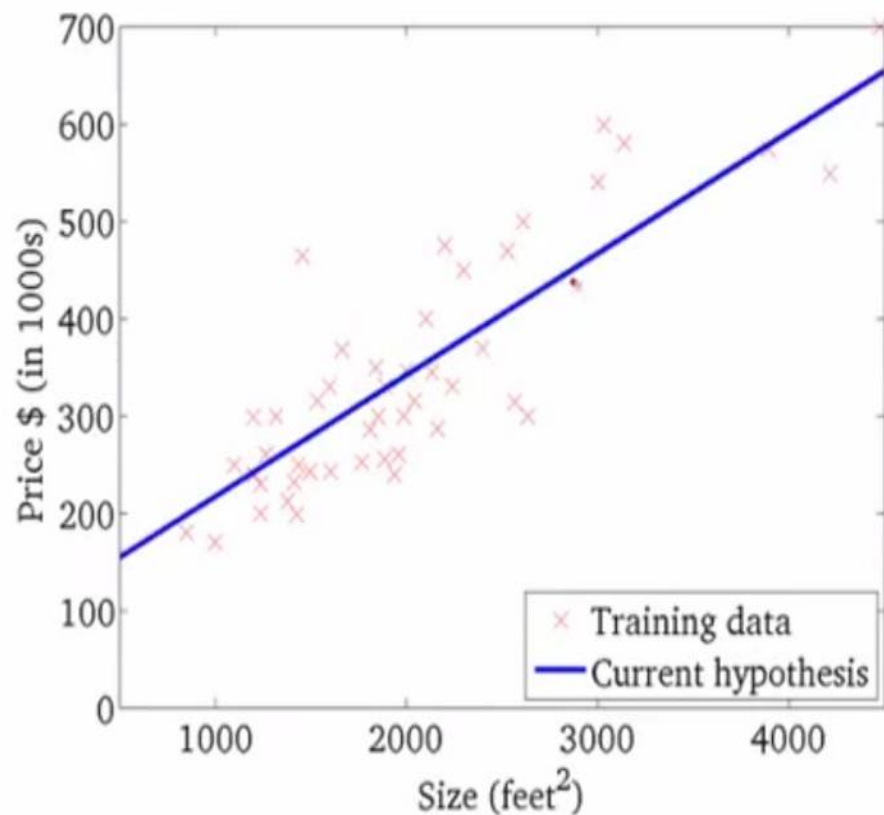
$$J(\theta_0, \theta_1)$$

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



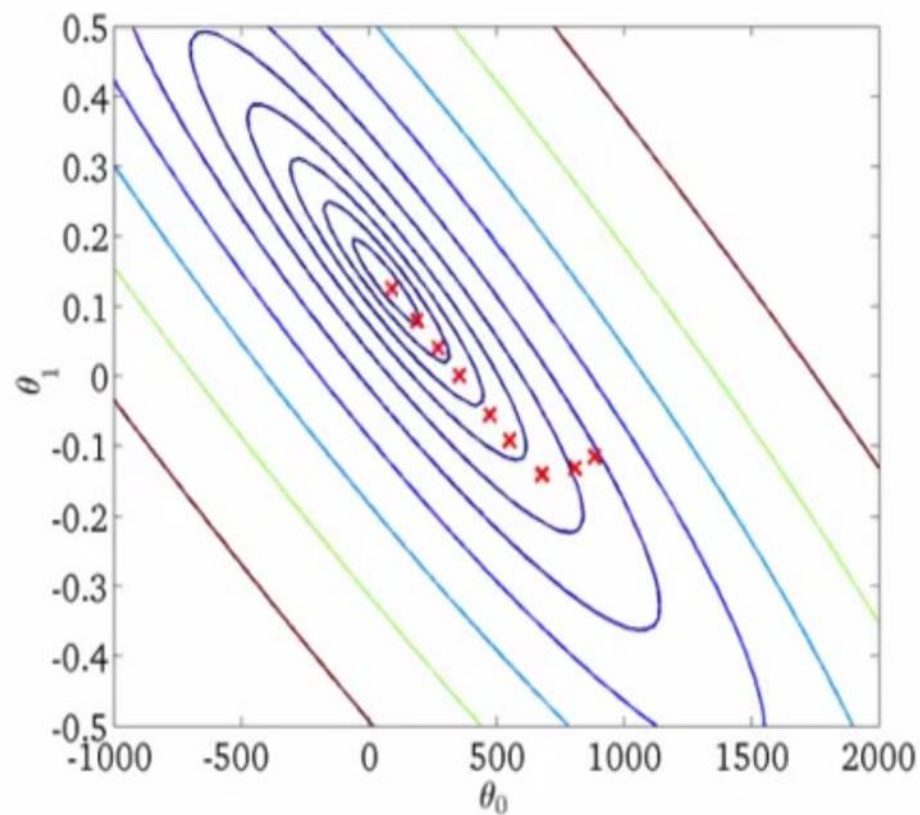
$$h_{\theta}(x)$$

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



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

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



# Gradient Descent (Summary)

**Aim:** To Find  $(\Theta_0, \Theta_1)$  that minimizes height of the surface  $J(\Theta)$ .

- Start at some point.
- Find gradient and go down in the direction (steepest descent direction chosen)
- **Note:** For non-linear cost function with multiple minima points, it may get stuck at local minima
- For linear regression, there is no local optima. The cost function for linear regression is a Convex Function (bowl shaped).  
Convergence guaranteed.

## “Batch” Gradient Descent

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

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

If  $m$  is large (e.g., 1M examples), then computing the derivative becomes very slow, just to make a tiny update.

**Stochastic Gradient Descent (SGD):** Use only one random sample to compute the derivative and then update the parameters.

- It takes a noisy (random) path, but on average is headed towards the global minima.
- Used in practice for large datasets. Much **faster** in practice. May not quite converge to the exact global minima.



# Multivariate LR: Using Multiple Features (Variables)

## Multiple features (variables).

Size (feet <sup>2</sup> ) $x_1$	Number of bedrooms $x_2$	Number of floors $x_3$	Age of home (years) $x_4$	Price (\$1000) $y$
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...

Notation:

- $n$  = number of features  $n=4$
- $x^{(i)}$  = input (features) of  $i^{th}$  training example.
- $x_j^{(i)}$  = value of feature  $j$  in  $i^{th}$  training example.

$$\underline{x^{(2)}} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$
$$x_3^{(2)} = 2$$

# Hypothesis (Multiple Features)

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$\text{e.g. } \underline{h_{\theta}(x)} = \underline{80} + \underline{0.1}x_1 + \underline{0.01}x_2 + 3x_3 - 2x_4$$

↑            ↑            ↑  
                                 age

$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

For convenience of notation, define  $x_0 = 1$ . ( $x_0^{(i)} = 1$ )

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\underbrace{[\theta_0 \ \theta_1 \ \dots \ \theta_n]}_{\theta^T} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

1 x (n + 1)  
matrix

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$\downarrow = 1$

$$= \theta^T x$$

Multivariate linear regression.  $\leftarrow$

# Gradient Descent for Multiple Variables

Hypothesis:  $\overset{x_0=1}{h_{\theta}(x)} = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters:  $\theta_0, \theta_1, \dots, \theta_n$   $\underline{\theta}$   $n+1$ -dimensional vector

Cost function:

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

Gradient descent:

Repeat {

$$\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \underset{J(\theta)}{J(\theta_0, \dots, \theta_n)}$$

(simultaneously update for every  $j = 0, \dots, n$ )

# Gradient Descent

Previously (n=1):

Repeat {

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

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \underbrace{x^{(i)}}_{x_1^{(i)}}$$

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm (n ≥ 2):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

}

# Gradient Descent

Previously (n=1):

Repeat {

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

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \underbrace{x_1^{(i)}}_{x_1^{(i)}}$$

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm (n ≥ 2):

Repeat {

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

(simultaneously update  $\theta_j$  for  $j = 0, \dots, n$ )

}

$$\rightarrow \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\rightarrow \theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\rightarrow \theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

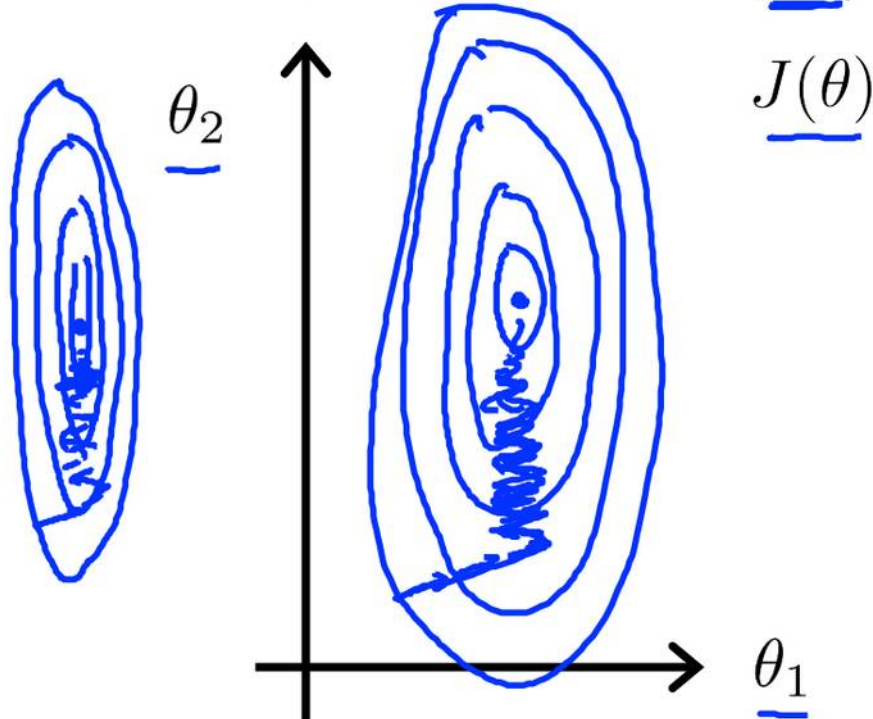


# Feature Scaling

Idea: Make sure features are on a similar scale.

E.g.  $x_1 = \text{size (0-2000 feet}^2\text{)}$  ←

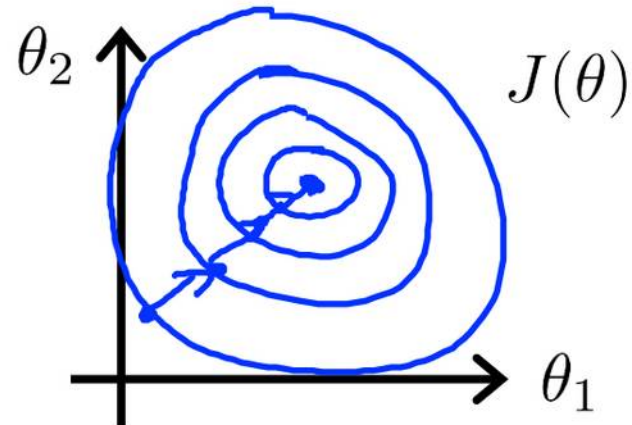
$x_2 = \text{number of bedrooms (1-5)}$  ←



$$\rightarrow x_1 = \frac{\text{size (feet}^2\text{)}}{2000} \quad \checkmark$$

$$\rightarrow x_2 = \frac{\text{number of bedrooms}}{5} \quad \checkmark$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



## Mean normalization

Replace  $x_i$  with  $x_i - \mu_i$  to make features have approximately zero mean  
(Do not apply to  $x_0 = 1$ ).

E.g.  $x_1 = \frac{\text{size} - 1000}{2000}$

$$x_2 = \frac{\text{\#bedrooms} - 2}{5}$$

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_1 \leftarrow \frac{x_1 - \mu_1}{S_1} \quad \left| \quad x_2 \leftarrow \frac{x_2 - \mu_2}{S_2}$$

Handwritten notes for  $x_1$ :

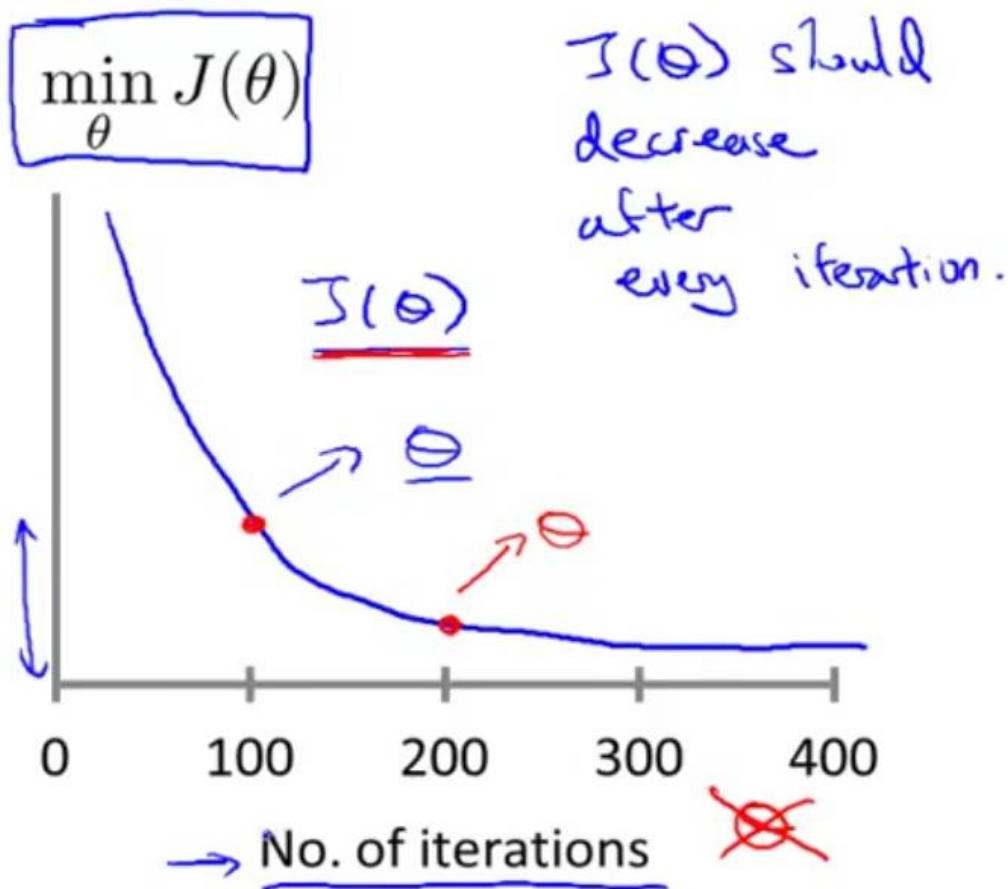
- $\mu_1$  is circled in red. An arrow points to it from the text "avg value of  $x_1$  in training set".
- $S_1$  is circled in red. An arrow points to it from the text "range (max-min) (or standard deviation)".

## Gradient descent

$$\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate  $\alpha$ .

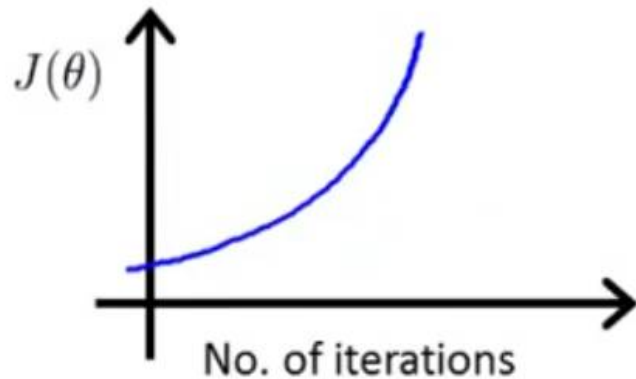
## Making sure gradient descent is working correctly.



Example automatic convergence test:

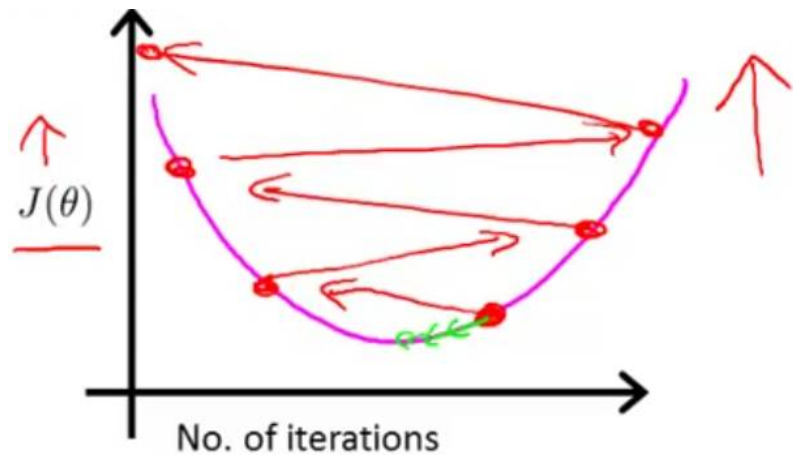
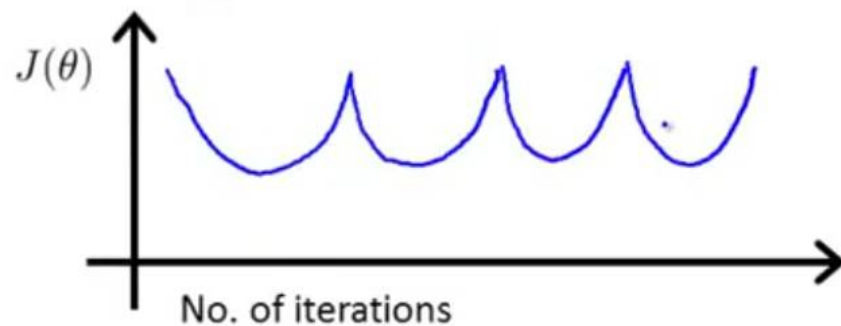
Declare convergence if  $J(\theta)$  decreases by less than  $10^{-3}$  in one iteration.

## Making sure gradient descent is working correctly.



Gradient descent not working.

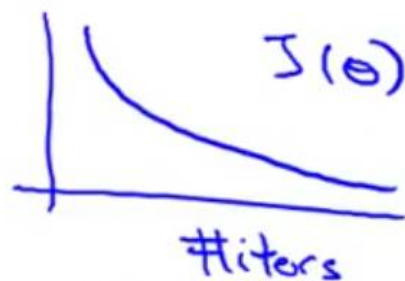
Use smaller  $\alpha$ .



- For sufficiently small  $\alpha$ ,  $J(\theta)$  should decrease on every iteration.
- But if  $\alpha$  is too small, gradient descent can be slow to converge.

## Summary:

- If  $\alpha$  is too small: slow convergence.
- If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration; may not converge. (Slow converge also possible.)



To choose  $\alpha$ , try

$\dots, 0.001, \quad , 0.01, \quad , 0.1, \quad , 1, \dots$



# Features and Polynomial Regression

## Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

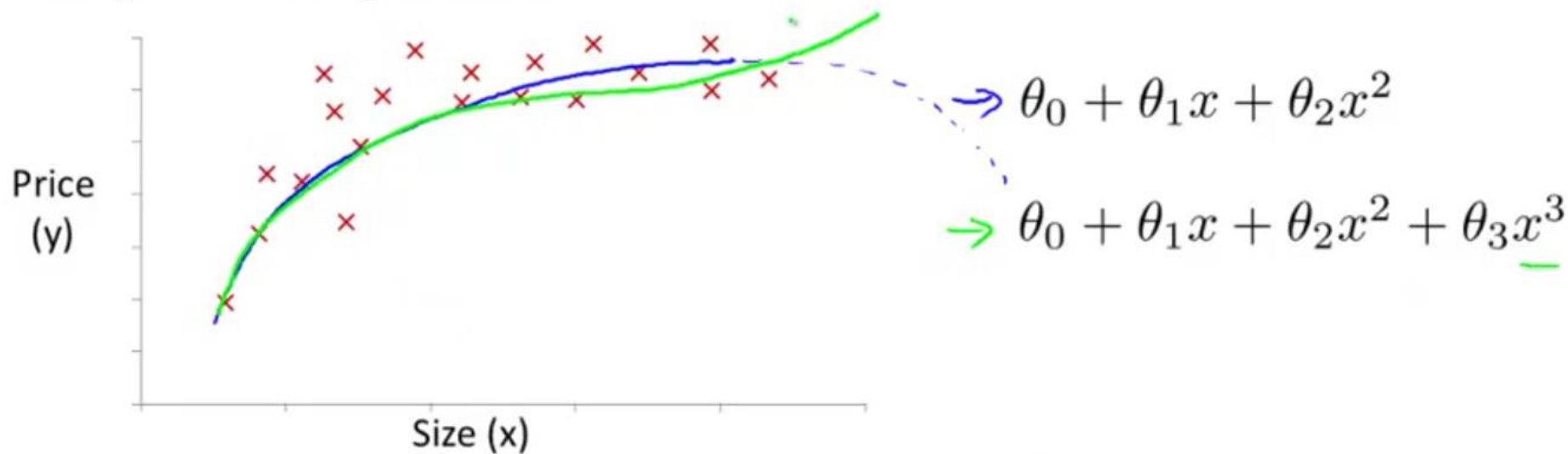
$$x = \underline{\text{frontage} \times \text{depth}}$$

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

↖ land area



## Polynomial regression

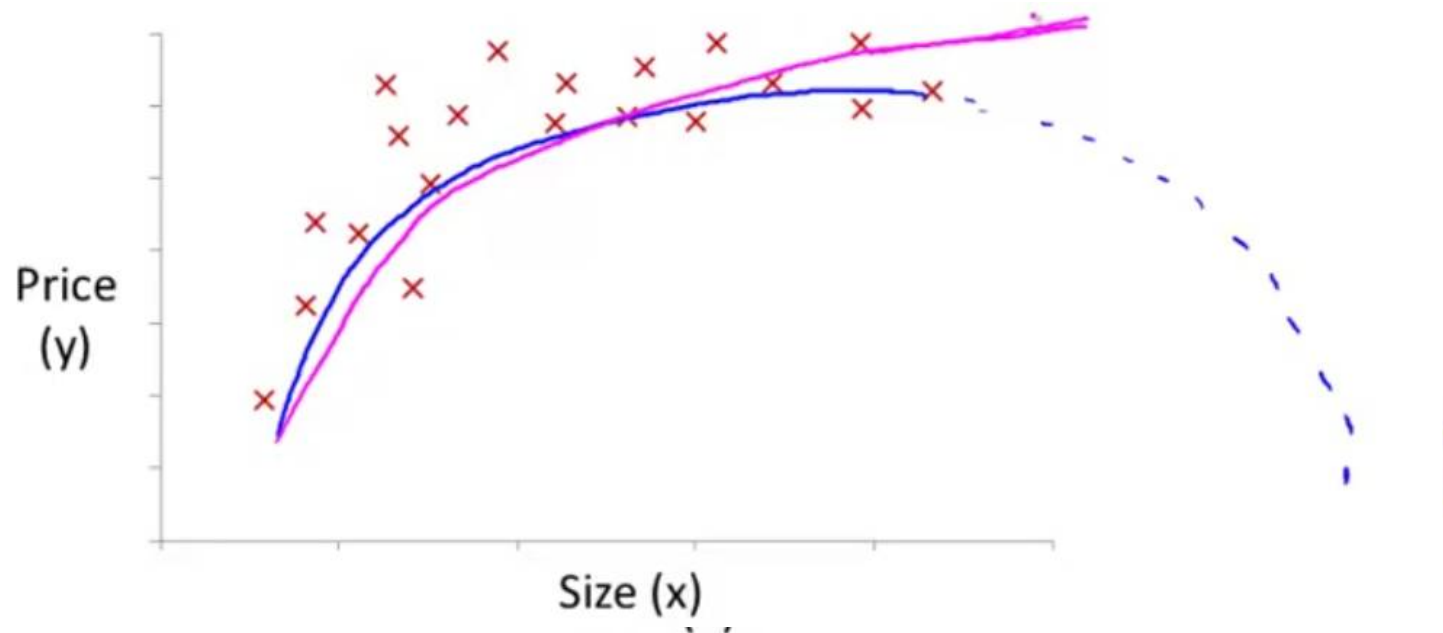


$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3 \end{aligned}$$

$$\begin{aligned} \rightarrow x_1 &= (\text{size}) \\ \rightarrow x_2 &= (\text{size})^2 \\ \rightarrow x_3 &= (\text{size})^3 \end{aligned}$$

$$\begin{aligned} \text{Size:} & 1 - 1000 \\ \text{Size}^2: & 1 - 1,000,000 \\ \text{Size}^3: & 1 - 10^9 \end{aligned}$$

## Choice of features



$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

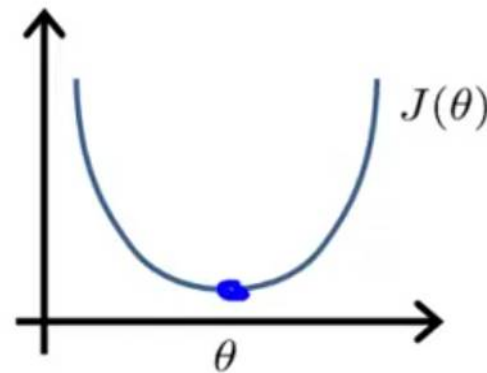
$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$



# Normal Equation: Method to solve for $\Theta$ analytically

Intuition: If 1D ( $\theta \in \mathbb{R}$ )

$\rightarrow \underline{J(\theta) = a\theta^2 + b\theta + c}$   
 $\frac{d}{d\theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$   
Solve for  $\Theta$



---

$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_n) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

Solve for  $\theta_0, \theta_1, \dots, \theta_n$

Examples:  $m = 4$ .

	Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$y$
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$   
 $m \times (n+1)$

$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$   
 $m$ -dimensional vector

$$\theta = (X^T X)^{-1} X^T y$$

$m$  examples  $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$  ;  $n$  features.

$$\underline{x^{(i)}} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad \left| \quad \begin{array}{c} \times \\ \text{(design} \\ \text{matrix)} \end{array} \right. = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix}$$

$m \times (n+1)$

$$\theta = (X^T X)^{-1} X^T y$$



# Question

- Which of the following is/are True w.r.t. normal equation?
  - It may not have a solution.
  - – It may have a unique solution.
  - – It may have infinite solutions.
  - – It is the least squares solution of  $Ax = b$

# Homework

- Derivation of Normal Equation.
- In which cases will a unique solution be not available using the Normal Equation (for linear regression).
- Is feature scaling necessary when using a Normal Equation to solve the linear regression task. (Hint: Are iterations needed.)

$m$  training examples,  $n$  features.

### Gradient Descent

- • Need to choose  $\alpha$ .
- • Needs many iterations.
- Works well even when  $n$  is large.

↗  
 $n = 10^6$

← -

### Normal Equation

- • No need to choose  $\alpha$ .
- • Don't need to iterate.
- Need to compute
- •  $\boxed{(X^T X)^{-1}}$   $n \times n$   $O(n^3)$
- Slow if  $n$  is very large.

$n = 100$

$n = 1000$

- - -  $n = 10000$

# References

- Gradient Descent Playground  
<https://uclaacm.github.io/gradient-descent-visualiser/>
- Stanford CS229 Machine Learning Course by Andrew Ng.

# End of Lecture