

---

# Linear Regression

Eun Yi Kim

---



Artificial Intelligence  
& Computer Vision  
L a b o r a t o r y

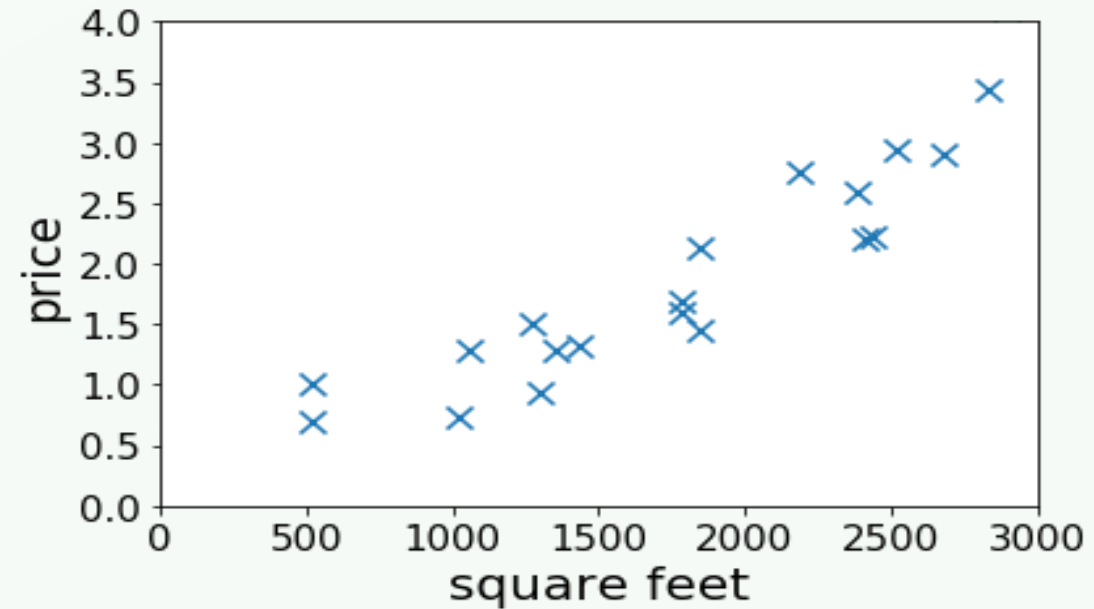


- Linear regression with one variable
  - ✓ Model representation
  - ✓ Cost function
  - ✓ Gradient decent
  - ✓ Gradient decent for linear regression
- Linear regression with multiple variable
  - ✓ Multiple features
  - ✓ Gradient decent for multiple variables
  - ✓ Gradient decent in practice
  - ✓ Features and Polynomial regression

Linear regression with one variable



## Housing Prices (in 100,000s of dollars)



- Supervised Learning
  - : Given the “right answers” for each example in the data
- Regression
  - : Predict real-valued output (price)
- Classification
  - : Discrete-valued output



## Training set of housing prices

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

Notation:

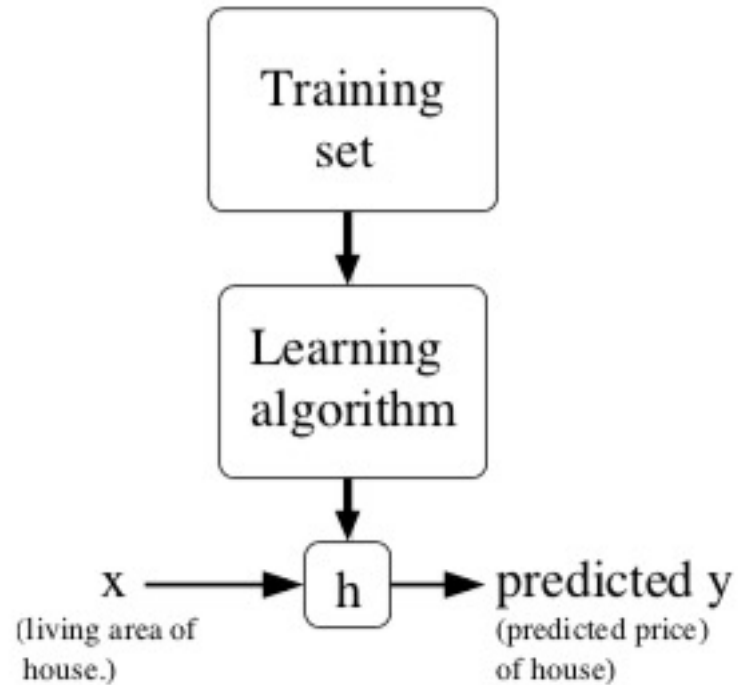
- $m$  = Number of training examples
- $x$ 's = “input” variable / features
- $y$ 's = “output” variable / “target” variable

# Model Representation



Artificial Intelligence  
& Computer Vision  
Laboratory

How do we represent  $h$ ?



Linear regression with one variable  
: Univariate linear regression

# Cost Function



Training set of housing prices

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

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

$\theta_i$ 's :  
Parameters

How to choose  $\theta_i$ 's ?

# Cost Function



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



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



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



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

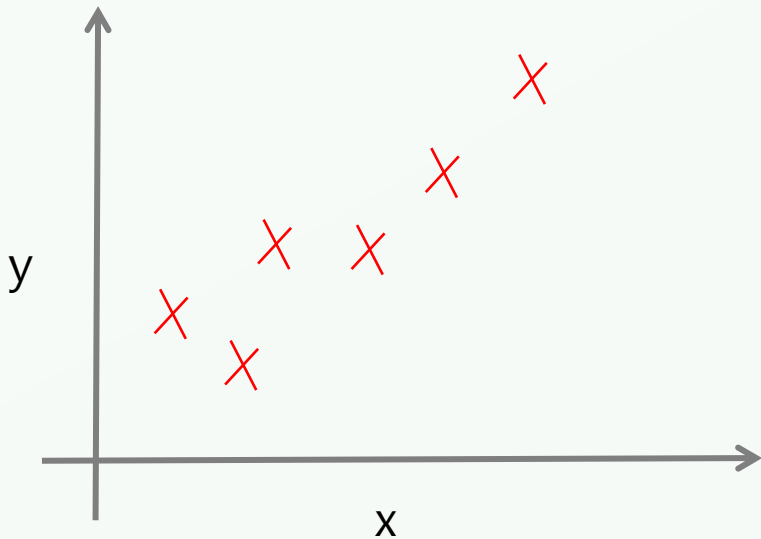


# Cost Function



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

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



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

# Cost Function (intuition)



Artificial Intelligence  
& Computer Vision  
Laboratory

Simplified

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1$$

Cost function:

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

Goal:

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

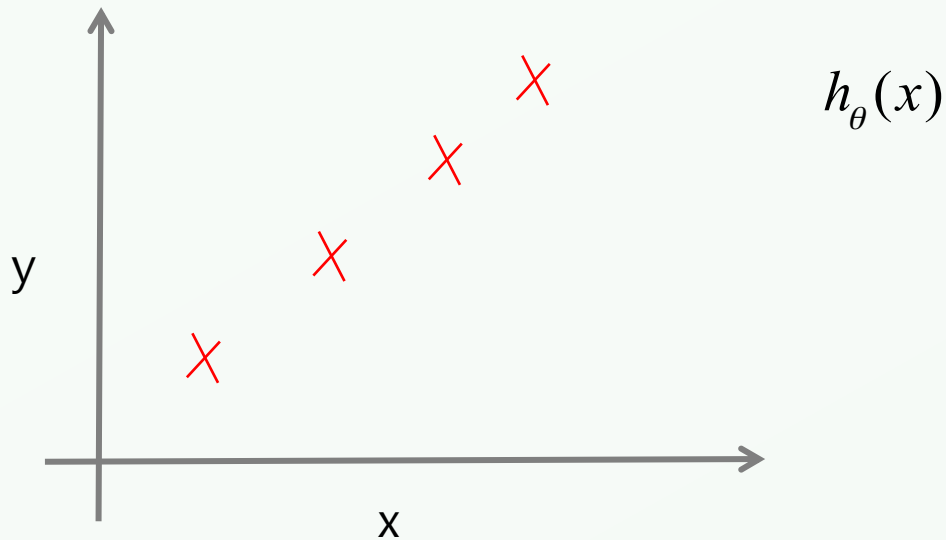
$$h_{\theta}(x) = \theta_1 x$$

# Cost Function (intuition)



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

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



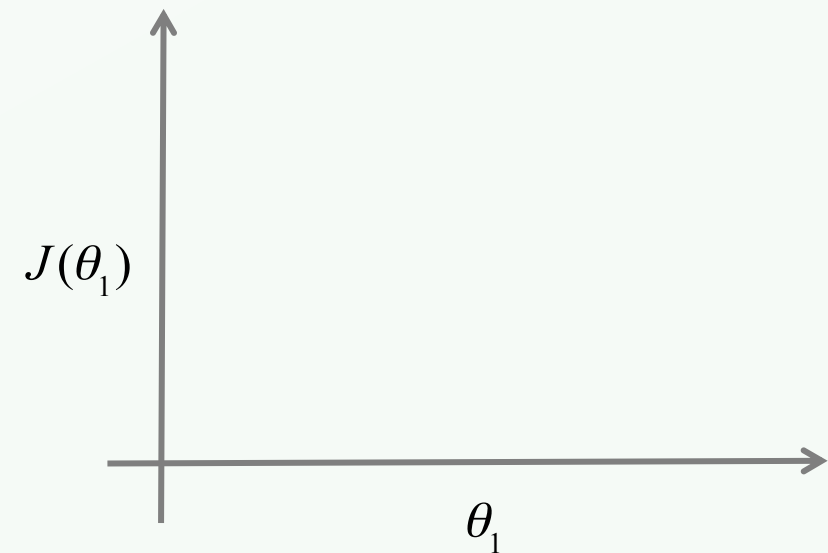
Idea: Choose  $\theta_0$  and  $\theta_1$  so that

$h_{\theta}(x)$  is close to  $y$

for our training example  $(x, y)$

$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )



# Cost Function (intuition)



Artificial Intelligence  
& Computer Vision  
Laboratory

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

Parameters:  $\theta_0, \theta_1$

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

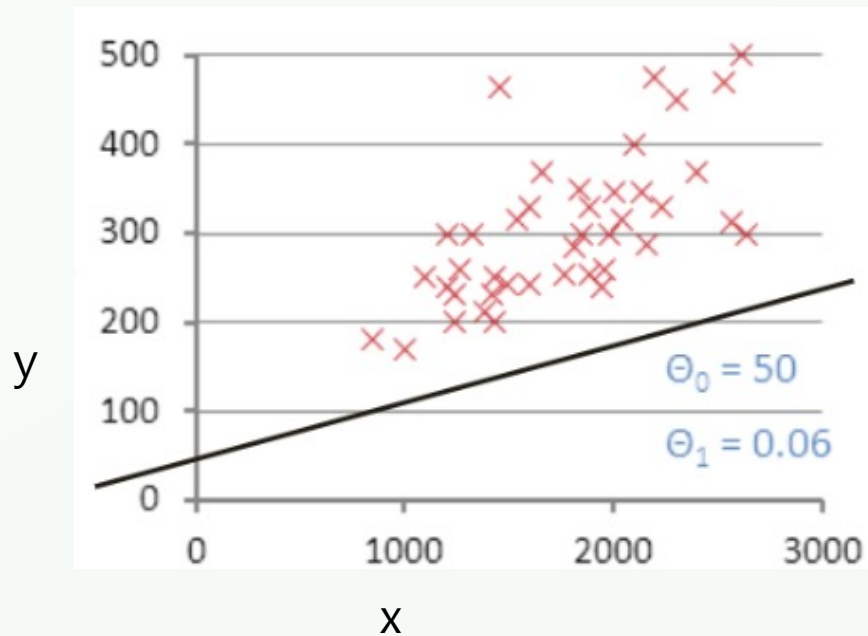
Goal:  $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

# Cost Function (intuition)



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

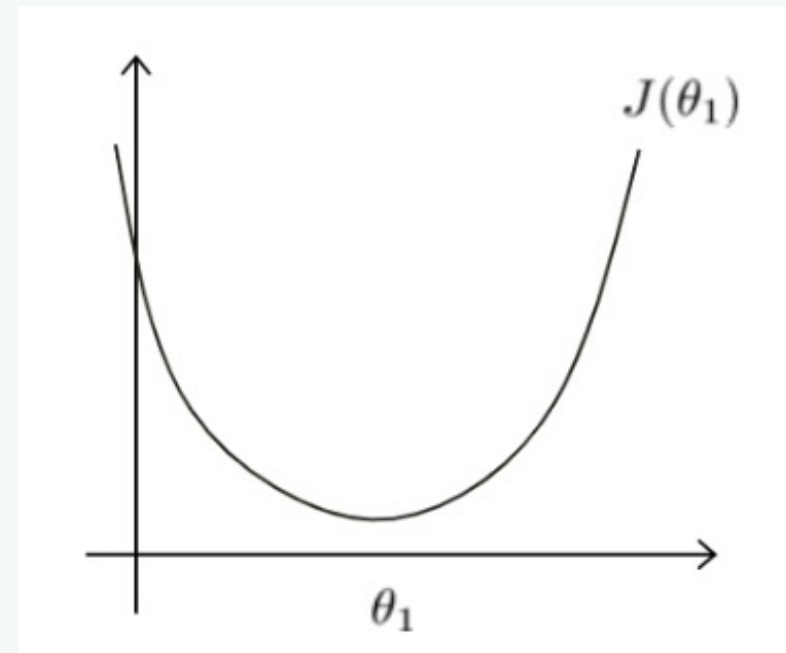
(for fixed  $\theta_1$ , this is a function of  $x$ )



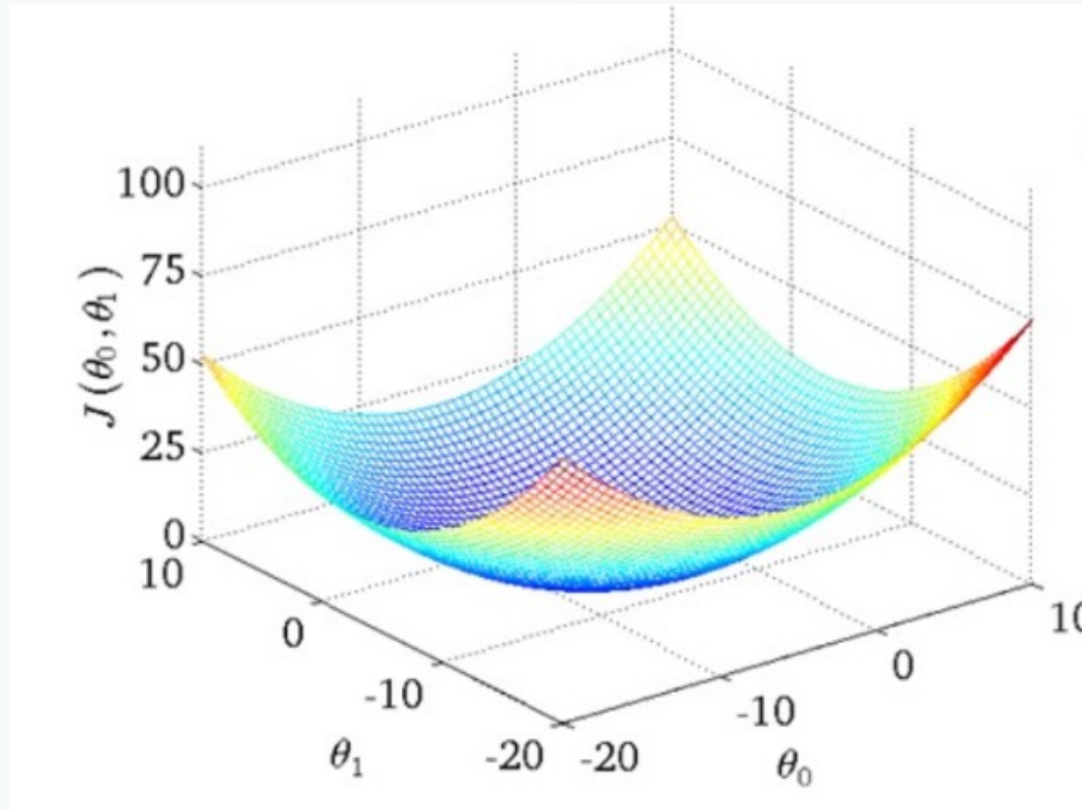
$$h_{\theta}(x) = 50 + 0.06x$$

$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )



# Cost Function (intuition)



Contour plots

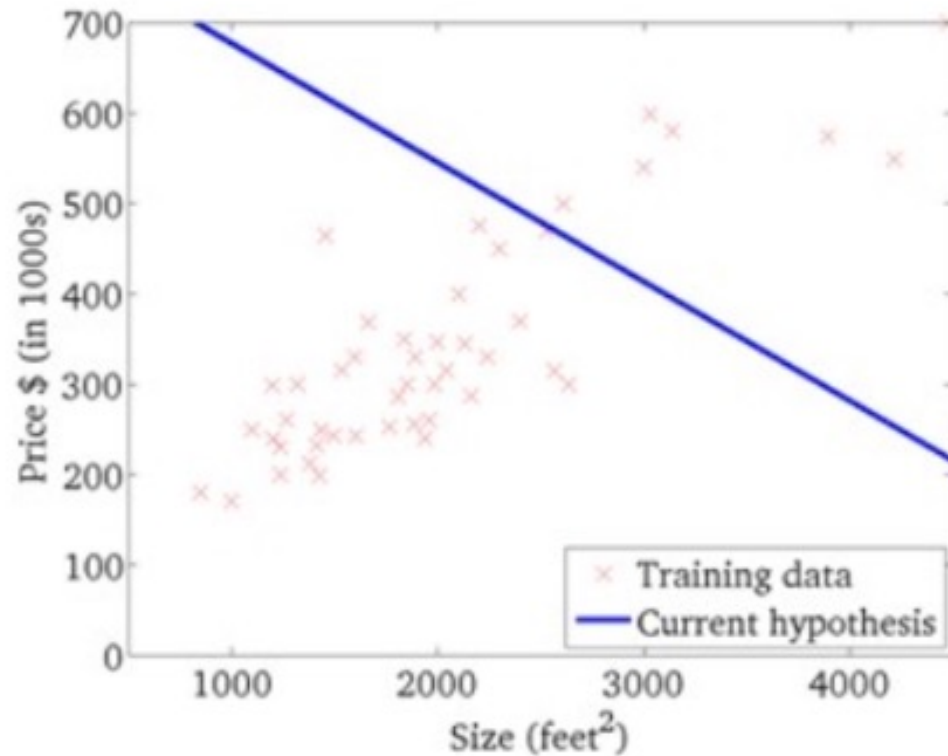
# Cost Function (intuition)



Artificial Intelligence  
& Computer Vision  
Laboratory

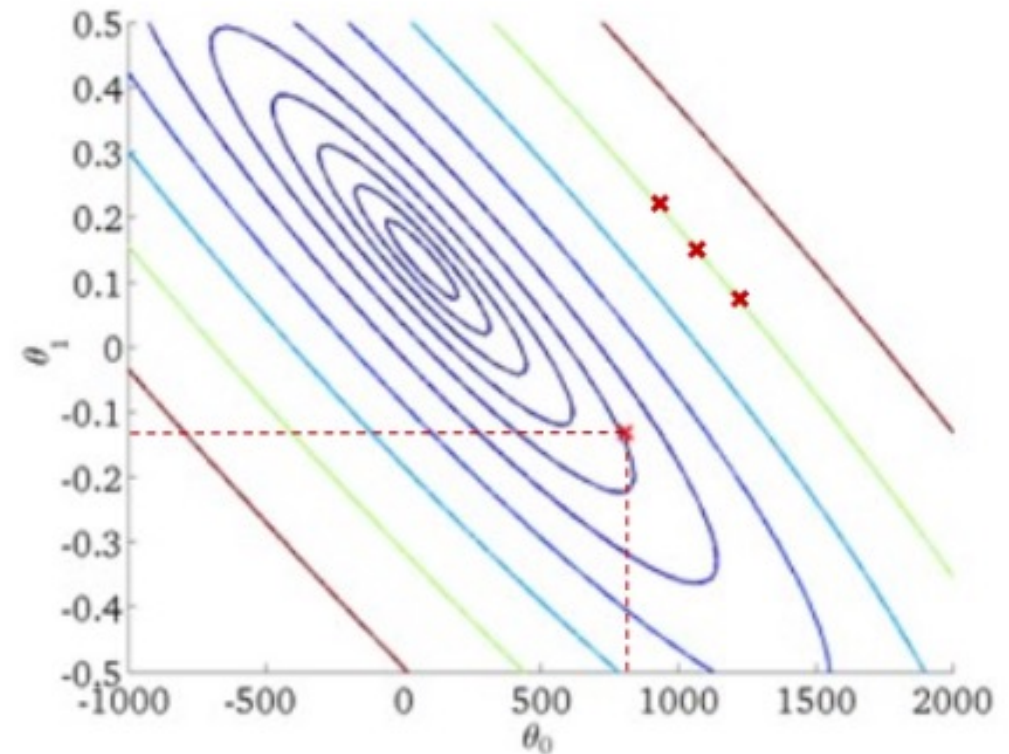
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

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

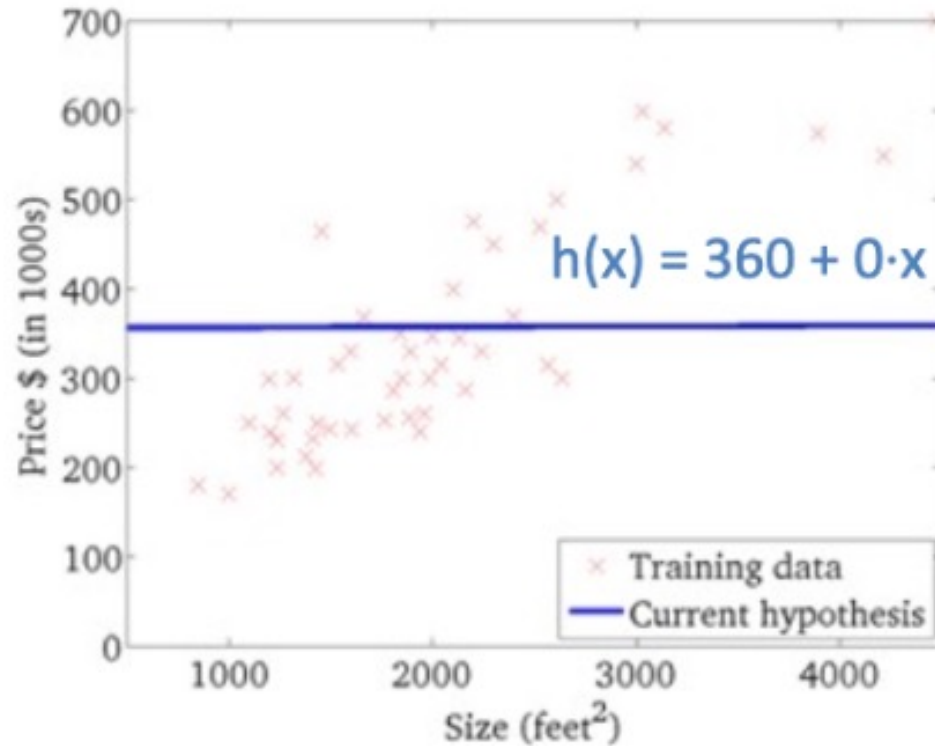


# Cost Function (intuition)



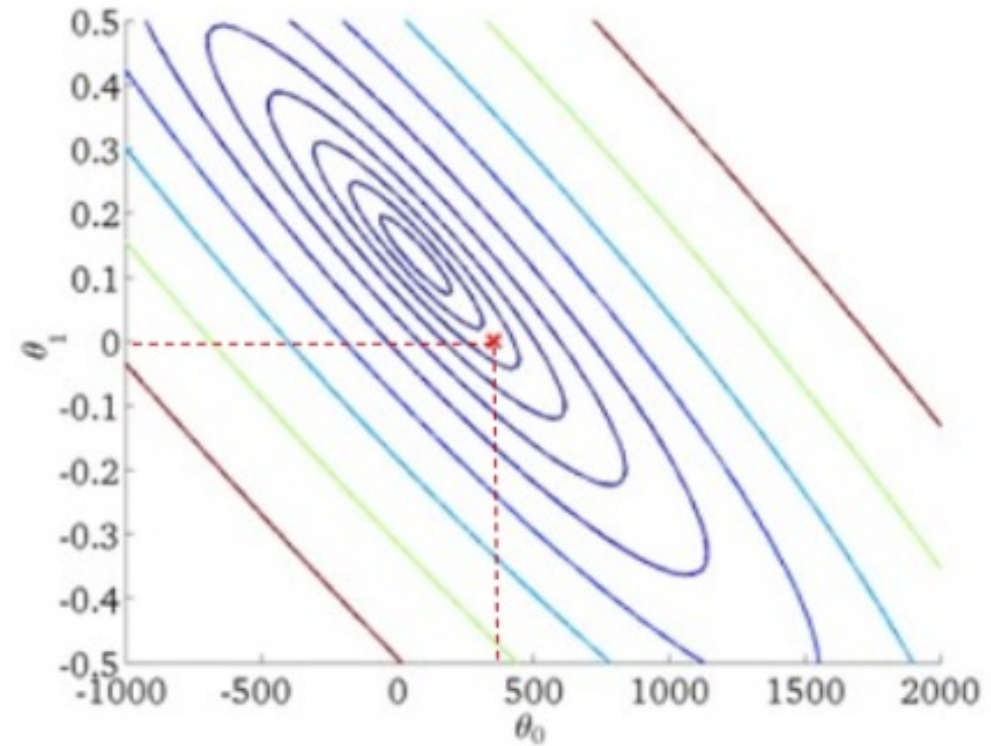
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

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



$$\Theta_0 = 360$$

$$\Theta_1 = 0$$

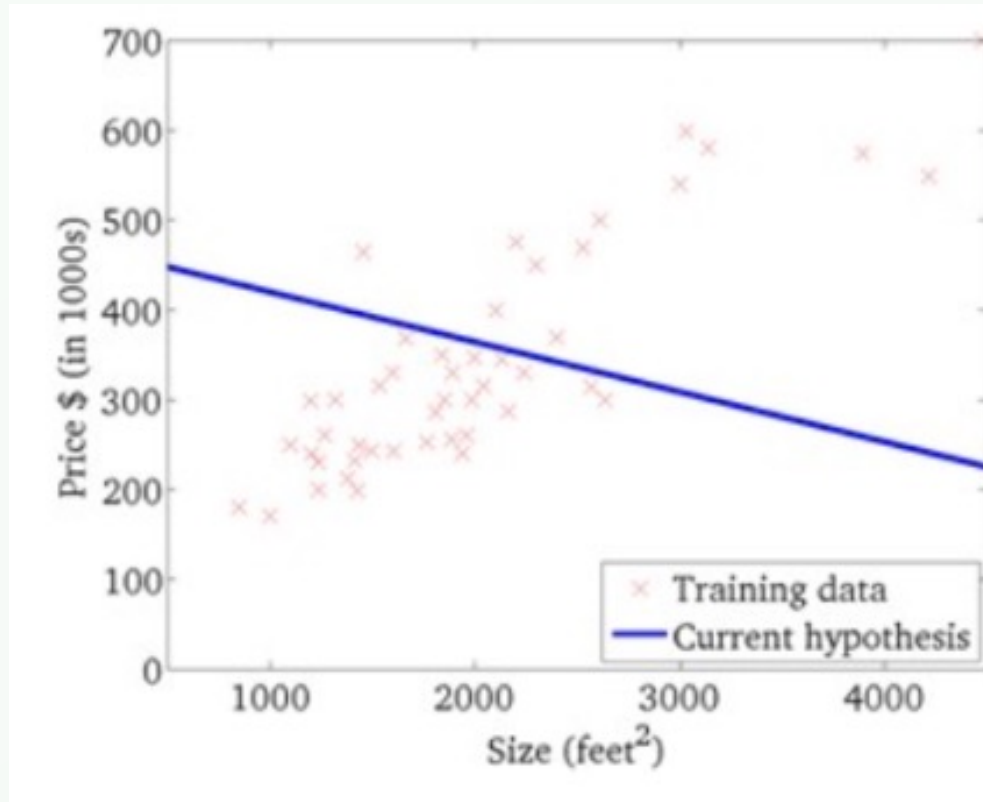


# Cost Function (intuition)



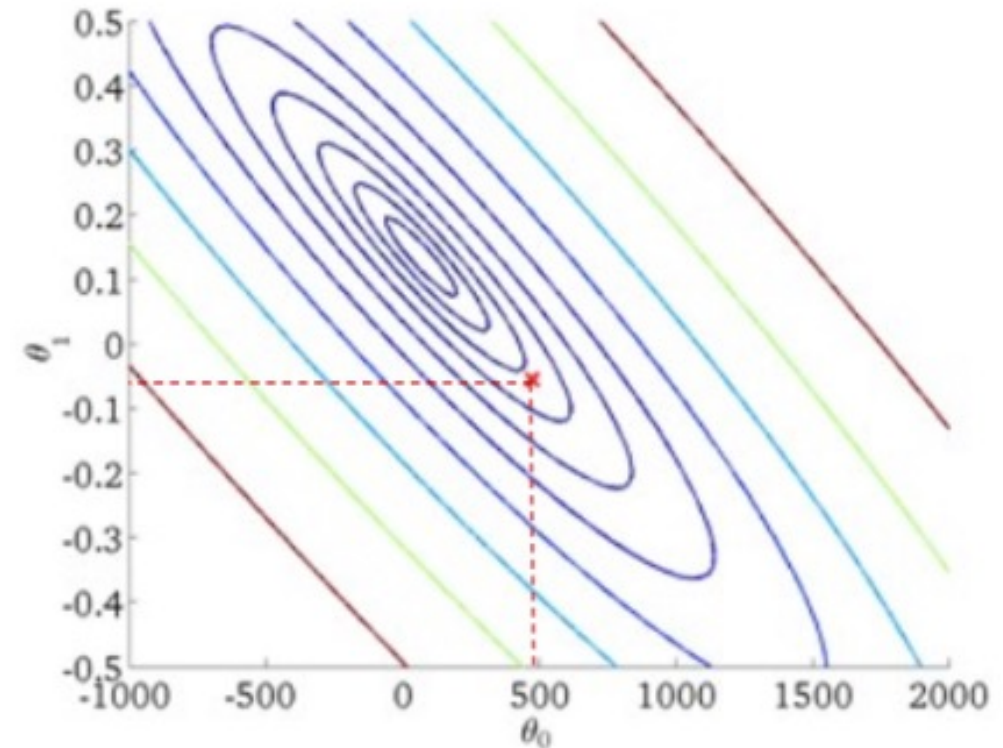
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

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

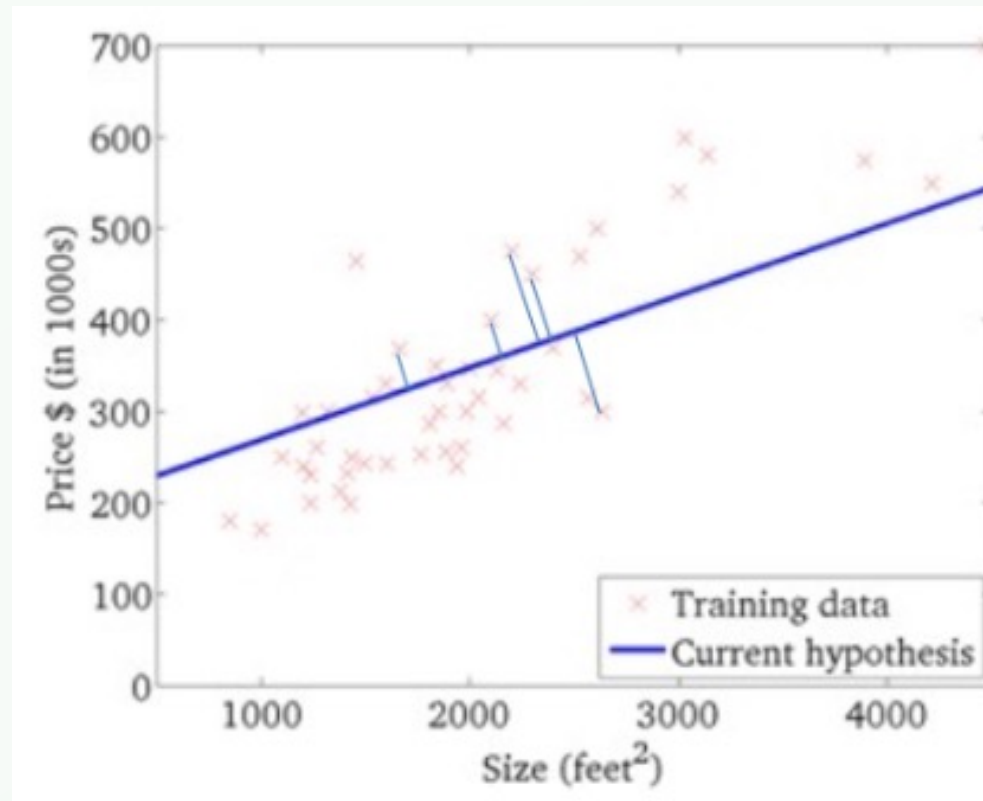


# Cost Function (intuition)



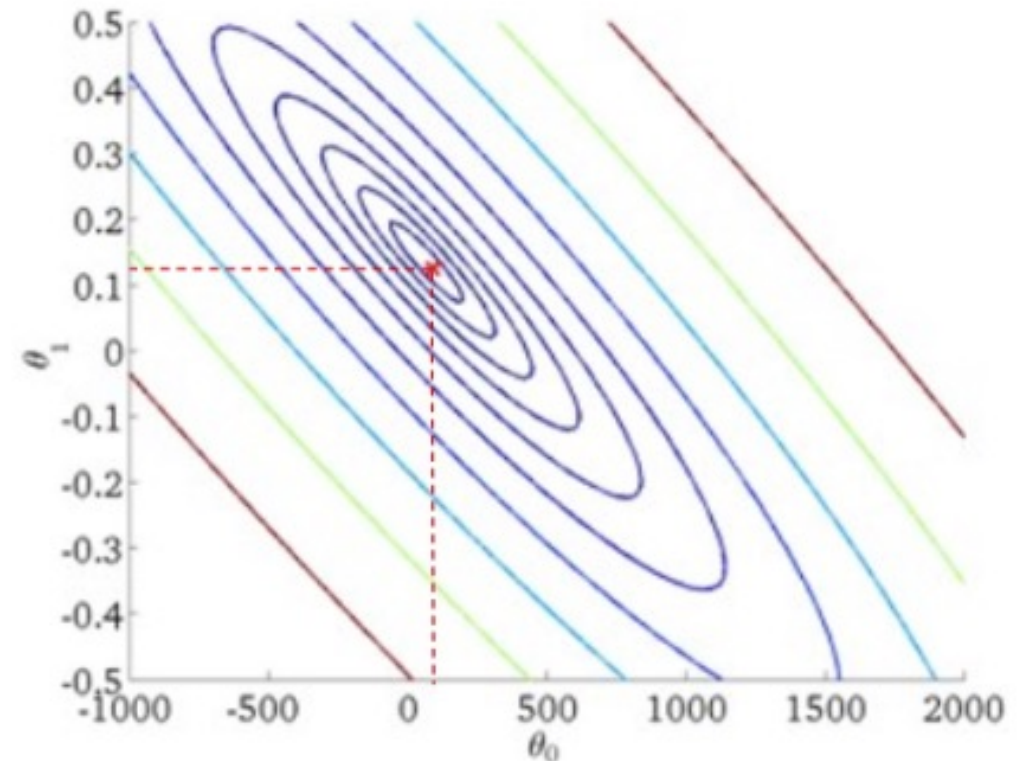
$$h_{\theta}(x)$$

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



$$J(\theta_1)$$

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





Have some function  $J(\theta_0, \theta_1)$

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

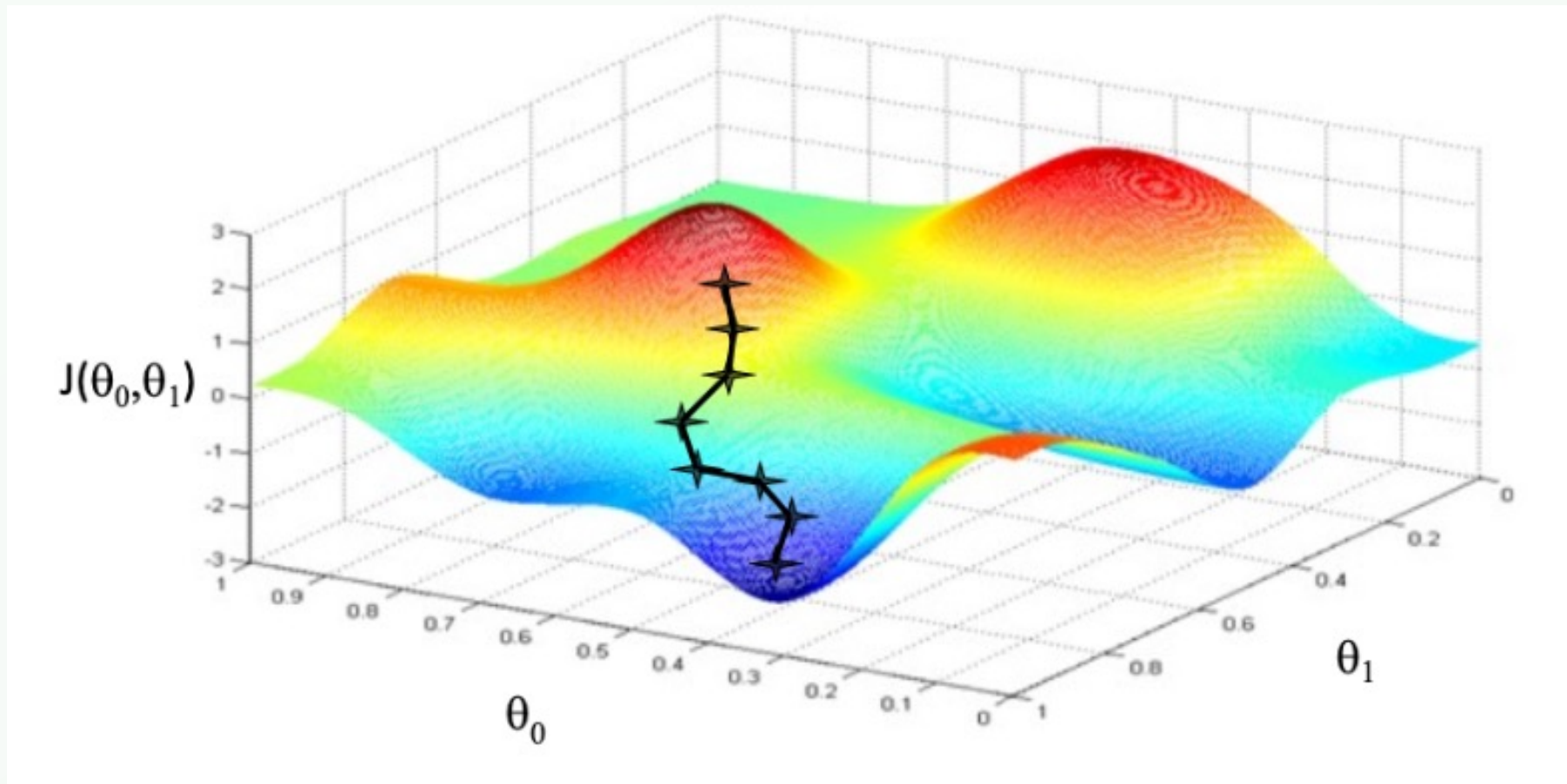
Outline:

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

# Gradient Descent



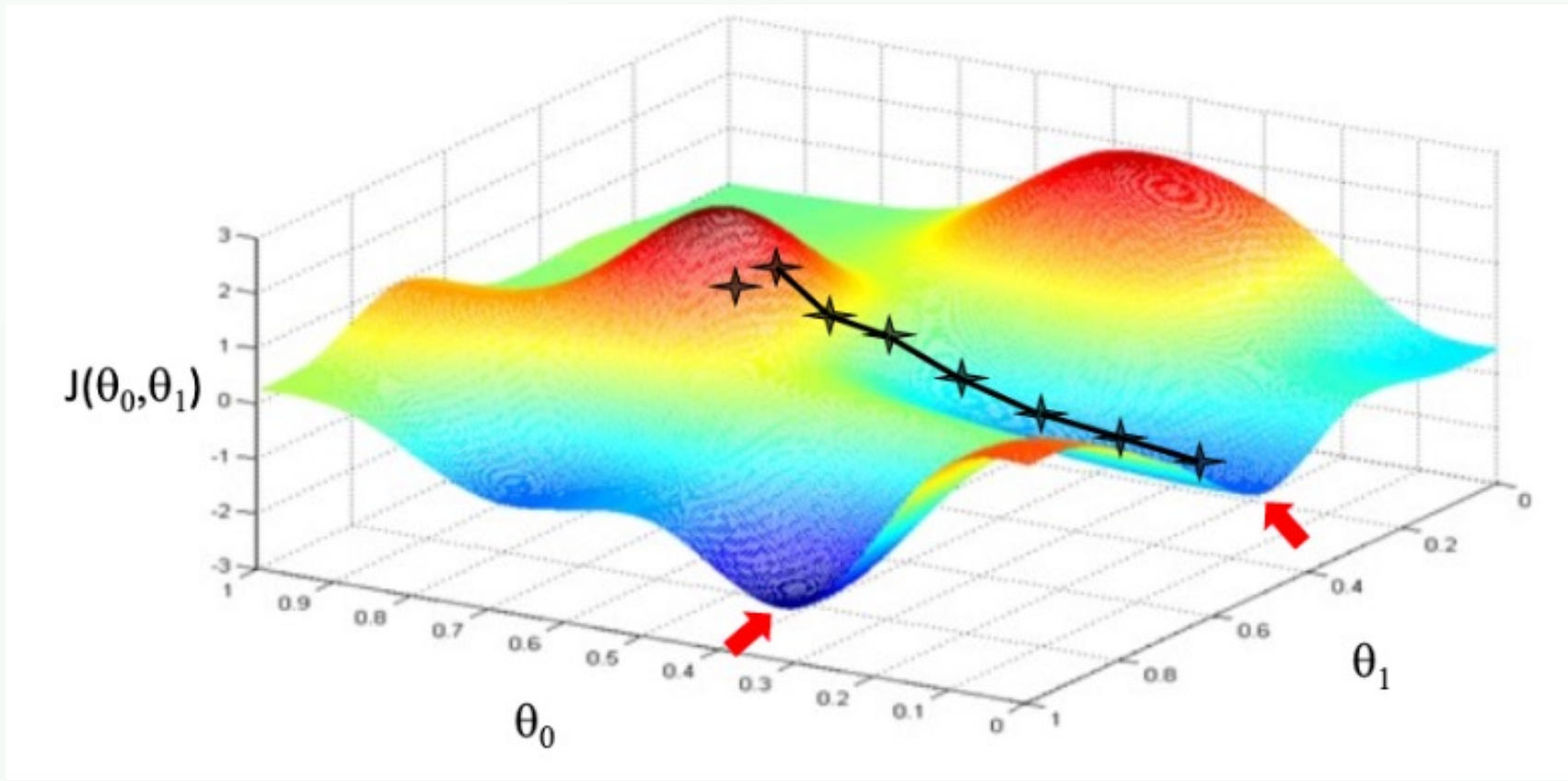
Artificial Intelligence  
& Computer Vision  
Laboratory



# Gradient Descent



Artificial Intelligence  
& Computer Vision  
Laboratory







## Gradient descent algorithm

repeat until convergence {

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

}

Learning rate

Simultaneously update  
 $\theta_0$  &  $\theta_1$

assignment  
 $a := b$

Correct: Simultaneous update

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

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

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

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

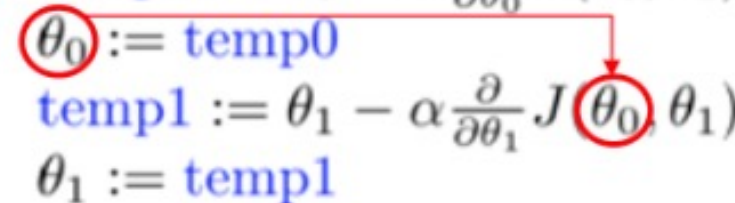
Incorrect:

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

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

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

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





## Gradient descent algorithm

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

Learning rate

derivative

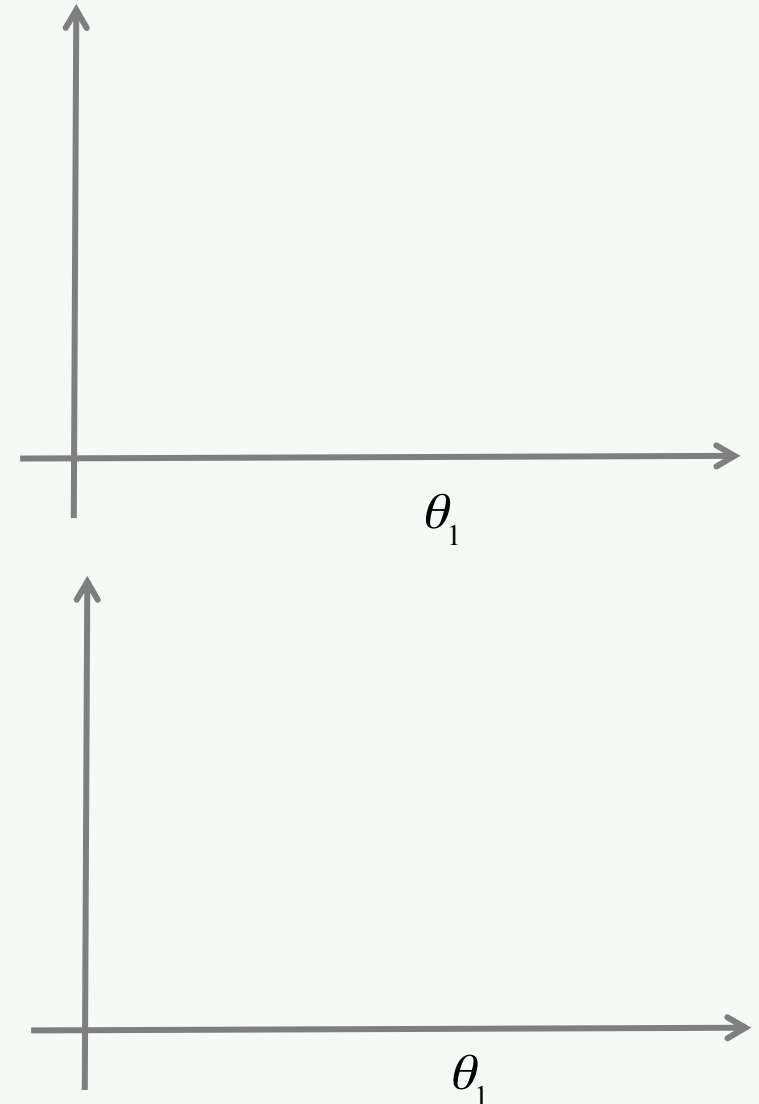
(simultaneously update  
 $j = 0$  and  $j = 1$ )

# Gradient Descent



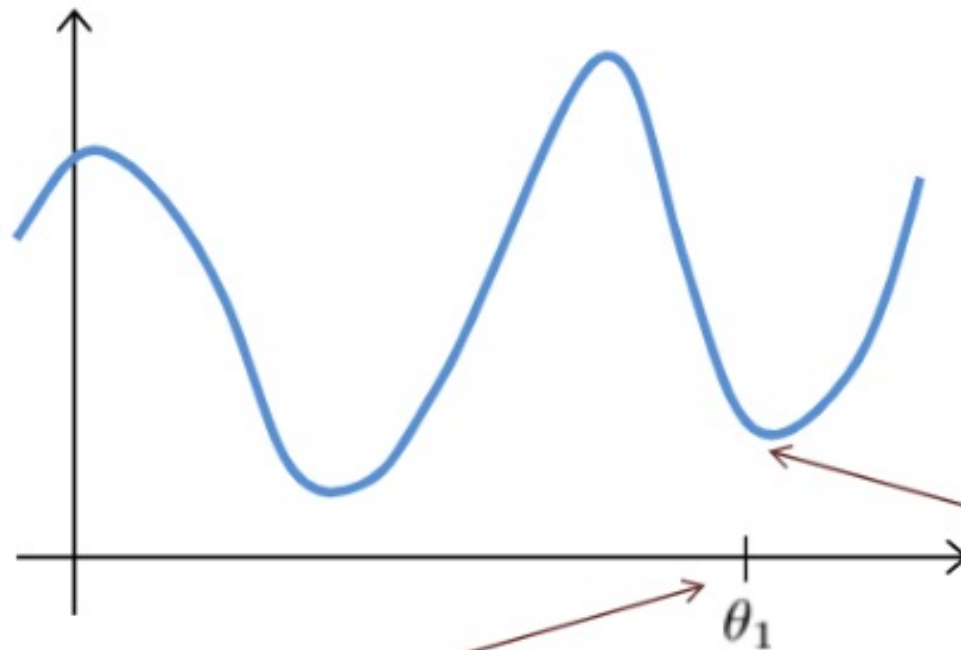
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

- 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





# Gradient Descent



Current value of  $\theta_1$

$\theta_1$  at local optima

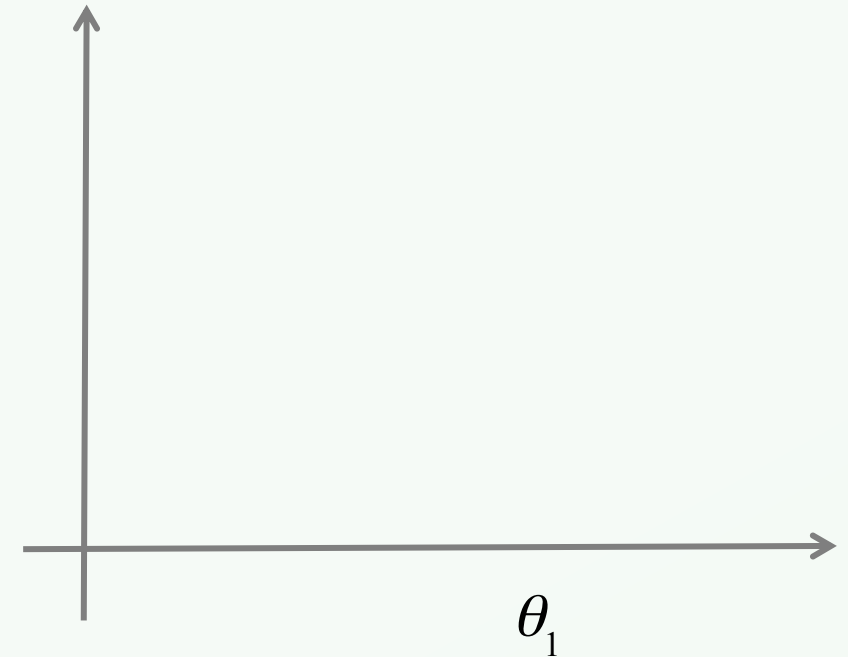
$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$



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

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

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



Artificial Intelligence  
& Computer Vision  
Laboratory

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

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

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

# Gradient Descent for linear regression



Artificial Intelligence  
& Computer Vision  
Laboratory

## Gradient descent algorithm

repeat until convergence {

$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$	}	update $\theta_0$ and $\theta_1$ simultaneously
$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$		

}

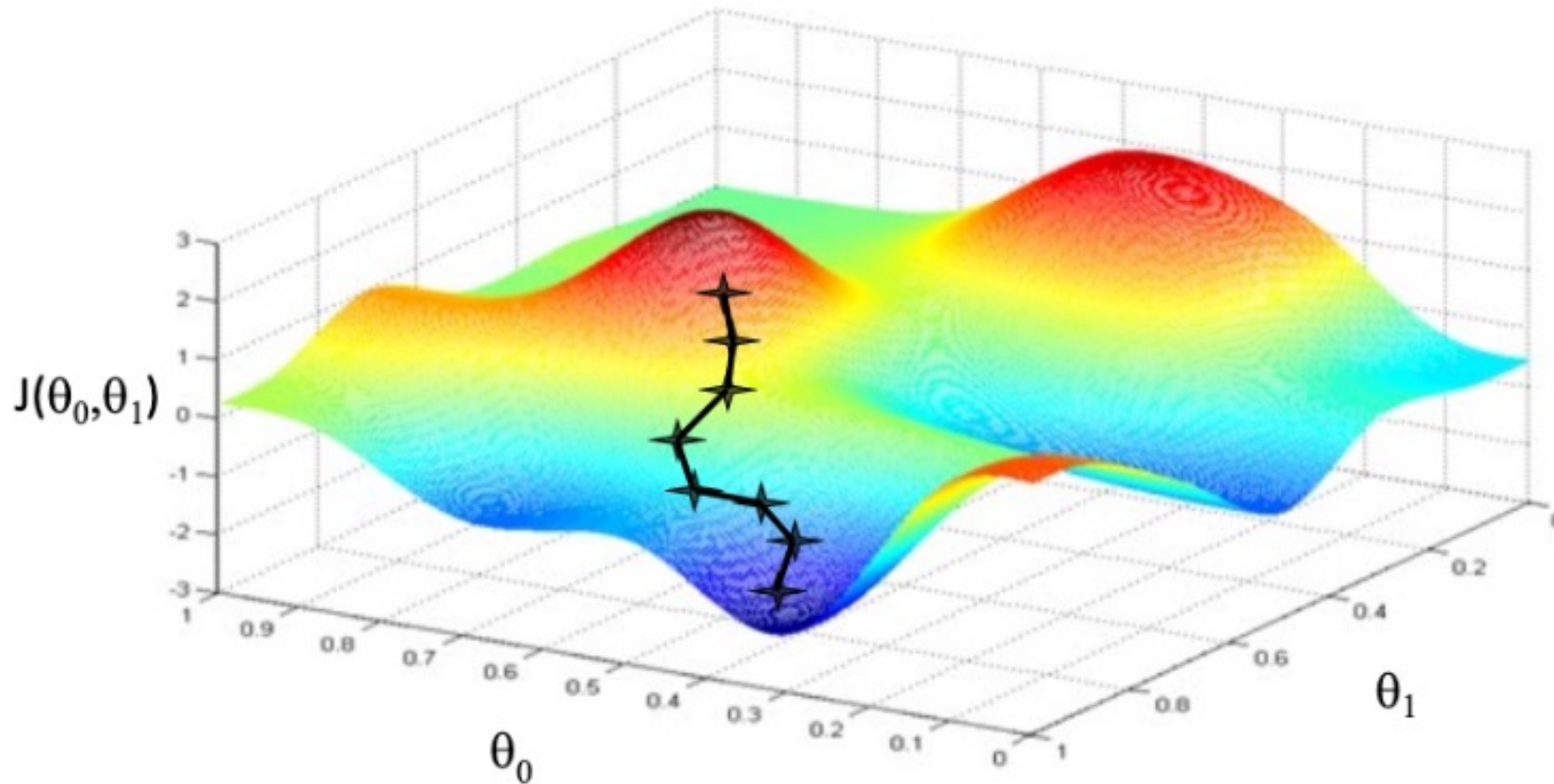
$\frac{d}{d\theta_0} \cdot J(\theta_0, \theta_1)$

$\frac{d}{d\theta_1} \cdot J(\theta_0, \theta_1)$

# Gradient Descent for linear regression



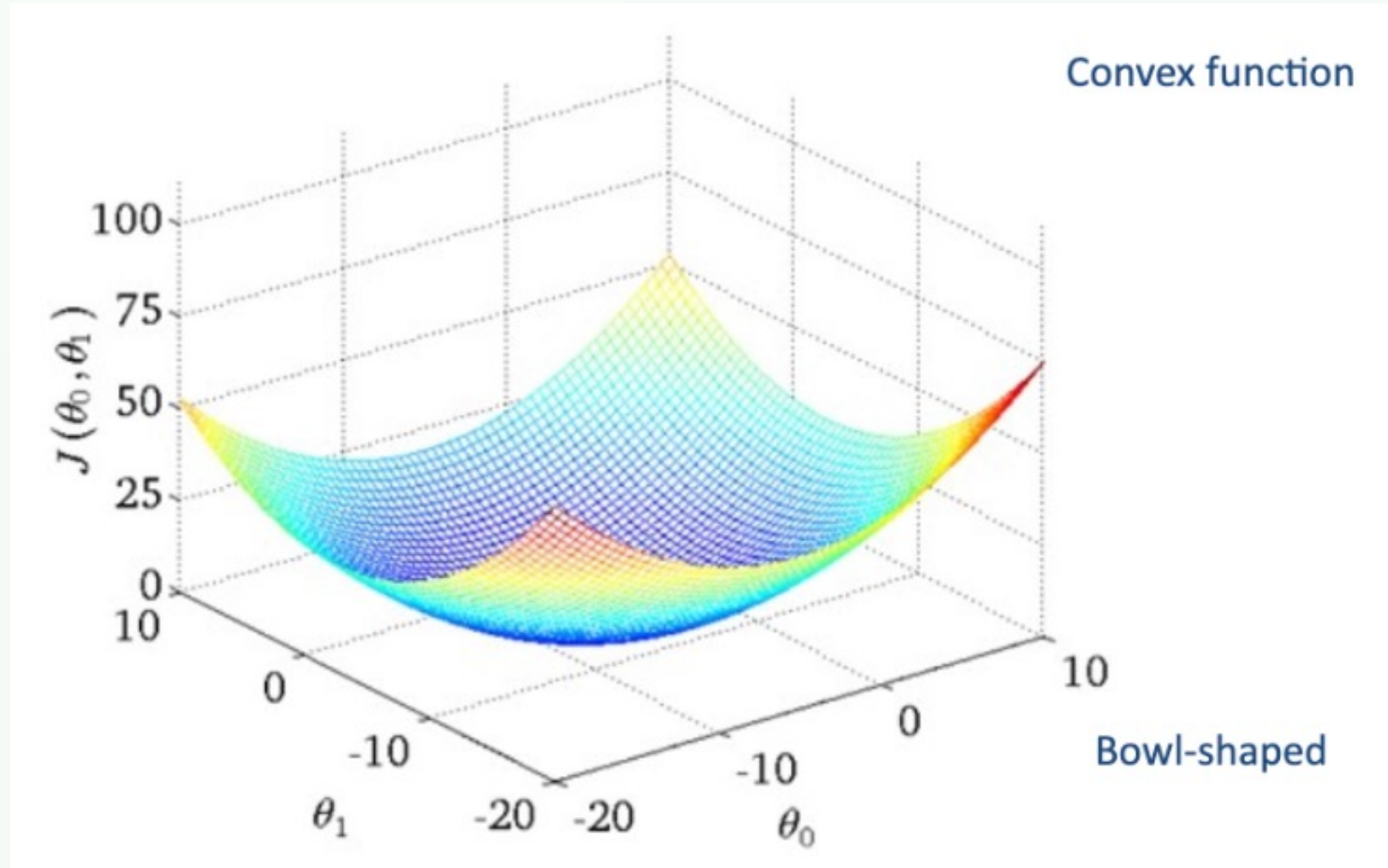
Artificial Intelligence  
& Computer Vision  
Laboratory



# Gradient Descent for linear regression



Artificial Intelligence  
& Computer Vision  
Laboratory





# Gradient Descent for linear regression

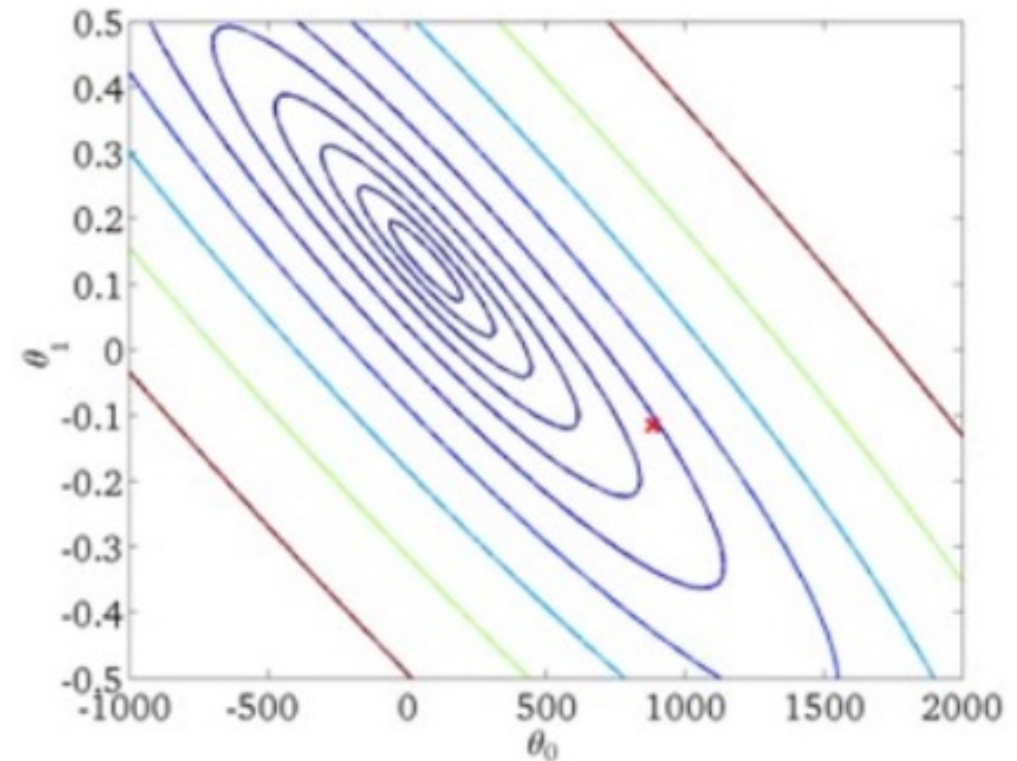
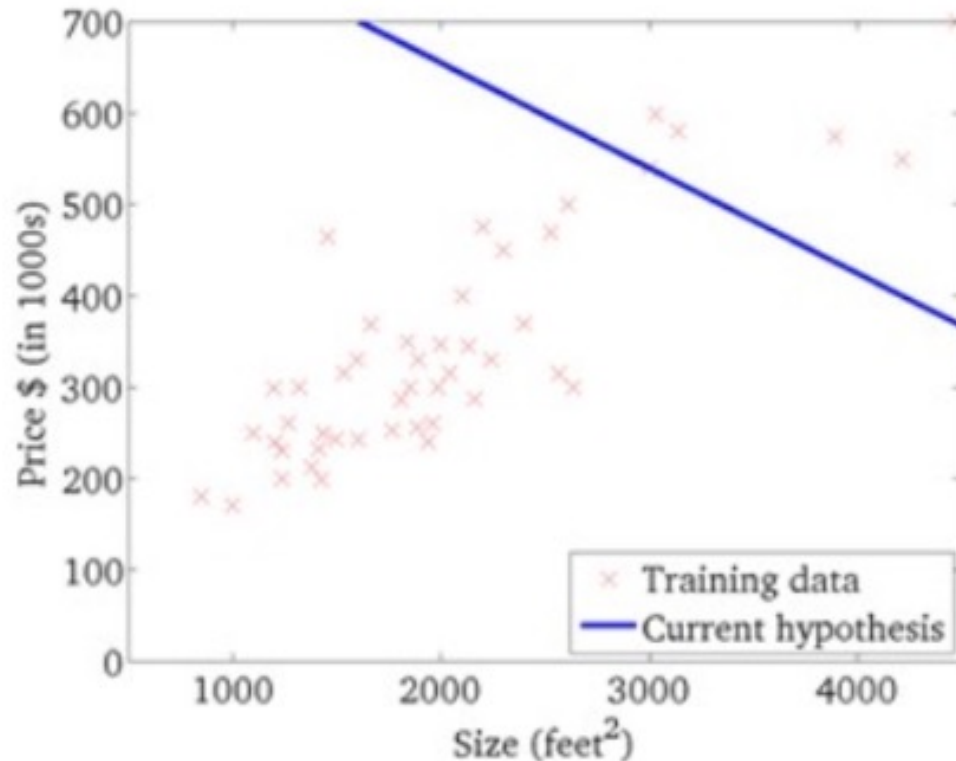


Artificial Intelligence  
& Computer Vision  
Laboratory

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

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

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



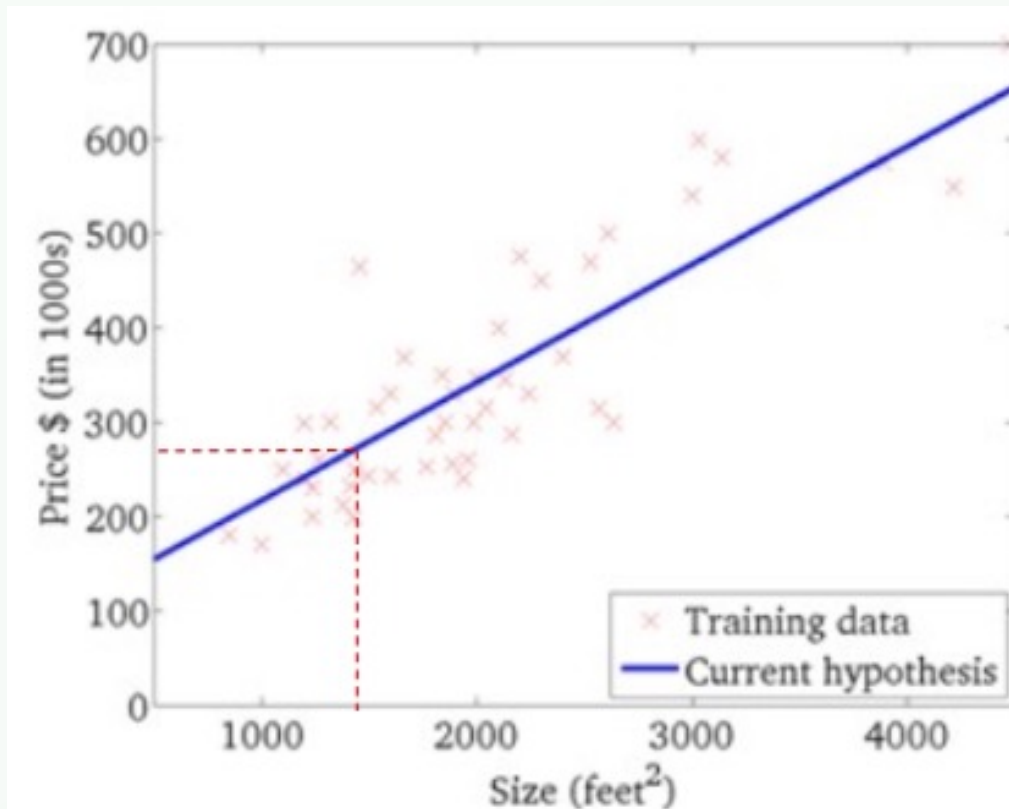
# Gradient Descent for linear regression



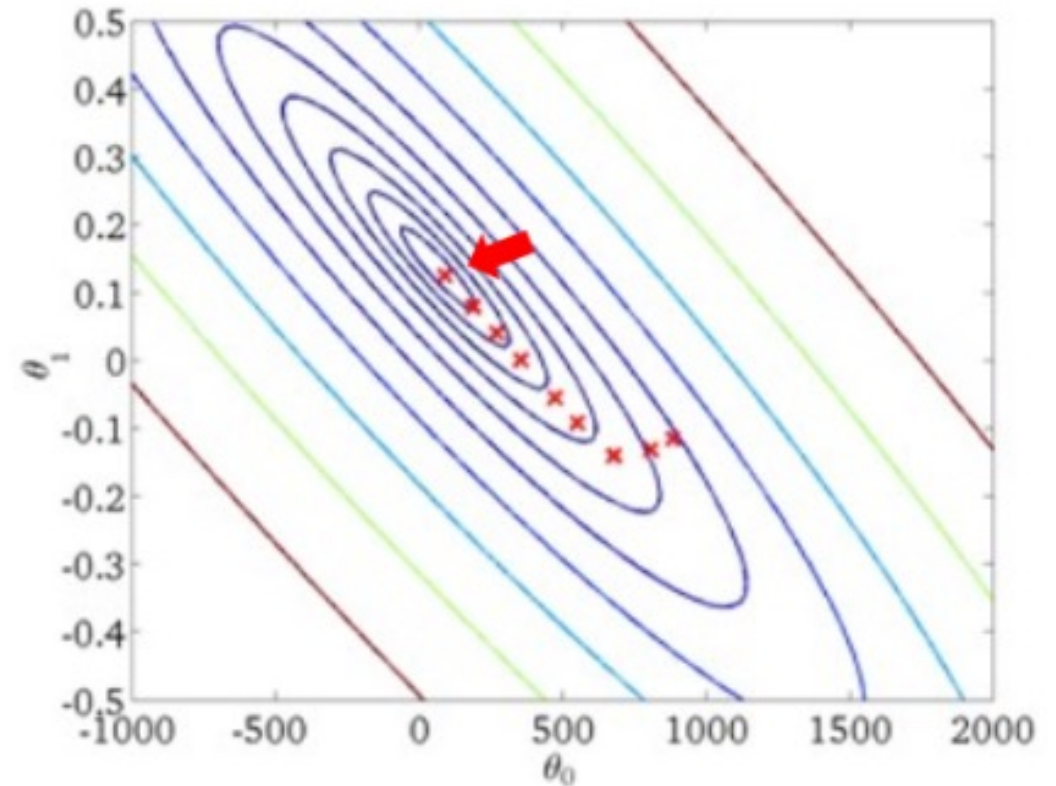
Artificial Intelligence  
& Computer Vision  
Laboratory

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

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



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





# “Batch” Gradient Descent



Artificial Intelligence  
& Computer Vision  
Laboratory

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

Linear regression with multiple variable

# Multiple Features (variables)



## Training set of housing prices

Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
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
- $\mathbf{x}^{(i)}$  = input (features) of  $i^{\text{th}}$  training example
- $x_j^{(i)}$  = value of feature  $j$  in  $i^{\text{th}}$  training example

# Hypothesis



Artificial Intelligence  
& Computer Vision  
Laboratory

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



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

For convenience of notation, define  $x_0 = 1$

Multivariate linear regression

# Gradient Descent for multiple variables



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

Parameters:  $\theta_0, \theta_1, \dots, \theta_n$

Cost function:

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

Gradient descent:

Repeat {

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

}

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

# Gradient Descent for multiple variables



## Gradient Descent

Previously ( $n=1$ ):

Repeat {

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

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

(simultaneously update  $\theta_0, \theta_1$ )

}

New algorithm ( $n \geq 1$ ):

Repeat {

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

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

}

---

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

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

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

...

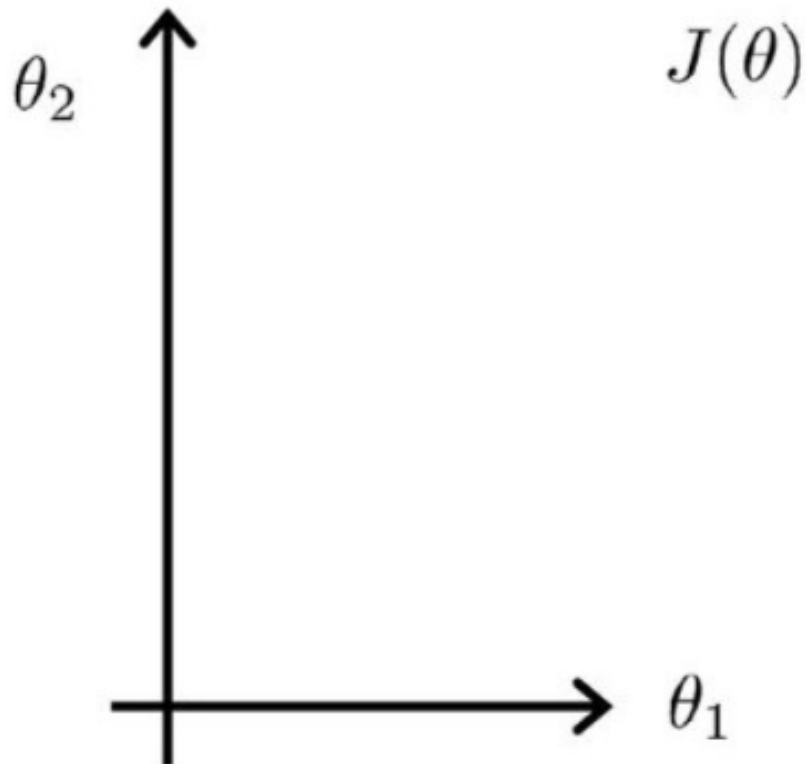
# Gradient Descent in practice : Feature scaling



Idea: Make sure features are on a similar scale

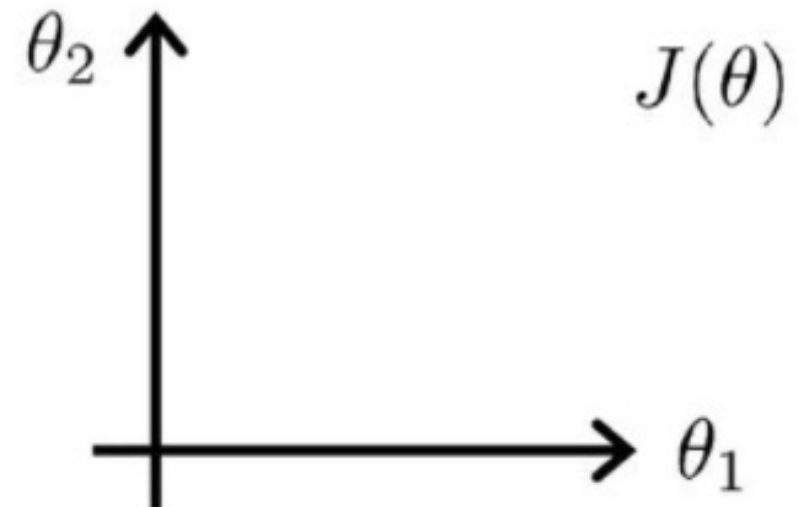
E.g.  $x_1 = \text{size (0-2000 feet}^2\text{)}$

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



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

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





# Gradient Descent in practice

## : Feature scaling



Artificial Intelligence  
& Computer Vision  
Laboratory

Get every feature into approximately a  $-1 \leq x_i \leq 1$  range

# Gradient Descent in practice

## : Feature scaling



Artificial Intelligence  
& Computer Vision  
Laboratory

### 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.

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

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

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

# Gradient Descent in practice : Learning rate



Artificial Intelligence  
& Computer Vision  
Laboratory

## Gradient decent

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

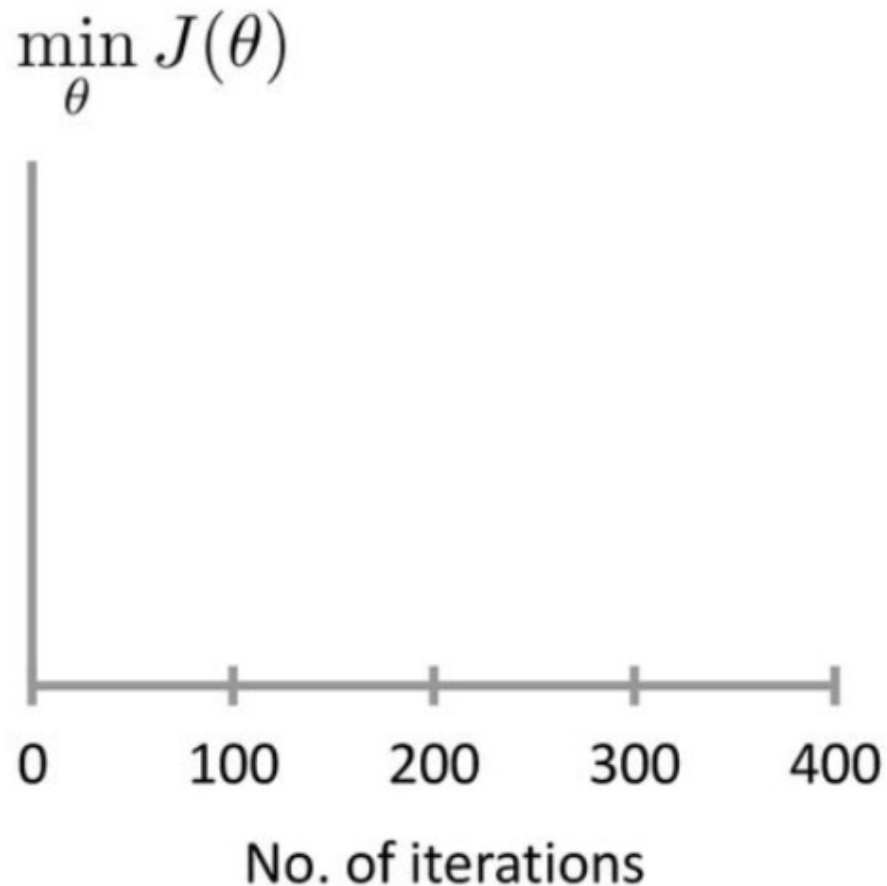
- “Debugging”: How to make sure gradient descent is working correctly
- How to choose learning rate  $\alpha$

# Gradient Descent in practice : Learning rate



Artificial Intelligence  
& Computer Vision  
Laboratory

**Make sure gradient descent is working correctly**



Example automatic  
convergence test:

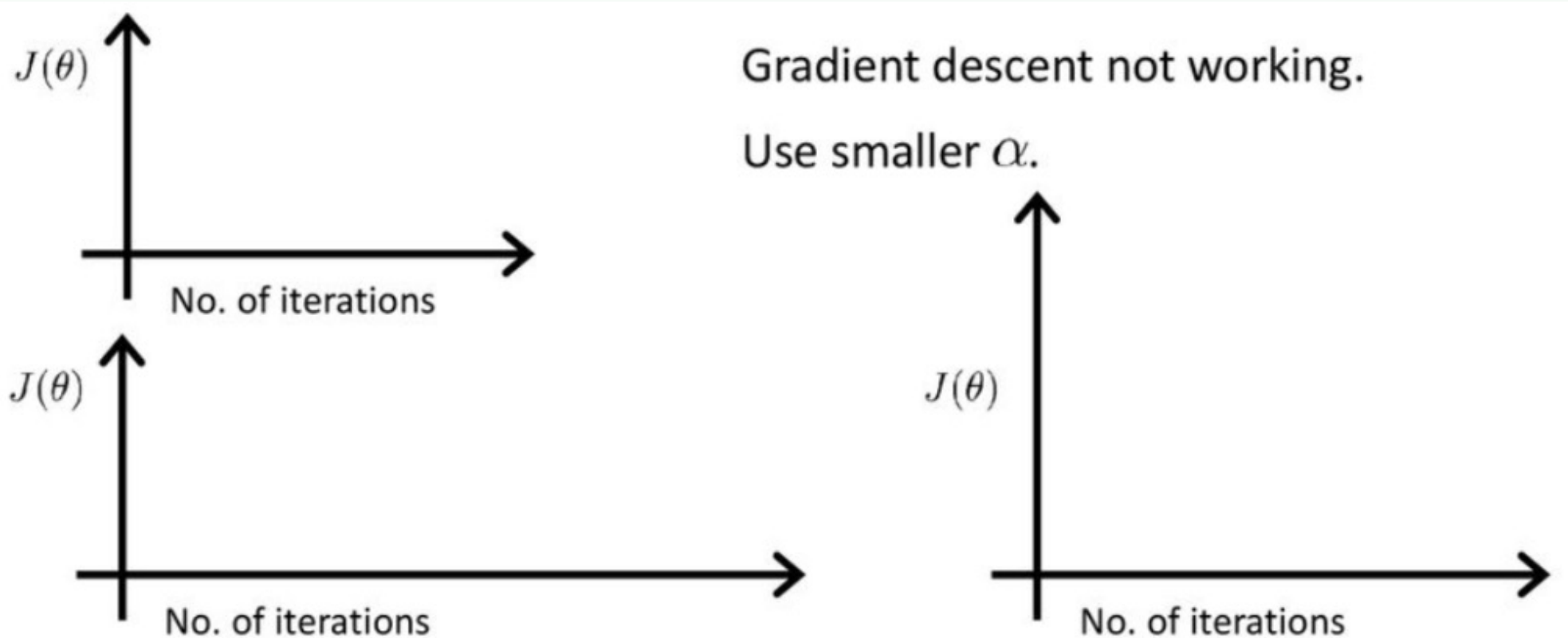
Declare convergence if  $J(\theta)$   
decreases by less than  $10^{-3}$   
in one iteration.

# Gradient Descent in practice

## : Learning rate



**Make sure gradient descent is working correctly**



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

# Gradient Descent in practice

## : Learning rate



Artificial Intelligence  
& Computer Vision  
Laboratory

### Summary

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

To choose learning rate  $\alpha$ , try

..., 0.001, , 0.01, 0.1, 1, .....



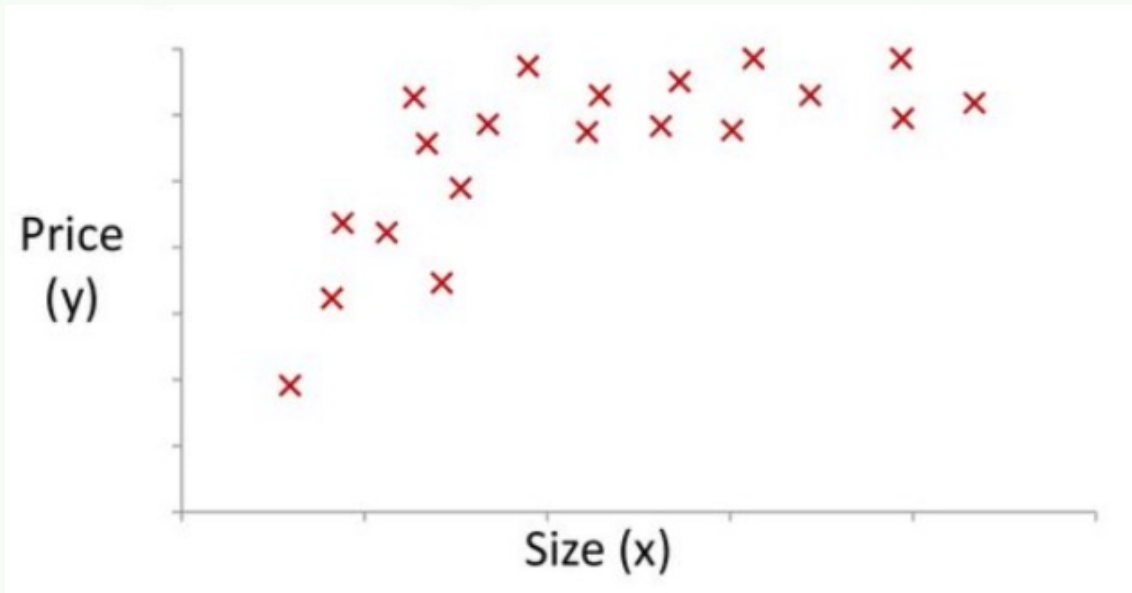
## Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \textit{frontage} + \theta_2 \times \textit{depth}$$





# Polynomial Regression



$$\theta_0 + \theta_1 x_1 + \theta_2 x_2$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

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

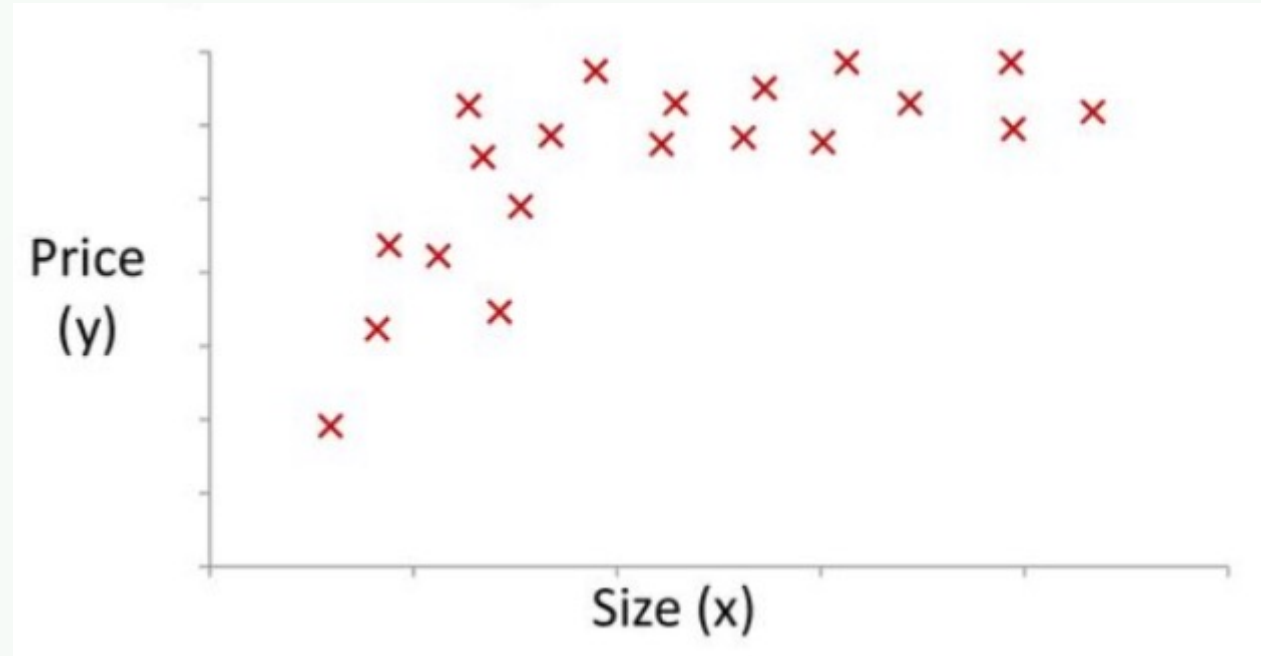
$$x_1 = (size)$$

$$x_2 = (size)^2$$

$$x_3 = (size)^3$$



## Choice of features



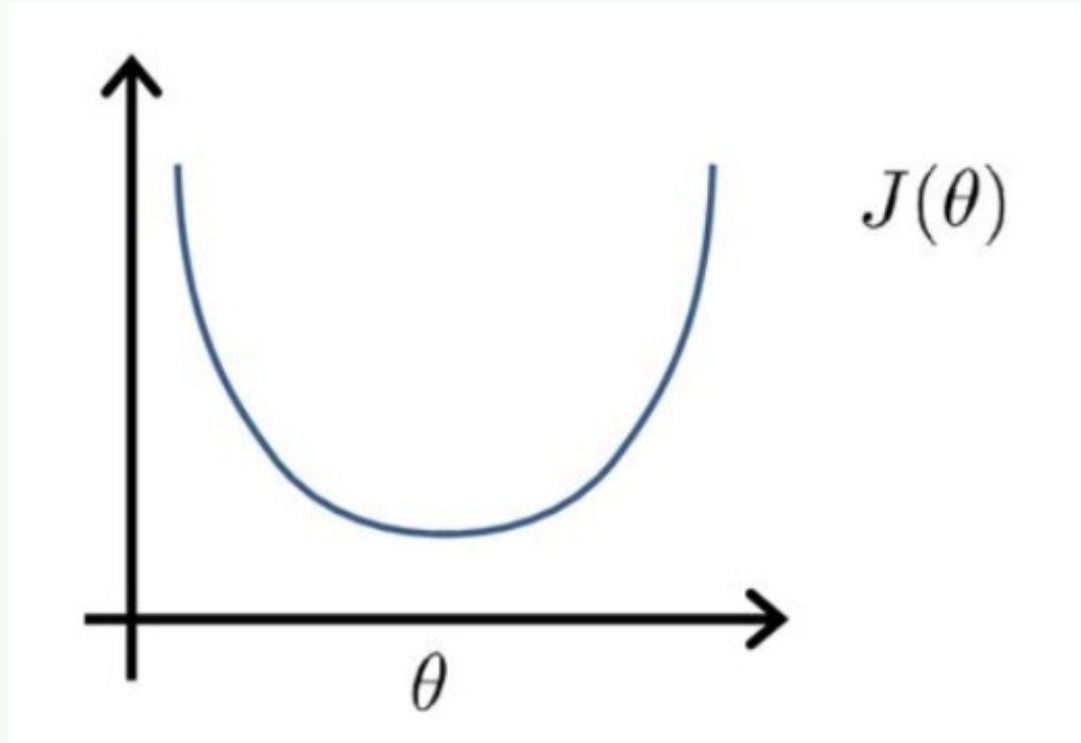
$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

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

# Normal Equation



## Gradient Descent



Normal equation:

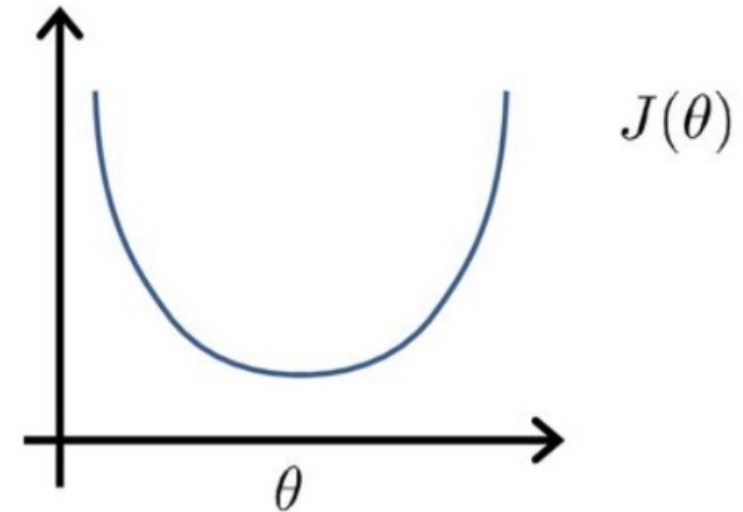
Method to solve for  $\theta$  analytically

# Normal Equation



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

$$J(\theta) = a\theta^2 + b\theta + c$$



---

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

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

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

# Normal Equation



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

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

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

# Normal Equation



$m=4$  examples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$ :  $n$  features

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

E.g. If  $\mathbf{x}^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$

# Normal Equation



Artificial Intelligence  
& Computer Vision  
Laboratory

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

$\left( X^T X \right)^{-1}$  is inverse of matrix  $X^T X$

Pinv( )



# Normal Equation



$m=4$  examples  $(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})$ :  $n$  features

## Gradient Descent

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

## Normal Equation

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



## Normal equation

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

- What if  $X^T X$  is non-invertible?



What if  $X^T X$  is non-invertible?

- Redundant features (linearly dependent)  
E.g.  $x_1 = \text{size in feet}^2$   
 $x_2 = \text{size in } m^2$
- Too many features (e.g.  $m \leq n$ )  
→ Delete some features, or use regularization