

Machine Learning

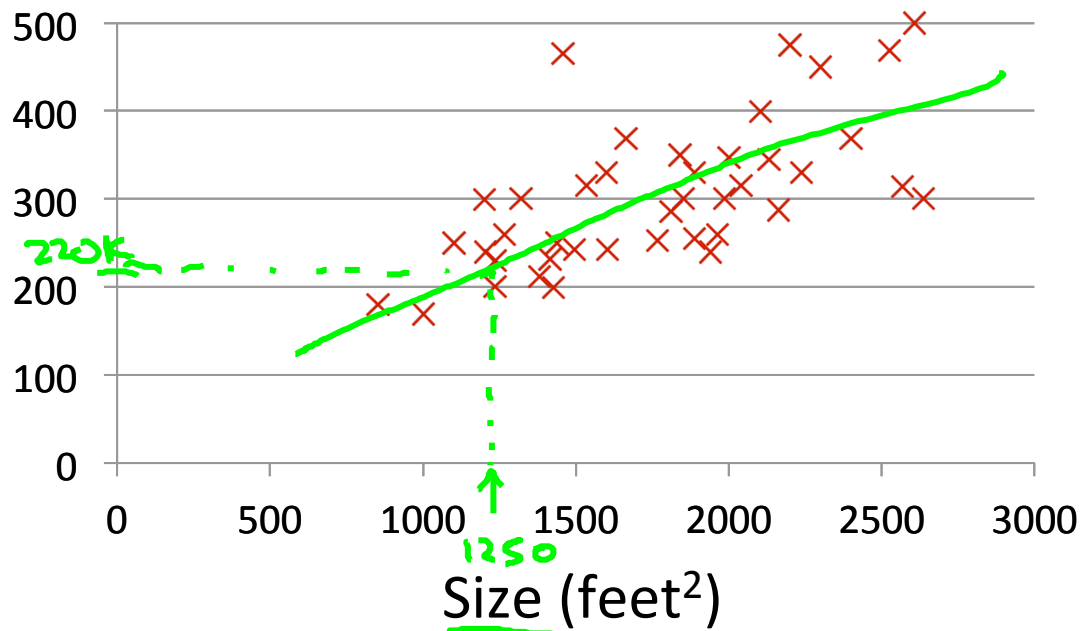
Linear regression  
with one variable

---

Model  
representation

# 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
...	...

$m = 47$

Notation:

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

$(x, y)$  - one training example

$(x^{(i)}, y^{(i)})$  -  $i$ th training example

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

Training Set

Learning Algorithm

Size of house

$h$

Estimated price

hypothesis

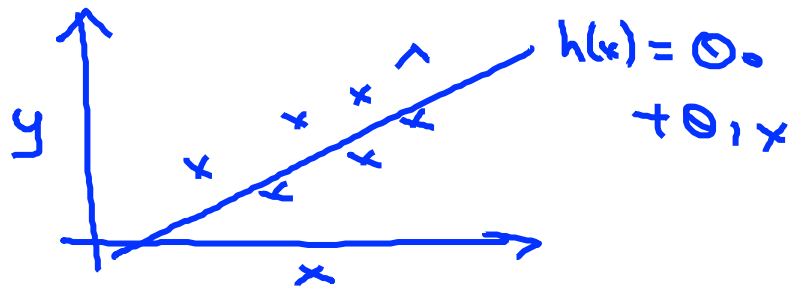
(estimated value of  $y$ )

$h$  maps from  $x$ 's to  $y$ 's.

How do we represent  $h$  ?

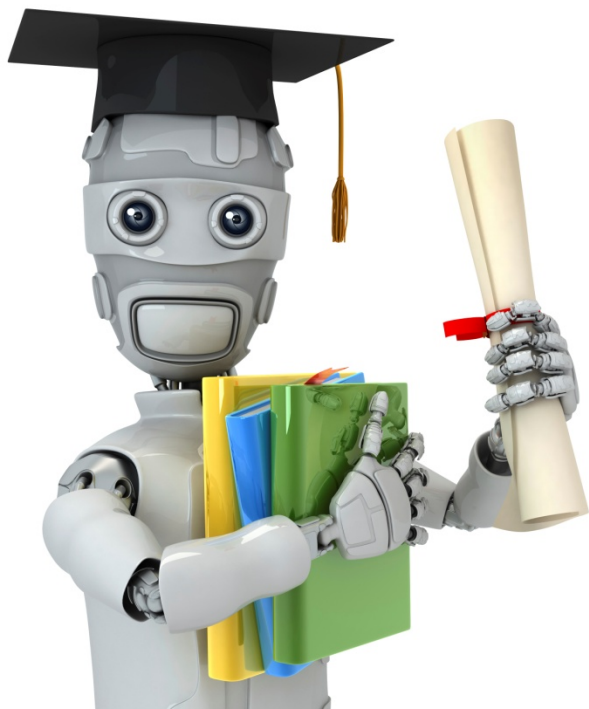
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

Shorthand:  $h(x)$



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

↳ one variable



Machine Learning

Linear regression  
with one variable

---

Cost function

# Training Set

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

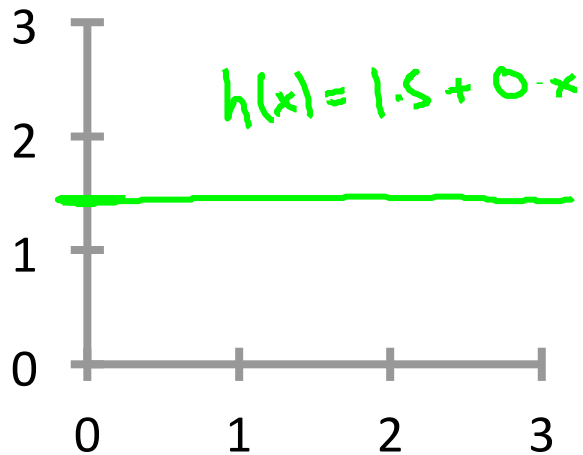
}  $n = 47$

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

$\theta_i$ 's: Parameters

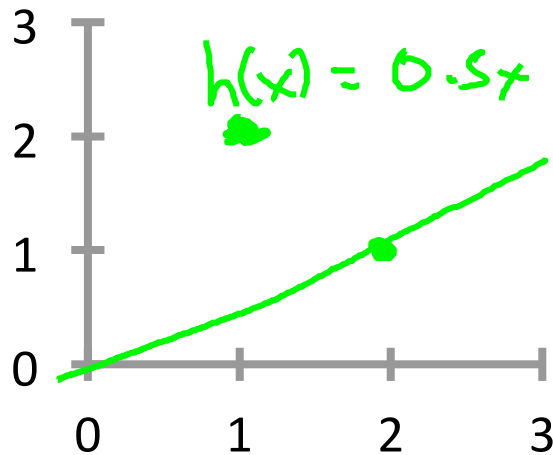
How to choose  $\theta_i$ 's ?

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



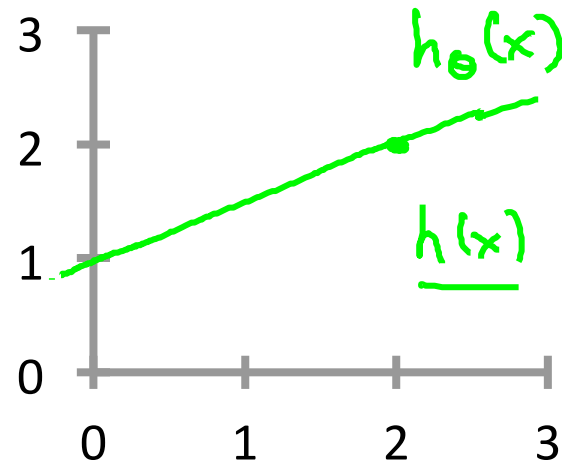
$$\rightarrow \theta_0 = 1.5$$

$$\rightarrow \theta_1 = 0$$



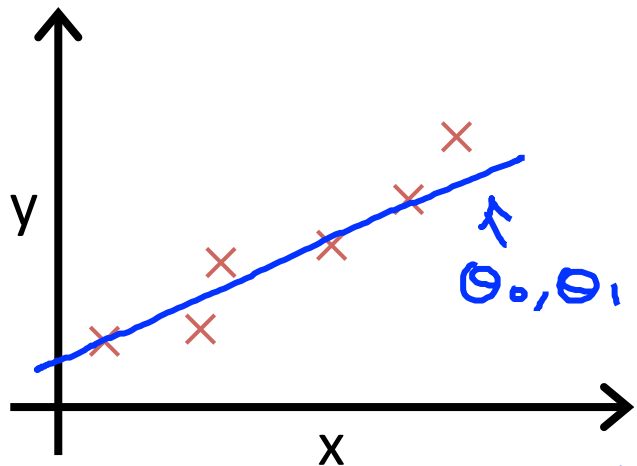
$$\rightarrow \theta_0 = 0$$

$$\rightarrow \theta_1 = 0.5$$



$$\rightarrow \theta_0 = 1$$

$$\rightarrow \theta_1 = 0.5$$



$(x^{(i)}, y^{(i)})$

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

$x, y$

minimize  $\theta_0, \theta_1$

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

# training examples

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

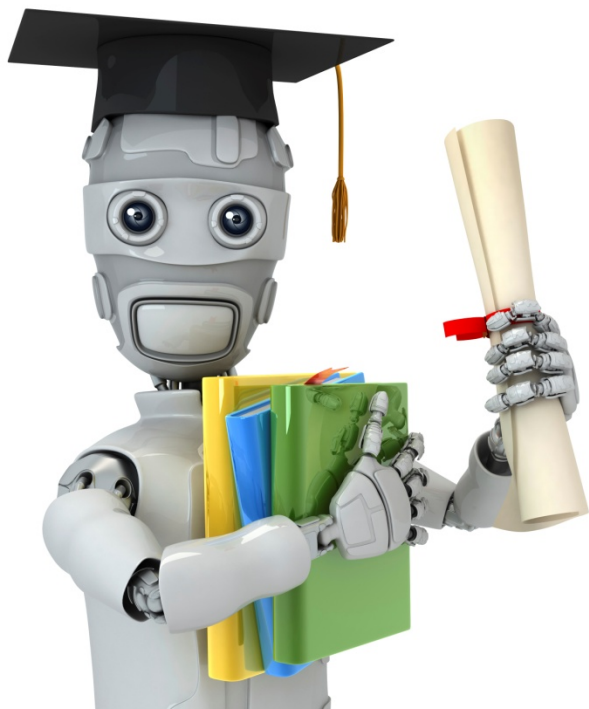
$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2$

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

Cost function

Squared error function





Machine Learning

Linear regression  
with one variable

---

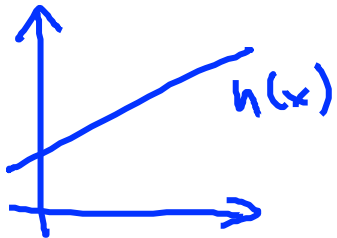
Cost function  
intuition I

Hypothesis:

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

Parameters:

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



Cost Function:

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

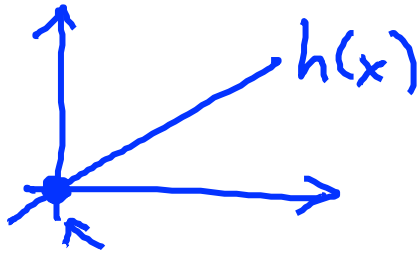
Goal: minimize  $J(\theta_0, \theta_1)$   
 $\nearrow \theta_0, \theta_1$

Simplified

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

$$\theta_0 = 0$$

$$\underline{\theta_1}$$



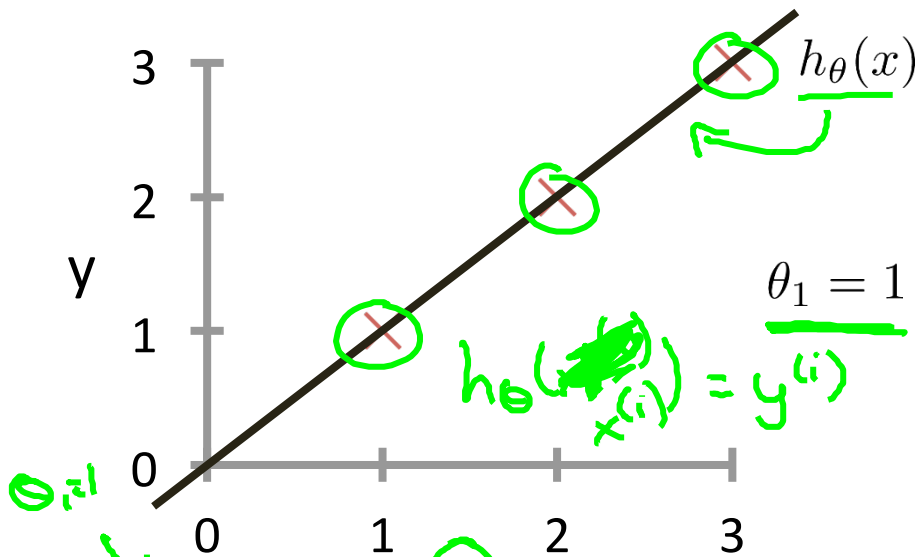
$$J(\theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

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

$$\theta_1, x^{(i)}$$

→  $h_{\theta}(x)$

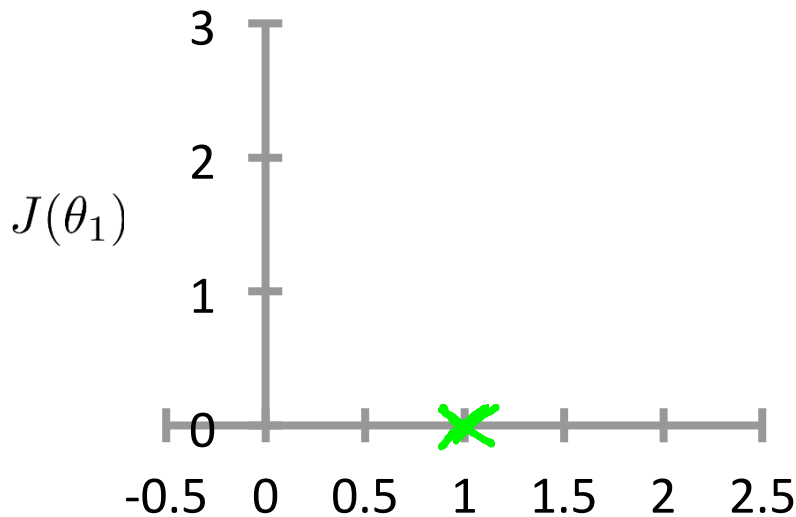
(for fixed  $\theta_1$ , this is a function of  $x$ )



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

→  $J(\theta_1)$

(function of the parameter  $\theta_1$ )

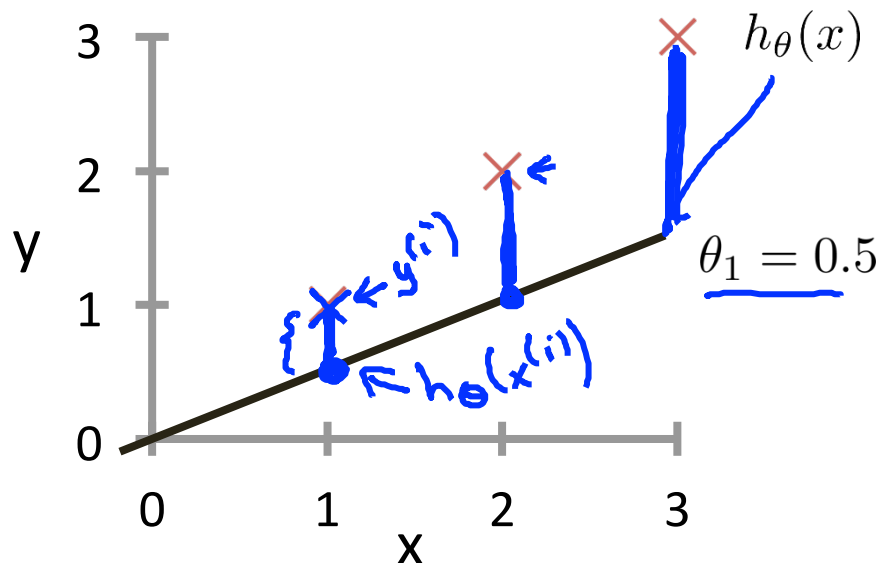


$\theta_1 = 0.5?$

$$\underline{J(1) = 0}$$

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

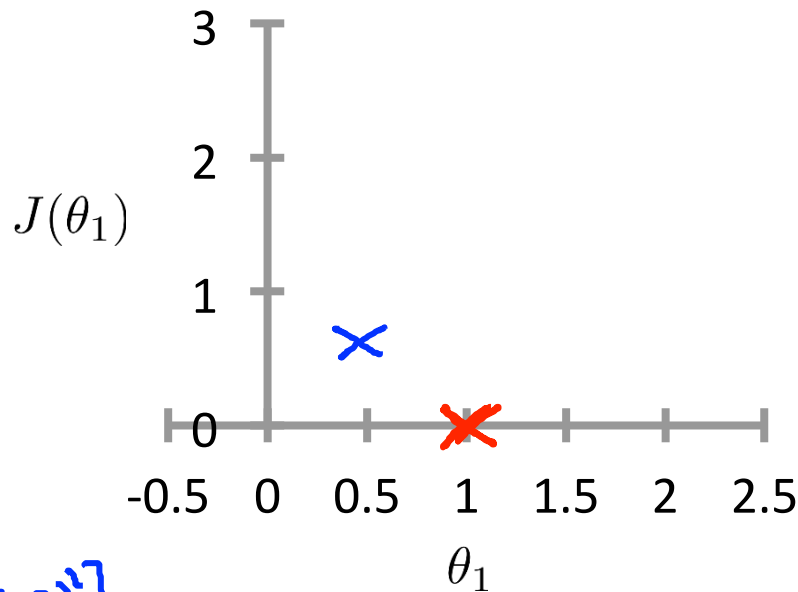
(for fixed  $\theta_1$ , this is a function of  $x$ )



$$\begin{aligned} J(0.5) &= \frac{1}{2M} [(0.5-1)^2 + (1-2)^2 + (1.5-3)^2] \\ &= \frac{1}{2 \times 3} (3.5) = \frac{3.5}{6} \approx \underline{0.58} \end{aligned}$$

$$J(\theta_1)$$

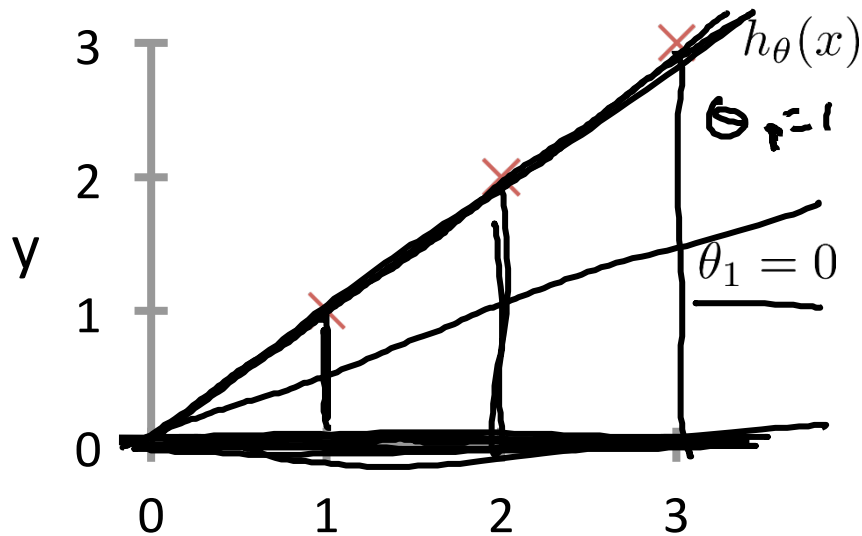
(function of the parameter  $\theta_1$ )



$$\begin{aligned} \theta_1 &= 0? \\ J(0) &=? \end{aligned}$$

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

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



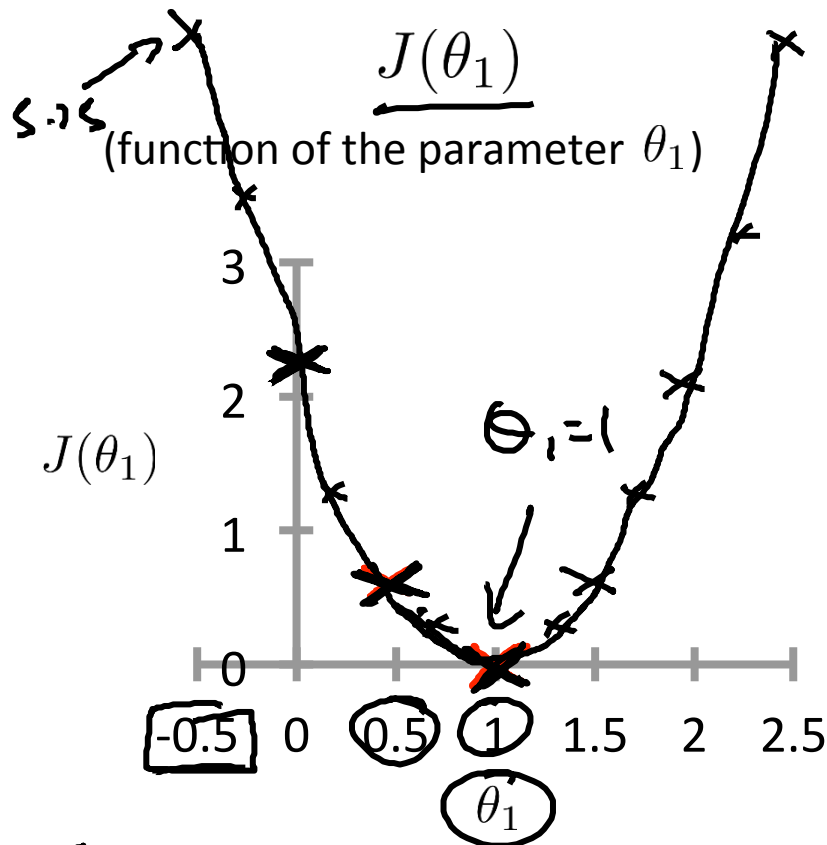
$$\begin{aligned} J(0) &= \frac{1}{2m} (1^2 + 2^2 + 3^2) \\ &= \frac{1}{6} \cdot 14 \approx 2.3 \end{aligned}$$

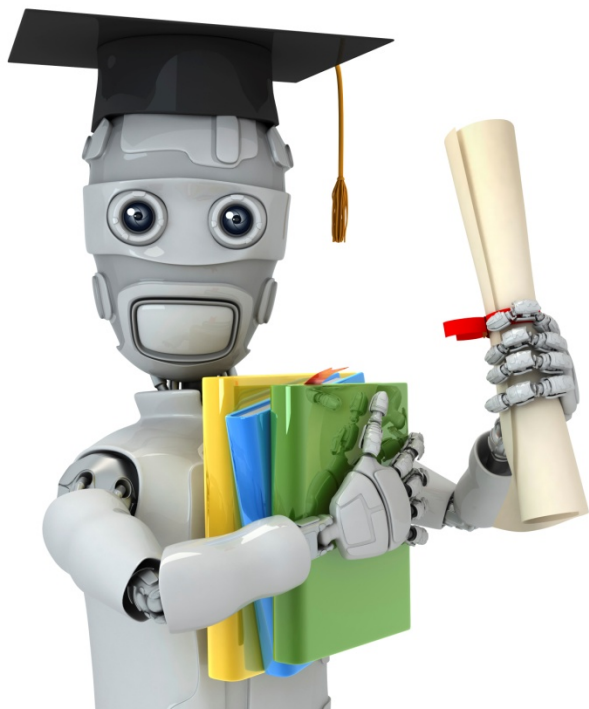
$$h(x) = -0.5x$$

minimize  $J(\theta_1)$

$$J(\theta_1)$$

(function of the parameter  $\theta_1$ )





Machine Learning

Linear regression  
with one variable

---

Cost function  
intuition II

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

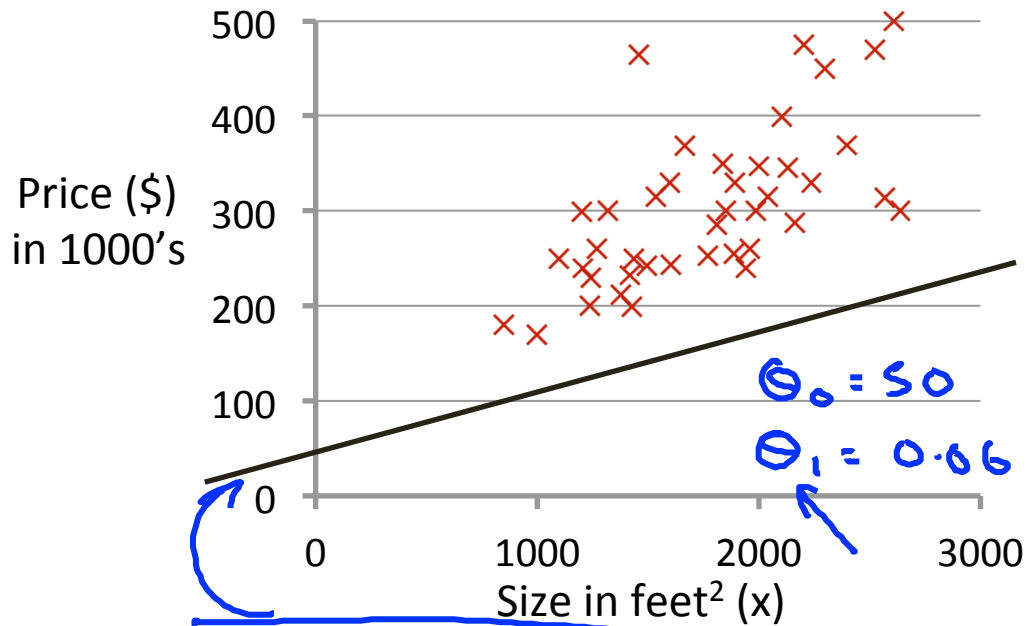
Parameters:  $\theta_0, \theta_1$

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

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

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

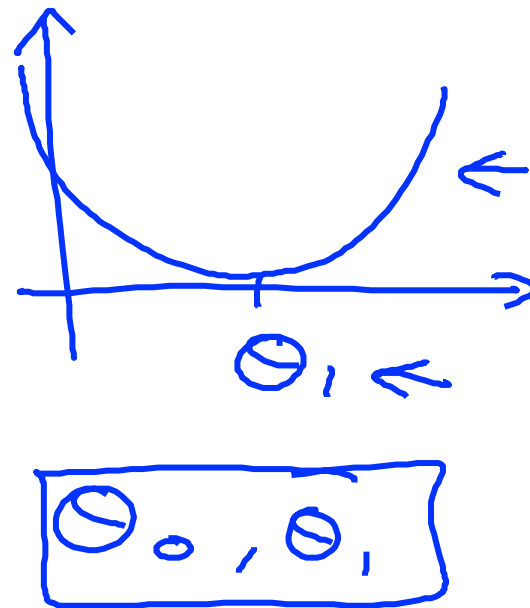
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

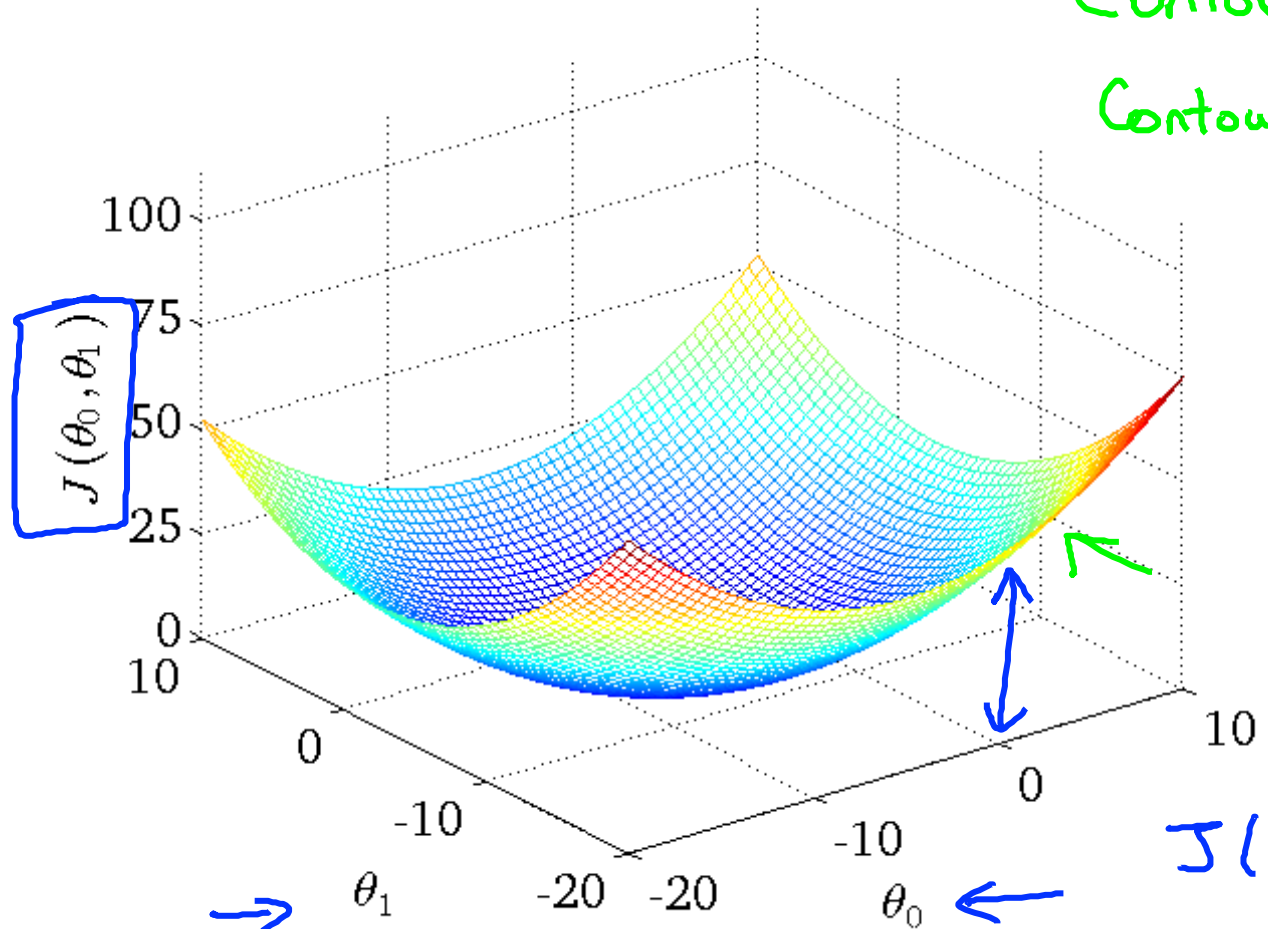
$$\underline{\underline{J(\theta_0, \theta_1)}}$$

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





Contour plots  
Contour figures -



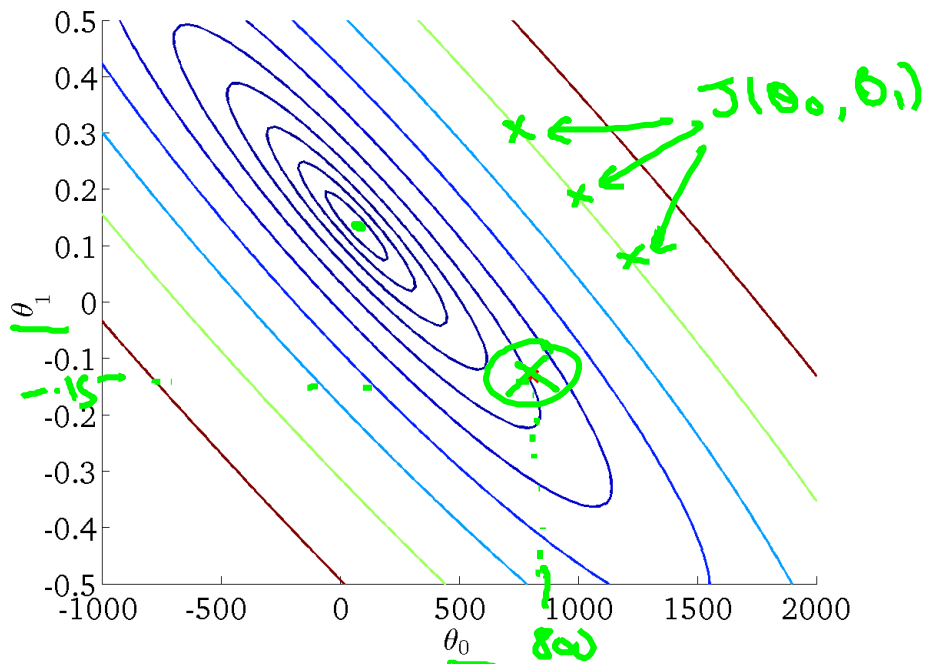
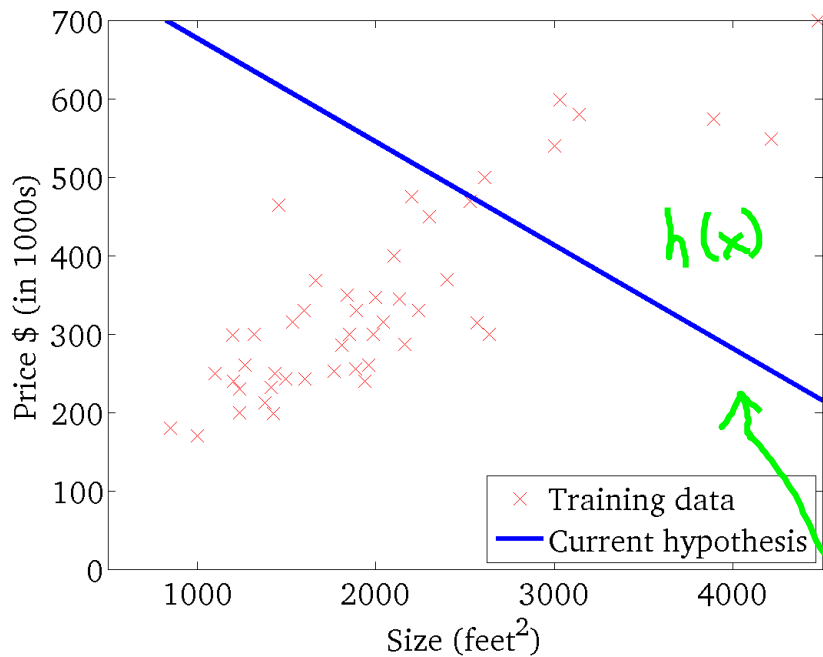
$J(\theta_0, \theta_1)$

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

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

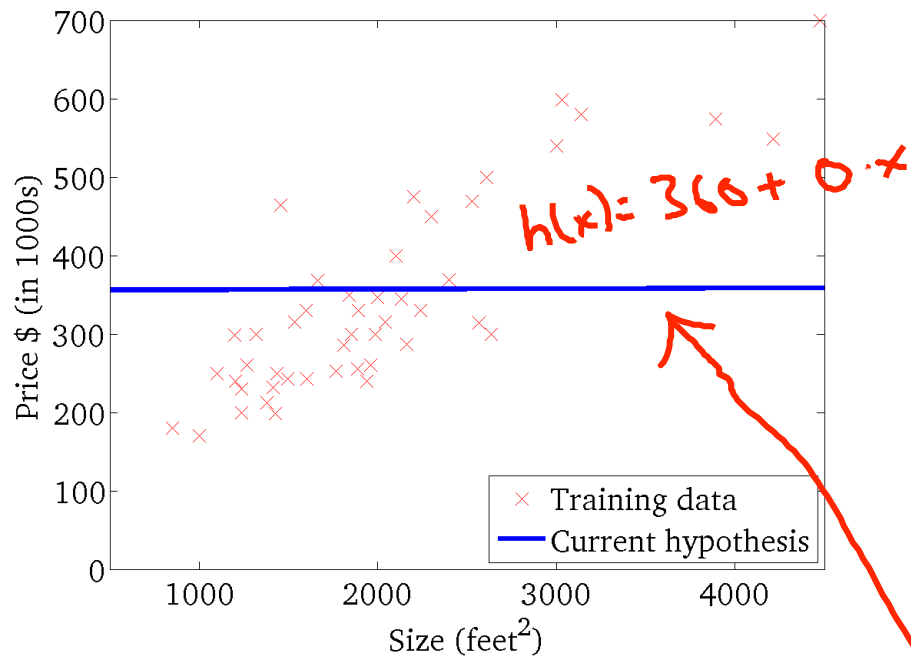
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )

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



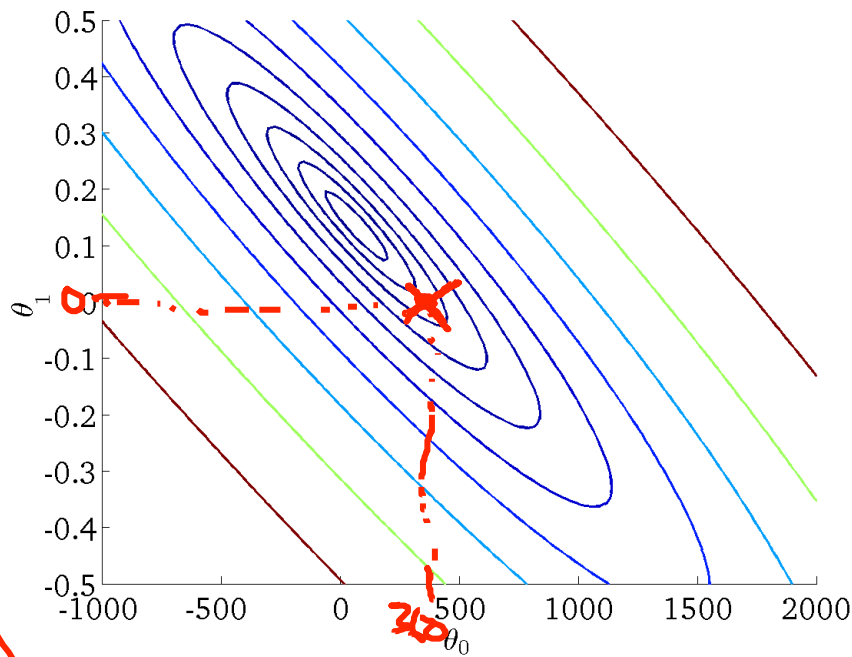
$$h_{\theta}(x)$$

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



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

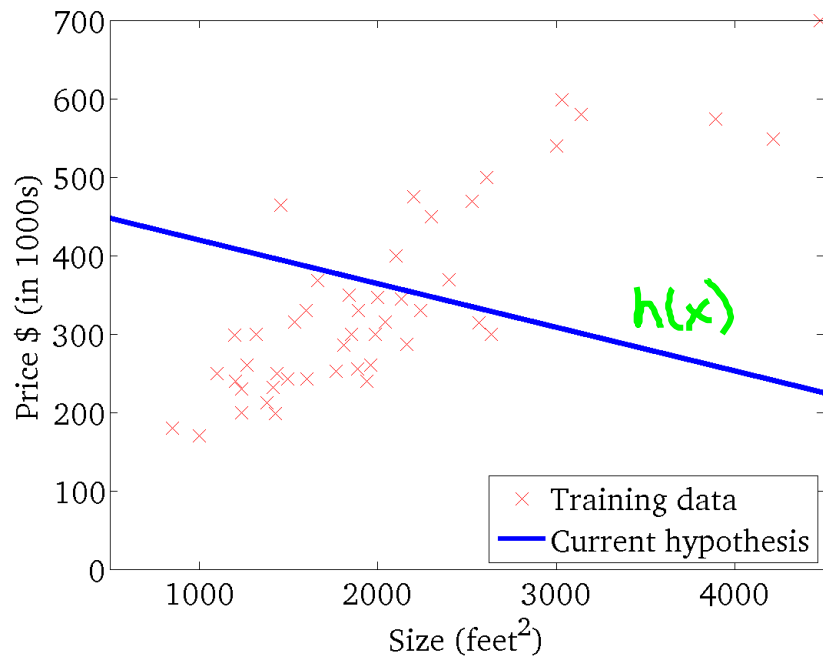
(function of the parameters  $\theta_0, \theta_1$ )



$$\begin{cases} \theta_0 = 360 \\ \theta_1 = 0 \end{cases}$$

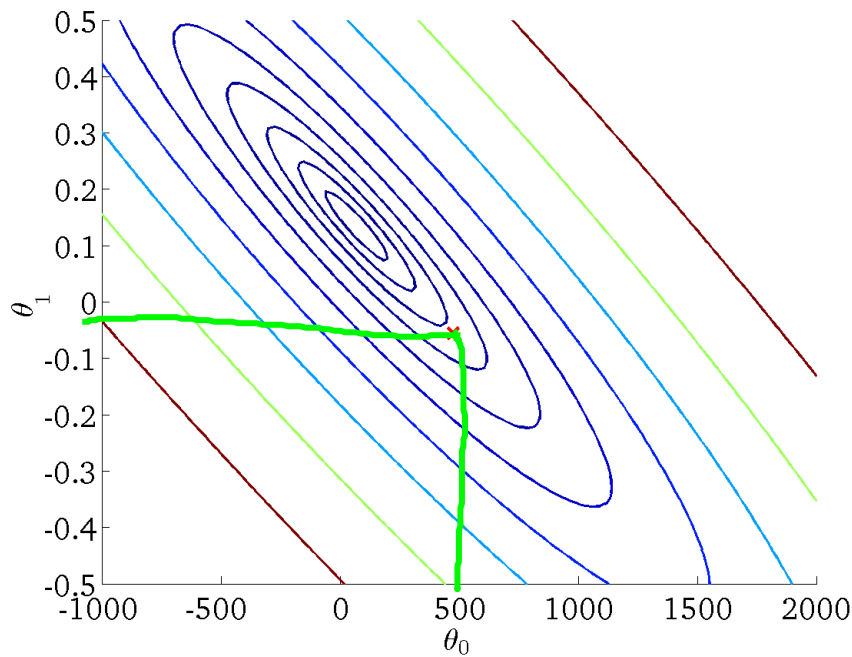
$$h_{\theta}(x)$$

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



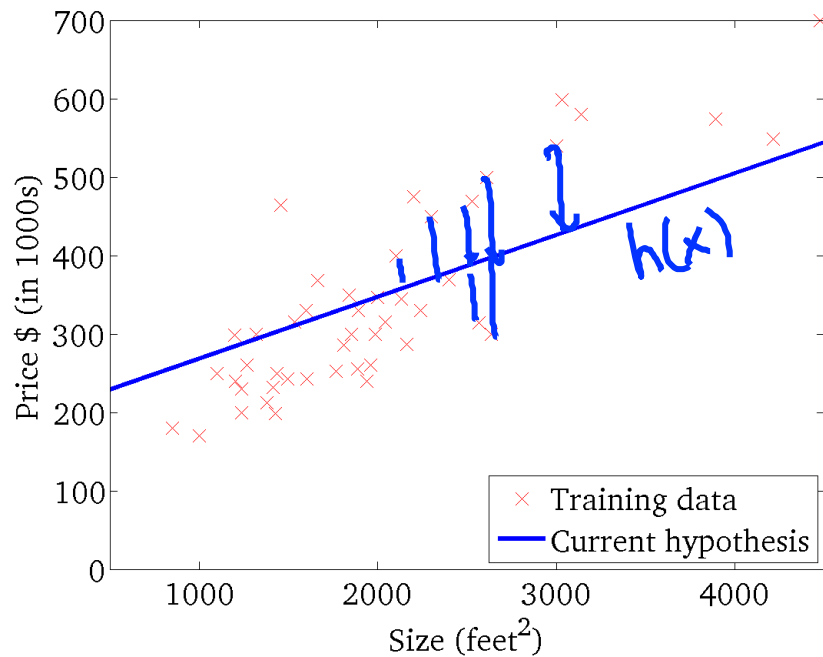
$$J(\theta_0, \theta_1)$$

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



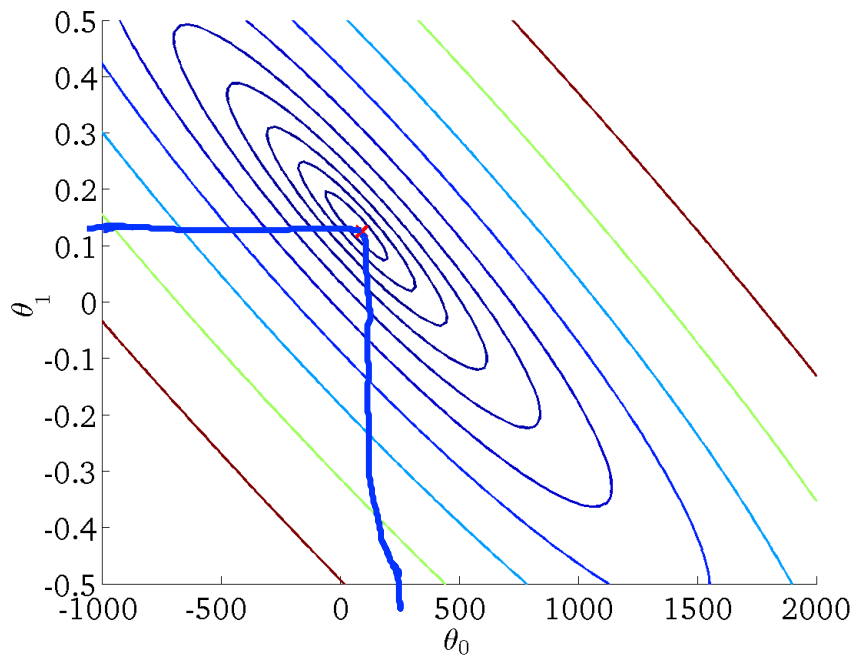
$$h_{\theta}(x)$$

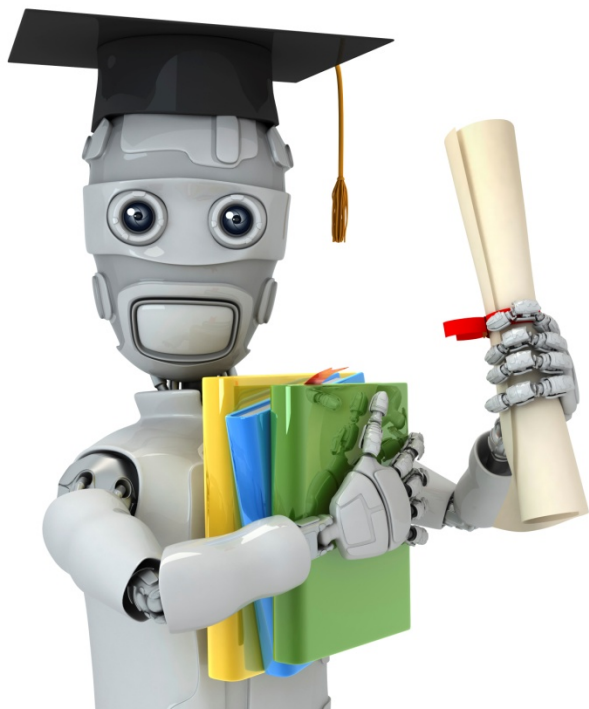
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

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





Machine Learning

Linear regression  
with one variable

---

Gradient  
descent

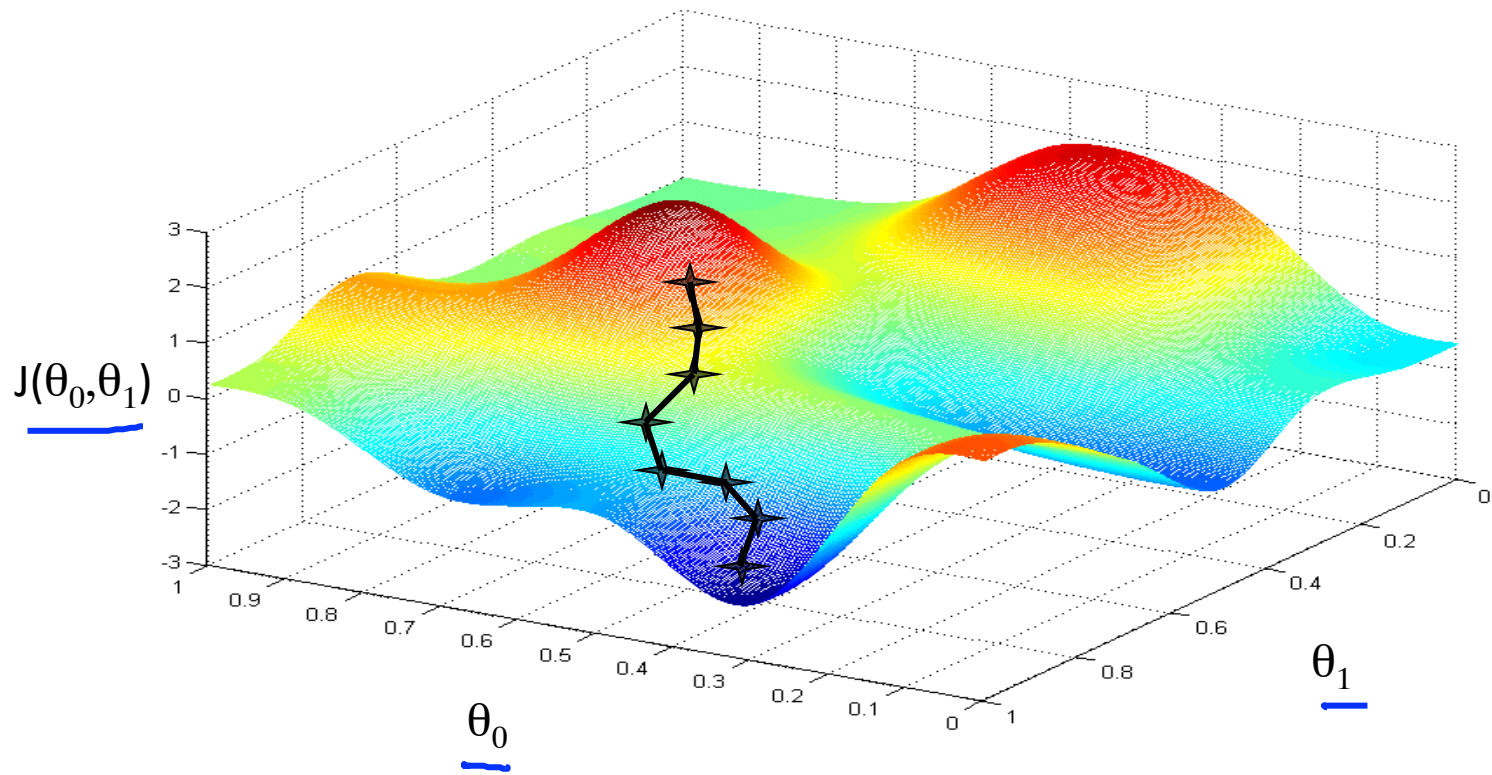
Have some function  $J(\theta_0, \theta_1)$   $J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)$

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

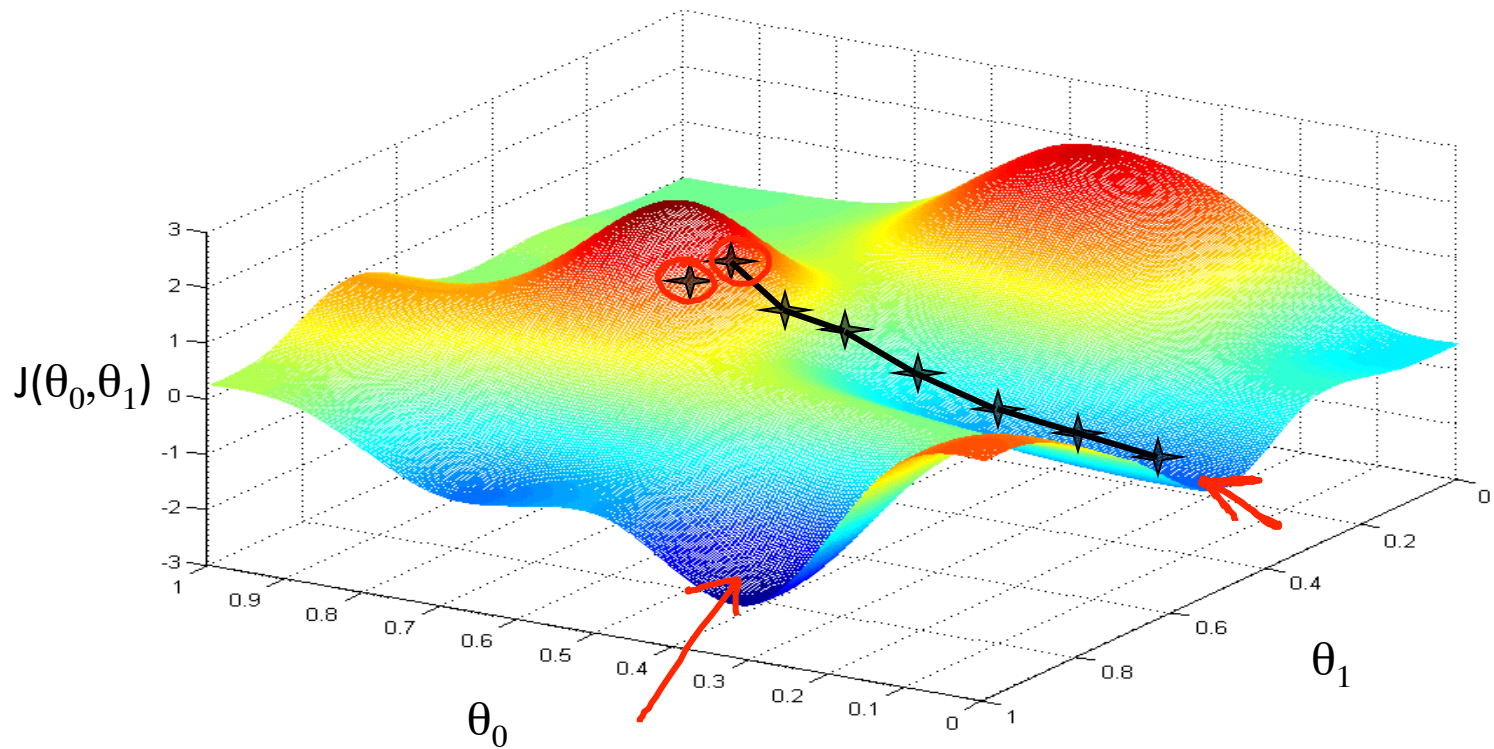
$\min_{\theta_0, \dots, \theta_n} J(\theta_0, \dots, \theta_n)$

## 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







# Gradient descent algorithm

$\theta_0, \theta_1$

repeat until convergence {

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

learning rate

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

Simultaneously update  
 $\theta_0$  and  $\theta_1$

Assignment

$$\begin{aligned} \rightarrow a &:= b \\ &\quad \uparrow \\ a &:= a + 1 \end{aligned}$$

Truth assertion

$$a = b \leftarrow$$

$$a = a + 1 \times$$

## Correct: Simultaneous update

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

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

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

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

↑

↑

## Incorrect:

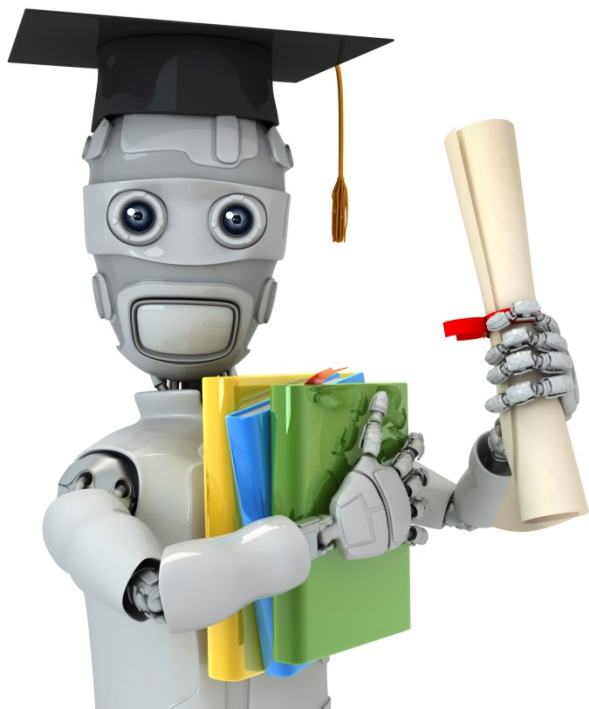
$$\rightarrow \text{temp0} := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

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

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

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

↑



Machine Learning

Linear regression  
with one variable

---

Gradient descent  
intuition

# Gradient descent algorithm

repeat until convergence {

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

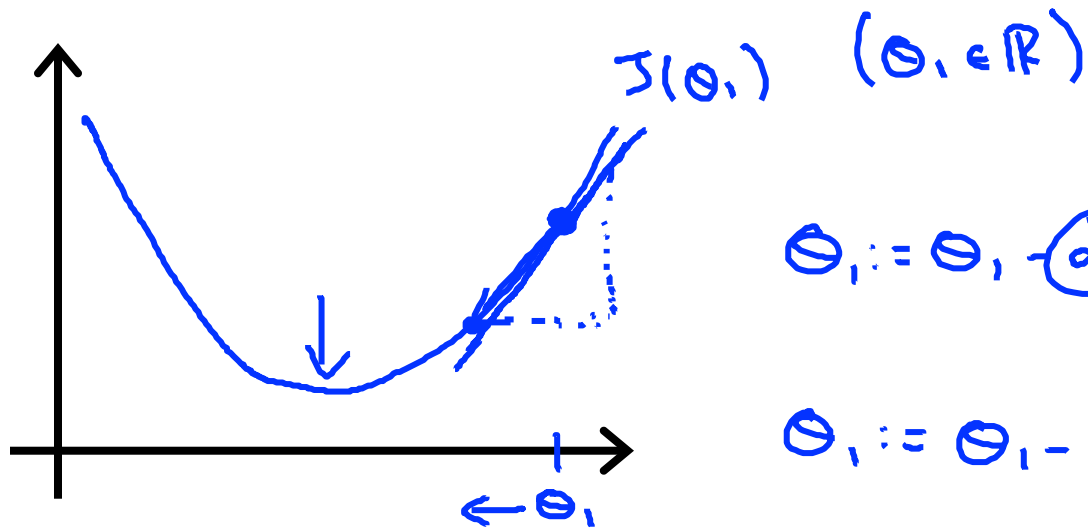
learning rate

derivative

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

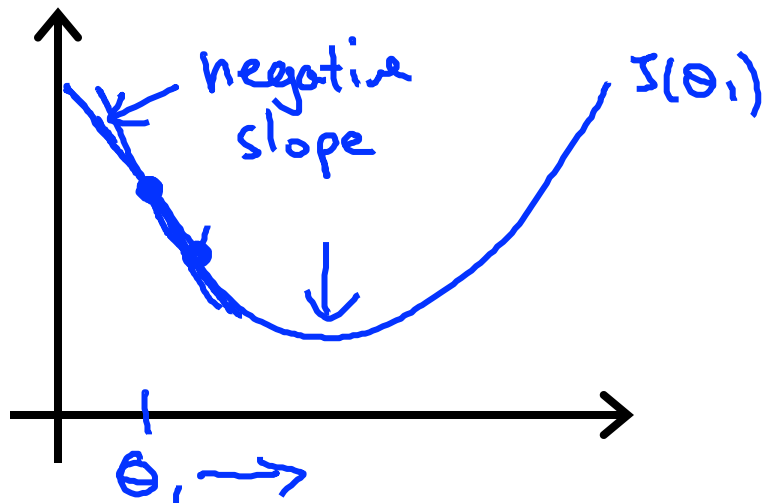
$$\min_{\theta_1} J(\theta_1)$$

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



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

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



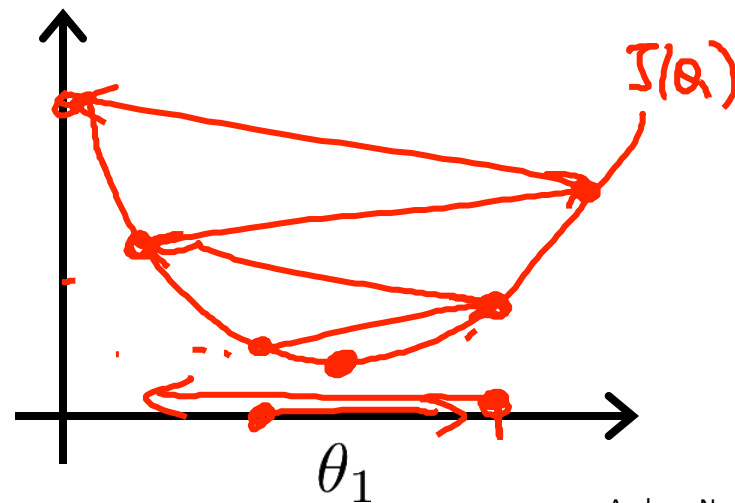
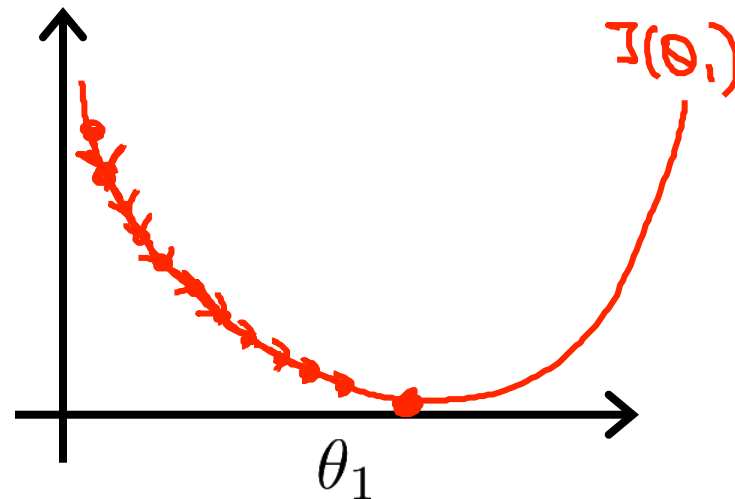
$$\frac{\partial}{\partial \theta_1} J(\theta_1) \leq 0$$

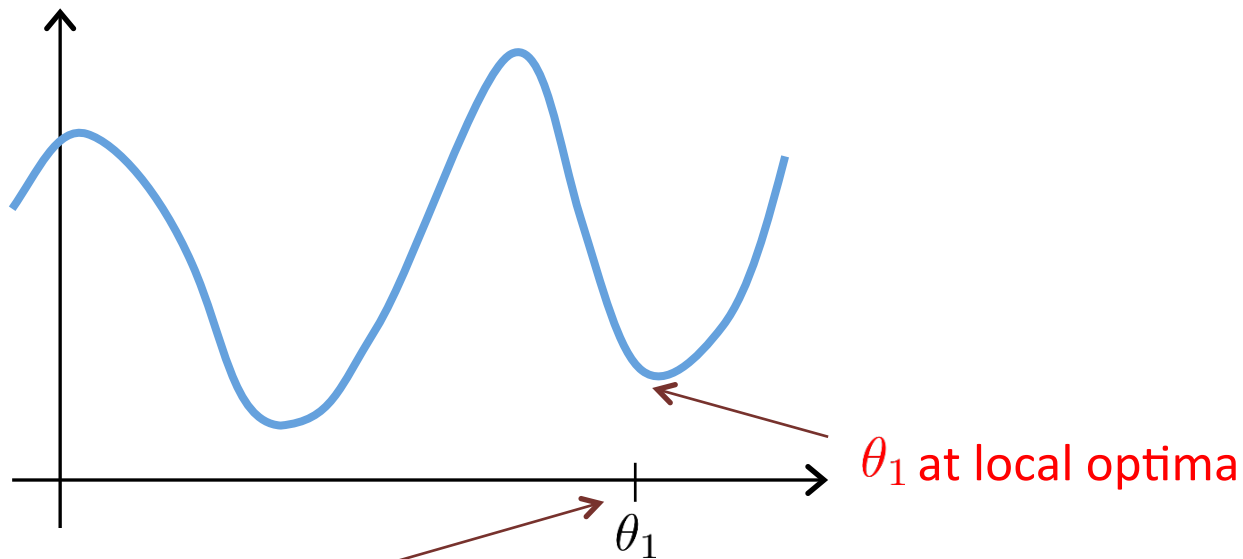
$$\theta_1 := \theta_1 - \alpha \cdot (\text{negative number})$$

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

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.





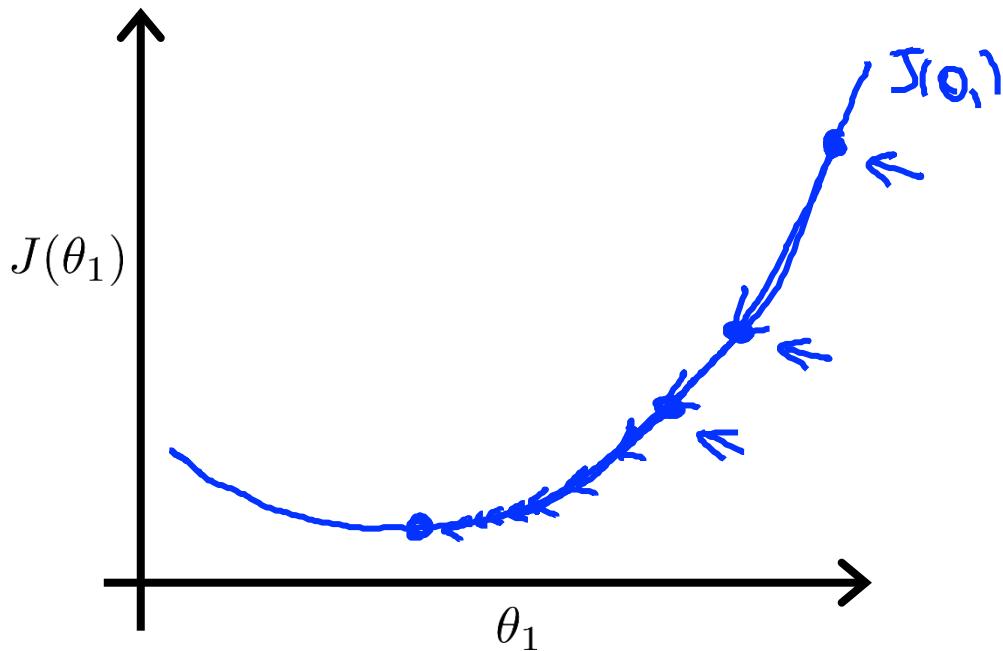
Current value of  $\theta_1$

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

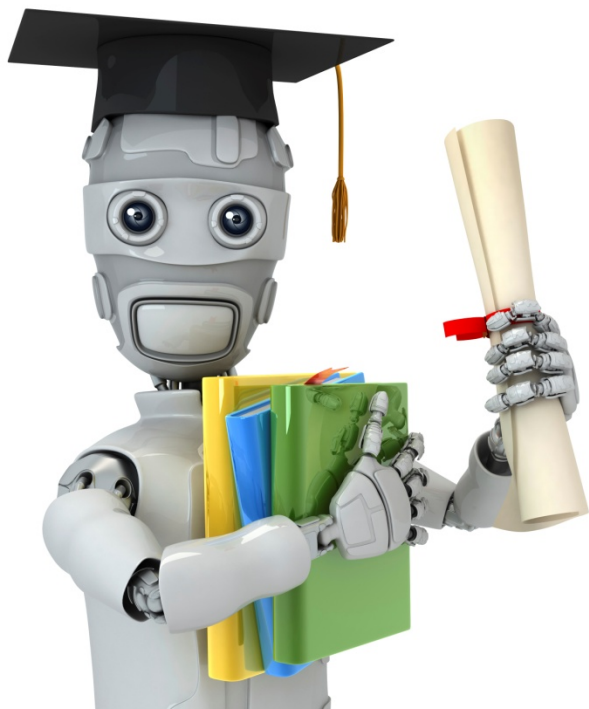
Gradient descent can converge to a local minimum, even with the learning rate  $\alpha$  fixed.

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

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







Machine Learning

Linear regression  
with one variable

---

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

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

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

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

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

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

# Gradient descent algorithm

repeat until convergence {

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

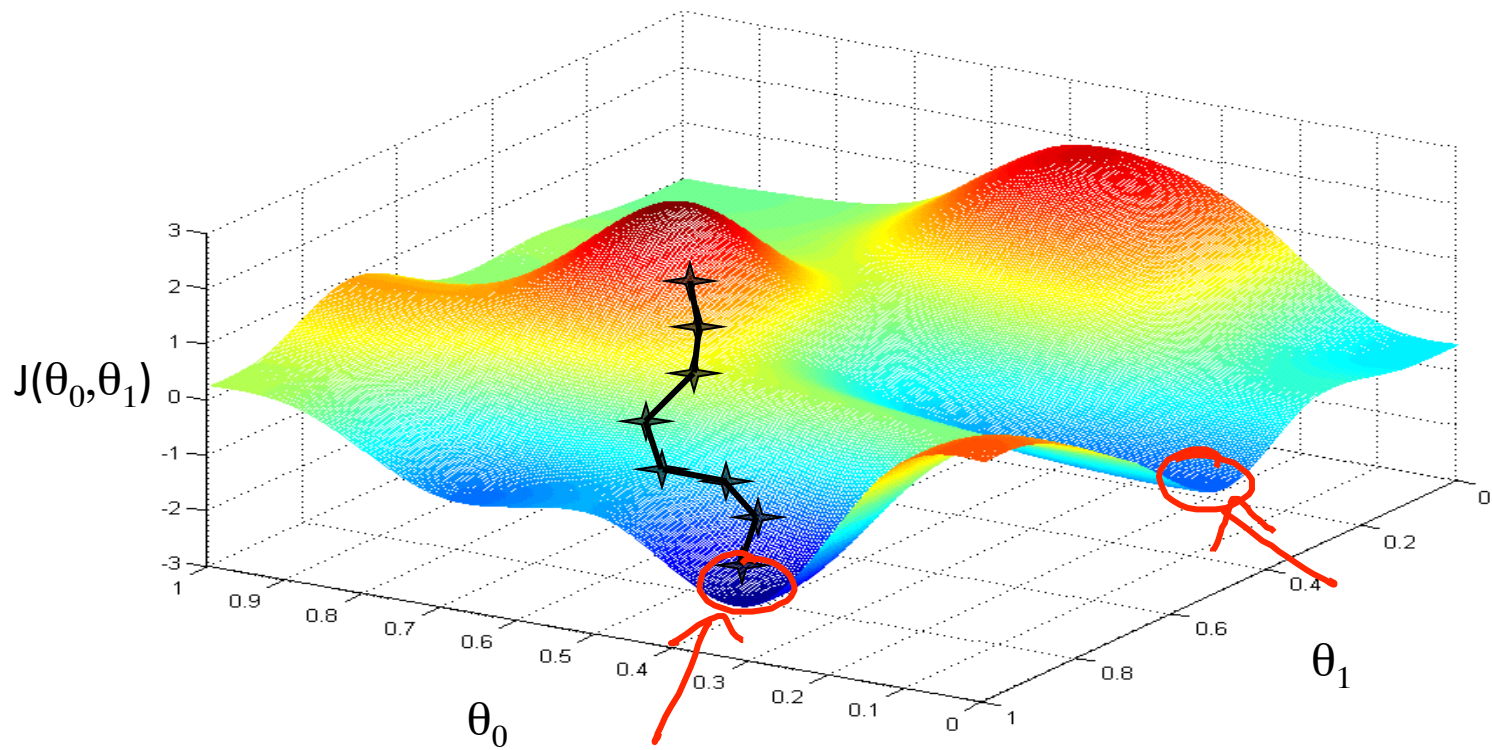
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

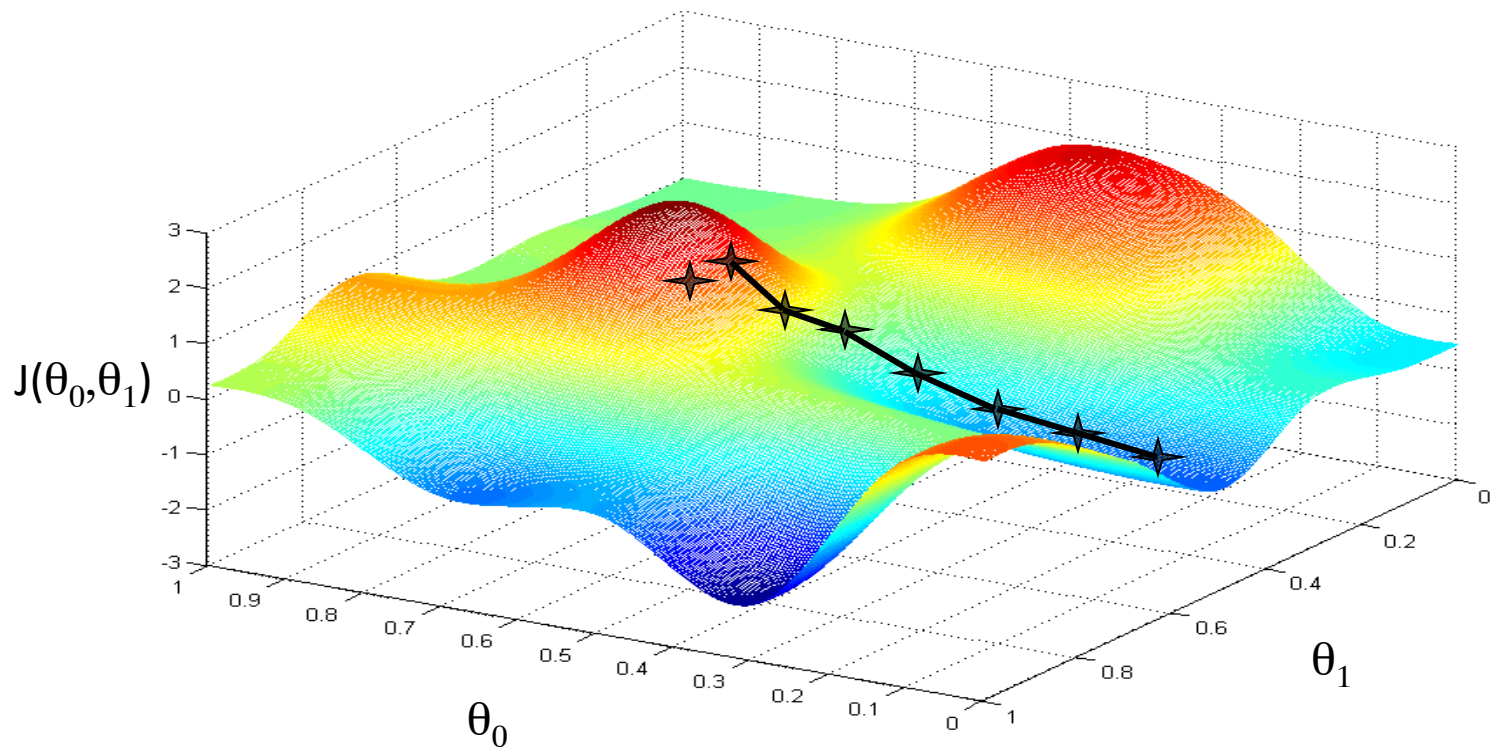
}

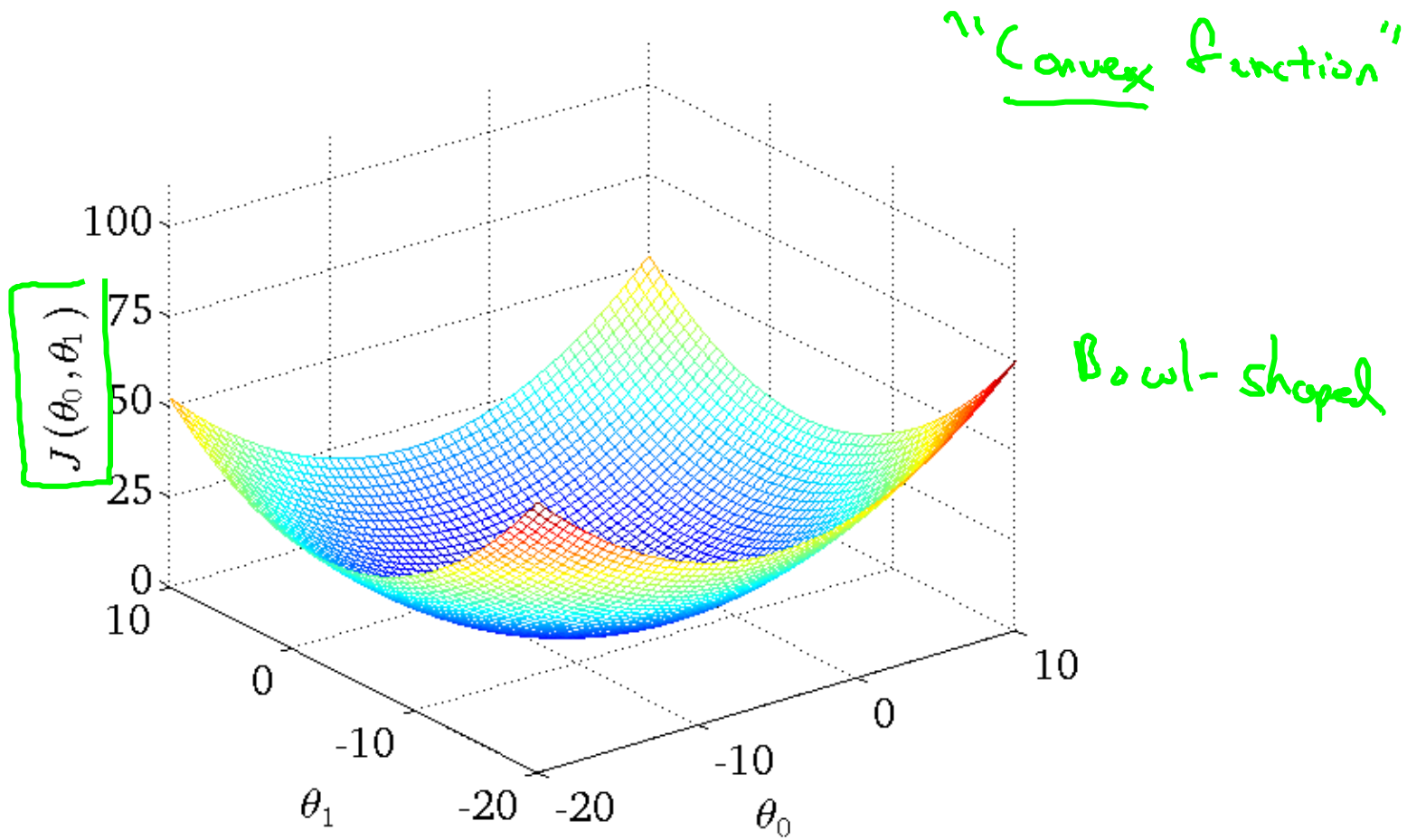
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

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

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

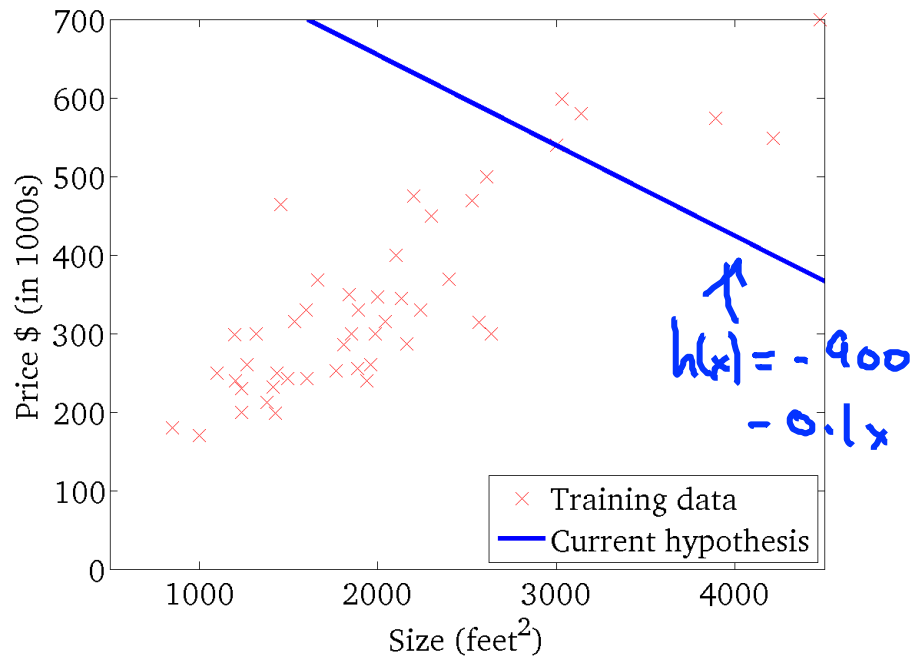






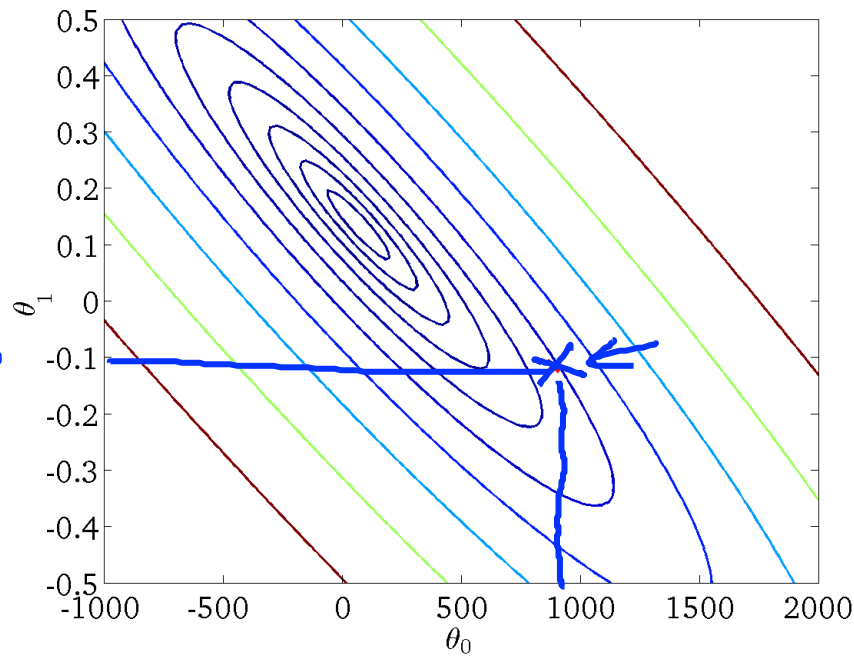
$$h_{\theta}(x)$$

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



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

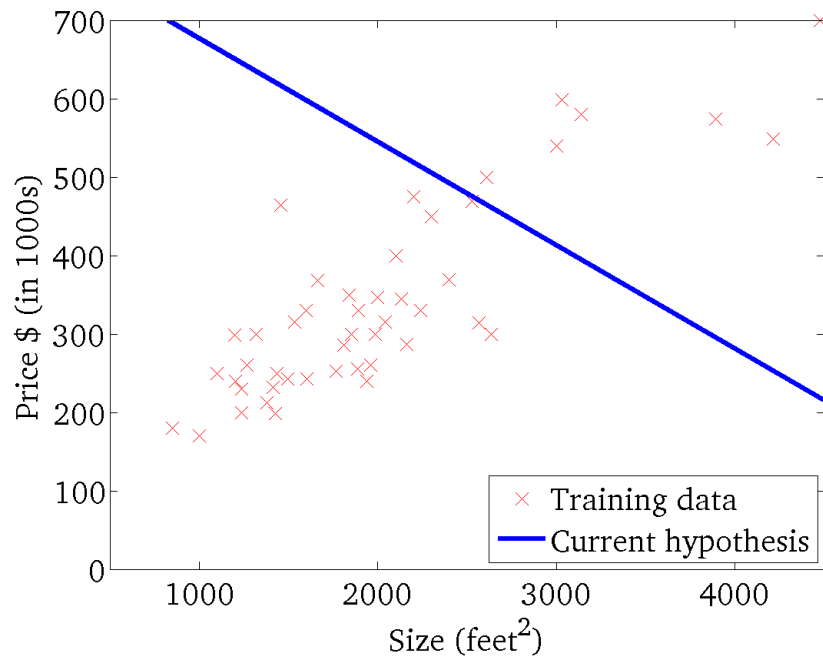
(function of the parameters  $\theta_0, \theta_1$ )





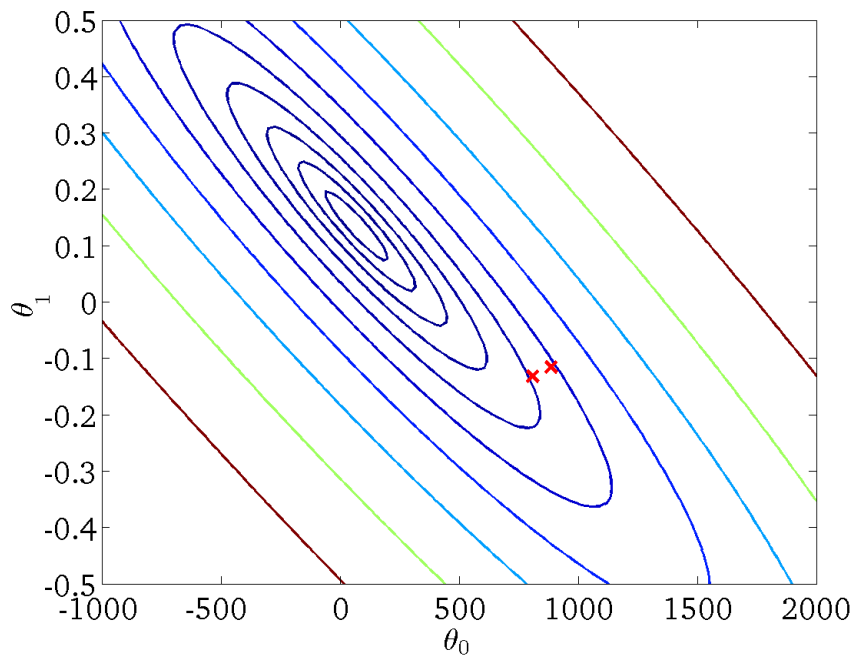
$$h_{\theta}(x)$$

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



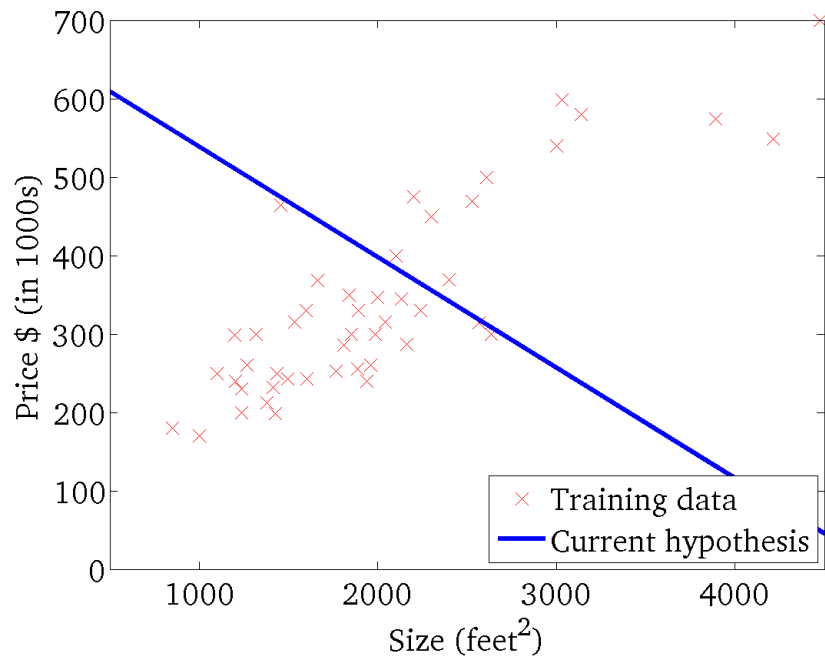
$$J(\theta_0, \theta_1)$$

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



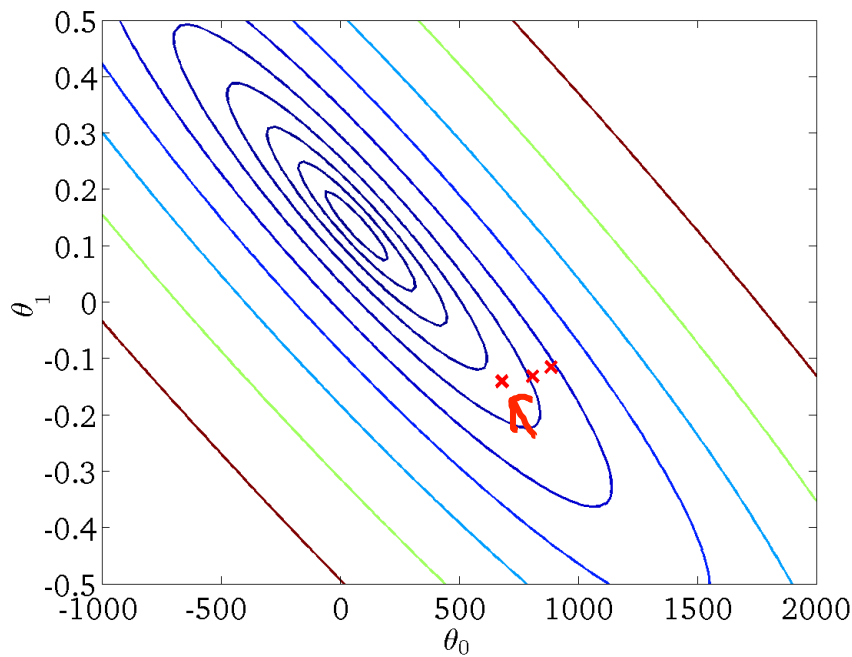
$$h_{\theta}(x)$$

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



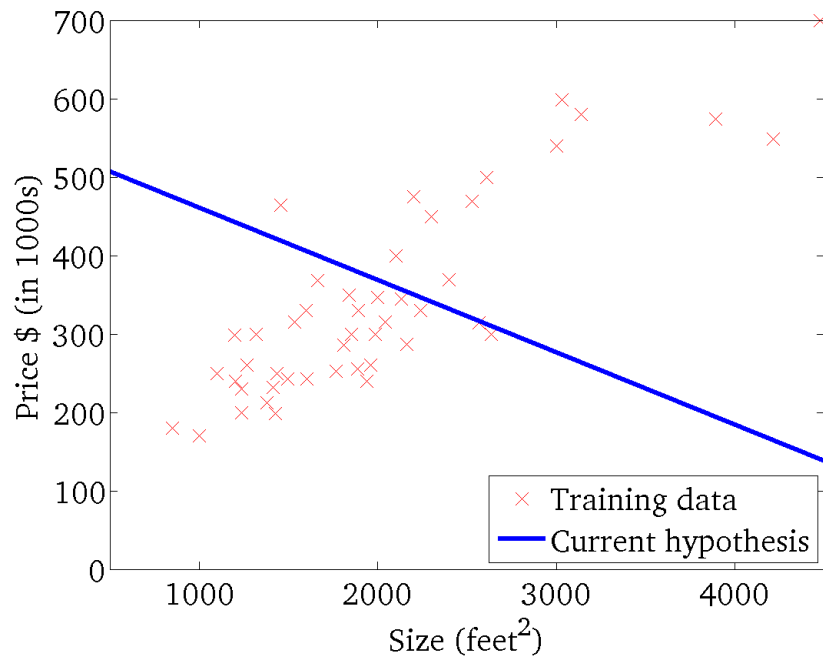
$$J(\theta_0, \theta_1)$$

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



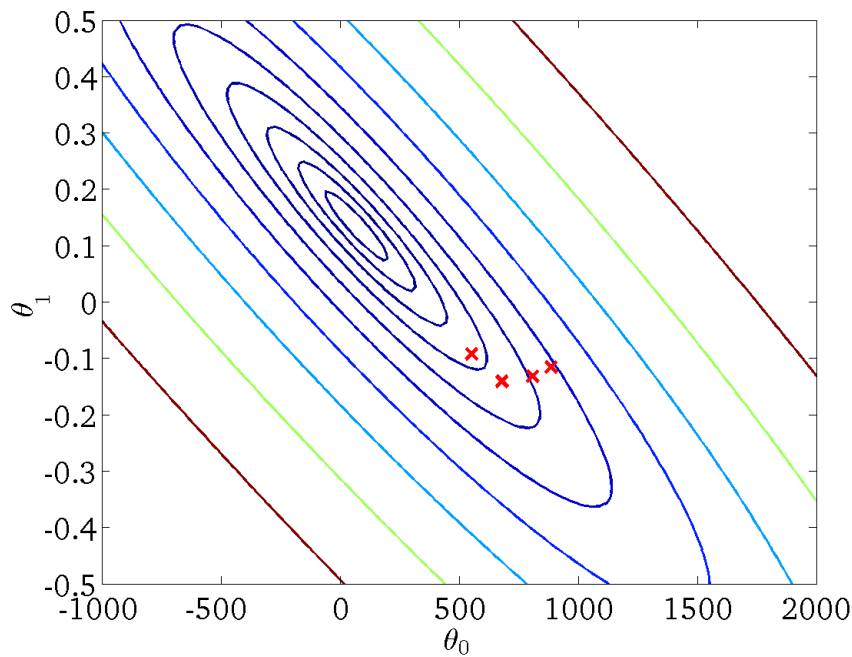
$$h_{\theta}(x)$$

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



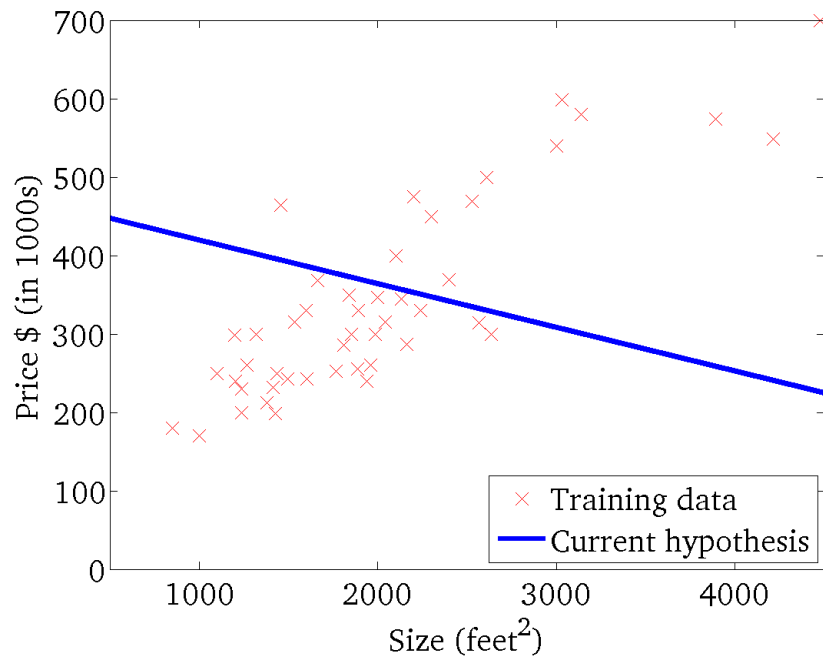
$$J(\theta_0, \theta_1)$$

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



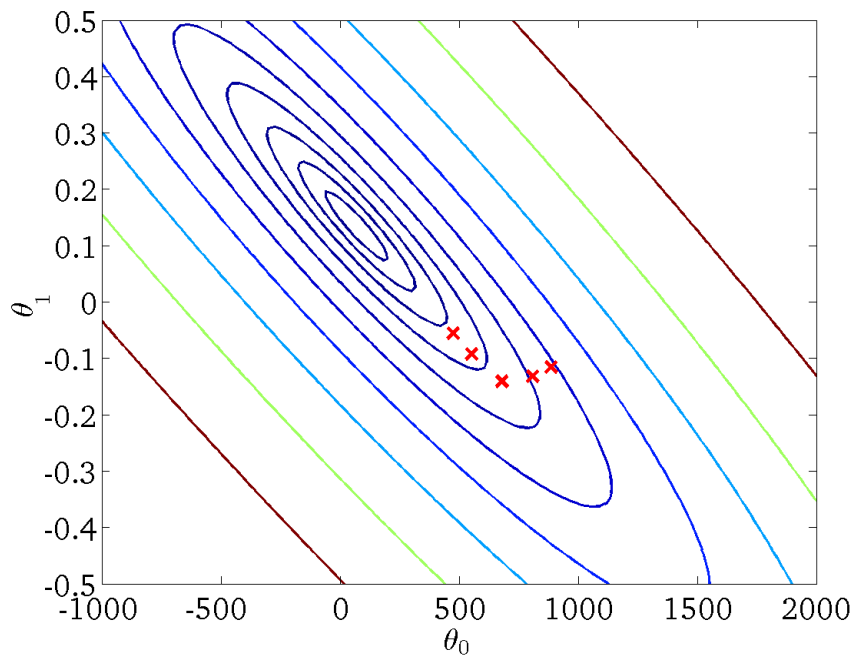
$$h_{\theta}(x)$$

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



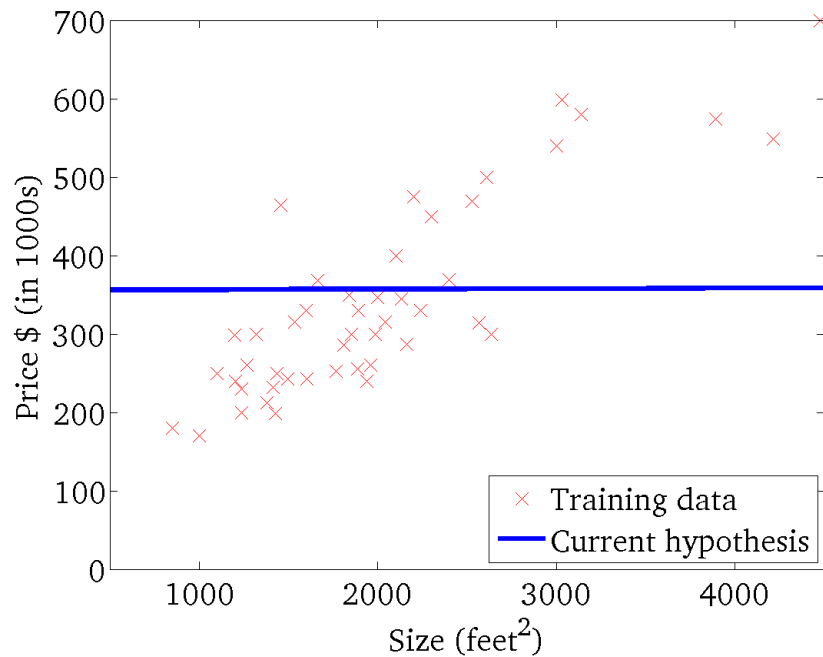
$$J(\theta_0, \theta_1)$$

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



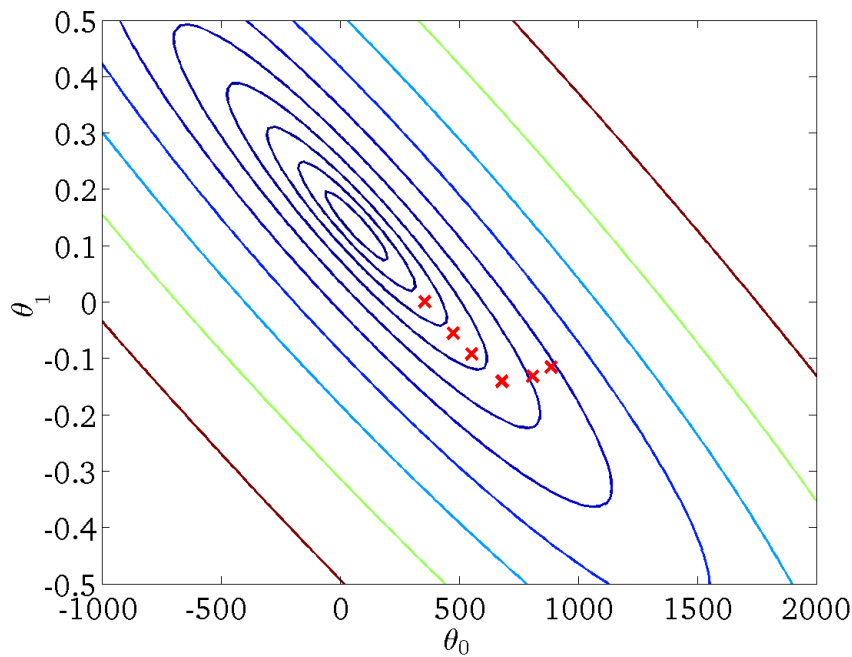
$$h_{\theta}(x)$$

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



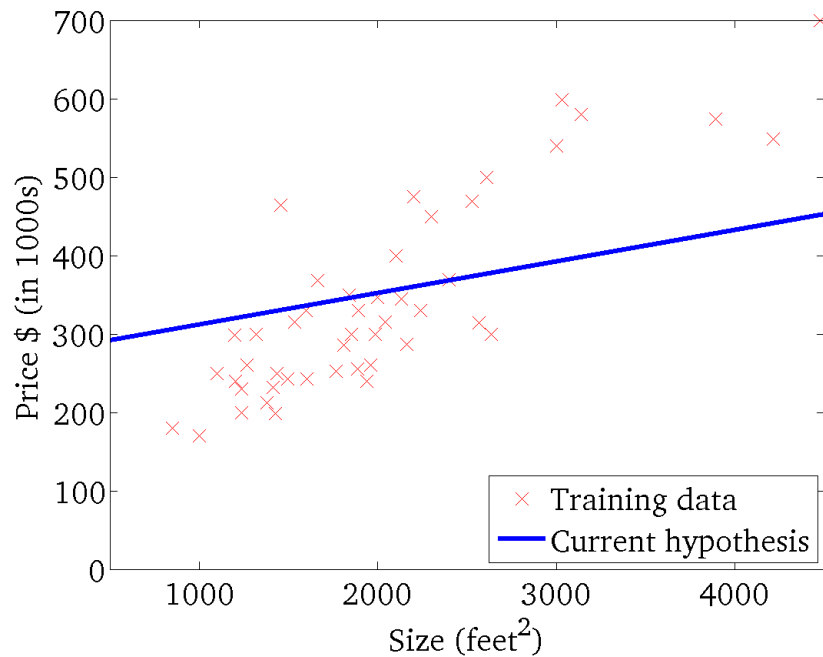
$$J(\theta_0, \theta_1)$$

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



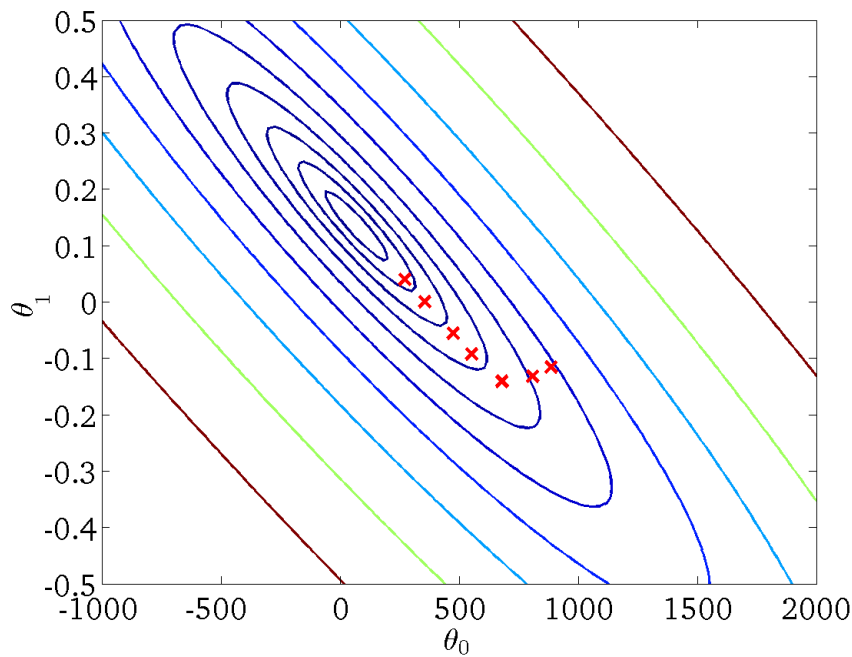
$$h_{\theta}(x)$$

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



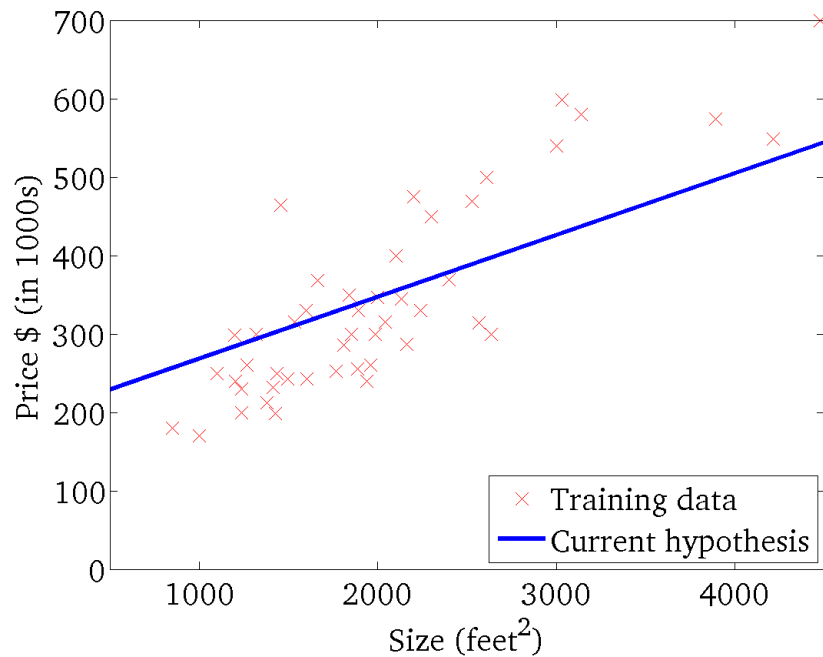
$$J(\theta_0, \theta_1)$$

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



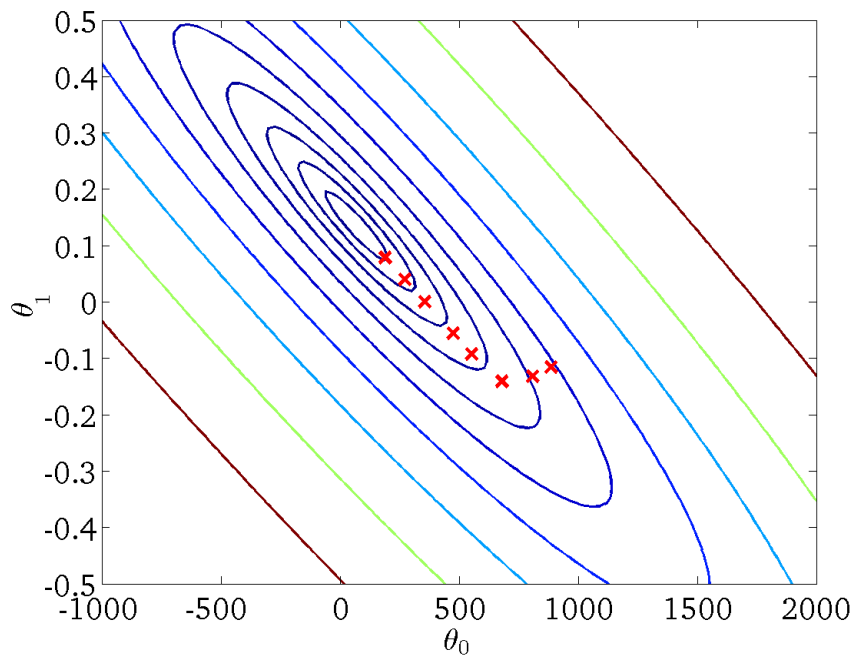
$$h_{\theta}(x)$$

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



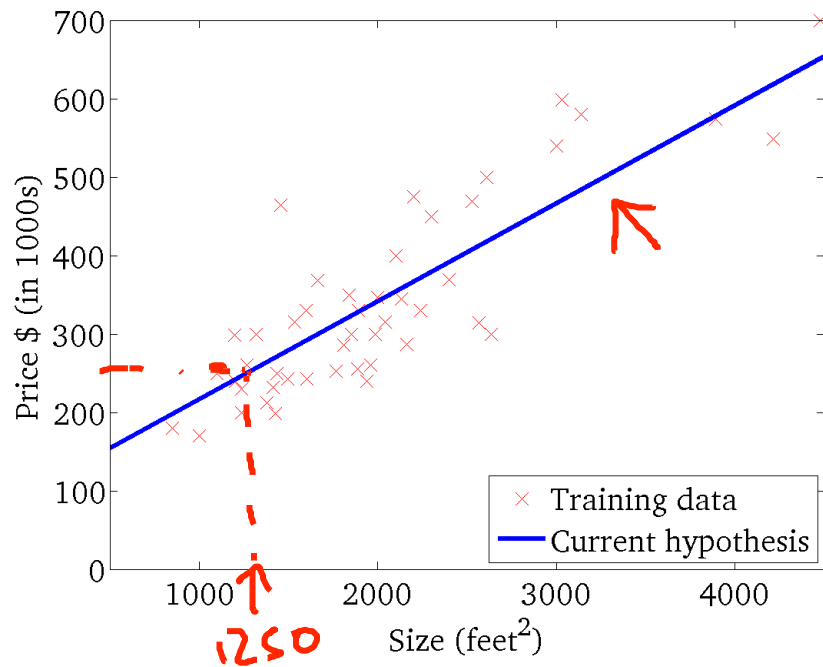
$$J(\theta_0, \theta_1)$$

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



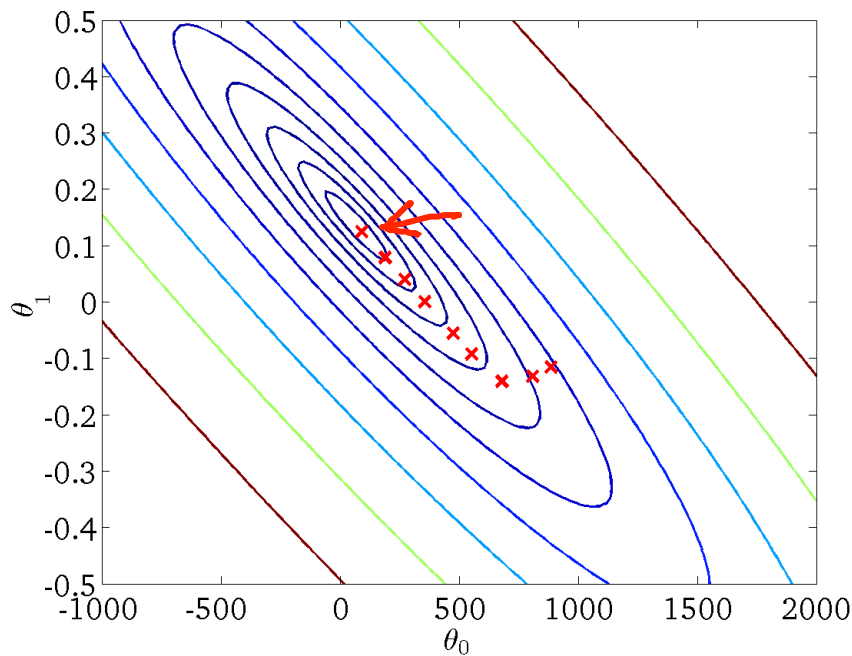
$$h_{\theta}(x)$$

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



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

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

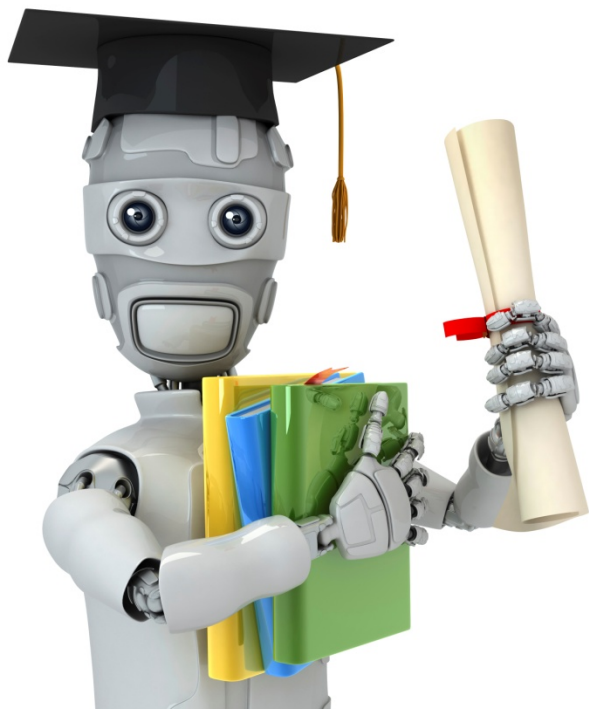




## “Batch” Gradient Descent

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

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



Machine Learning

Logistic  
Regression

---

Classification

# Classification

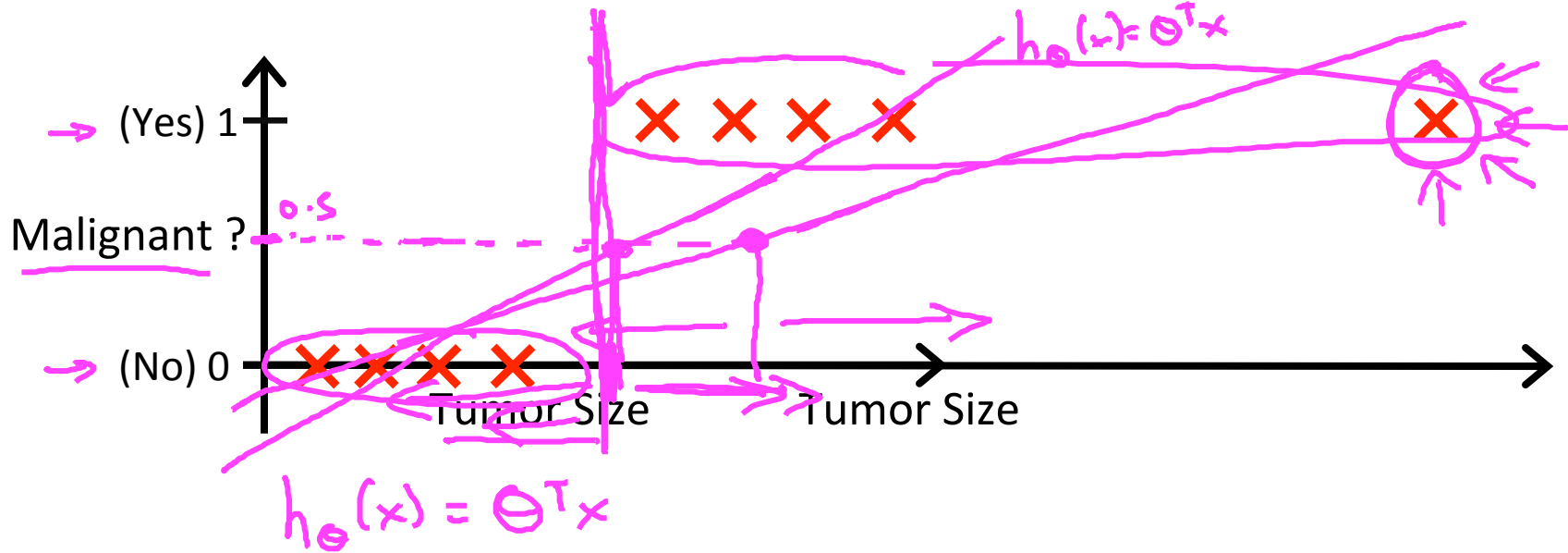
- Email: Spam / Not Spam?
- Online Transactions: Fraudulent (Yes / No)?
- Tumor: Malignant / Benign ?

→  $y \in \{0, 1\}$

0: "Negative Class" (e.g., benign tumor)

1: "Positive Class" (e.g., malignant tumor)

→  $y \in \{0, 1, 2, 3\}$



→ Threshold classifier output  $h_{\theta}(x)$  at 0.5:

→ If  $h_{\theta}(x) \geq 0.5$ , predict "y = 1"

If  $h_{\theta}(x) < 0.5$ , predict "y = 0"

Classification:

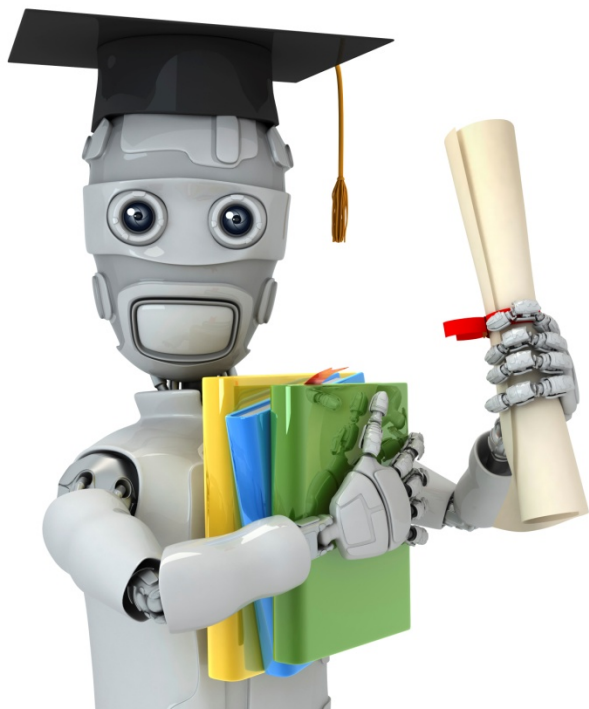
$$y = 0 \text{ or } 1$$

$h_{\theta}(x)$  can be  $> 1$  or  $< 0$

Logistic Regression:

$$0 \leq h_{\theta}(x) \leq 1$$

Classification



Machine Learning

Logistic

Regression

---

Hypothesis

Representation

# Logistic Regression Model

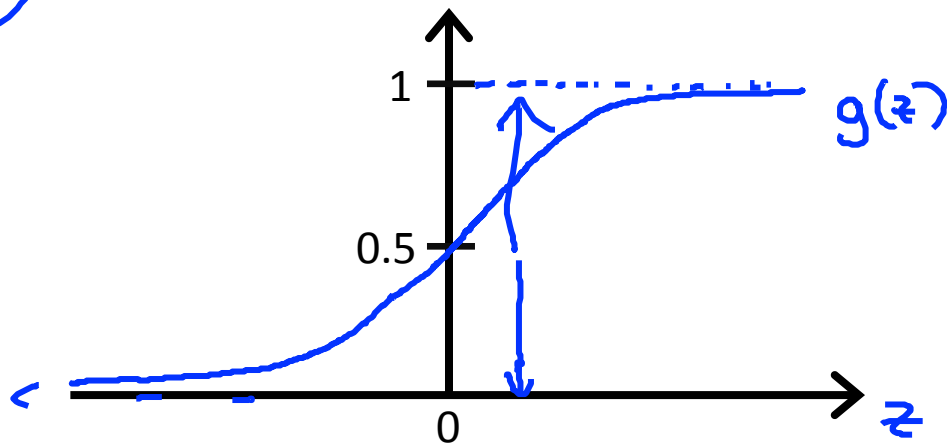
Want  $0 \leq h_{\theta}(x) \leq 1$

$$h_{\theta}(x) = g(\theta^T x)$$

$$\rightarrow g(z) = \frac{1}{1 + e^{-z}}$$

$\theta^T x$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



Parameters  $\underline{\theta}$ .

Sigmoid function  
Logistic function

## Interpretation of Hypothesis Output

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

$h_{\theta}(x)$  = estimated probability that  $y = 1$  on input  $x$  ←

Example: If  $\underline{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \leftarrow \\ \text{tumorSize} \leftarrow \end{bmatrix}$

$$\underline{h_{\theta}(x)} = \underline{0.7}$$

$$y = 1$$

Tell patient that 70% chance of tumor being malignant

$$\underline{h_{\theta}(x)} = \underline{P(y=1|x;\theta)}$$

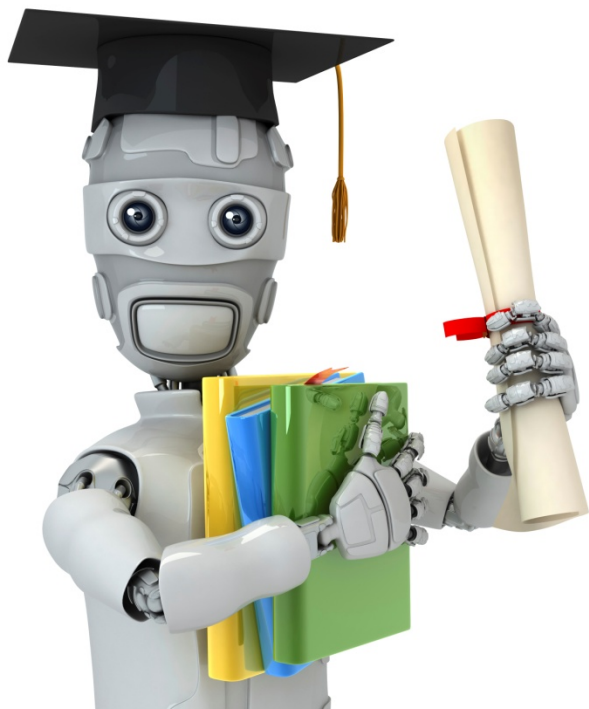
$$\underline{y = 0 \text{ or } 1}$$

“probability that  $y = 1$ , given  $x$ , parameterized by  $\theta$ ”

$$\rightarrow P(y = 0 | \theta) + P(y = 1 | \theta) = 1$$

$$\rightarrow P(y = 0 | x; \theta) = 1 - P(y = 1 | x; \theta)$$





Machine Learning

# Logistic Regression

---

## Decision boundary

## Logistic regression

$$h_{\theta}(x) = g(\theta^T x)$$

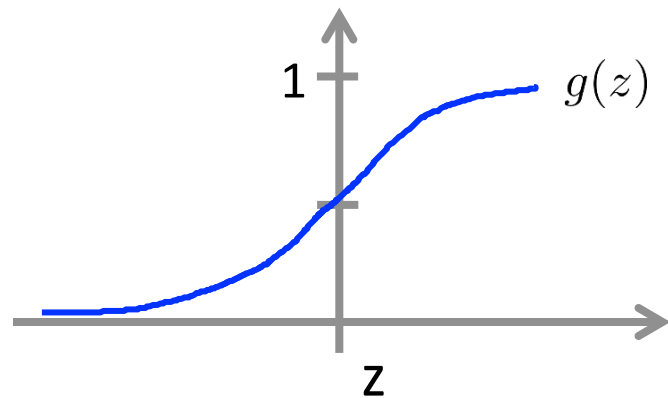
$$g(z) = \frac{1}{1+e^{-z}}$$

Suppose predict “ $y = 1$ ” if  $h_{\theta}(x) \geq 0.5$

$$\theta^T x \geq 0$$

predict “ $y = 0$ ” if  $h_{\theta}(x) < 0.5$

$$\theta^T x < 0$$

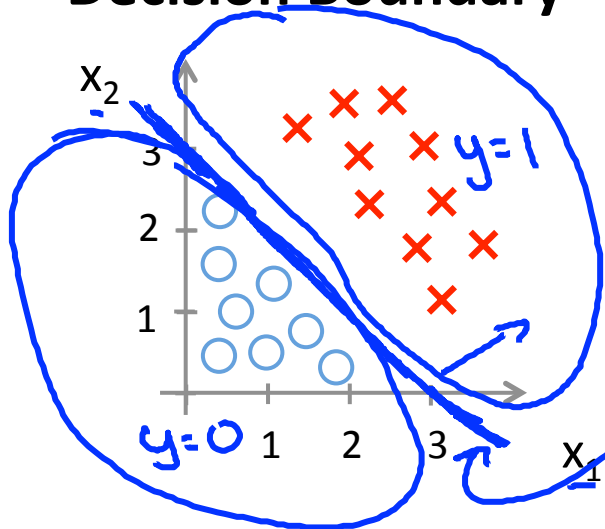


$$g(z) \geq 0.5 \\ \text{when } z \geq 0$$

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) < 0.5 \\ \text{when } z < 0$$

# Decision Boundary



$$\theta = \begin{bmatrix} -3 \\ 1 \\ 1 \end{bmatrix} \leftarrow$$

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

Decision boundary

Predict " $y = 1$ " if  $-3 + x_1 + x_2 \geq 0$

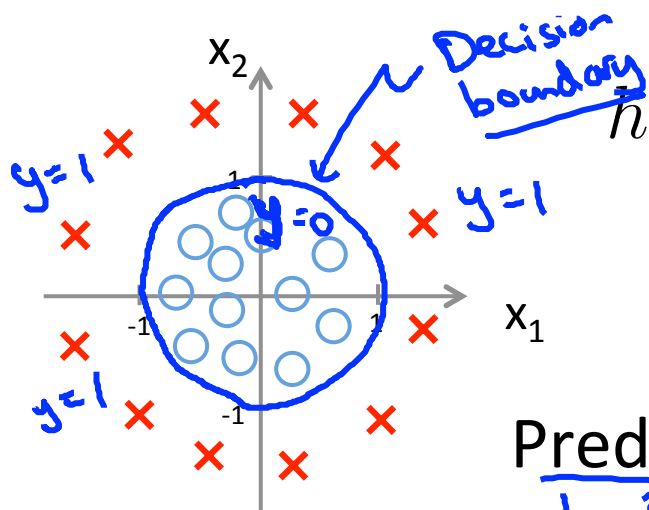
$\theta^T x$

$x_1, x_2$   
 $\rightarrow h_{\theta}(x) = 0.5$   
 $x_1 + x_2 = 3$

$x_1 + x_2 \geq 3$

$x_1 + x_2 < 3$   
 $y = 0$

# Non-linear decision boundaries



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_2^2)$$

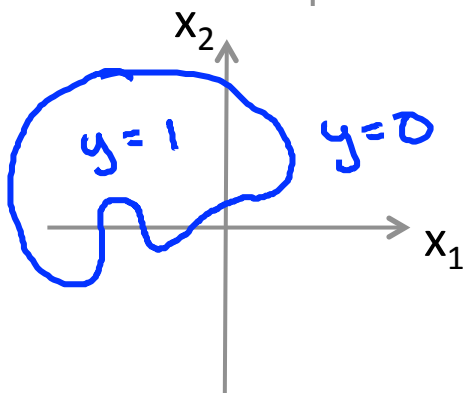
$\begin{matrix} \text{---} \\ \text{---} \\ \text{---} \end{matrix}$

$$\theta = \begin{bmatrix} -1 \\ 1 \\ 1 \\ 0 \\ 0 \end{bmatrix}$$

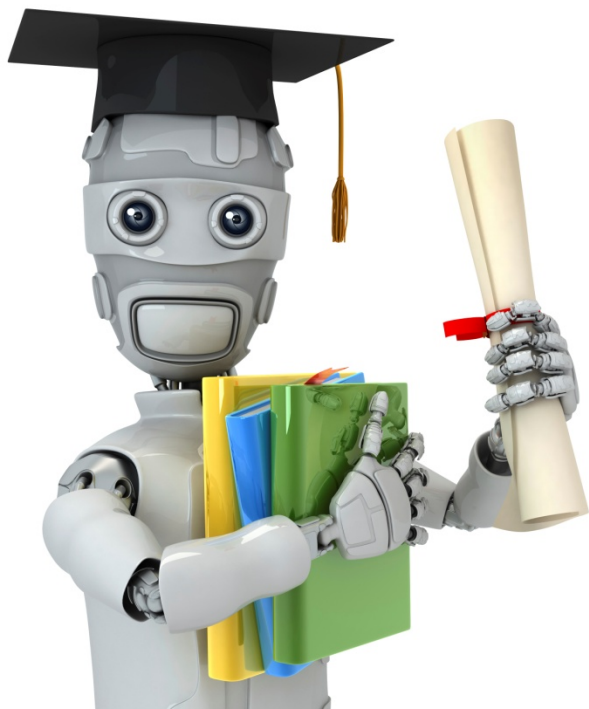
Predict "y = 1" if  $-1 + x_1^2 + x_2^2 \geq 0$

$$\boxed{x_1^2 + x_2^2 = 1}$$

$$x_1^2 + x_2^2 \geq 1$$



$$h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1^2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^2 x_2^2 + \theta_6 x_1^3 x_2 + \dots)$$



Machine Learning

# Logistic Regression

---

## Cost function

Training  
set:

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

m examples

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

$$x_0 = 1, y \in \{0, 1\}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

How to choose parameters  $\theta$  ?

# Cost function

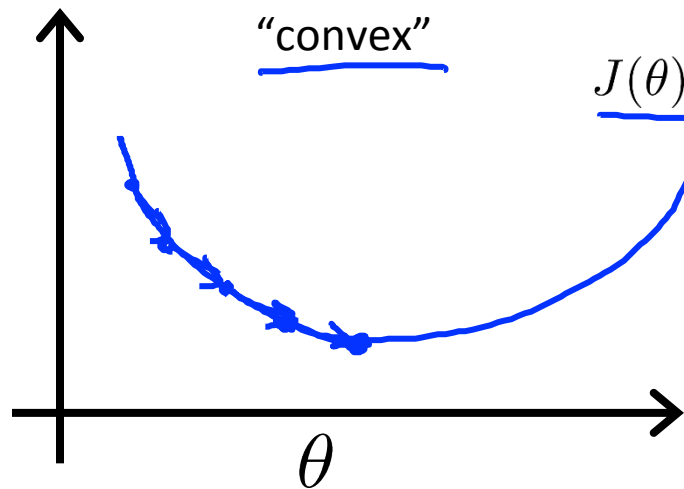
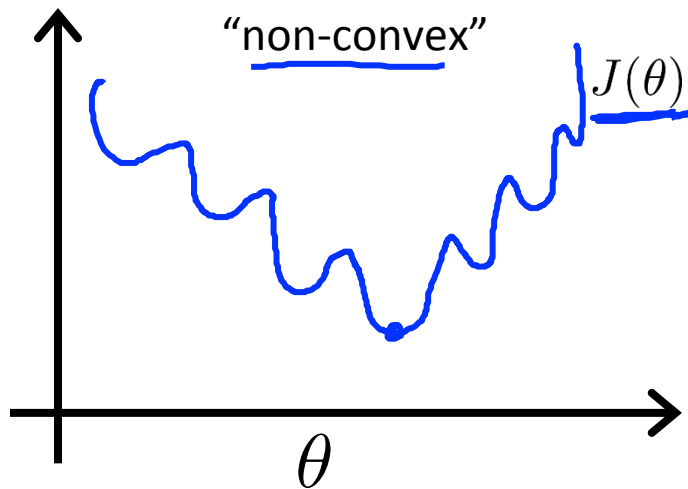
→ Linear regression:  
logistic

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

$$\text{cost}(h_{\theta}(x^{(i)}), y)$$

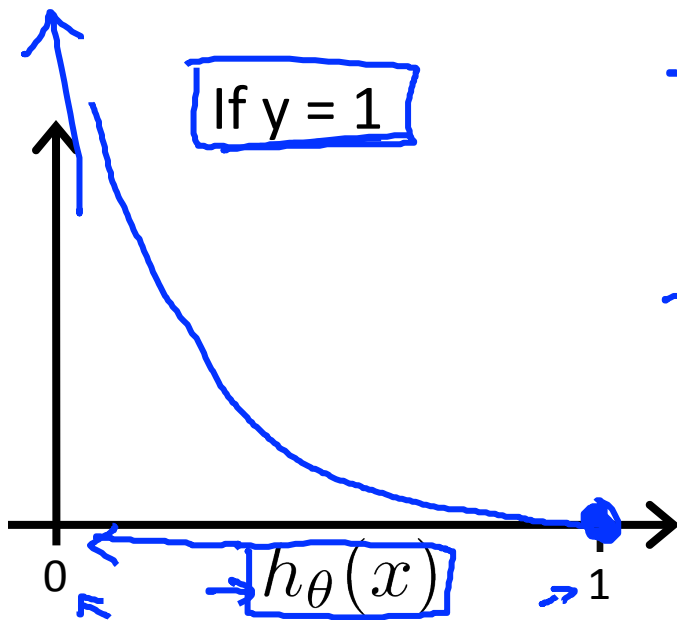
$$\text{Cost}(h_{\theta}(x), y) = \frac{1}{2} (h_{\theta}(x) - y)^2$$

$$\frac{1}{1 + e^{-\theta^T x}}$$



# Logistic regression cost function

$$\text{Cost}(\underbrace{h_\theta(x)}_{\uparrow}, y) = \begin{cases} \boxed{-\log(h_\theta(x))} & \text{if } y = 1 \\ \underline{-\log(1 - h_\theta(x))} & \text{if } y = 0 \end{cases}$$



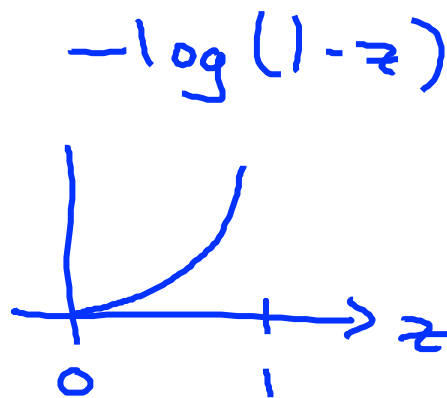
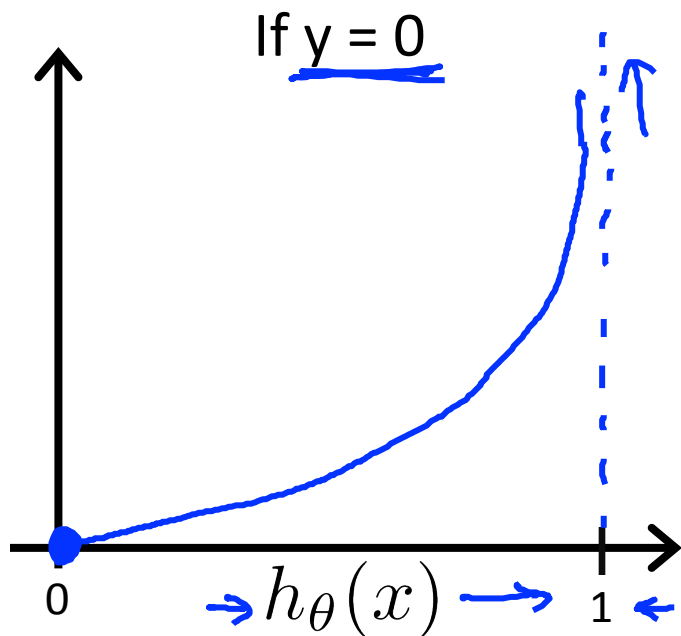
→ Cost = 0 if  $y = 1, h_\theta(x) = 1$   
But as  $h_\theta(x) \rightarrow 0$   
Cost  $\rightarrow \infty$

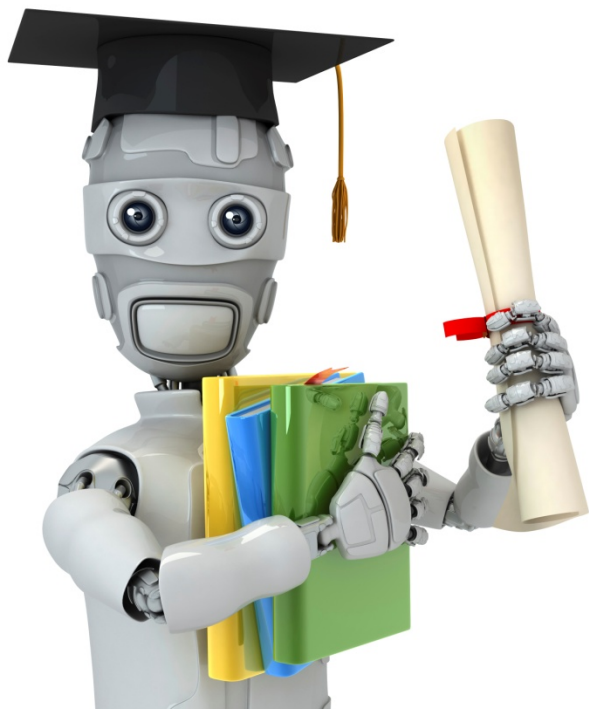
→ Captures intuition that if  $h_\theta(x) = 0$ ,  
(predict  $P(y = 1|x; \theta) = 0$ ), but  $y = 1$ ,  
we'll penalize learning algorithm by a very  
large cost.



## Logistic regression cost function

$$\text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$





Machine Learning

# Logistic Regression

---

Simplified cost function  
and gradient descent

## Logistic regression cost function

$$\rightarrow J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$

$$\rightarrow \text{Cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \leftarrow \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \leftarrow \end{cases}$$

Note:  $y = 0$  or  $1$  always

$$\rightarrow \text{Cost}(h_{\theta}(x), y) = \underbrace{-y \log(h_{\theta}(x))}_{=0} - \underbrace{(1-y) \log(1-h_{\theta}(x)})_{=1}$$

If  $y=1$ :  $\text{Cost}(h_{\theta}(x), y) = -\log h_{\theta}(x) \leftarrow$

If  $y=0$ :  $\text{Cost}(h_{\theta}(x), y) = \underline{-\log(1-h_{\theta}(x))}$

## Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(x^{(i)}), y^{(i)})$$
$$= -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

To fit parameters  $\theta$  :

$$\min_{\theta} J(\theta) \quad \text{Get } \underline{\theta}$$

To make a prediction given new  $\underline{x}$  :

$$\text{Output } \underline{h_{\theta}(x)} = \frac{1}{1 + e^{-\theta^T x}}$$

$$\underline{p(y=1 | x; \theta)}$$

## Gradient Descent

$$\rightarrow J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Want  $\min_{\theta} J(\theta)$ :

Repeat {

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

}

(simultaneously update all  $\theta_j$ )

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

# Gradient Descent

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

Want  $\min_{\theta} J(\theta)$ :

Repeat {

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

}

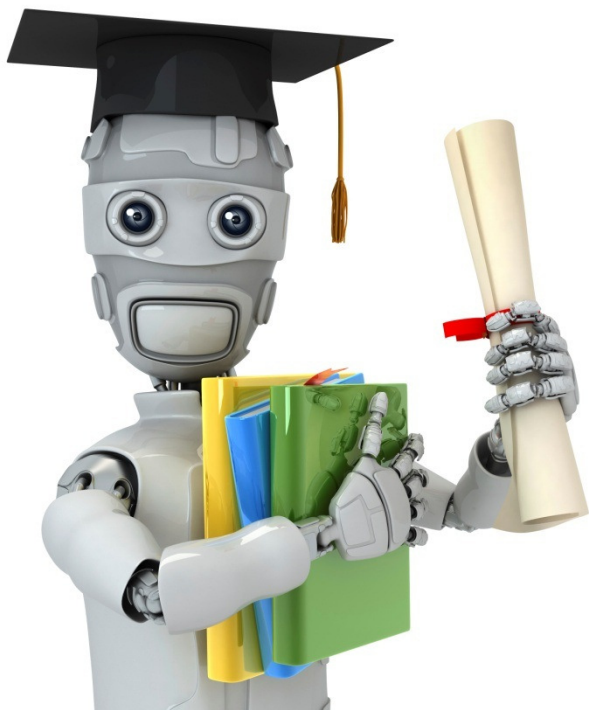
(simultaneously update all  $\theta_j$ )

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \begin{matrix} \uparrow \\ \uparrow \\ \uparrow \\ \uparrow \\ \uparrow \end{matrix} \text{ for } i=0 \text{ to } n$$

$$h_{\theta}(x) = \Theta^T x$$

$$\rightarrow h_{\theta}(x) = \frac{1}{1 + e^{-\Theta^T x}}$$

Algorithm looks identical to linear regression!



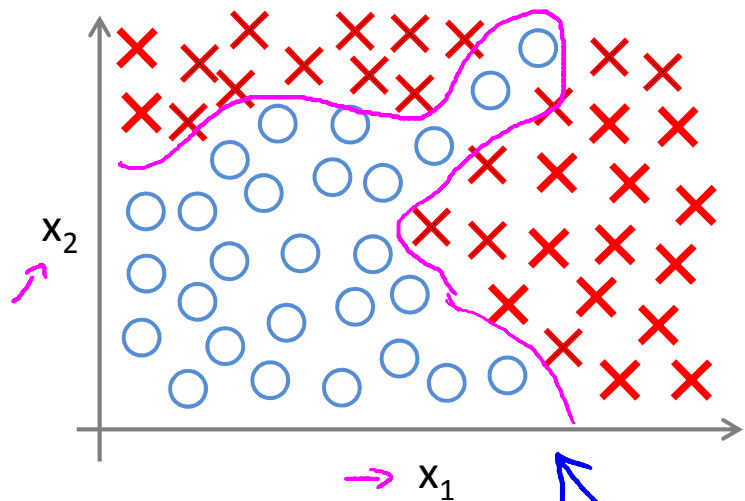
Machine Learning

# Neural Networks: Representation

---

## Non-linear hypotheses

# Non-linear Classification



$$g(\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_1 x_2 + \theta_4 x_1^2 x_2 + \theta_5 x_1^3 x_2 + \theta_6 x_1 x_2^2 + \dots)$$

$$\rightarrow x_1^2, x_1 x_2, x_1 x_3, x_1 x_4 \dots x_1 x_{100}$$

$$x_2^2, x_2 x_3 \dots$$

~ 5000 feature  $O(n^2)$

$$\rightarrow x_1^2, x_2^2, x_3^2, \dots, x_{100}^2$$

$\approx \frac{n^2}{2}$

$$\rightarrow x_1 x_2 x_3, x_1^2 x_2, x_{10} x_{11} x_{17}, \dots$$

$O(n^3)$

170,000

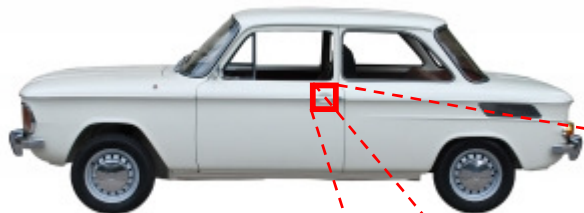
- $x_1$  = size
- $x_2$  = # bedrooms
- $x_3$  = # floors
- $x_4$  = age
- ...
- $x_{100}$

$n=100$



# What is this?

You see this:



But the camera sees this:

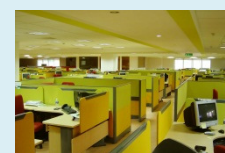
194	210	201	212	199	213	215	195	178	158	182	209
180	189	190	221	209	205	191	167	147	115	129	163
114	126	140	188	176	165	152	140	170	106	78	88
87	103	115	154	143	142	149	153	173	101	57	57
102	112	106	131	122	138	152	147	128	84	58	66
94	95	79	104	105	124	129	113	107	87	69	67
68	71	69	98	89	92	98	95	89	88	76	67
41	56	68	99	63	45	60	82	58	76	75	65
20	43	69	75	56	41	51	73	55	70	63	44
50	50	57	69	75	75	73	74	53	68	59	37
72	59	53	66	84	92	84	74	57	72	63	42
67	61	58	65	75	78	76	73	59	75	69	50



# Computer Vision: Car detection



Cars

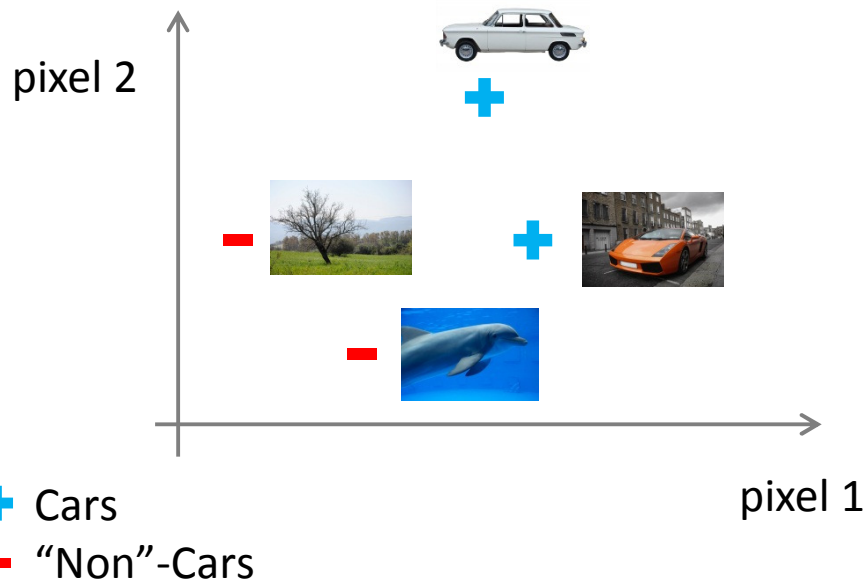
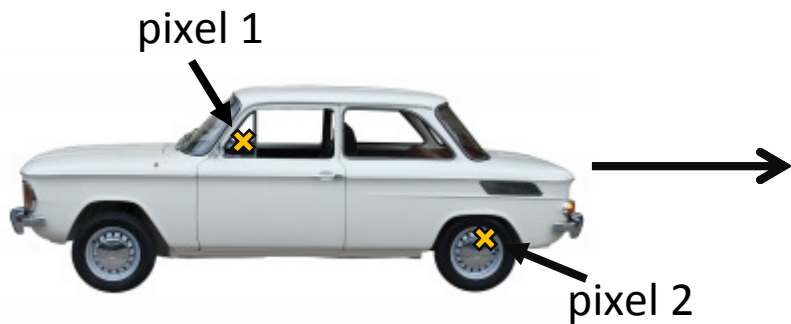


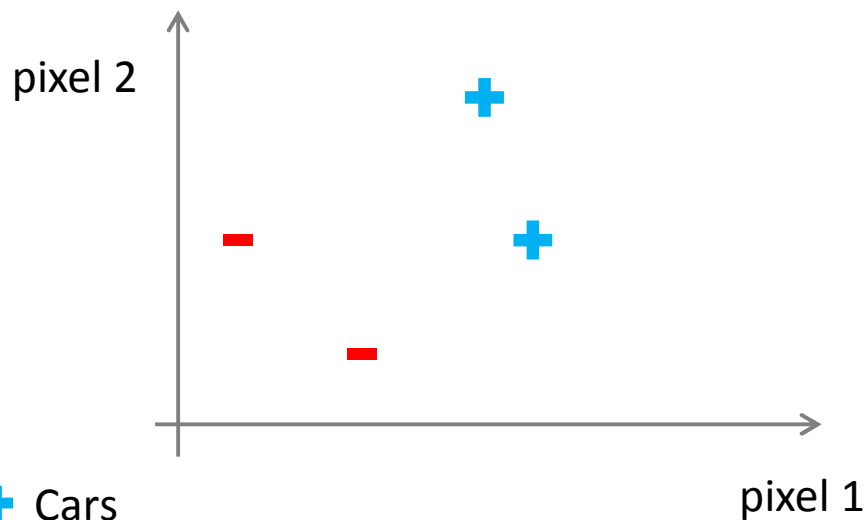
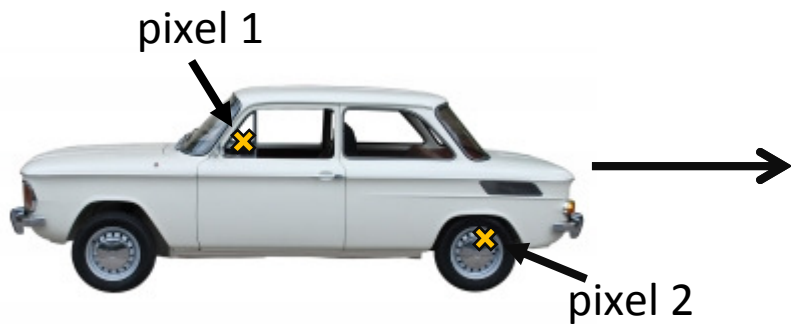
Not a car

Testing:

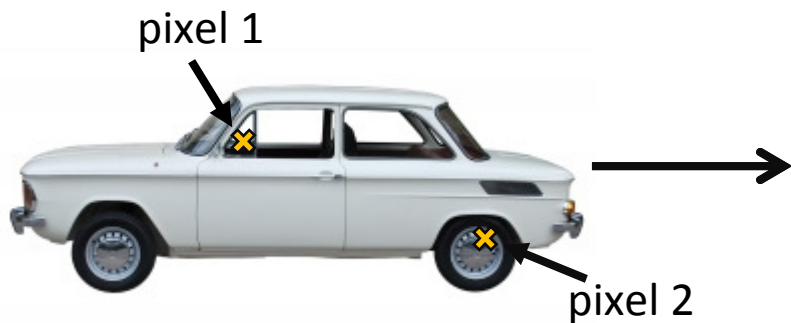


What is this?

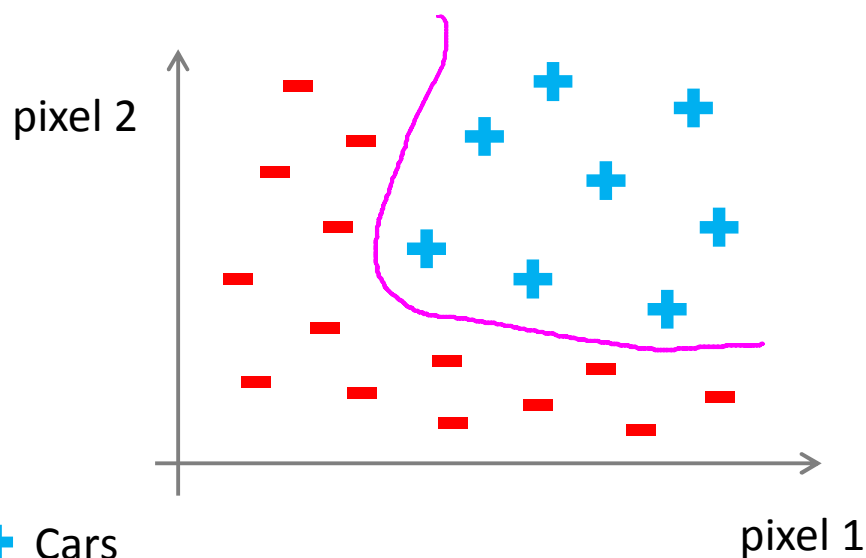




+ Cars  
- "Non"-Cars



Learning  
Algorithm



+ Cars  
- "Non"-Cars

50 x 50 pixel images  $\rightarrow$  2500 pixels

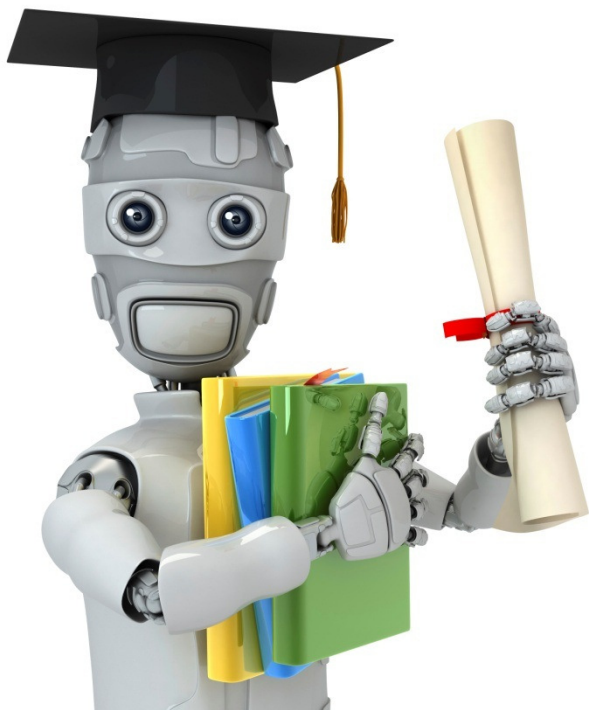
$n = 2500$  (7500 if RGB)

$$\rightarrow x = \begin{bmatrix} \text{pixel 1 intensity} \\ \text{pixel 2 intensity} \\ \vdots \\ \text{pixel 2500 intensity} \end{bmatrix}$$

0-255

Quadratic features ( $x_i \times x_j$ ):  $\approx$  3 million features





Machine Learning

# Neural Networks: Representation

---

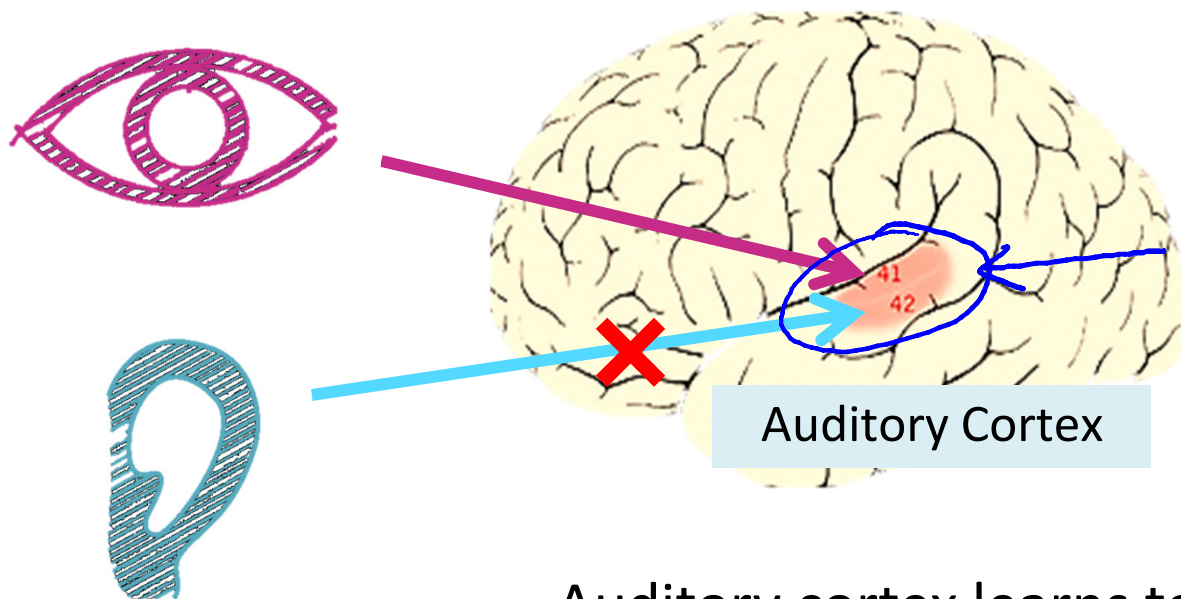
## Neurons and the brain

# Neural Networks

- Origins: Algorithms that try to mimic the brain.
- Was very widely used in 80s and early 90s; popularity diminished in late 90s.
- Recent resurgence: State-of-the-art technique for many applications

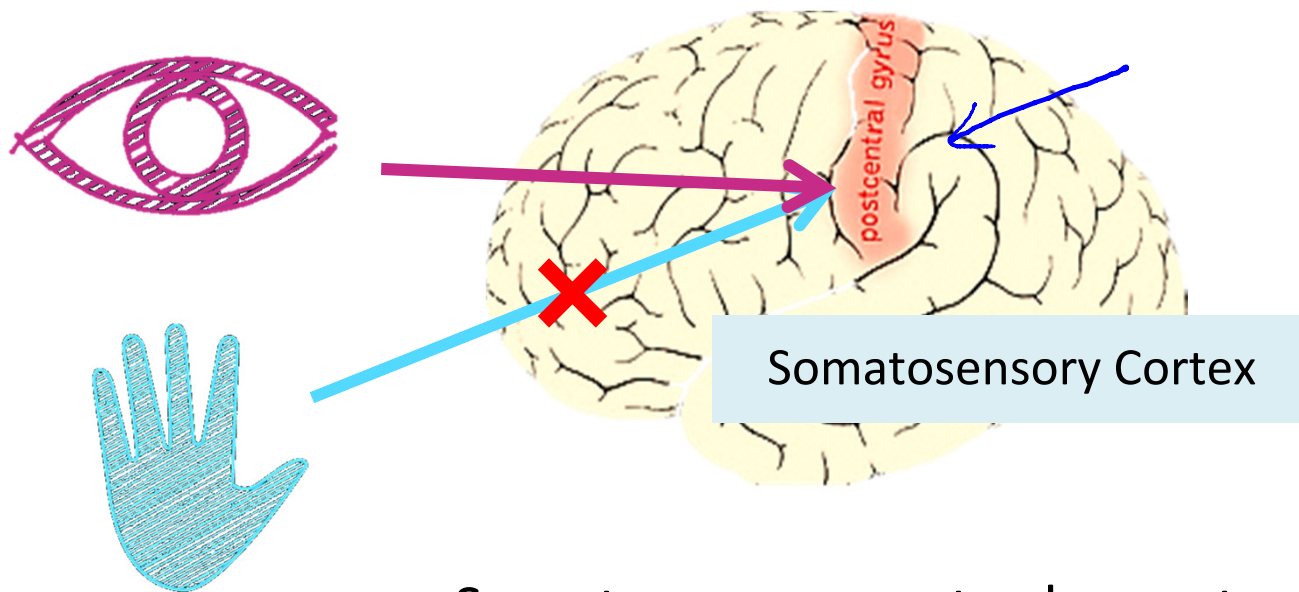


# The “one learning algorithm” hypothesis



Auditory cortex learns to see

# The “one learning algorithm” hypothesis



Somatosensory cortex learns to see

pressure, pain, warm

# Sensor representations in the brain



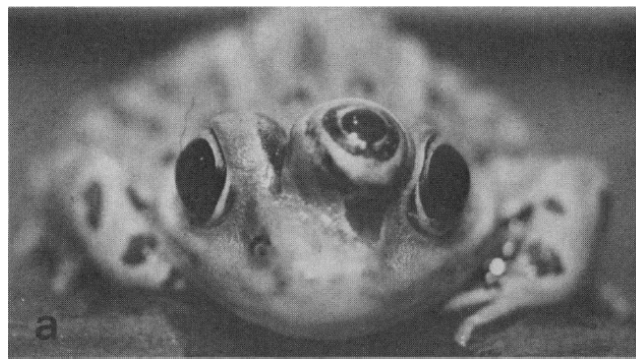
Seeing with your tongue



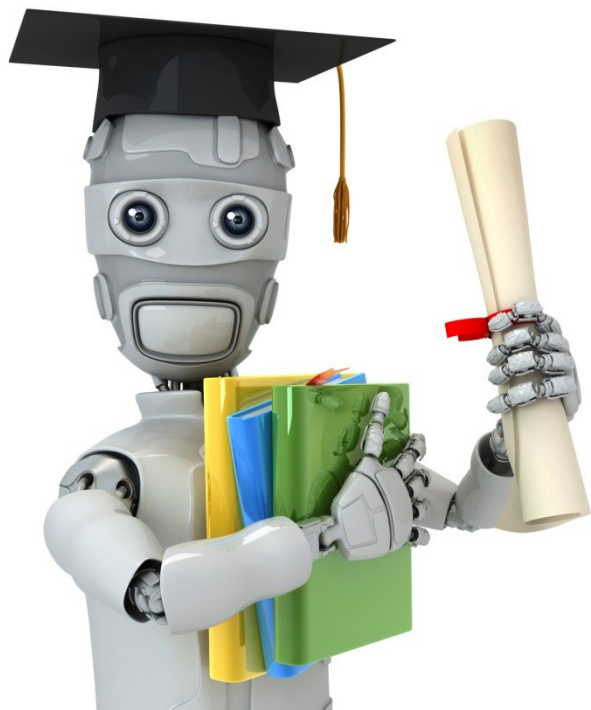
Human echolocation (sonar)



Haptic belt: Direction sense



Implanting a 3<sup>rd</sup> eye



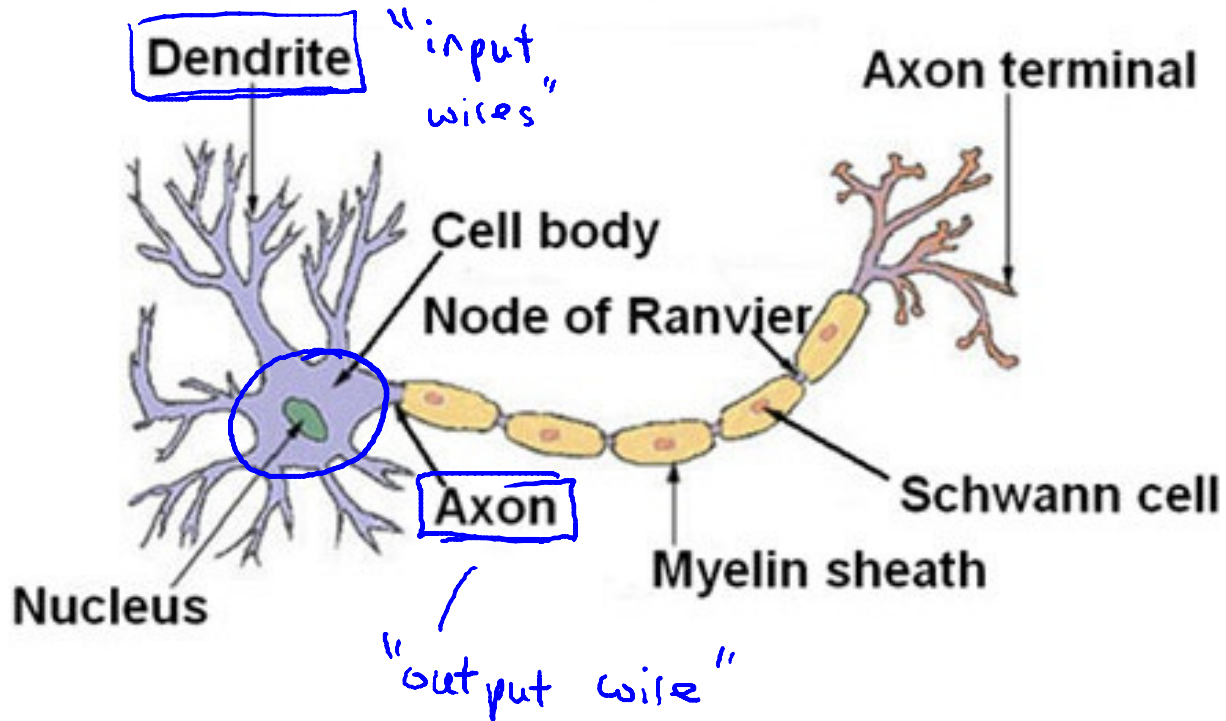
Machine Learning

# Neural Networks: Representation

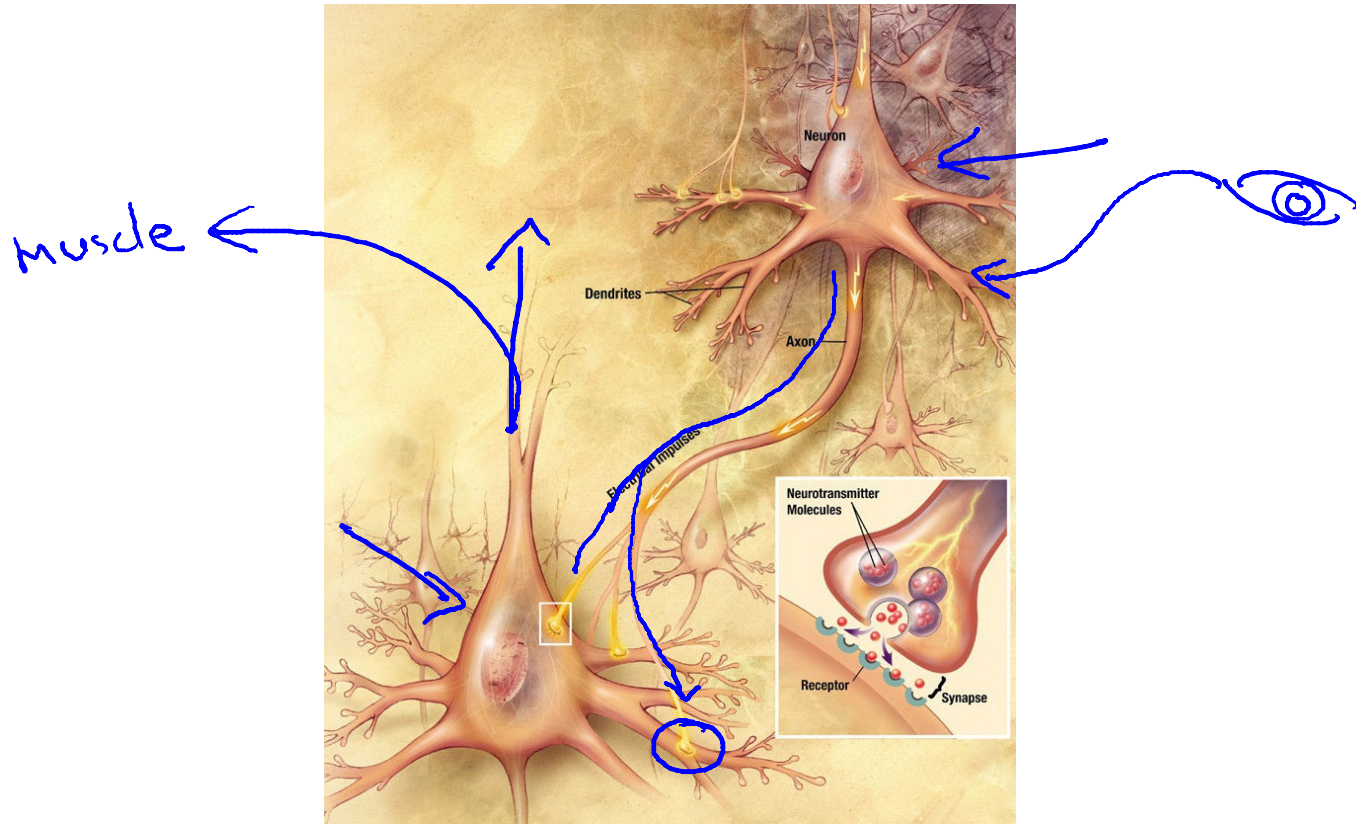
---

## Model representation I

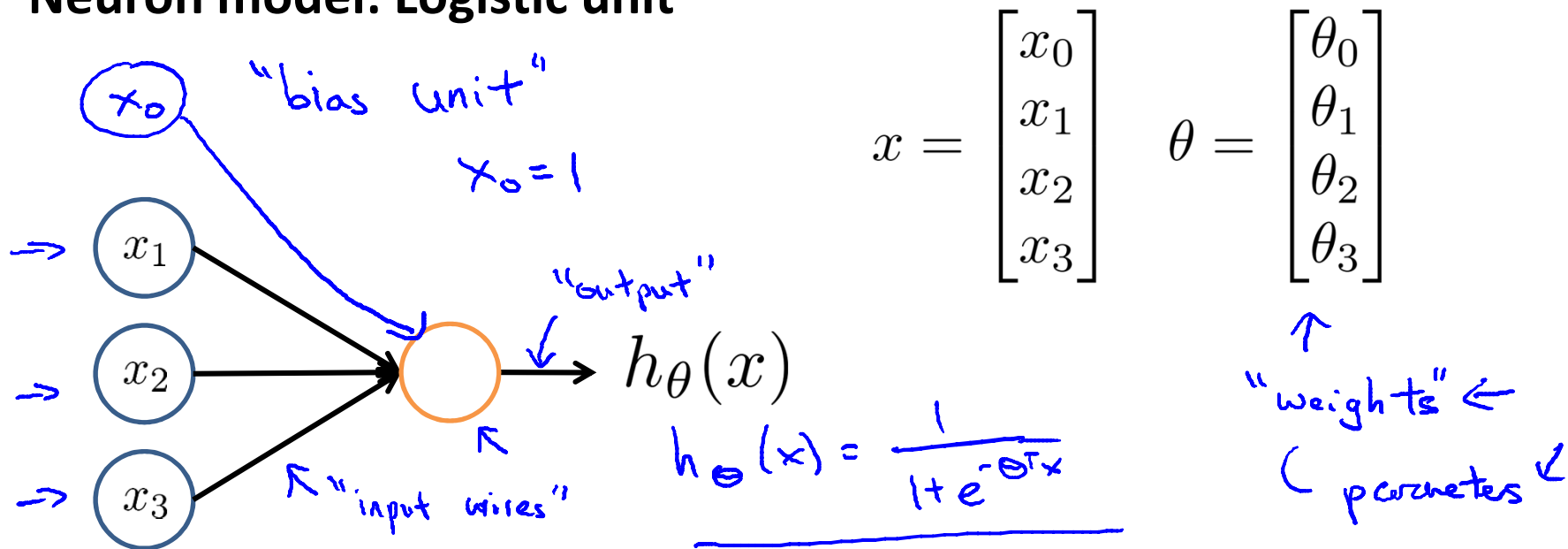
# Neuron in the brain



# Neurons in the brain



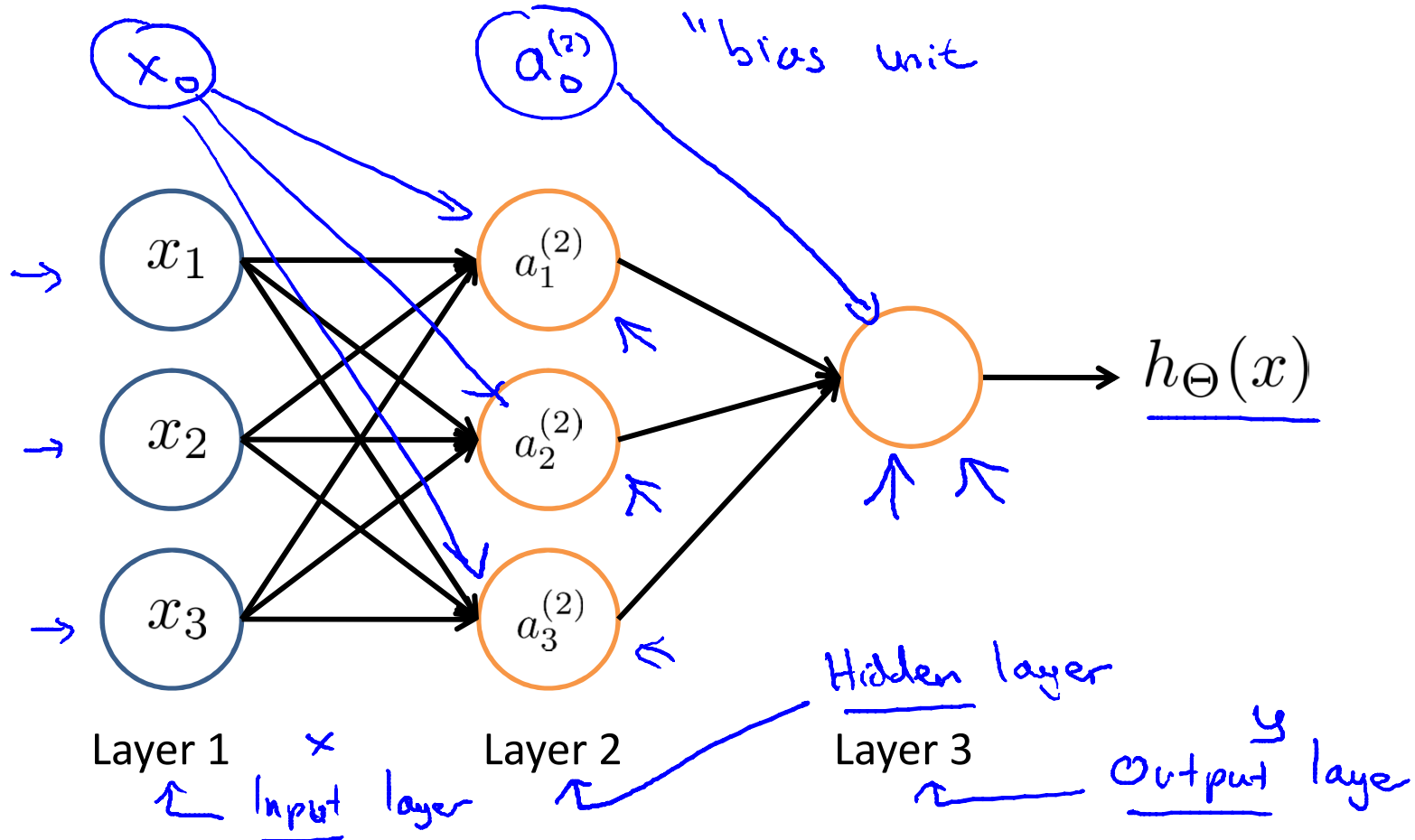
## Neuron model: Logistic unit



Sigmoid (logistic) activation function.

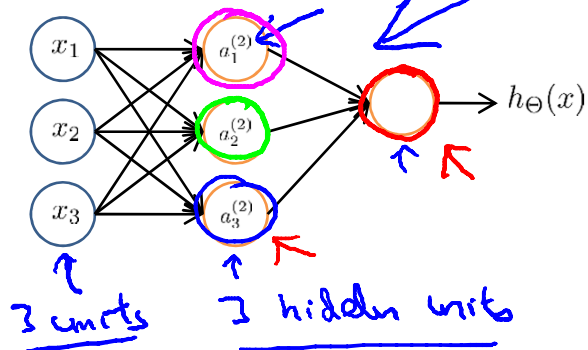
$$g(z) = \frac{1}{1 + e^{-z}}$$

# Neural Network





# Neural Network



$\rightarrow a_i^{(j)}$  = “activation” of unit  $i$  in layer  $j$

$\rightarrow \Theta^{(j)}$  = matrix of weights controlling function mapping from layer  $j$  to layer  $j + 1$

$$\Theta^{(1)} \in \mathbb{R}^{3 \times 4}$$

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

$$\rightarrow a_1^{(2)} = g(\Theta_{10}^{(1)} x_0 + \Theta_{11}^{(1)} x_1 + \Theta_{12}^{(1)} x_2 + \Theta_{13}^{(1)} x_3)$$

$$\rightarrow a_2^{(2)} = g(\Theta_{20}^{(1)} x_0 + \Theta_{21}^{(1)} x_1 + \Theta_{22}^{(1)} x_2 + \Theta_{23}^{(1)} x_3)$$

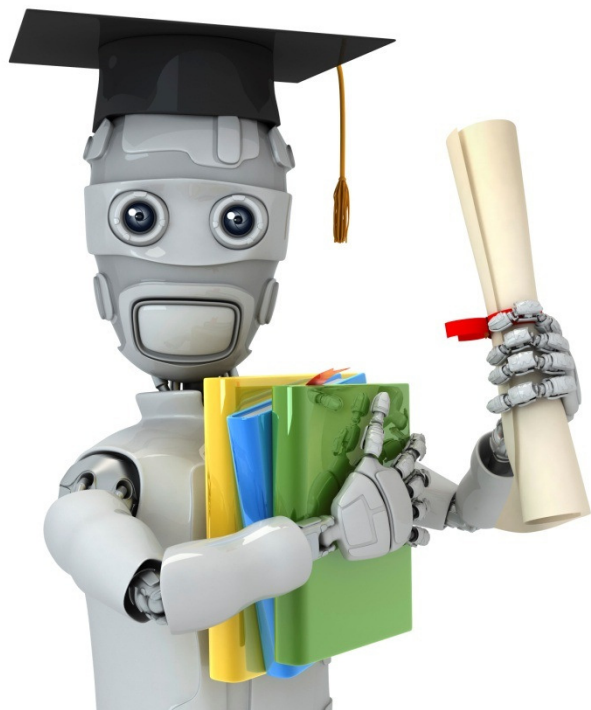
$$\Theta^{(2)}$$

$$\rightarrow a_3^{(2)} = g(\Theta_{30}^{(1)} x_0 + \Theta_{31}^{(1)} x_1 + \Theta_{32}^{(1)} x_2 + \Theta_{33}^{(1)} x_3)$$

$$\rightarrow h_{\Theta}(x) = a_1^{(3)} = g(\Theta_{10}^{(2)} a_0^{(2)} + \Theta_{11}^{(2)} a_1^{(2)} + \Theta_{12}^{(2)} a_2^{(2)} + \Theta_{13}^{(2)} a_3^{(2)})$$

$\rightarrow$  If network has  $s_j$  units in layer  $j$ ,  $s_{j+1}$  units in layer  $j + 1$ , then  $\Theta^{(j)}$  will be of dimension  $s_{j+1} \times (s_j + 1)$ .

$$s_{j+1} \times (s_j + 1)$$



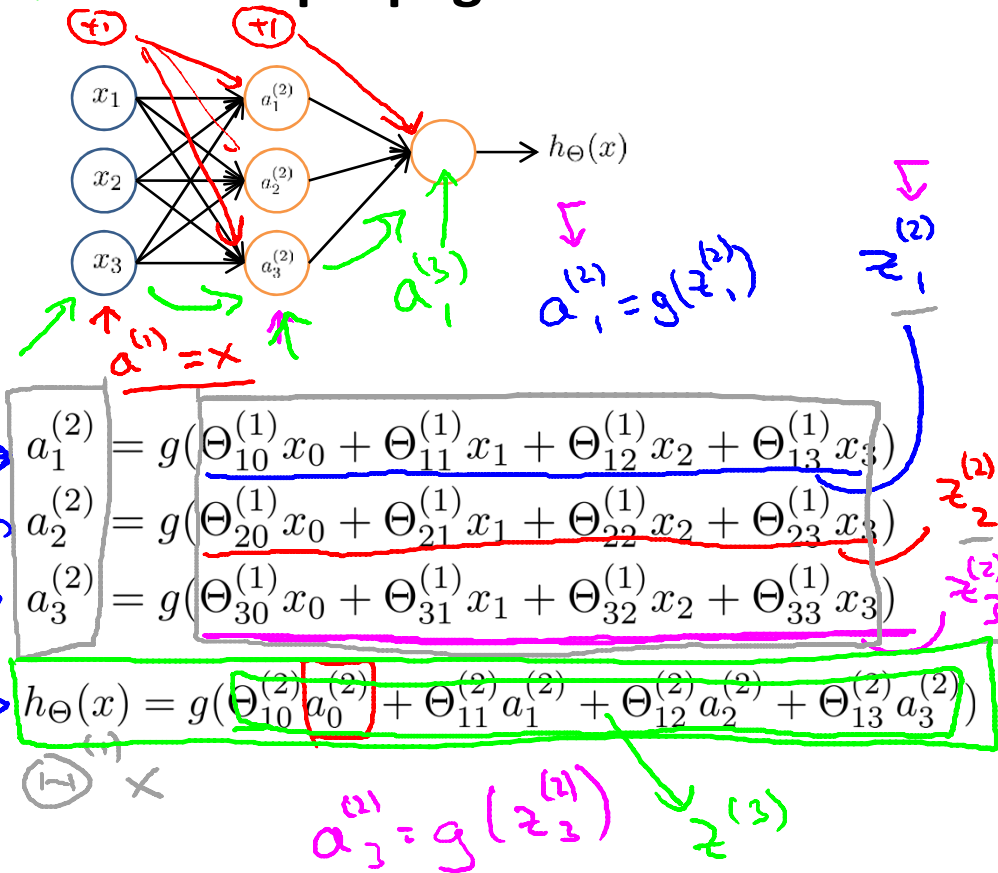
Machine Learning

# Neural Networks: Representation

---

## Model representation II

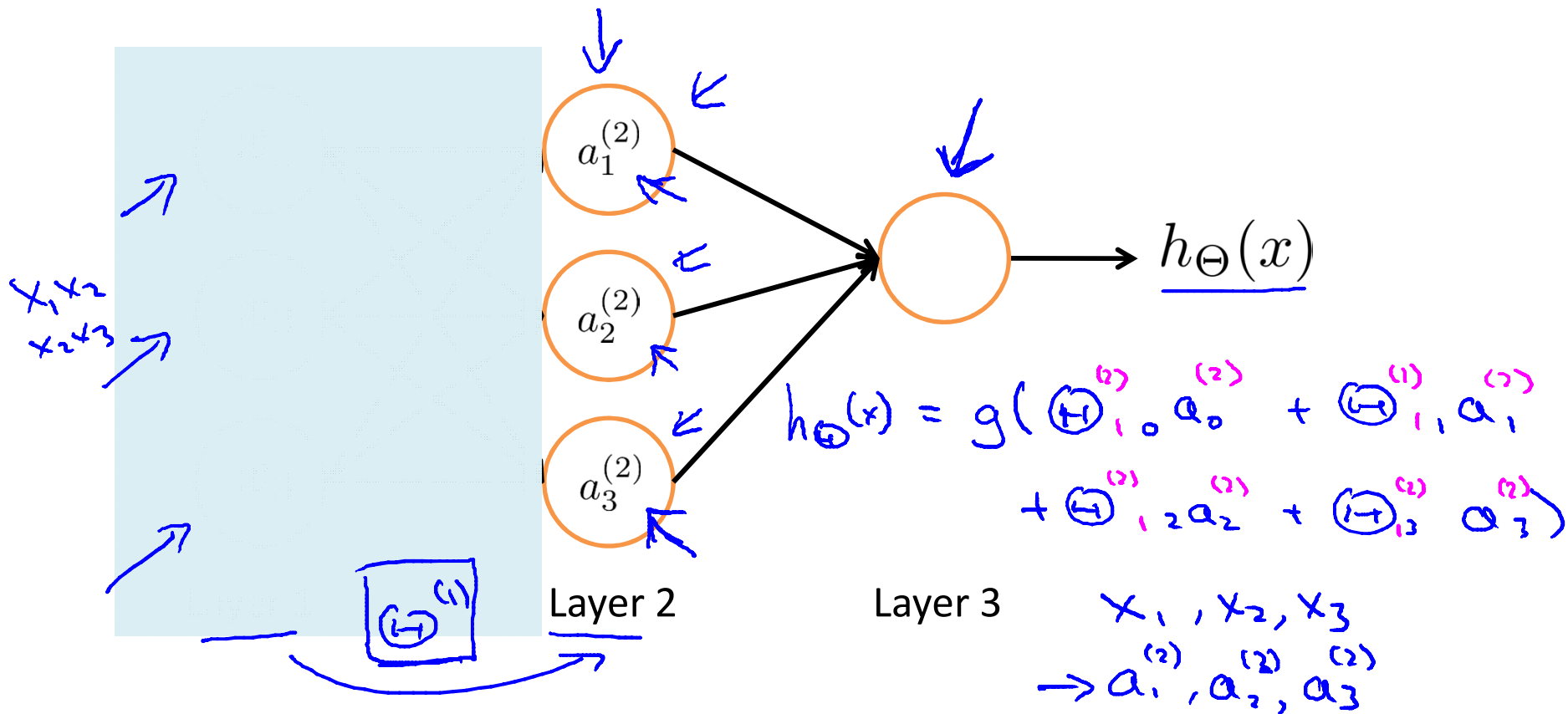
# Forward propagation: Vectorized implementation



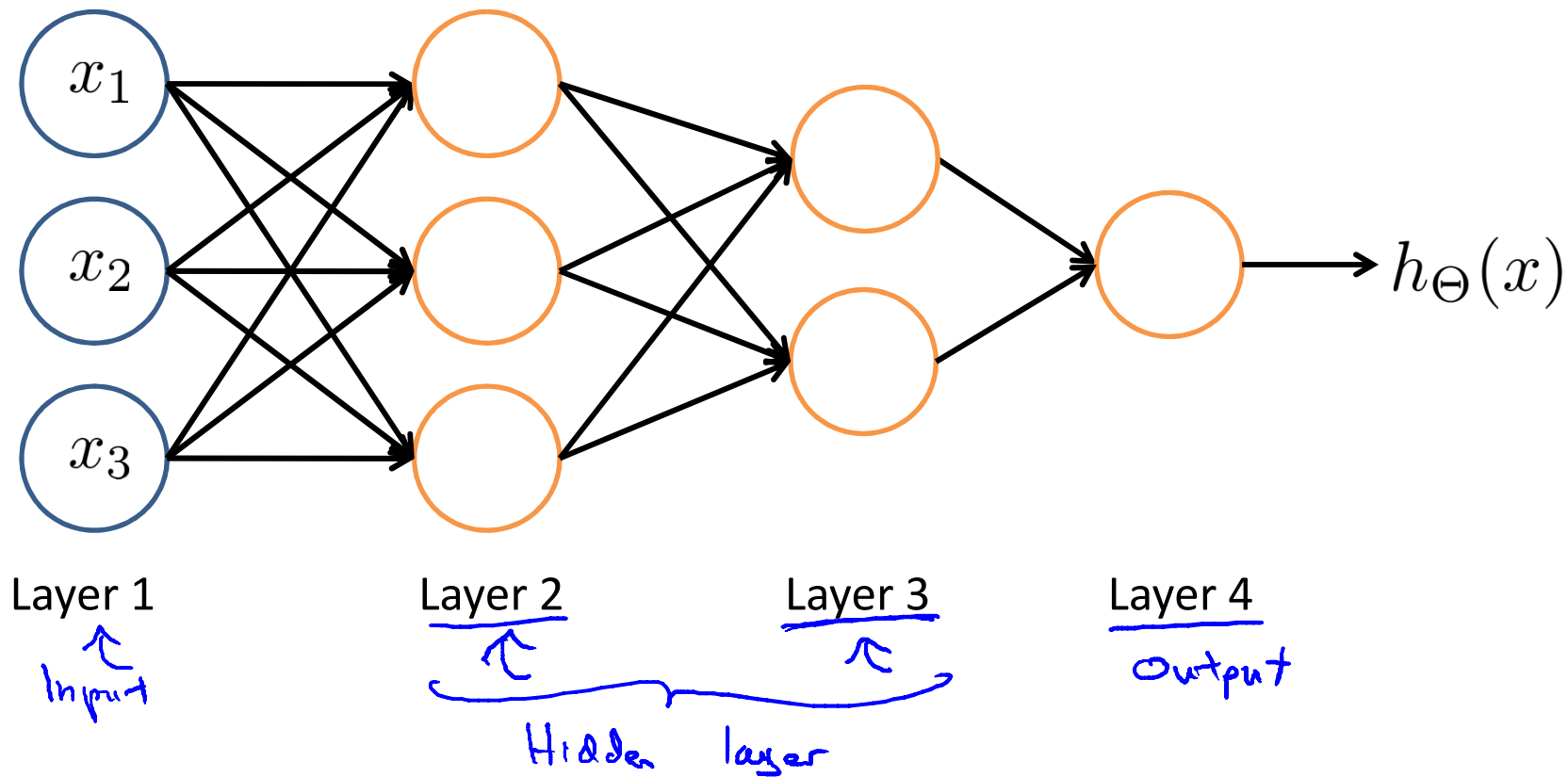
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{bmatrix} \quad z^{(2)} = \begin{bmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{bmatrix}$$

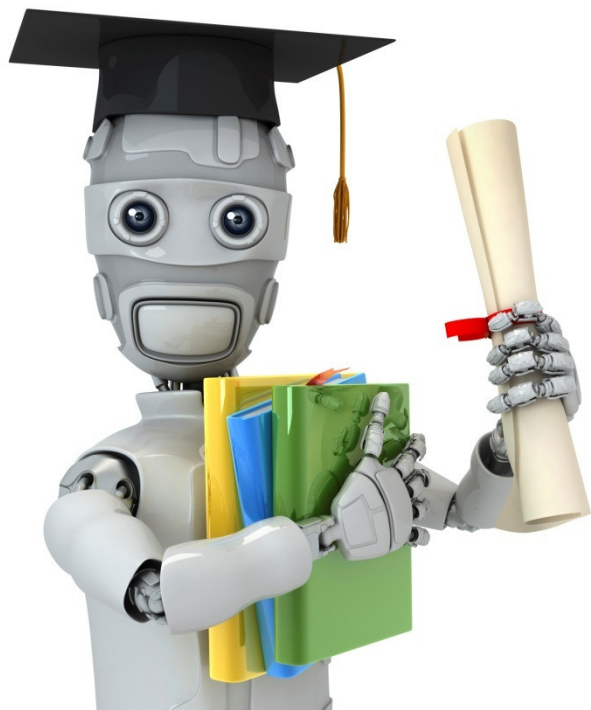
$z^{(2)} = \Theta^{(1)} a^{(1)}$   
 $a^{(2)} = g(z^{(2)})$   
 Add  $a_0^{(2)} = 1$ .  $\rightarrow a^{(2)} \in \mathbb{R}^4$   
 $z^{(3)} = \Theta^{(2)} a^{(2)}$   
 $h_{\Theta}(x) = a^{(3)} = g(z^{(3)})$

# Neural Network learning its own features



## Other network architectures





Machine Learning

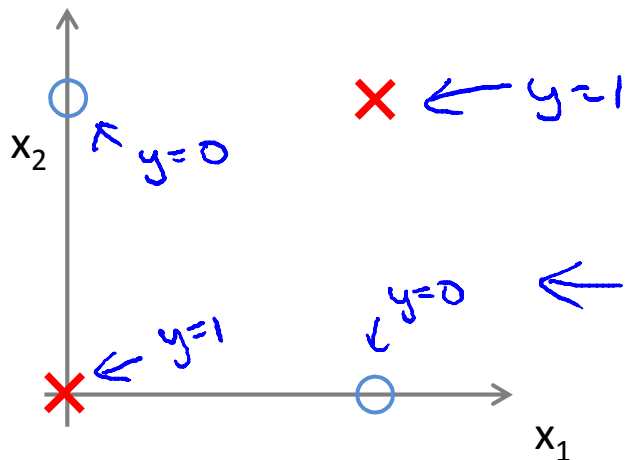
# Neural Networks: Representation

---

## Examples and intuitions I

# Non-linear classification example: XOR/XNOR

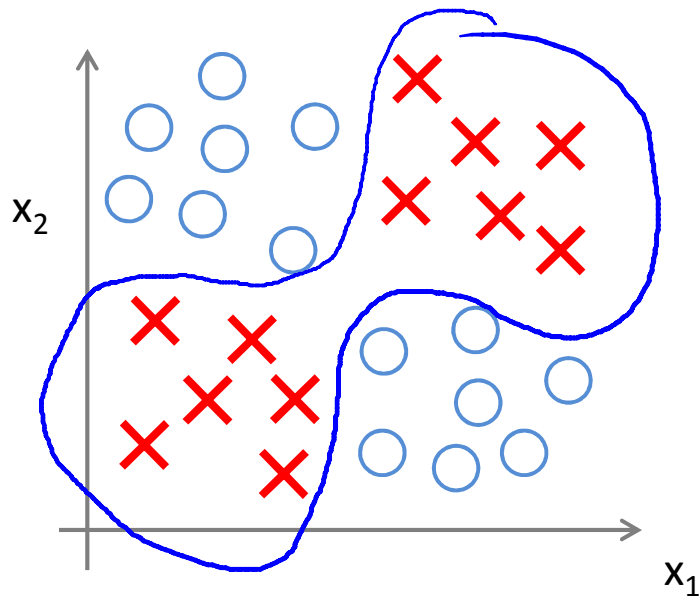
→  $x_1, x_2$  are binary (0 or 1).



$$y = \underline{x_1 \text{ XOR } x_2}$$

$$\underbrace{x_1 \text{ XNOR } x_2}_{\leftarrow}$$

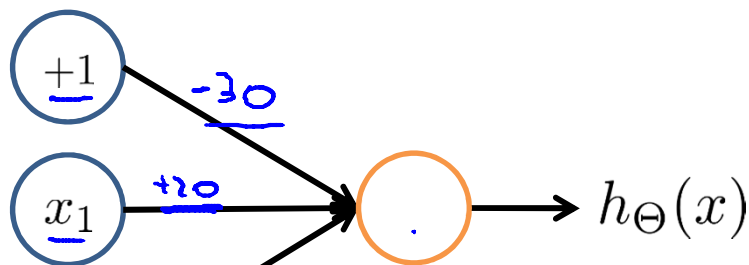
$$\underbrace{\text{NOT } (x_1 \text{ XOR } x_2)}$$



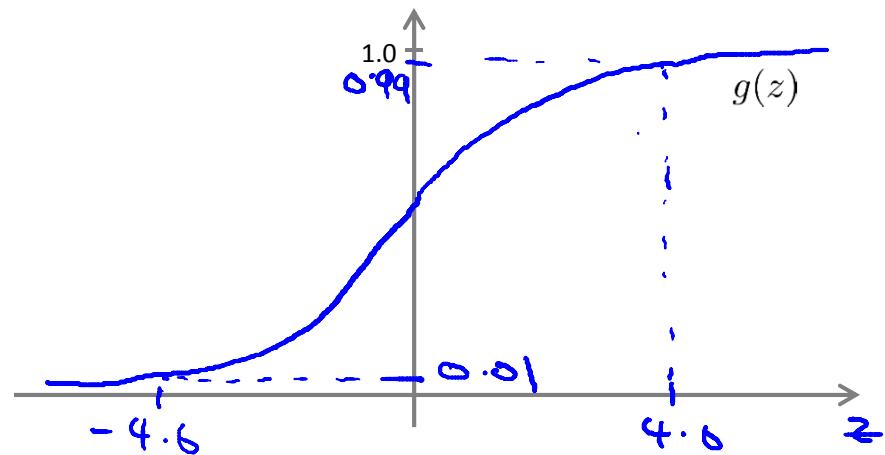
# Simple example: AND

→  $x_1, x_2 \in \{0, 1\}$

→  $y = x_1 \text{ AND } x_2$



→ 
$$h_{\Theta}(x) = g\left(\underbrace{-30}_{w_{10}} + \underbrace{20}_{w_{11}}x_1 + \underbrace{20}_{w_{12}}x_2\right)$$

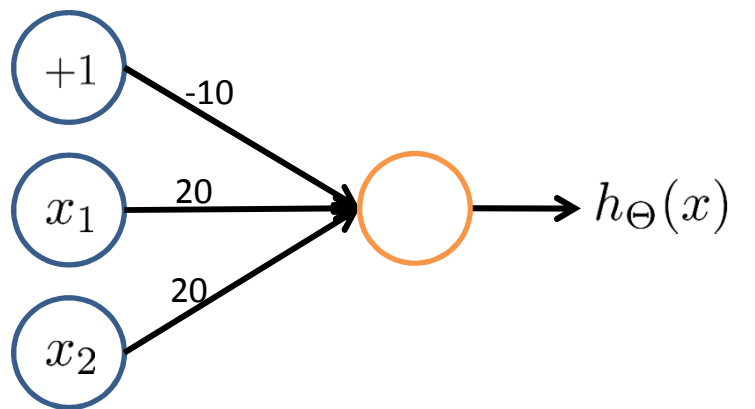


$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	$g(-30) \approx 0$
→ 0	1	$g(-10) \approx 0$
1	0	$g(-10) \approx 0$
→ 1	1	$g(10) \approx 1$

$h_{\Theta}(x) \approx x_1 \text{ AND } x_2$



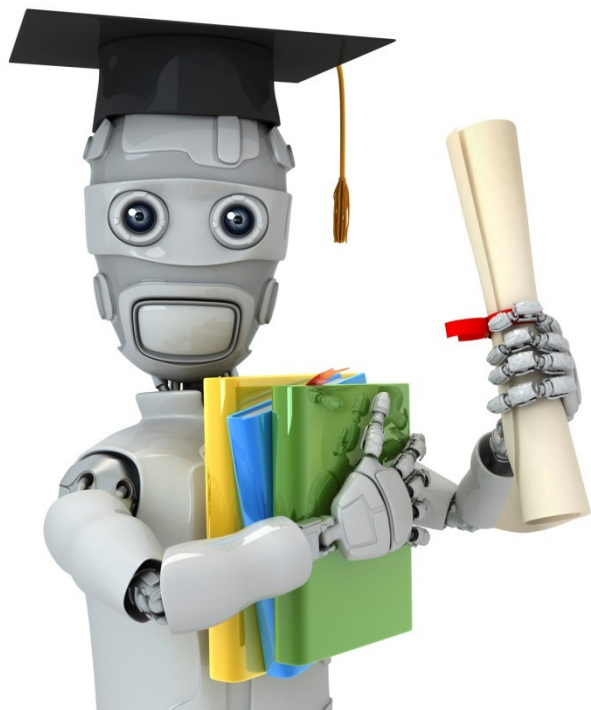
## Example: OR function



$$g(-10 + 20x_1 + 20x_2)$$

$x_1$	$x_2$	$h_{\Theta}(x)$
0	0	$g(-10) \approx 0$
0	1	$g(10) \approx 1$
1	0	$\approx 1$
1	1	$\approx 1$





Machine Learning

# Neural Networks: Representation

---

## Examples and intuitions II

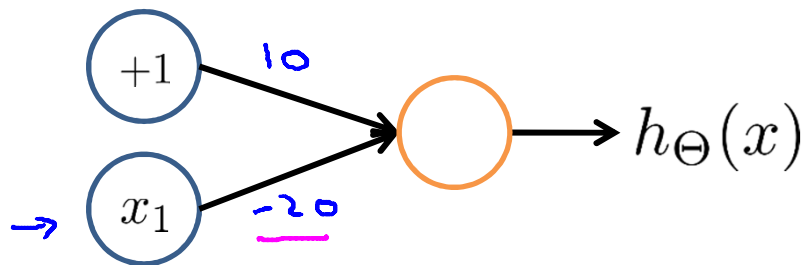
$\rightarrow x_1$  AND  $x_2$

$\rightarrow x_1$  OR  $x_2$

$\{0,1\}$ .

**Negation:**

NOT  $x_1$

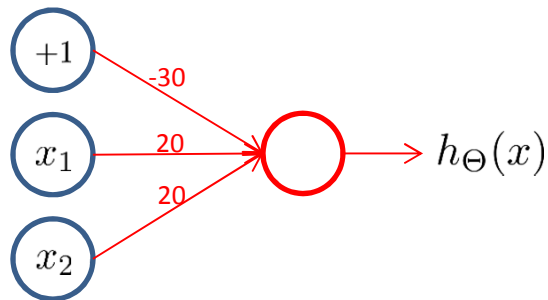
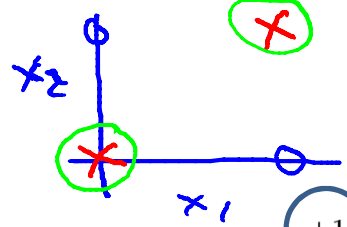


$x_1$	$h_{\Theta}(x)$
0	$g(10) \approx 1$
1	$g(-10) \approx 0$

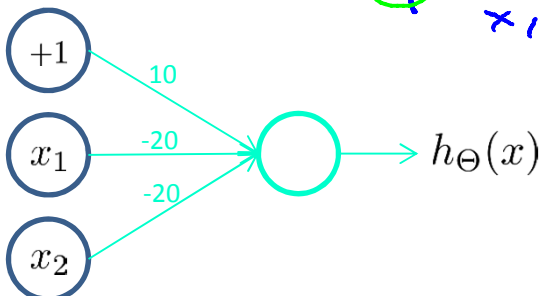
$$h_{\Theta}(x) = g(10 - 20x_1)$$

$\rightarrow$  (NOT  $x_1$ ) AND (NOT  $x_2$ )  
(= 1 if and only if  
 $\rightarrow x_1 = x_2 = 0$ )

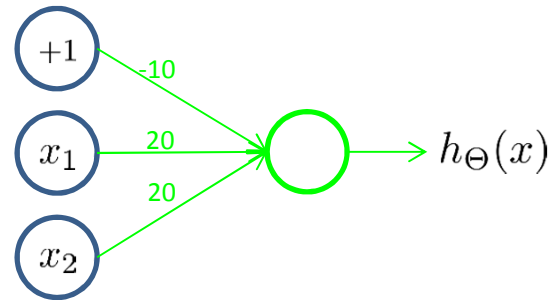
# Putting it together: $x_1$ XNOR $x_2$



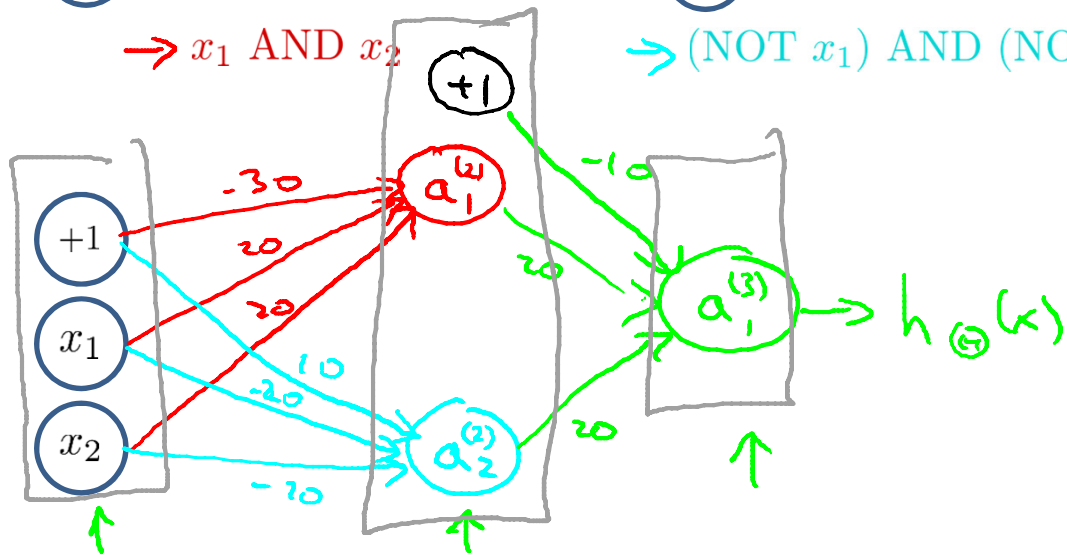
$\rightarrow x_1$  AND  $x_2$



$\rightarrow$  (NOT  $x_1$ ) AND (NOT  $x_2$ )

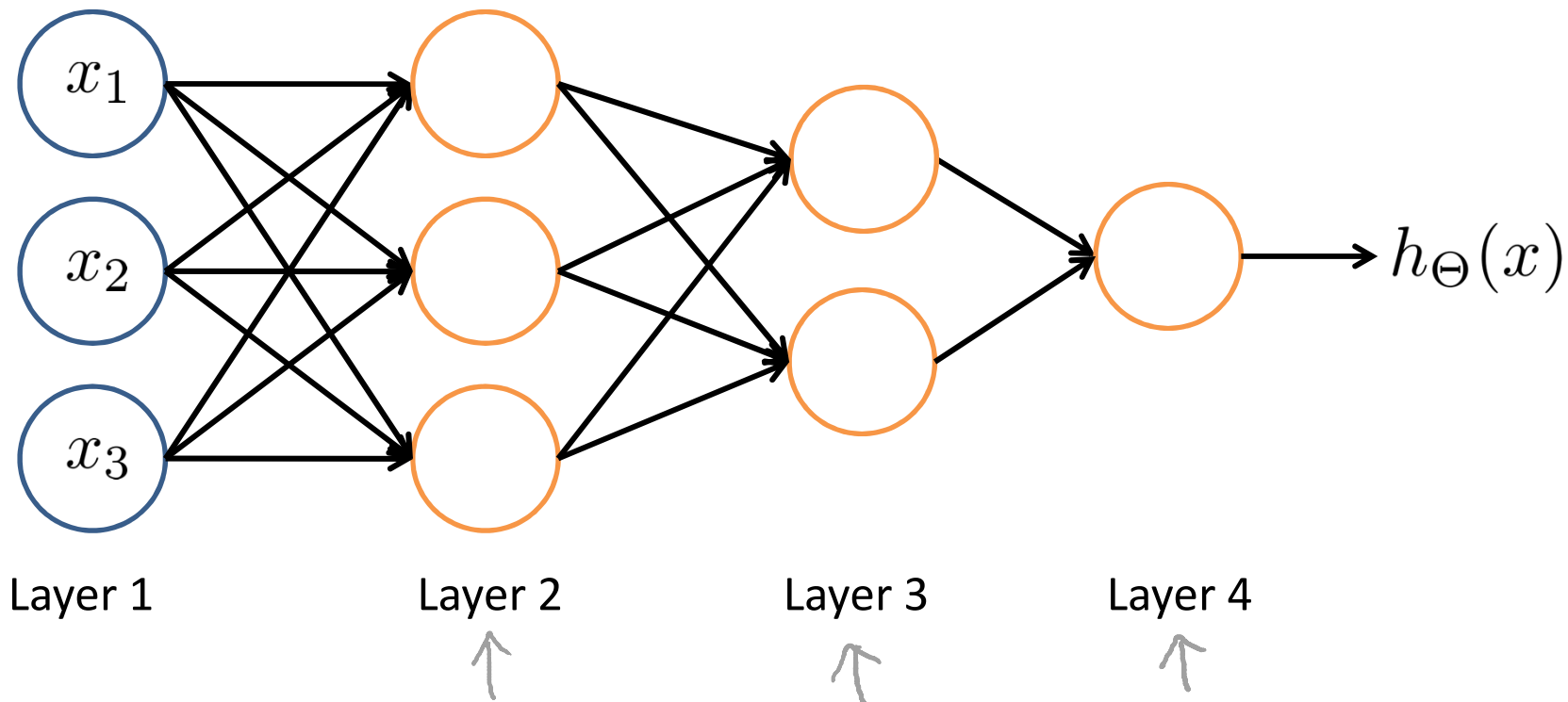


$\rightarrow x_1$  OR  $x_2$

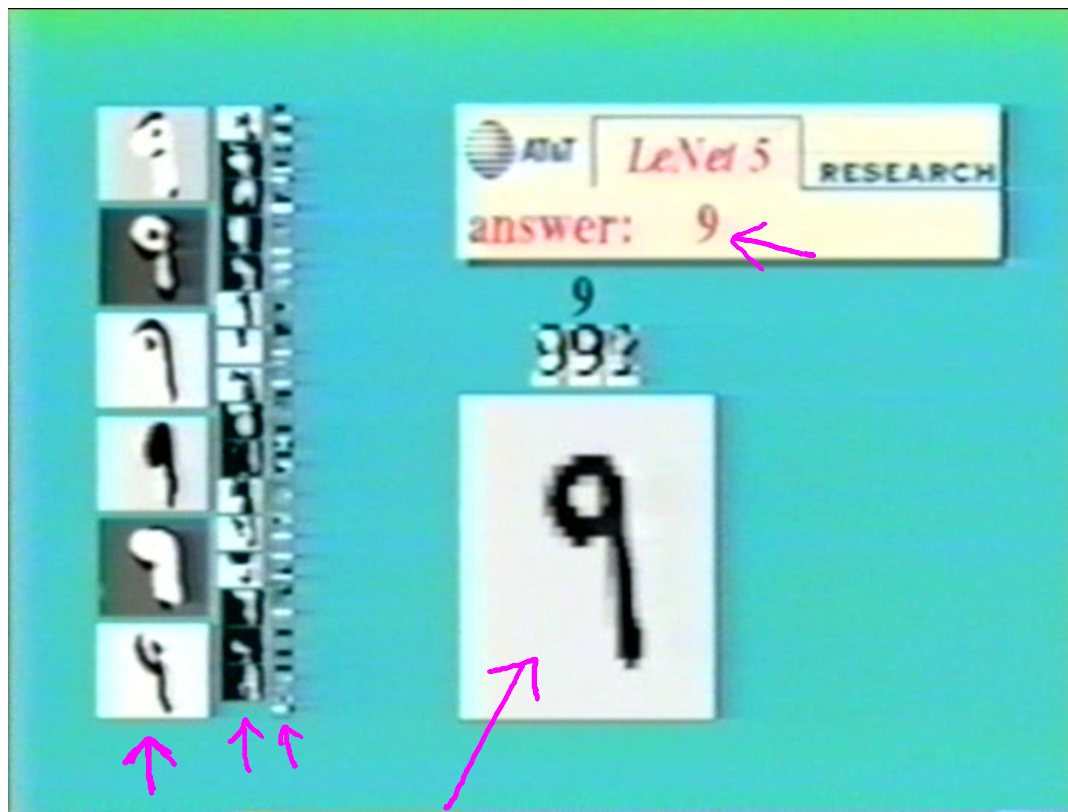


$x_1$	$x_2$	$a_1^{(2)}$	$a_2^{(2)}$	$h_{\Theta}(x)$
0	0	0	1	1
0	1	0	0	0
1	0	0	0	0
1	1	1	0	1

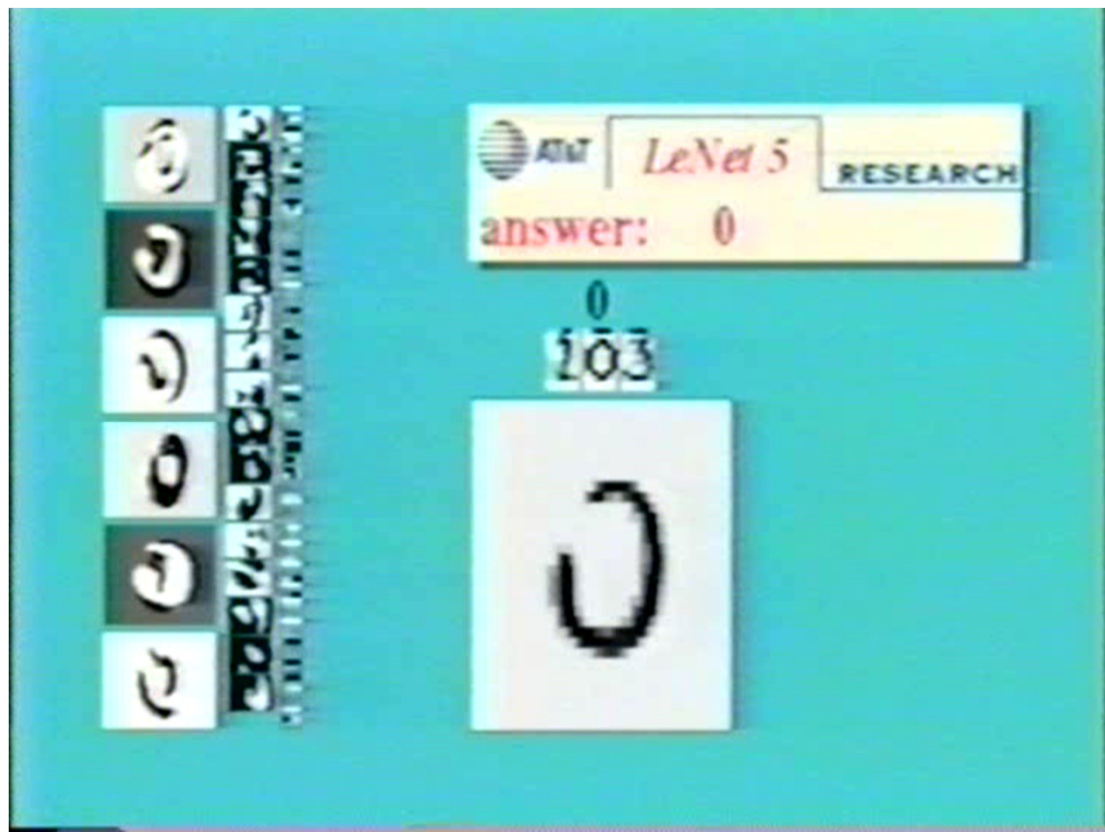
# Neural Network intuition



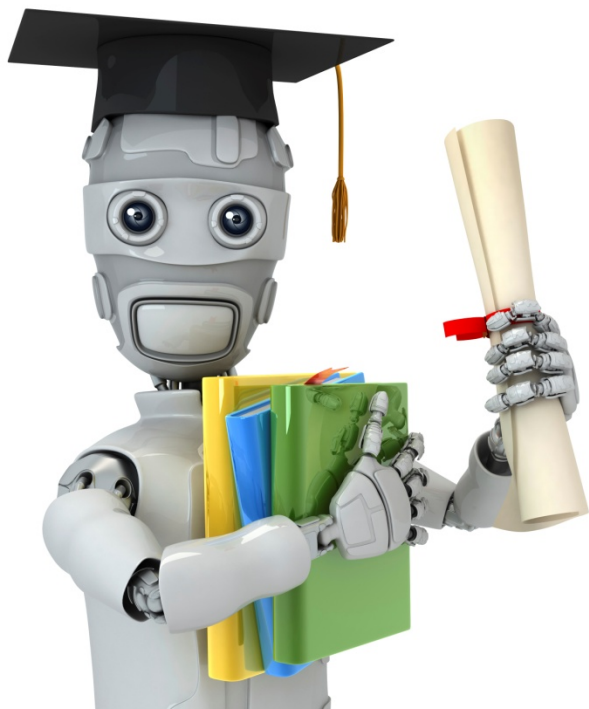
# Handwritten digit classification



# Handwritten digit classification







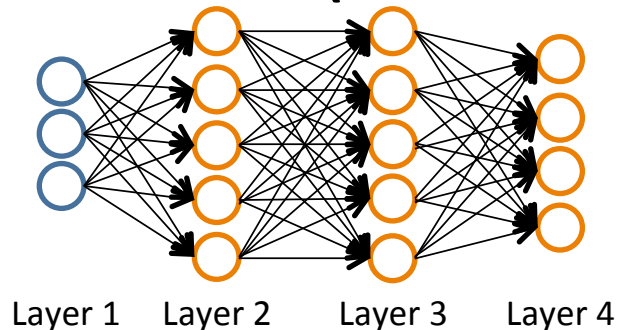
Machine Learning

# Neural Networks: Learning

---

## Cost function

# Neural Network (Classification)



$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$$

$L$  = total no. of layers in network

$s_l$  = no. of units (not counting bias unit) in layer  $l$

## Binary classification

$y = 0$  or  $1$

1 output unit

## Multi-class classification (K classes)

$y \in \mathbb{R}^K$  E.g.  $\begin{bmatrix} 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}$ ,  $\begin{bmatrix} 0 \\ 0 \\ 0 \\ 1 \end{bmatrix}$   
pedestrian car motorcycle truck

K output units

## Cost function

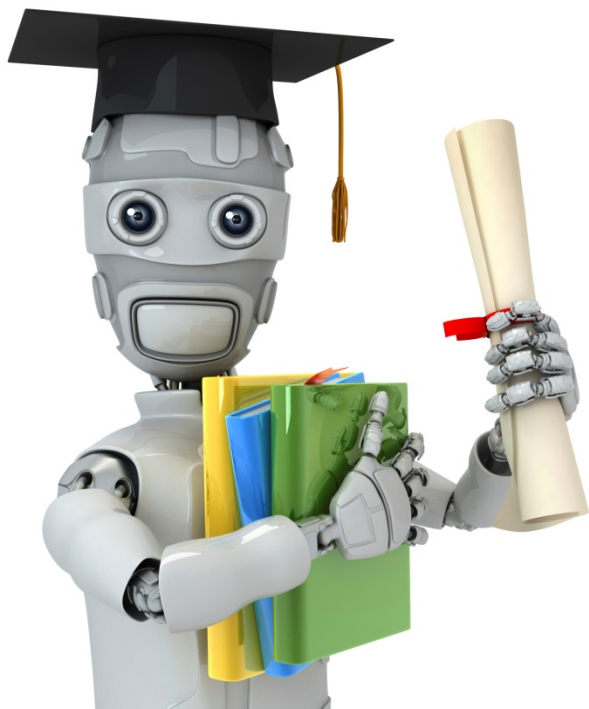
Logistic regression:

$$J(\theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(x^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

Neural network:

$$h_{\Theta}(x) \in \mathbb{R}^K \quad (h_{\Theta}(x))_i = i^{th} \text{ output}$$

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log(h_{\Theta}(x^{(i)}))_k + (1 - y_k^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$



Machine Learning

# Neural Networks: Learning

---

## Backpropagation algorithm

## Gradient computation

$$\begin{aligned} \rightarrow \underline{J(\Theta)} &= -\frac{1}{m} \left[ \sum_{i=1}^m \sum_{k=1}^K y_k^{(i)} \log h_{\theta}(x^{(i)})_k + (1 - y_k^{(i)}) \log(1 - h_{\theta}(x^{(i)})_k) \right] \\ &+ \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_j^{(l)})^2 \end{aligned}$$

$$\rightarrow \min_{\Theta} J(\Theta)$$

Need code to compute:

$$\rightarrow - \underline{J(\Theta)}$$

$$\rightarrow - \underline{\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)} \leftarrow$$

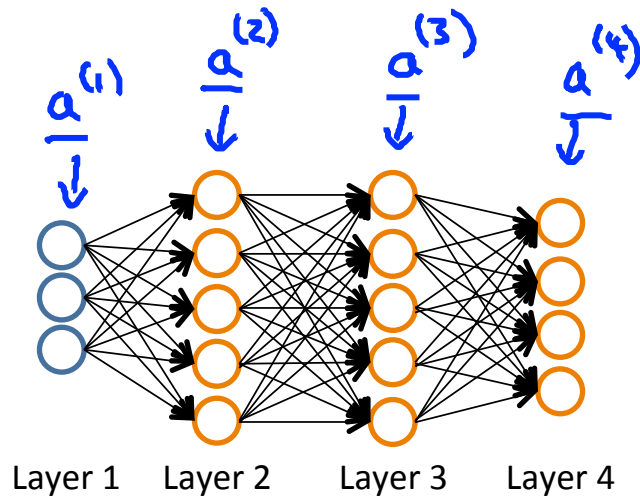
$$\Theta_{ij}^{(l)} \in \mathbb{R}$$

# Gradient computation

Given one training example  $(x, y)$ :

Forward propagation:

- $\rightarrow \underline{a^{(1)}} = \underline{x}$
- $\rightarrow z^{(2)} = \Theta^{(1)} a^{(1)}$
- $\rightarrow a^{(2)} = g(z^{(2)})$  (add  $\underline{a_0^{(2)}}$ )
- $\rightarrow z^{(3)} = \Theta^{(2)} a^{(2)}$
- $\rightarrow a^{(3)} = g(z^{(3)})$  (add  $a_0^{(3)}$ )
- $\rightarrow z^{(4)} = \Theta^{(3)} a^{(3)}$
- $\rightarrow \underline{a^{(4)}} = \underline{h_{\Theta}(x)} = g(z^{(4)})$



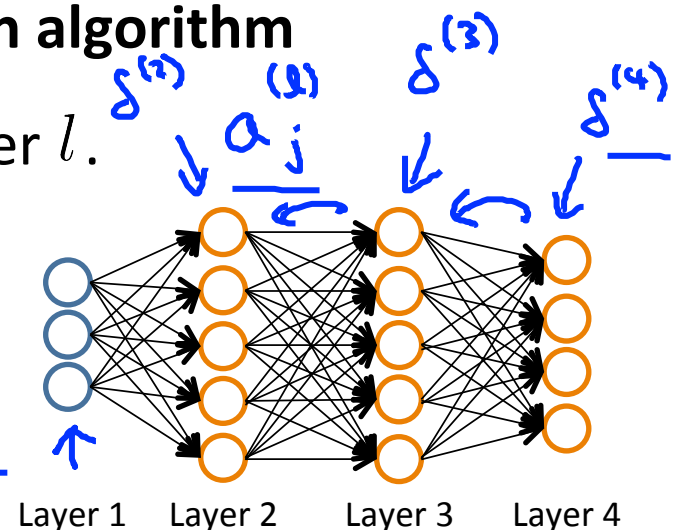
# Gradient computation: Backpropagation algorithm

Intuition:  $\delta_j^{(l)}$  = "error" of node  $j$  in layer  $l$ .

For each output unit (layer  $L = 4$ )

$$\delta_j^{(4)} = a_j^{(4)} - y_j$$

$$(h_{\Theta^{(4)}})_j \quad \delta_j^{(4)} = a_j^{(4)} - y_j$$



$$\delta^{(3)} = (\Theta^{(3)})^T \delta^{(4)} \cdot * g'(z^{(3)})$$

$$\delta^{(2)} = (\Theta^{(2)})^T \delta^{(3)} \cdot * g'(z^{(2)})$$

$$\frac{a^{(3)}}{a^{(2)}} \cdot * (1 - a^{(3)})$$

$$\frac{a^{(3)}}{a^{(2)}} \cdot * (1 - a^{(2)})$$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = a_j^{(l)} \delta_i^{(l+1)}$$

(ignoring  $\lambda$ ; if  $\lambda = 0$ )

# Backpropagation algorithm

→ Training set  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

Set  $\Delta_{ij}^{(l)} = 0$  (for all  $l, i, j$ ).

(used to compute  $\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta)$ )

For  $i = 1$  to  $m$  ←  $(\underline{x^{(i)}}, \underline{y^{(i)}})$

Set  $\underline{a^{(1)}} = \underline{x^{(i)}}$

→ Perform forward propagation to compute  $\underline{a^{(l)}}$  for  $l = 2, 3, \dots, L$

→ Using  $\underline{y^{(i)}}$ , compute  $\underline{\delta^{(L)}} = \underline{a^{(L)}} - \underline{y^{(i)}}$

→ Compute  $\underline{\delta^{(L-1)}}, \underline{\delta^{(L-2)}}, \dots, \underline{\delta^{(2)}}$  ~~use~~

→  $\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$  ←

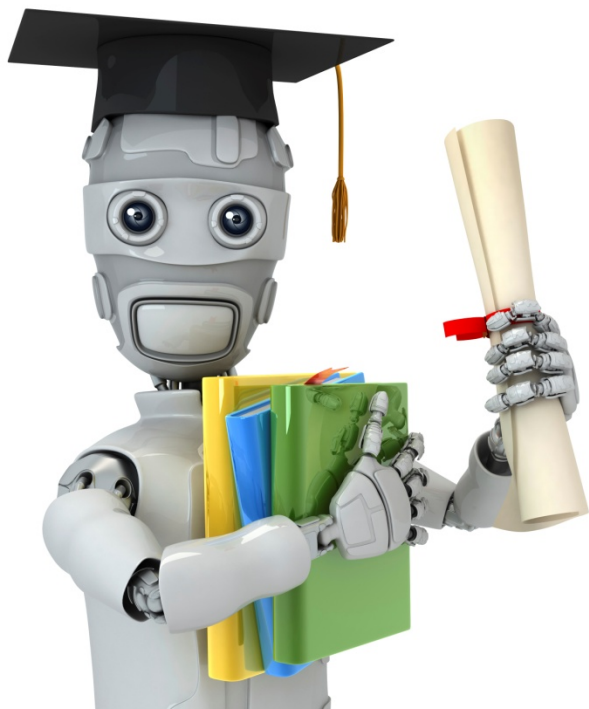
$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + \delta^{(l+1)} (a^{(l)})^T$

→  $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}$  if  $j \neq 0$

→  $D_{ij}^{(l)} := \frac{1}{m} \Delta_{ij}^{(l)}$  if  $j = 0$

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} J(\Theta) = D_{ij}^{(l)}$$





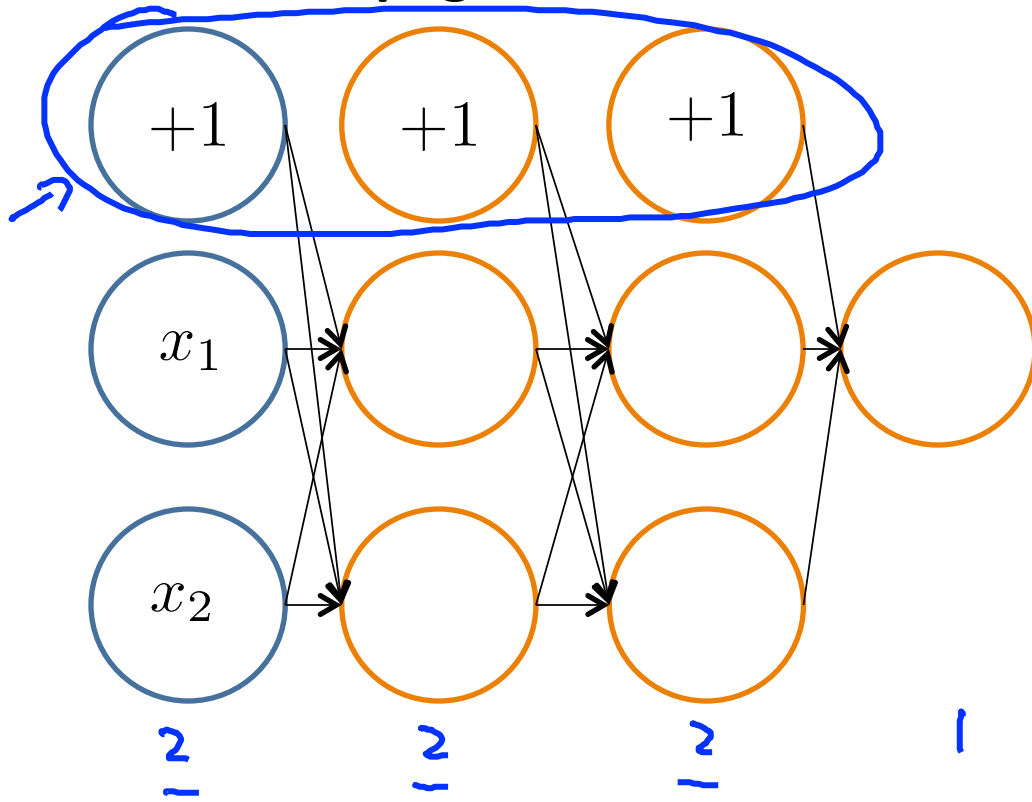
Machine Learning

# Neural Networks: Learning

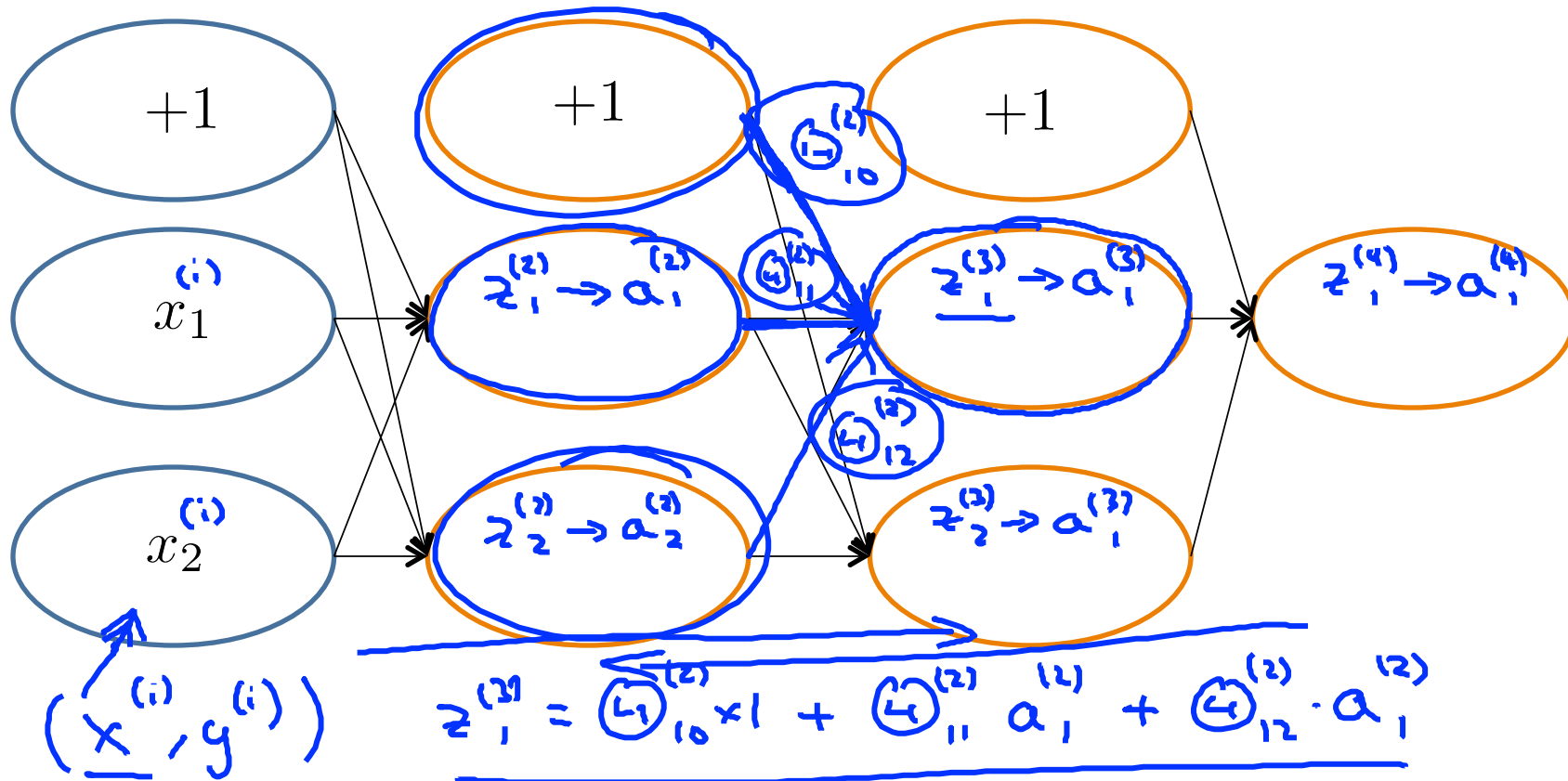
---

## Backpropagation intuition

# Forward Propagation



# Forward Propagation



# What is backpropagation doing?

$$J(\Theta) = -\frac{1}{m} \left[ \sum_{i=1}^m y^{(i)} \log(h_{\Theta}(x^{(i)})) + (1 - y^{(i)}) \log(1 - (h_{\Theta}(x^{(i)}))) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2$$

$(x^{(i)}, y^{(i)})$

Focusing on a single example  $x^{(i)}, y^{(i)}$ , the case of 1 output unit, and ignoring regularization ( $\lambda = 0$ ),

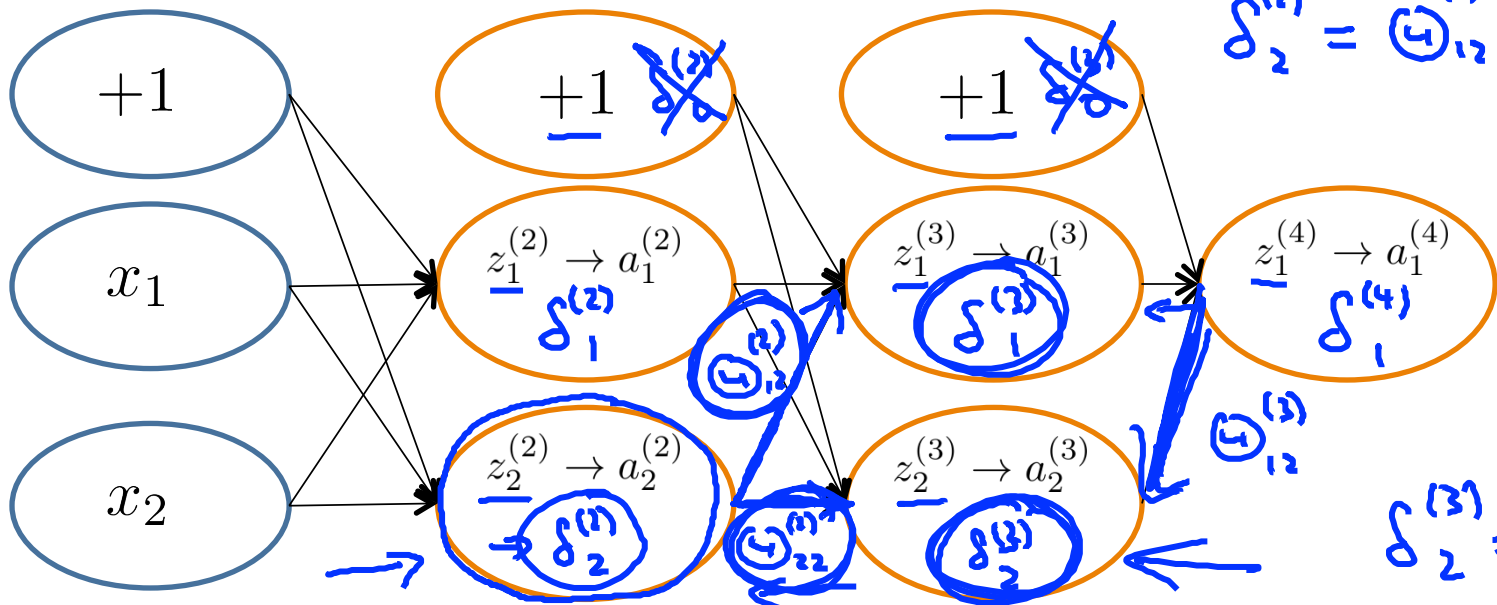
Note: Mistake on lecture, it is supposed to be  $1-h(x)$ .

$$\text{cost}(i) = y^{(i)} \log h_{\Theta}(x^{(i)}) + (1 - y^{(i)}) \log h_{\Theta}(x^{(i)})$$

(Think of  $\text{cost}(i) \approx (h_{\Theta}(x^{(i)}) - y^{(i)})^2$ )

I.e. how well is the network doing on example  $i$ ?

# Forward Propagation



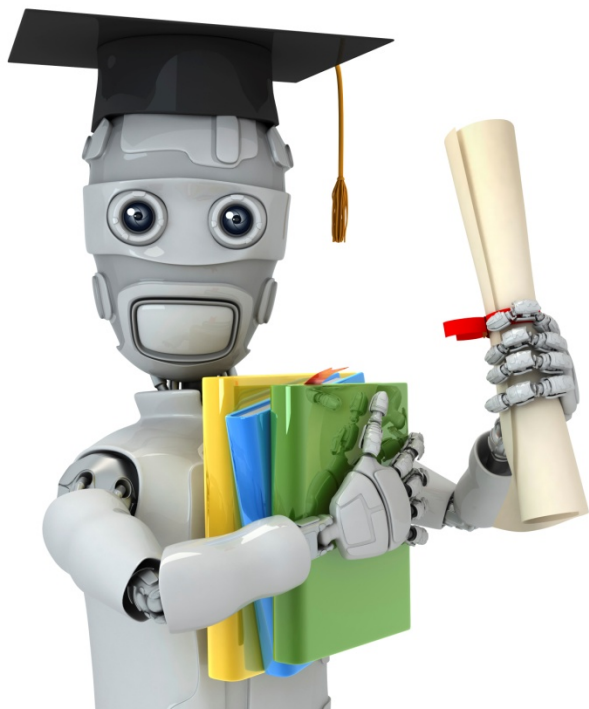
$$\delta_2^{(2)} = \delta_{12}^{(3)} \delta_1^{(3)} + \delta_{22}^{(3)} \delta_2^{(3)}$$

$$\delta_2^{(3)} = \delta_{12}^{(3)} \delta_1^{(4)}$$

→  $\delta_j^{(l)}$  = "error" of cost for  $a_j^{(l)}$  (unit  $j$  in layer  $l$ ).

Formally,  $\delta_j^{(l)} = \frac{\partial \text{cost}(i)}{\partial z_j^{(l)}}$  (for  $j \geq 0$ ), where

$$\text{cost}(i) = y^{(i)} \log h_{\Theta}(x^{(i)}) + (1 - y^{(i)}) \log h_{\Theta}(x^{(i)})$$



Machine Learning

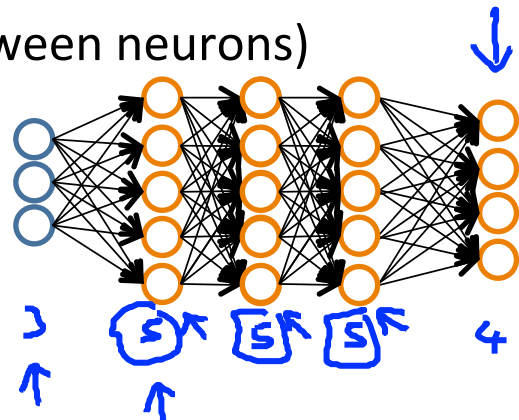
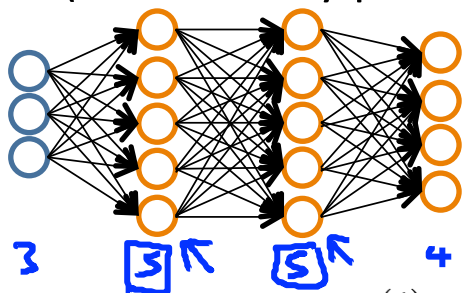
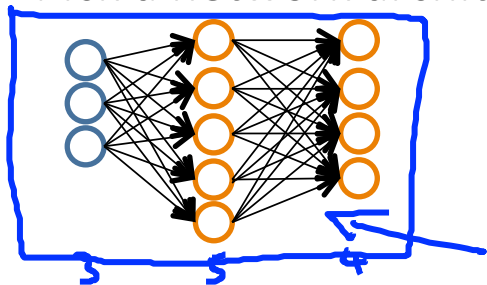
# Neural Networks: Learning

---

# Putting it together

# Training a neural network

Pick a network architecture (connectivity pattern between neurons)



→ No. of input units: Dimension of features  $x^{(i)}$

→ No. output units: Number of classes

Reasonable default: 1 hidden layer, or if >1 hidden layer, have same no. of hidden units in every layer (usually the more the better)

$y \in \{1, 2, 3, \dots, 10\}$   
 ~~$y \in \{1, 2, 3, \dots, 10\}$~~

$$y = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \text{ or } \begin{bmatrix} 0 \\ 0 \\ 0 \\ \vdots \\ 0 \\ 1 \end{bmatrix}$$

# Training a neural network

- 1. Randomly initialize weights
- 2. Implement forward propagation to get  $h_{\Theta}(x^{(i)})$  for any  $x^{(i)}$
- 3. Implement code to compute cost function  $J(\Theta)$
- 4. Implement backprop to compute partial derivatives  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$

$$\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$$

→ for  $i = 1:m$  {  $(x^{(1)}, y^{(1)})$   $(x^{(2)}, y^{(2)})$ , ...,  $(x^{(m)}, y^{(m)})$  }

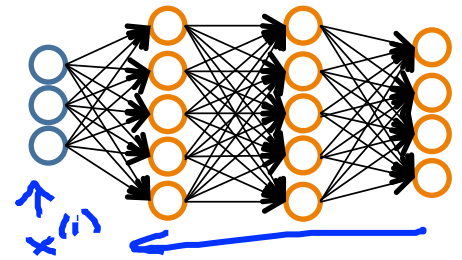
→ Perform forward propagation and backpropagation using example  $(x^{(i)}, y^{(i)})$

(Get activations  $a^{(l)}$  and delta terms  $\delta^{(l)}$  for  $l = 2, \dots, L$ ).

→  $\Delta^{(2)} := \Delta^{(2)} - \delta^{(2)} (a^{(1)})^T$

...

compute  $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ .





## Training a neural network

- 5. Use gradient checking to compare  $\frac{\partial}{\partial \Theta^{(l)}_{ik}} J(\Theta)$  computed using backpropagation vs. using numerical estimate of gradient of  $J(\Theta)$ .
- Then disable gradient checking code.
- 6. Use gradient descent or advanced optimization method with backpropagation to try to minimize  $J(\Theta)$  as a function of parameters  $\Theta$

$$\frac{\partial}{\partial \Theta^{(l)}_{ik}} J(\Theta)$$

$J(\Theta)$  — non-convex.

