

# INTRODUÇÃO AO ESTUDO DAS REDES NEURAIS ARTIFICIAIS

Laboratório de Conexionismo e Ciências

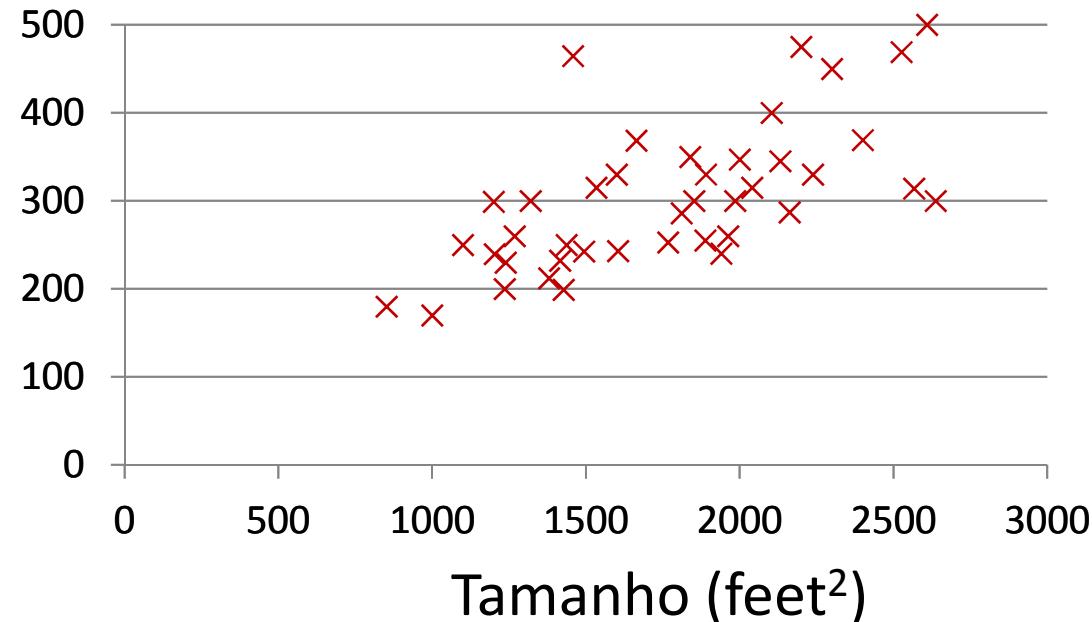
Cognitivas L3C  
Grupo SICRES  
INE - UFSC

## Regressão Linear com Uma Variável

## Introdução à Regressão Linear com Uma Variável – A questão do Modelo

### Preço dos Imóveis (Portland, OR)

Preço  
(em 1000s  
of dollars)



#### Aprendizado Supervisionado

Fornecer a “resposta correta” para cada exemplo visto nos dados.

#### Problema de Regressão

Predizer um valor contínuo de saída

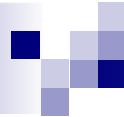
## Introdução à Regressão Linear com Uma Variável – A questão do Modelo

<b>Conjunto de treinamento</b>	<b>Tamanaho em feet<sup>2</sup> (x)</b>	<b>Preço (\$) em 1000's (y)</b>
Para o preço dos imóveis (Portland, OR)	2104	460
	1416	232
	1534	315
	852	178
Notation:	...	...

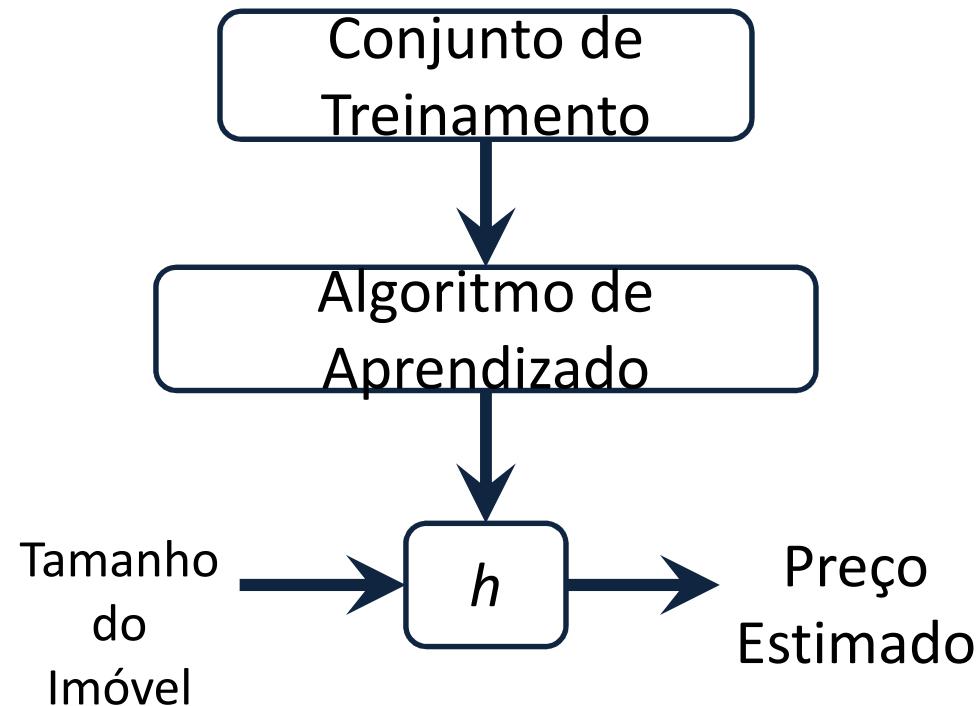
**m** = Number of training examples

**x**'s = “input” variable / features

**y**'s = “output” variable / “target” variable

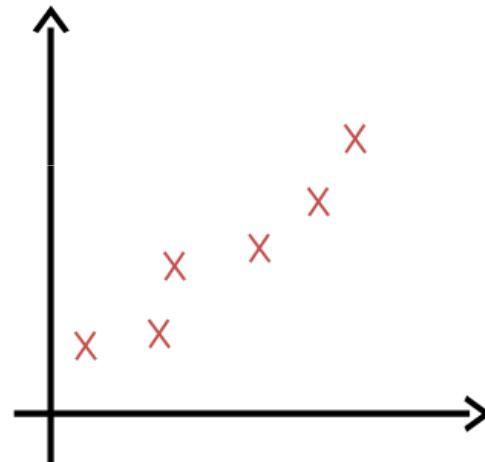


## Introdução à Regressão Linear com Uma Variável – a questão do Modelo



**Como representar a hipótese  $h$ ?**

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



Regressão Linear com uma variável.  
Regressão Linear Univariável.



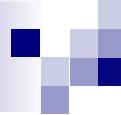
## Introdução à Regressão Linear com Uma Variável – a questão do Custo

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

Parâmetros:  $\theta_0, \theta_1$

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

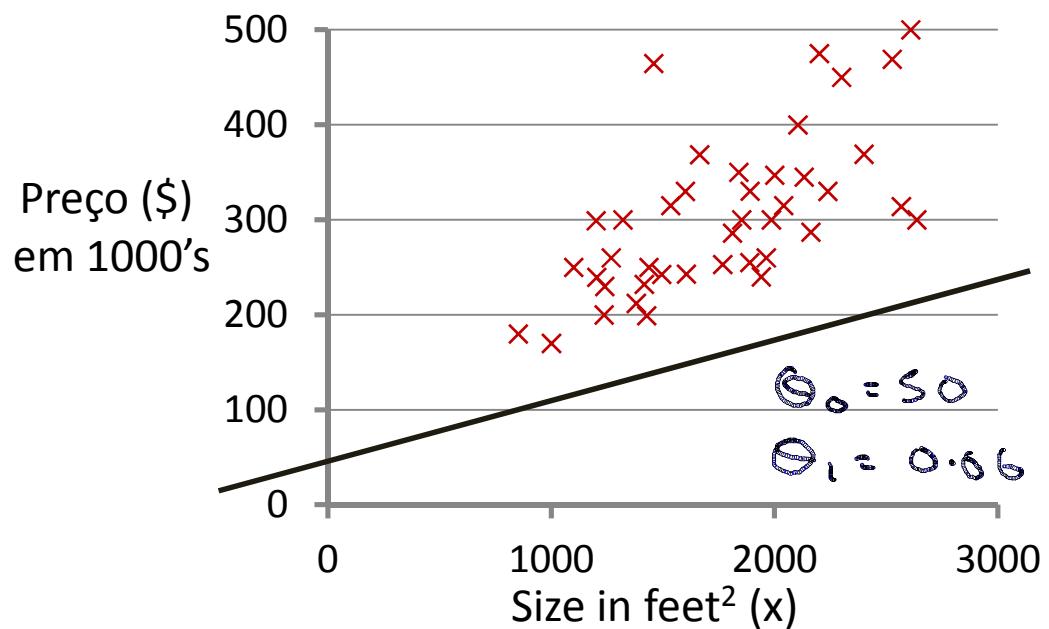
Objetivo:  $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$



## Introdução à Regressão Linear com Uma Variável – a questão do Custo

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

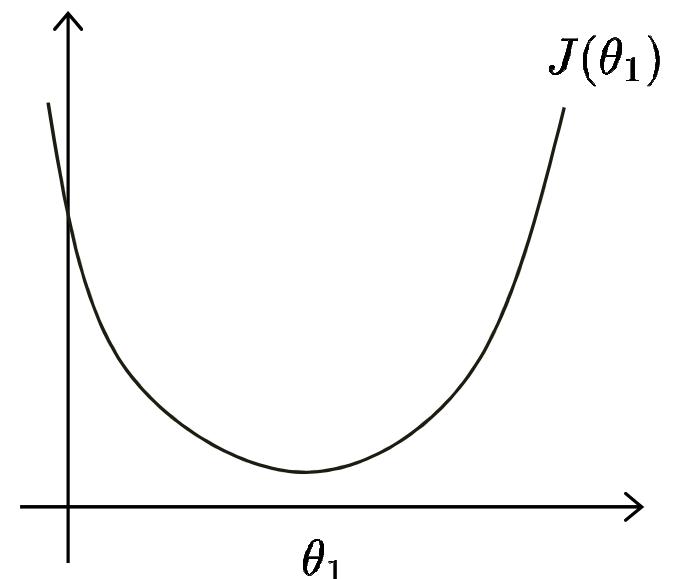
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )

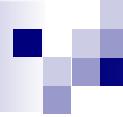


$$h_{\theta}(x) = 50 + 0.06x$$

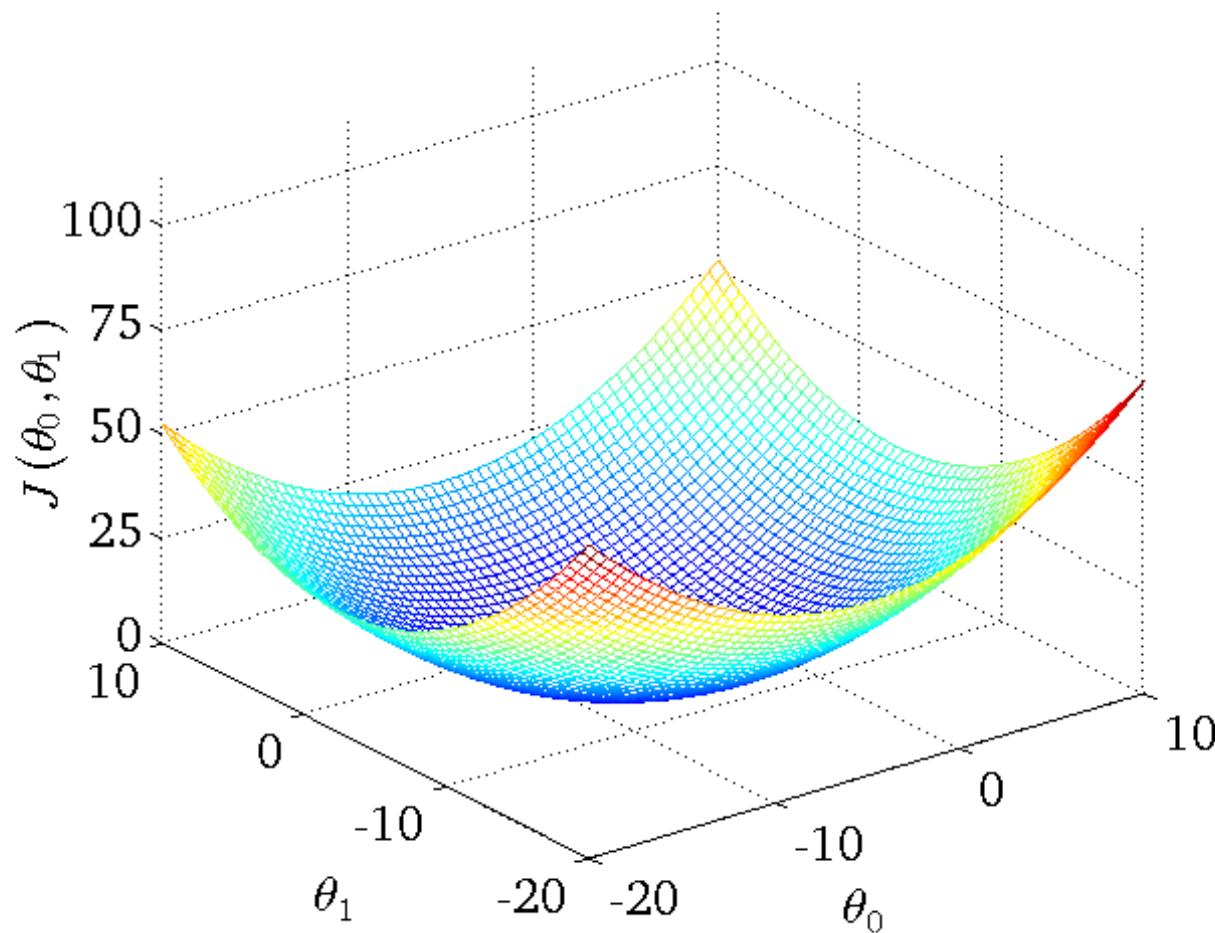
$$J(\theta_0, \theta_1)$$

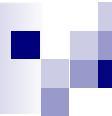
(function of the parameters  $\theta_0, \theta_1$ )





## Introdução à Regressão Linear com Uma Variável – a questão do Custo

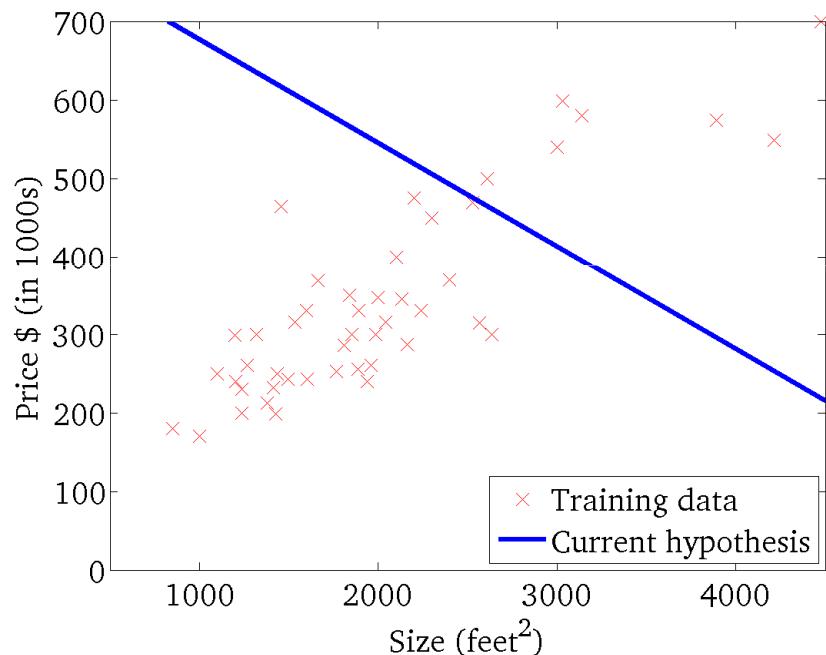




## Introdução à Regressão Linear com Uma Variável – a questão do Custo

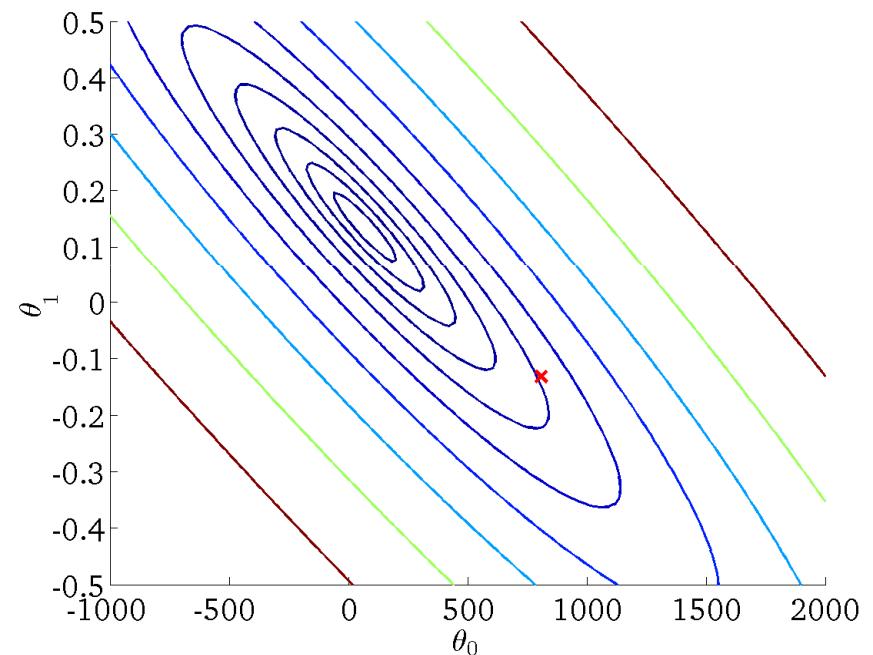
$$h_{\theta}(x)$$

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



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

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



## Introdução à Regressão Linear com Uma Variável – Minimizando a Função de Custo O MÉTODO DOS MÍNINOS QUADRÁTICOS

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \quad \leftarrow \text{Esta é a função que queremos minimizar}$$

$$h_\Theta(X) = X.\Theta \qquad \Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \quad X = \begin{bmatrix} 1 & x^1 \\ 1 & x^2 \\ 1 & x^2 \\ \dots & \dots \\ 1 & x^m \end{bmatrix}$$

$$e = Y - h_\Theta(X) = Y - X.\Theta$$

$$J(\Theta) = e^2$$

Na forma vetorial ela fica assim

$$e^2 = e^T \cdot e = (Y - X\Theta)^T \cdot (Y - X\Theta)$$

$$e^2 = Y^T Y - Y^T X\Theta - \Theta^T X^T Y - \Theta^T X^T X\Theta$$

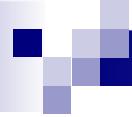
$$\frac{\partial J}{\partial \Theta} = -2X^T Y + 2X^T X\Theta = 0$$

$$-X^T Y + X^T X\Theta = 0$$

$$X^T X\Theta = X^T Y$$

$$\Theta = (X^T X)^{-1} \cdot X^T Y \quad \leftarrow$$

Esta é a fórmula analítica para calcular os valores de  $\Theta$  que minimizam a função de custo



## Introdução à Regressão Linear com Uma Variável – Minimizando a Função de Custo O MÉTODO DO GRADIENTE DESCENDENTE

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

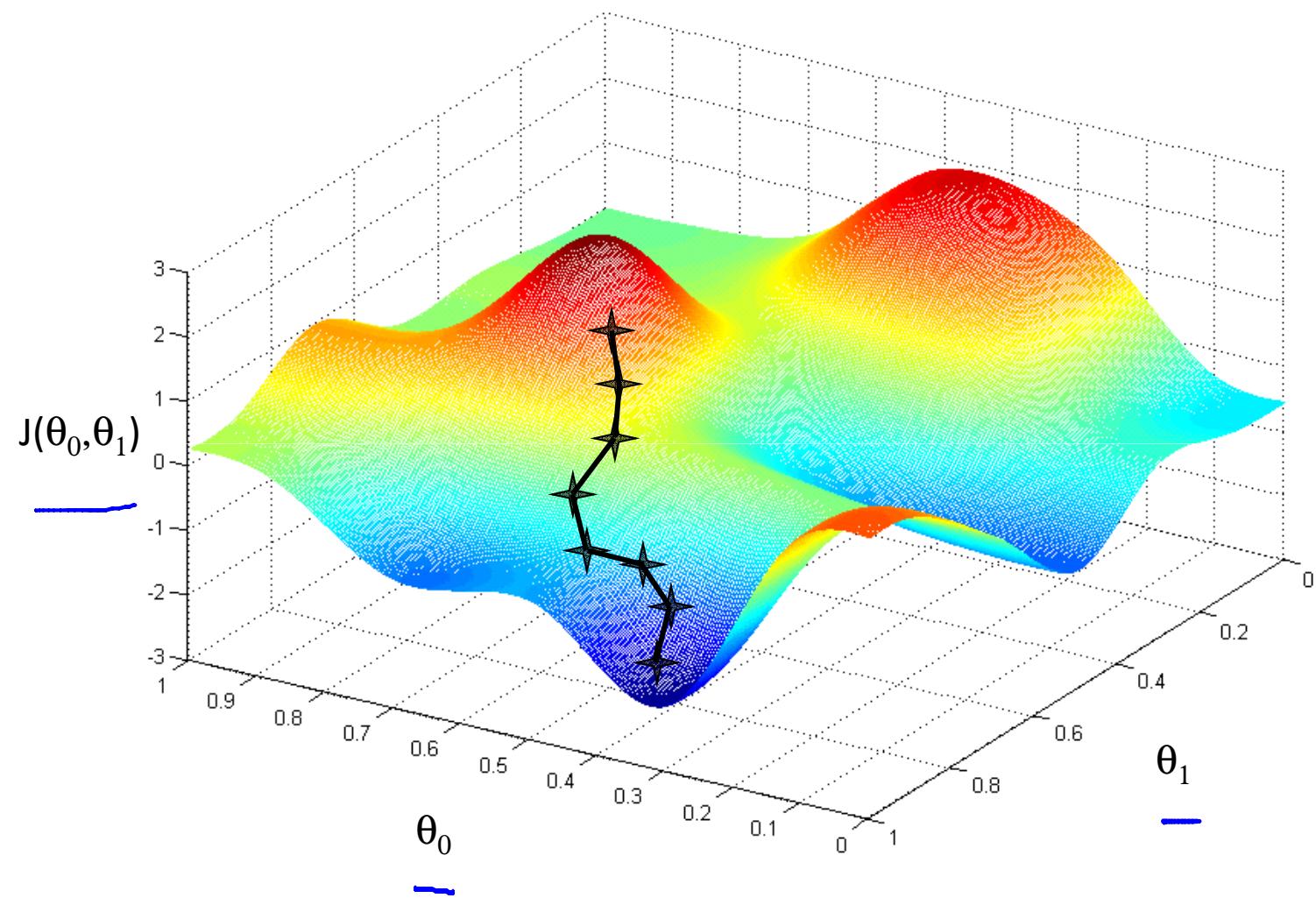
Esta é a função que queremos minimizar

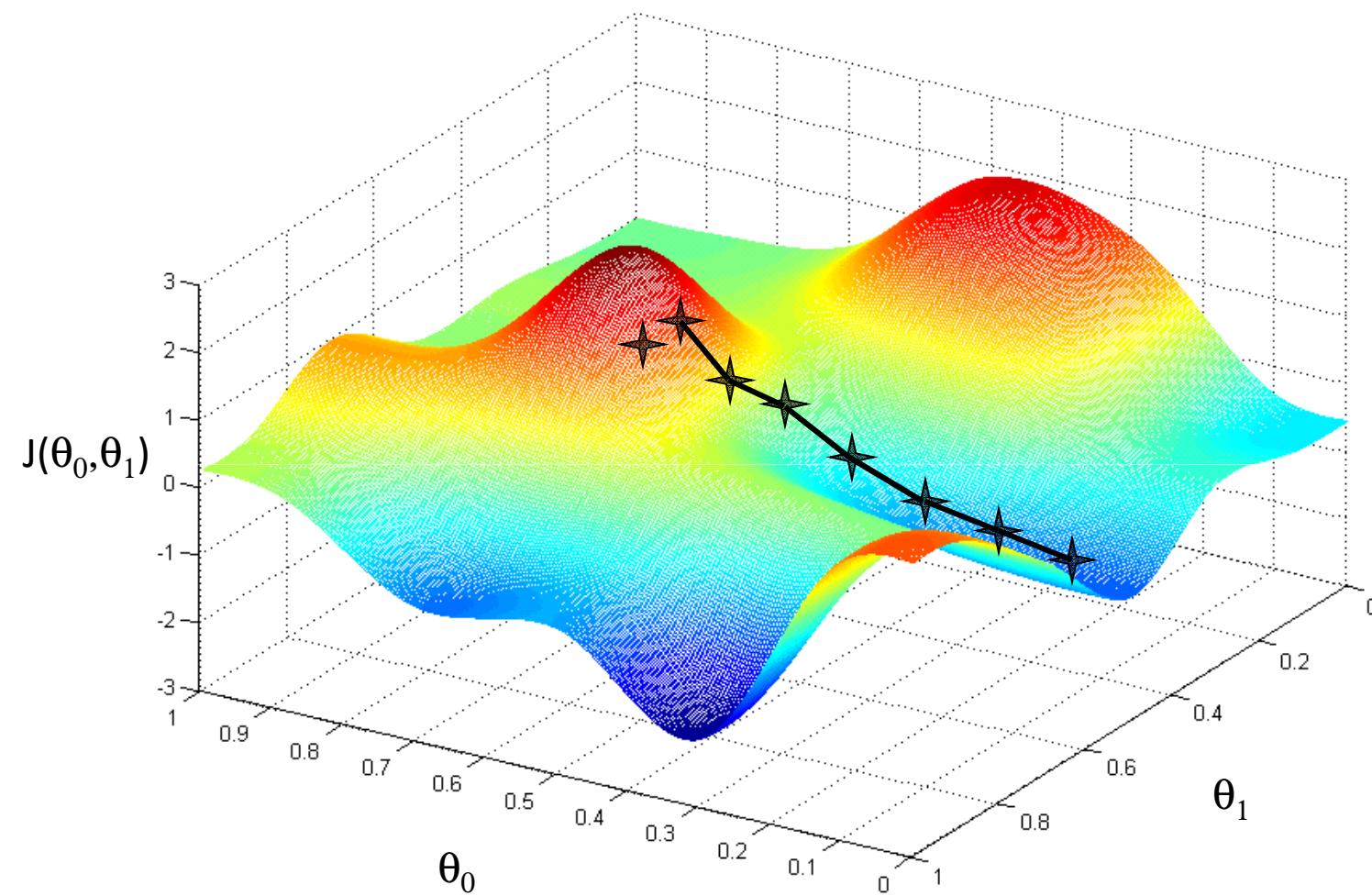
Dada uma função  $J(\theta_0, \theta_1)$

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

### Roteiro:

- Inicie com valores (aleatório ou não)  $\theta_0, \theta_1$
- Permaneça alterando  $\theta_0, \theta_1$  para reduzir  $J(\theta_0, \theta_1)$   
até que “esperançosamente” alcancemos  
um mínimo





# Gradient descent algorithm

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

---

Correct: Simultaneous update

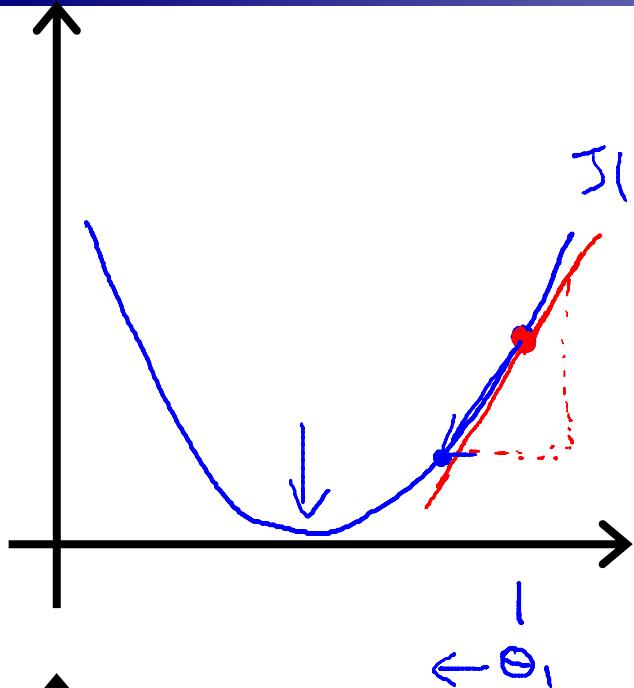
```
temp0 :=  $\theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$   
temp1 :=  $\theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$   
 $\theta_0 := \text{temp0}$   
 $\theta_1 := \text{temp1}$ 
```

Incorrect:

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

# Gradient descent algorithm

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



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

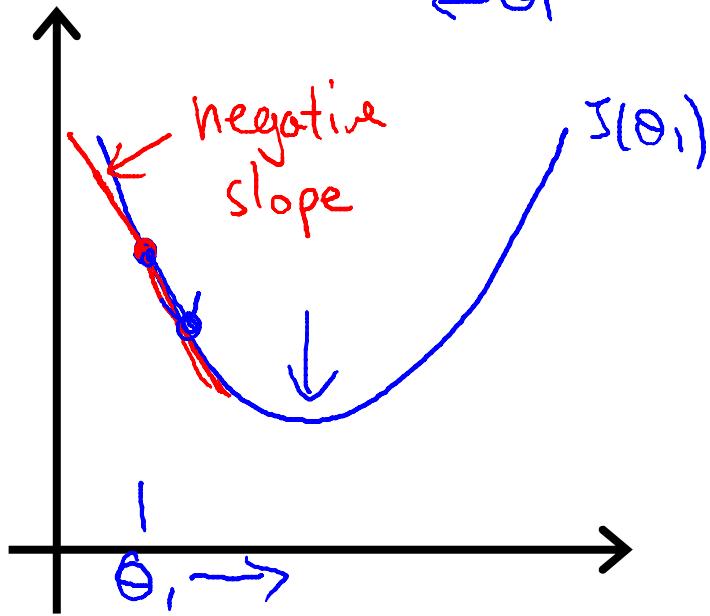
$$\theta_1 := \theta_1 - \frac{\alpha}{\frac{d}{d\theta_1} J(\theta_1)} \geq 0$$

Diagram illustrating the update rule for gradient descent:

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

The term  $\frac{d}{d\theta_1} J(\theta_1)$  is highlighted in a pink box. A red circle highlights the step size  $\alpha$ .

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



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

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

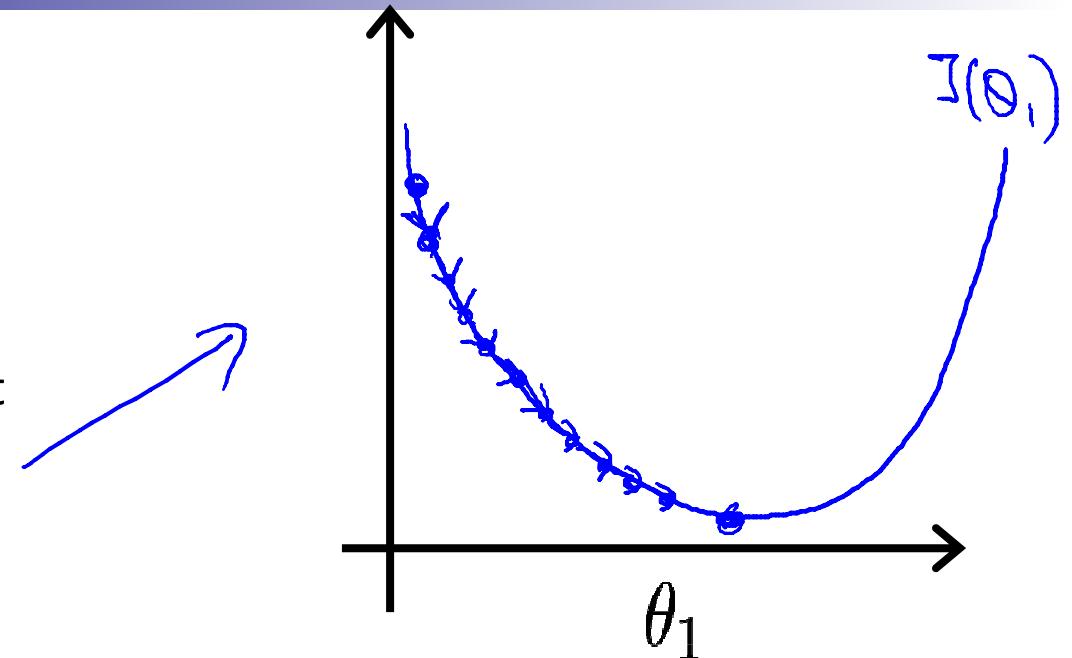
Diagram illustrating the update rule for gradient ascent:

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

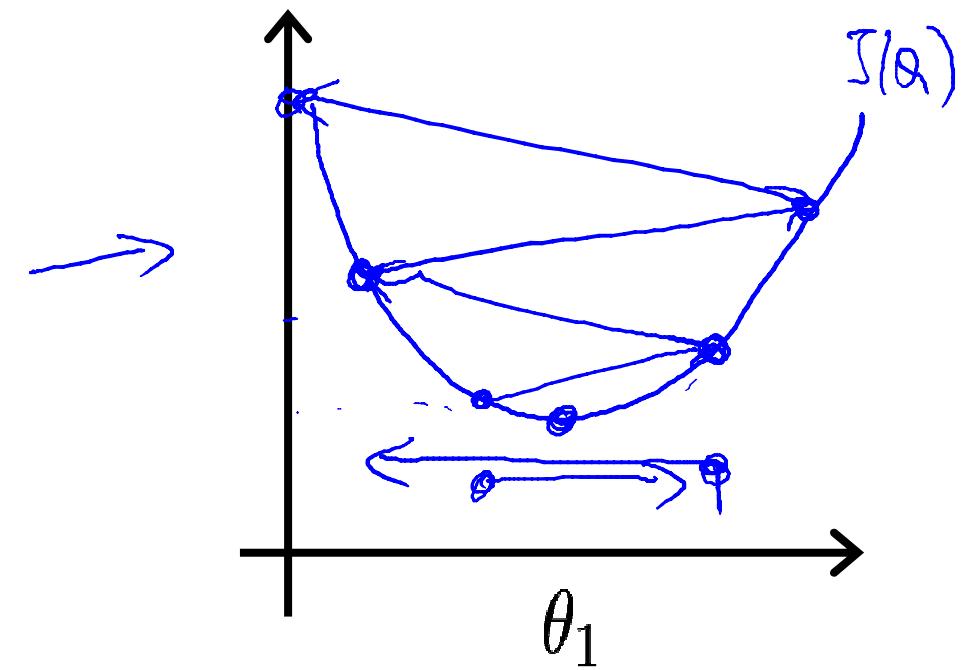
The term  $\frac{d}{d\theta_1} J(\theta_1)$  is highlighted in a pink box. A red circle highlights the step size  $\alpha$ .

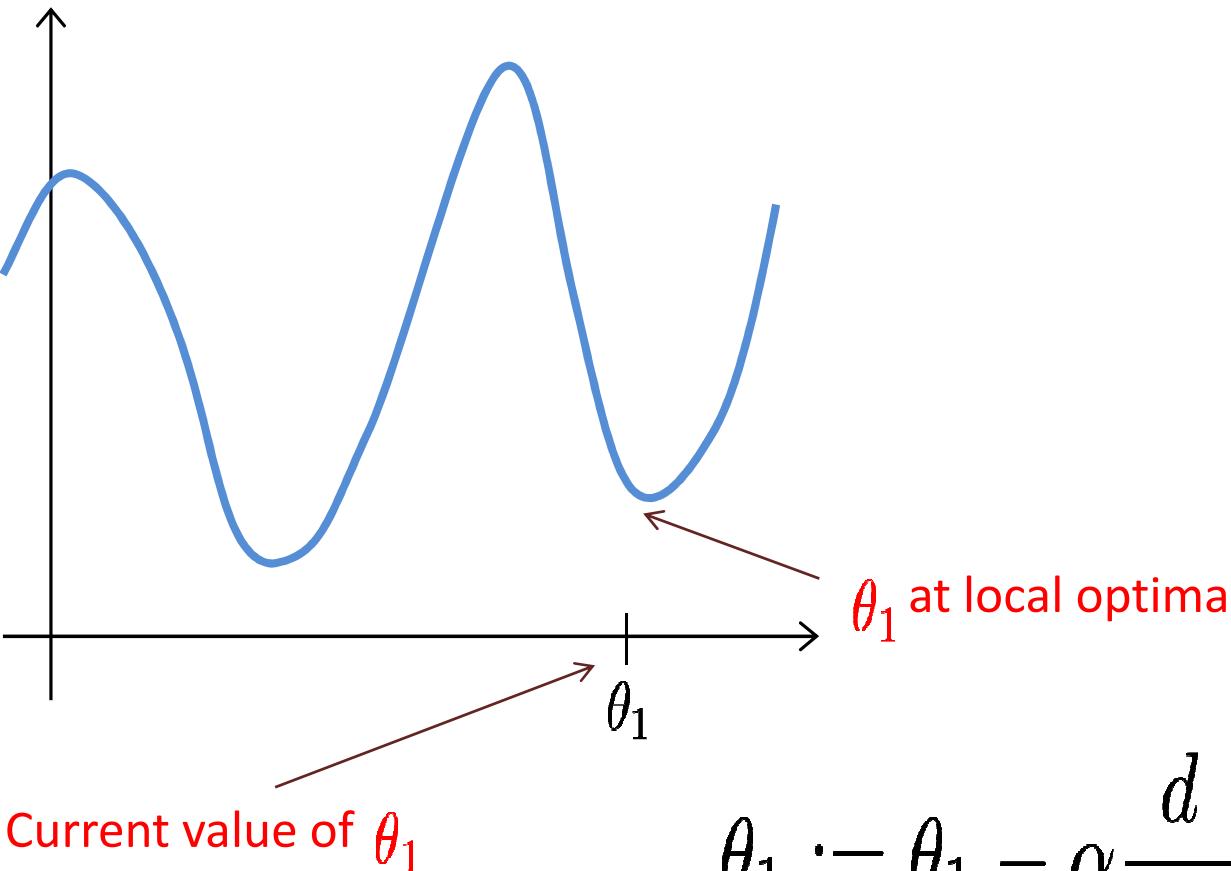
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

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



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



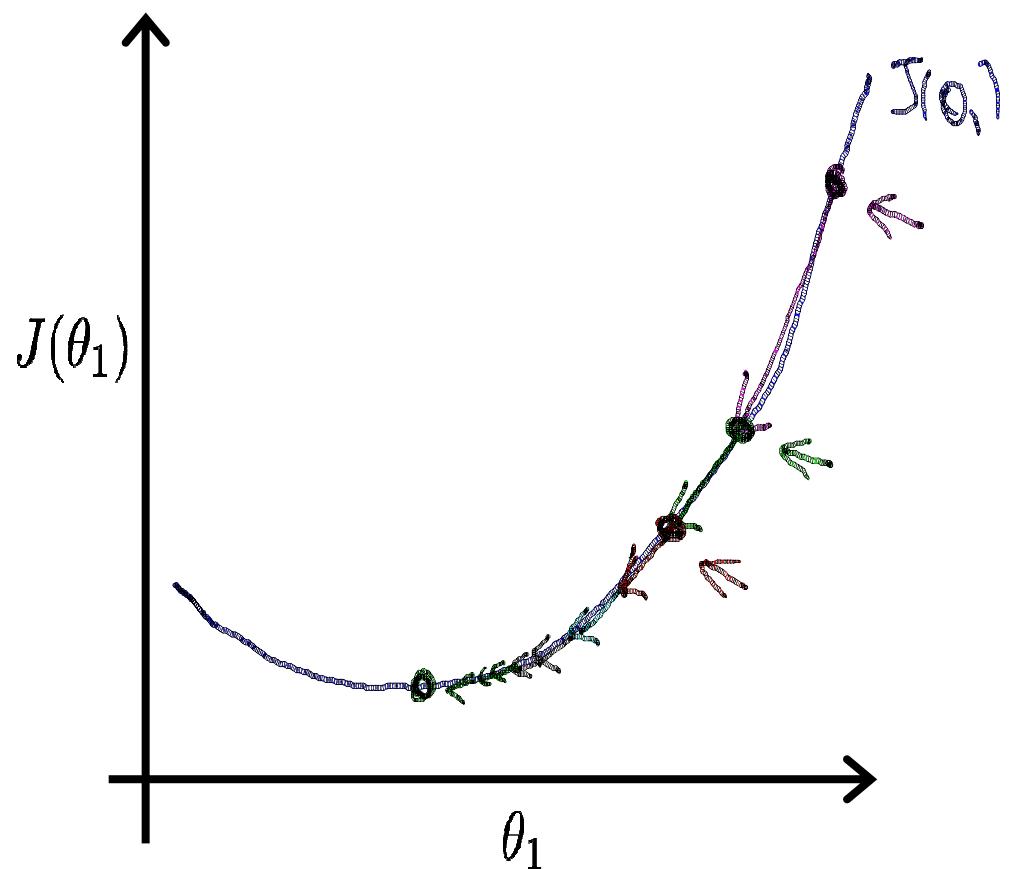


$$\theta_1 := \theta_1 - \alpha \frac{d}{d\theta_1} J(\theta_1)$$

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

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

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



## Gradient descent algorithm

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

## Linear Regression Model

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

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

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

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

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

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

## Gradient descent algorithm

repeat until convergence {

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

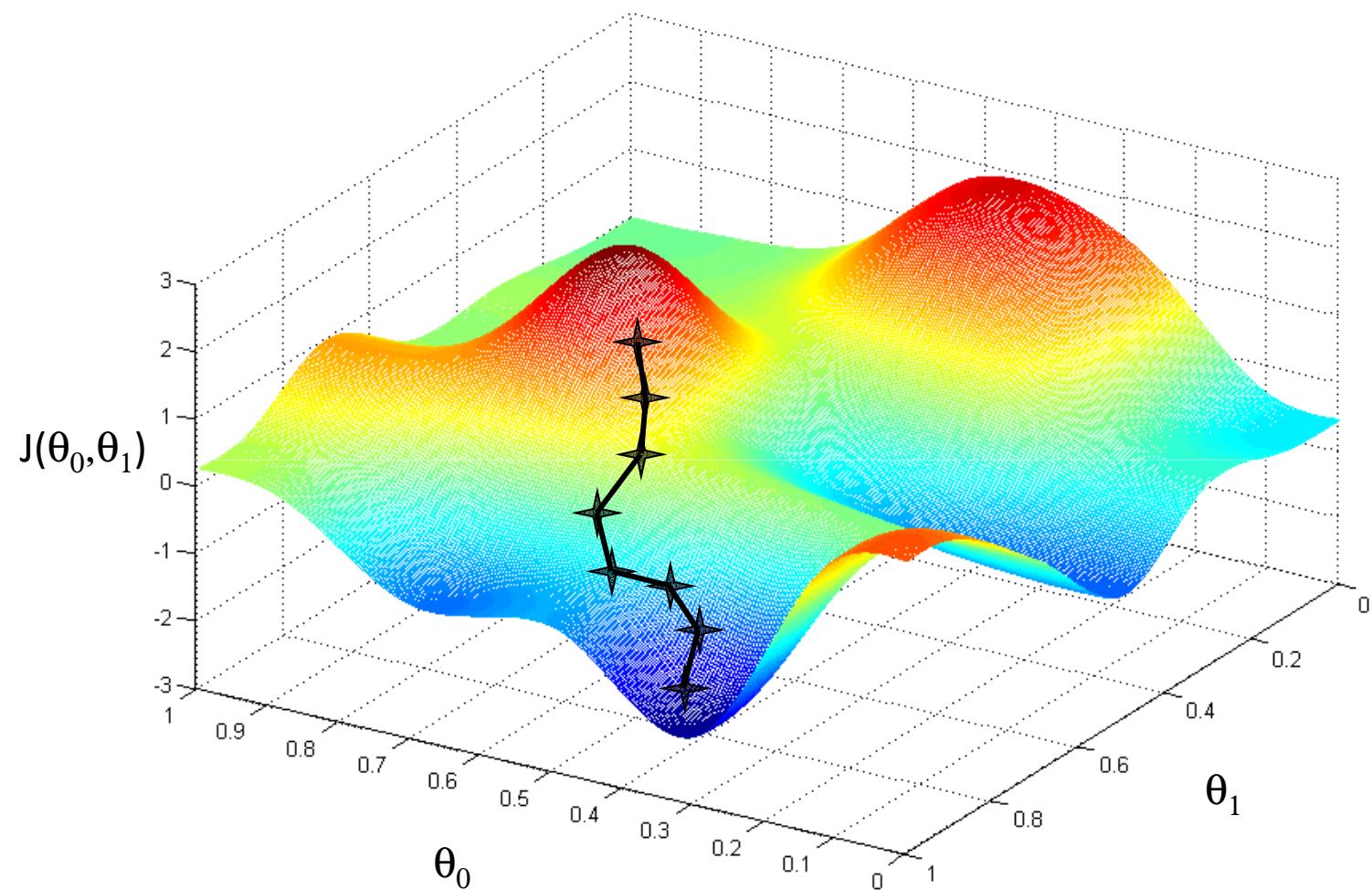
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

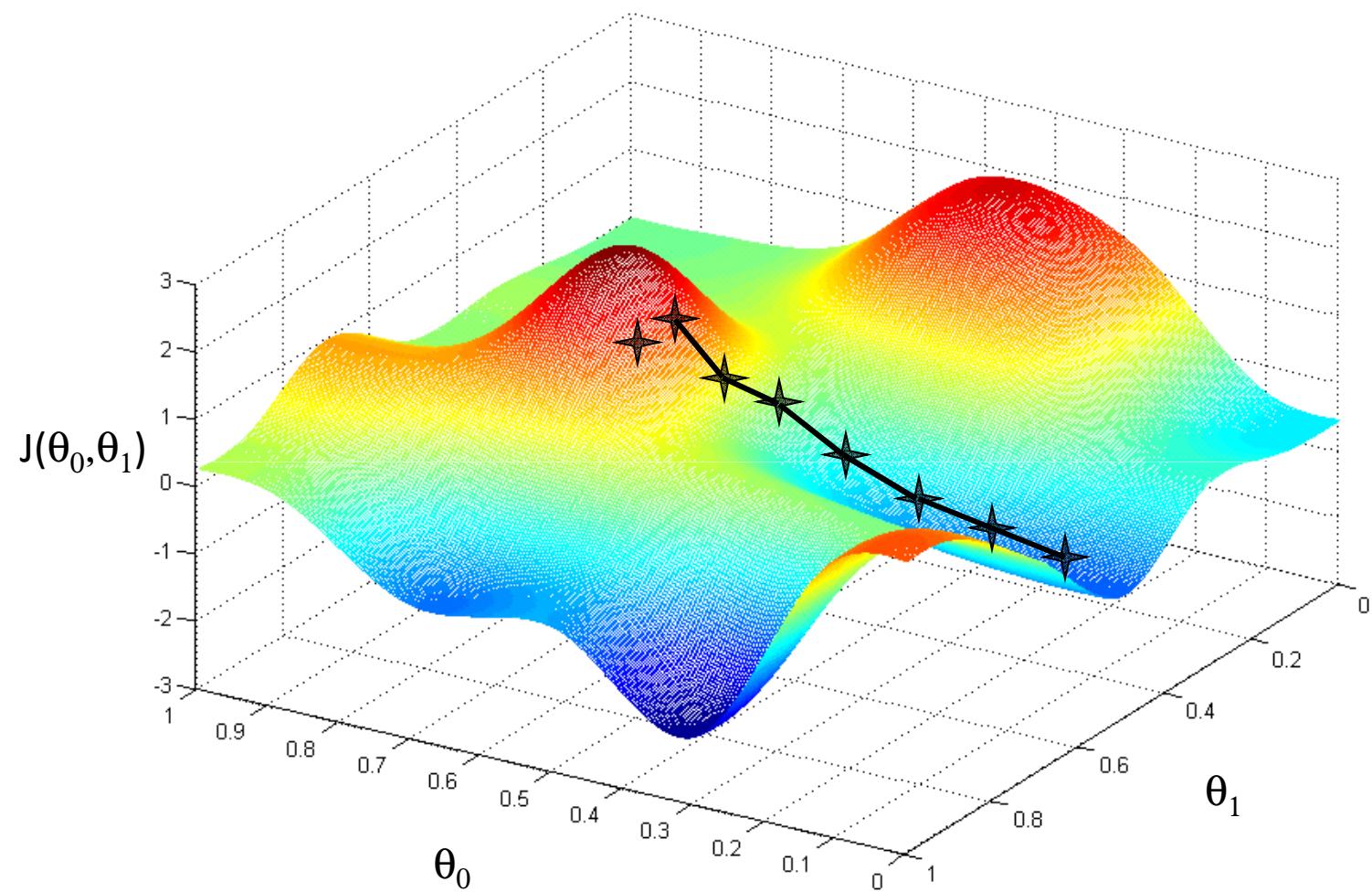
}

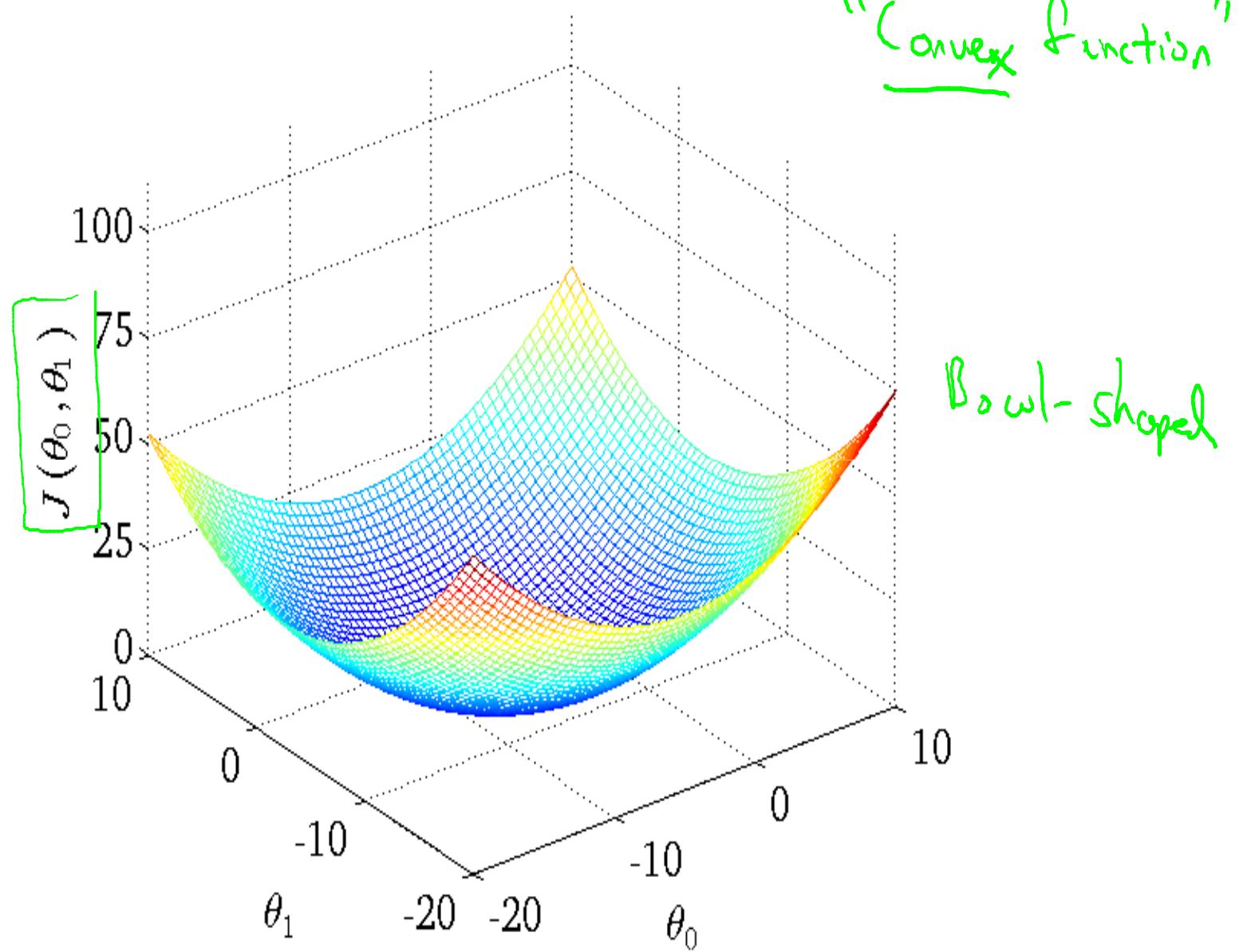
$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$$

update  
 $\theta_0$  and  $\theta_1$   
simultaneously

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

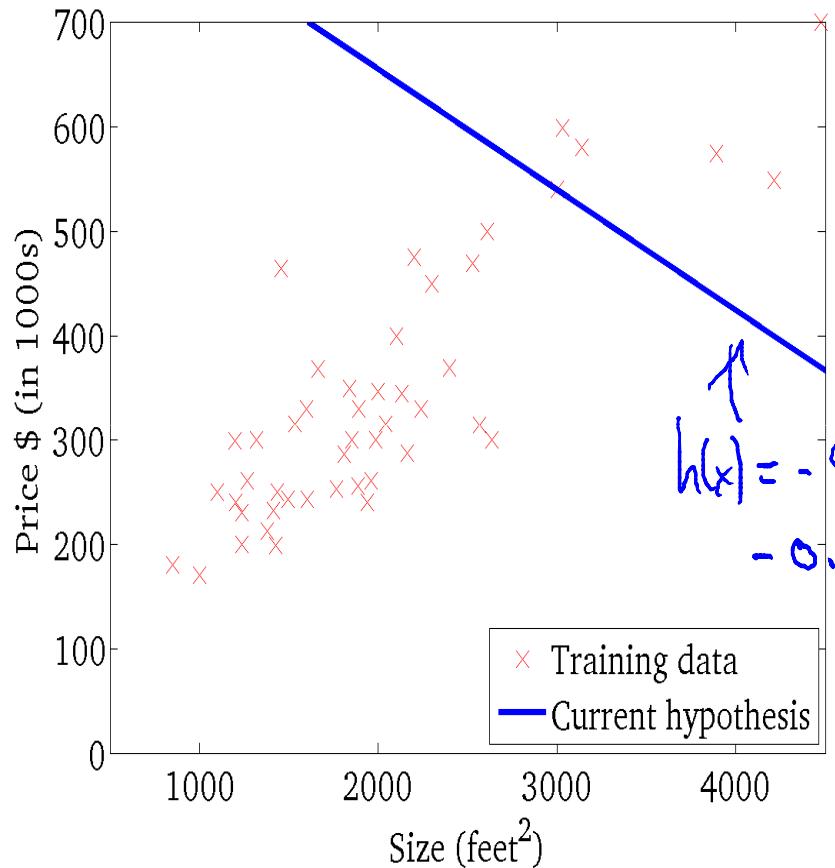






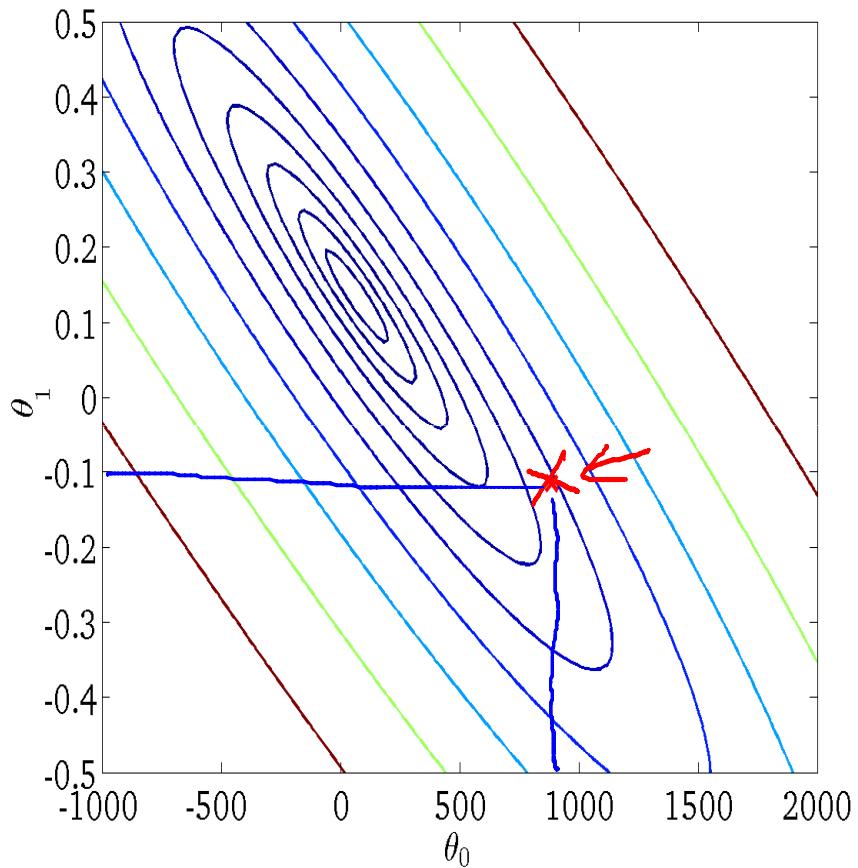
$$h_{\theta}(x)$$

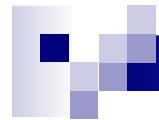
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

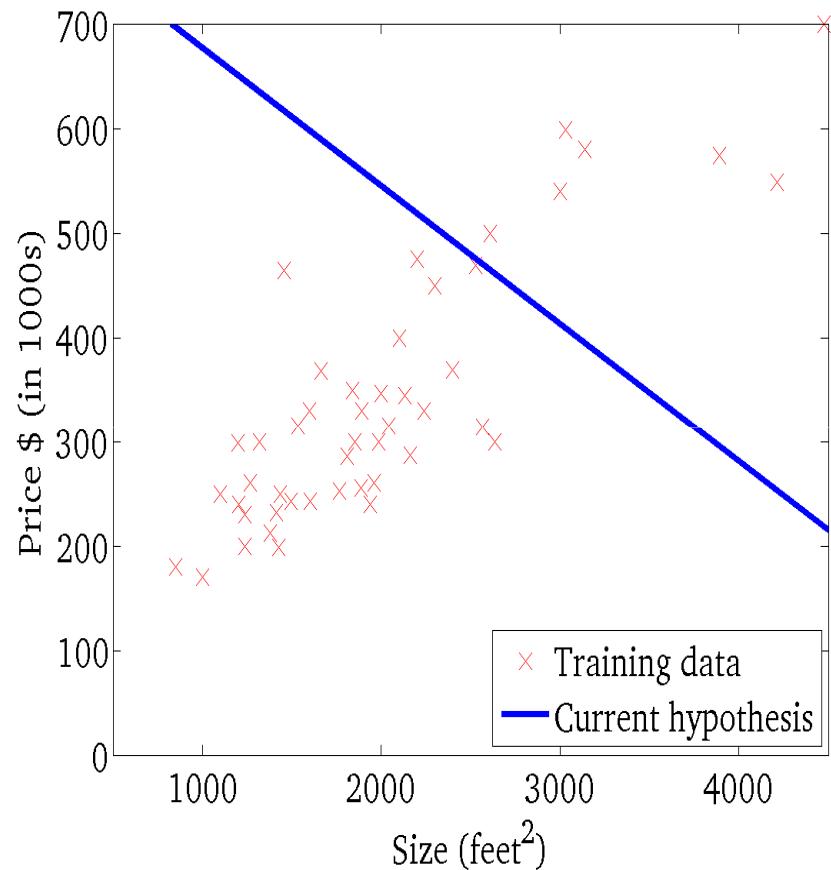
(function of the parameters  $\theta_0, \theta_1$ )





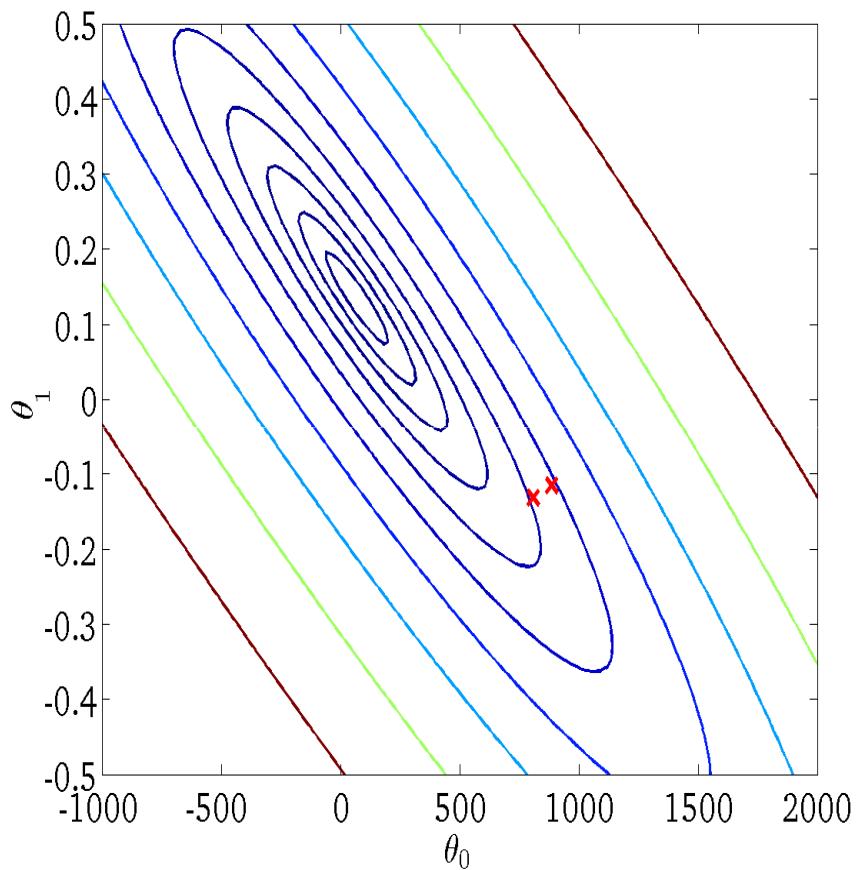
$$h_{\theta}(x)$$

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



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

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

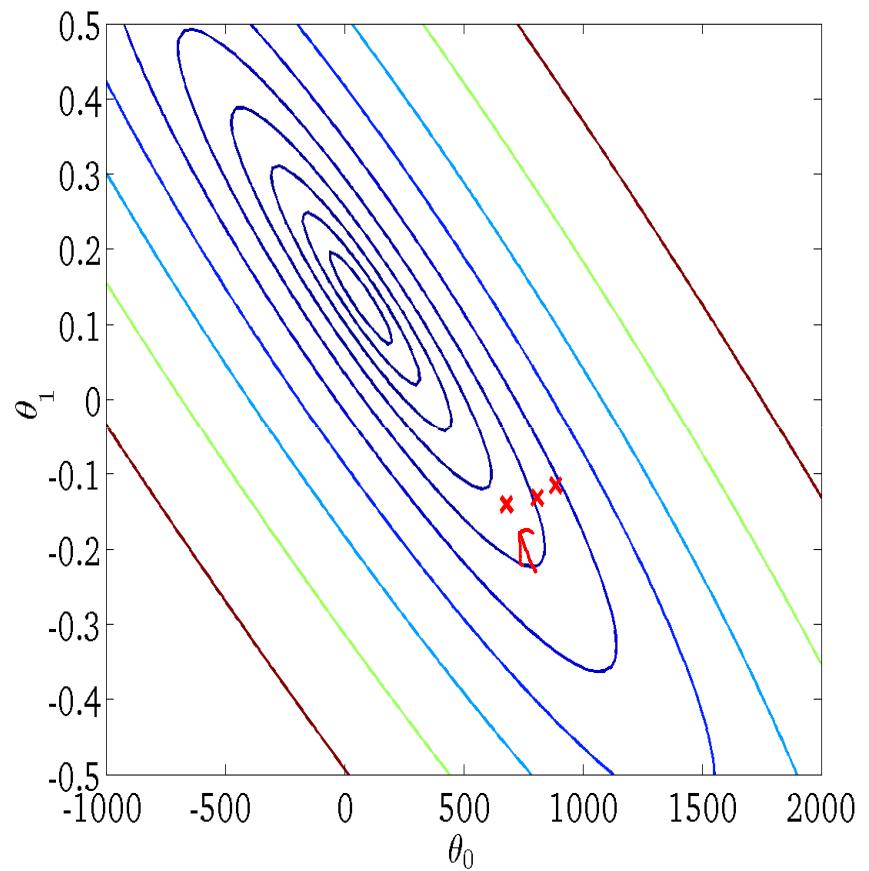
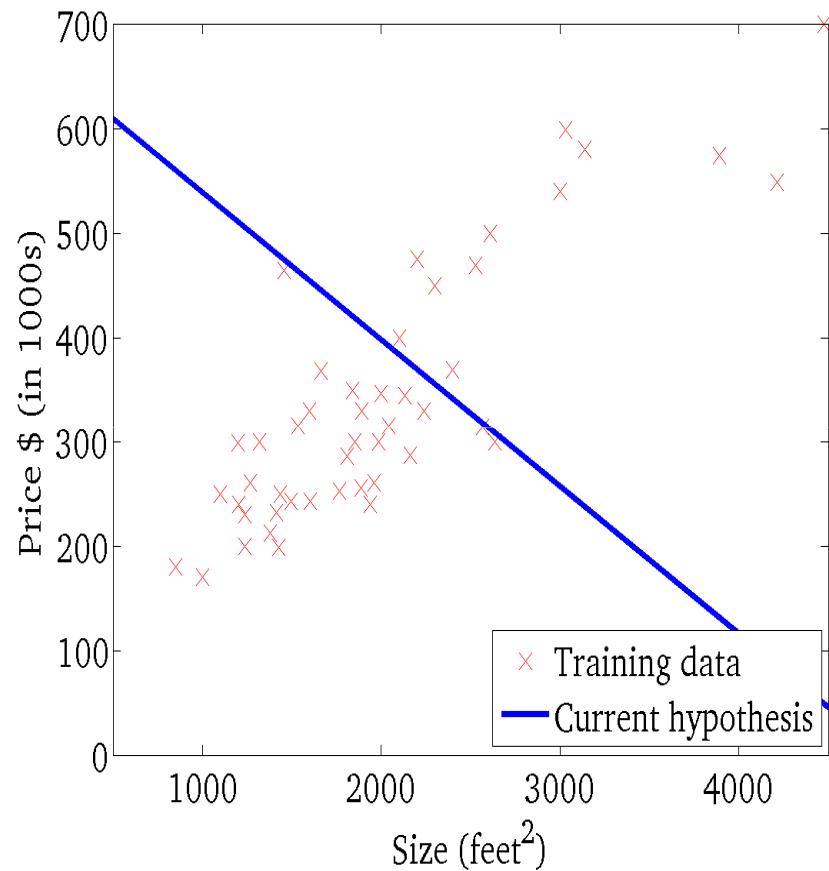


$$h_{\theta}(x)$$

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

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

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

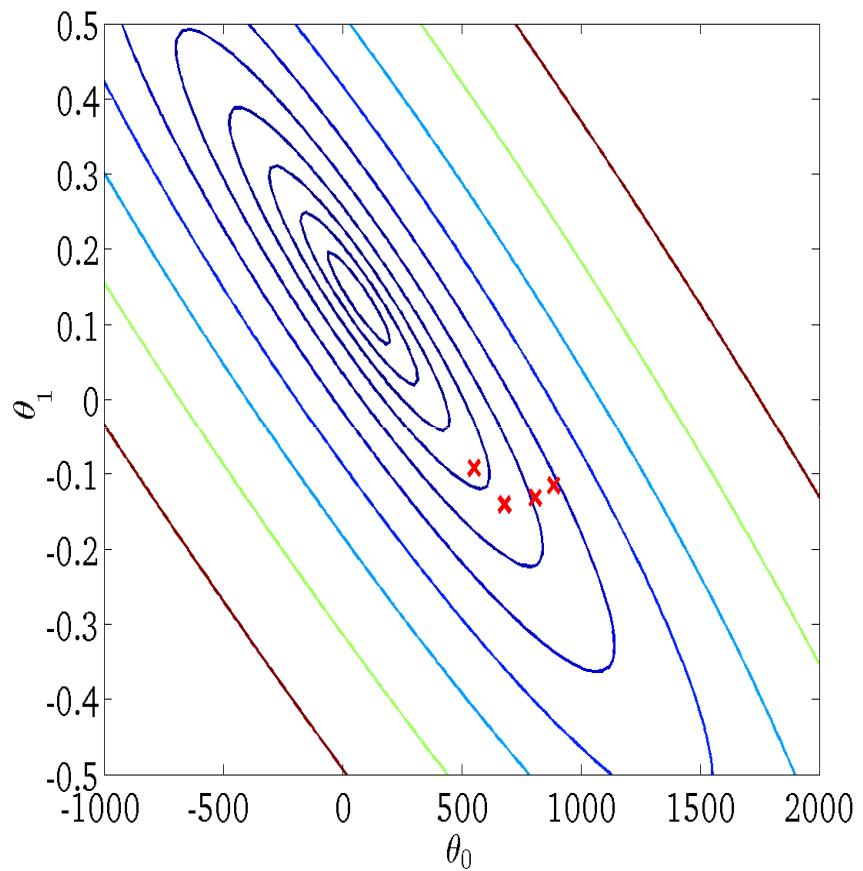
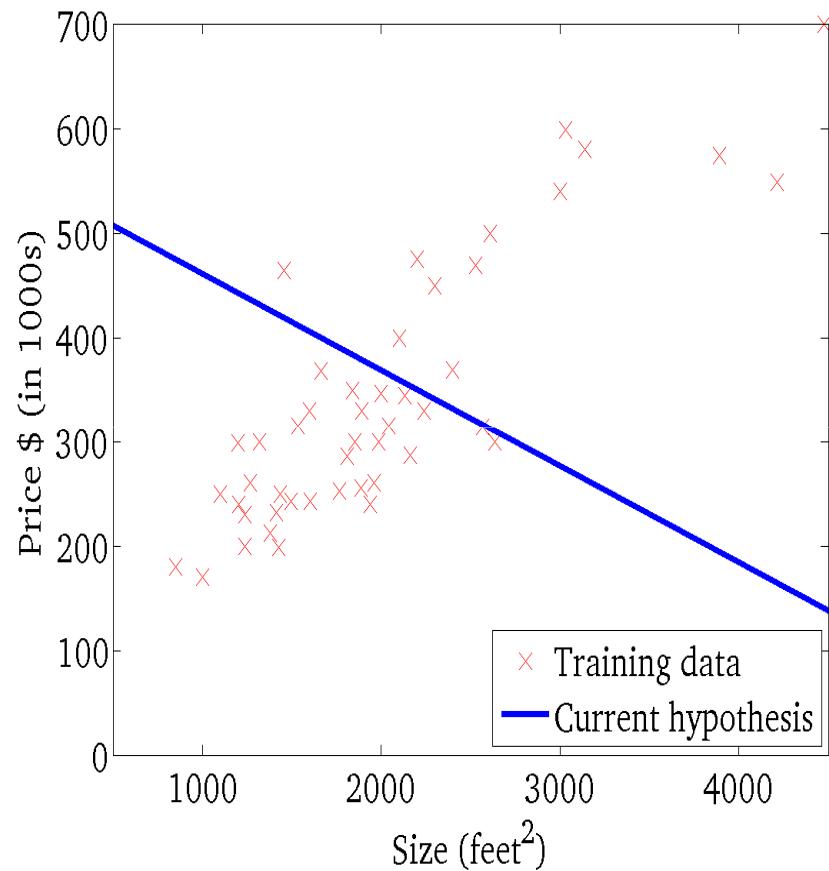


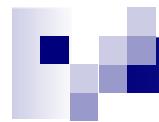
$$h_{\theta}(x)$$

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

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

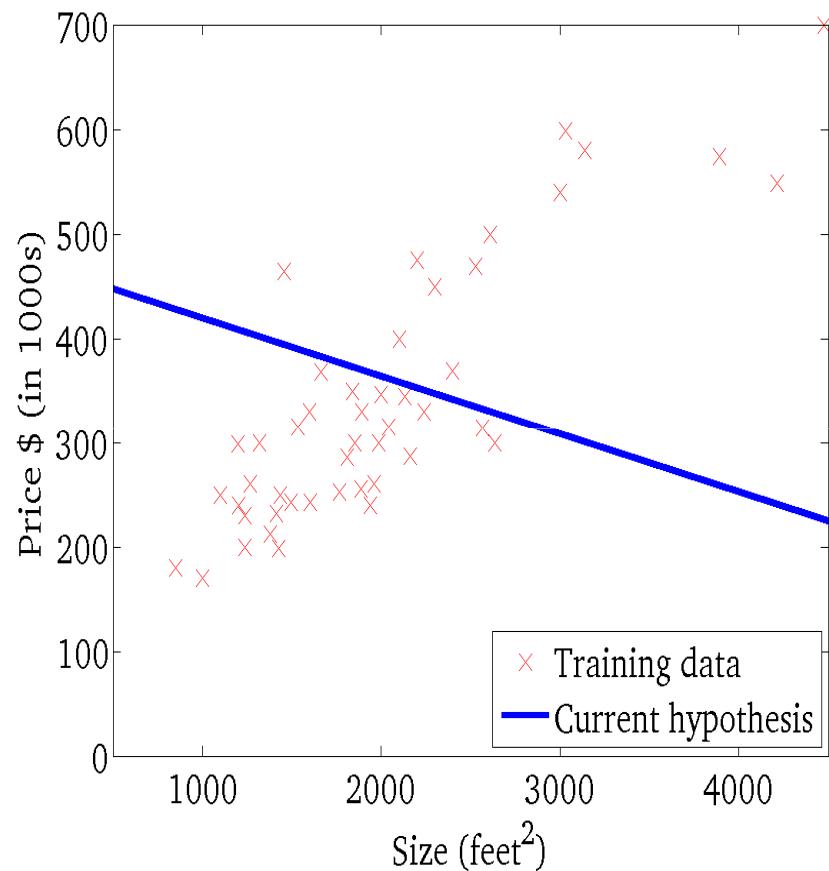
(function of the parameters  $\theta_0, \theta_1$ )





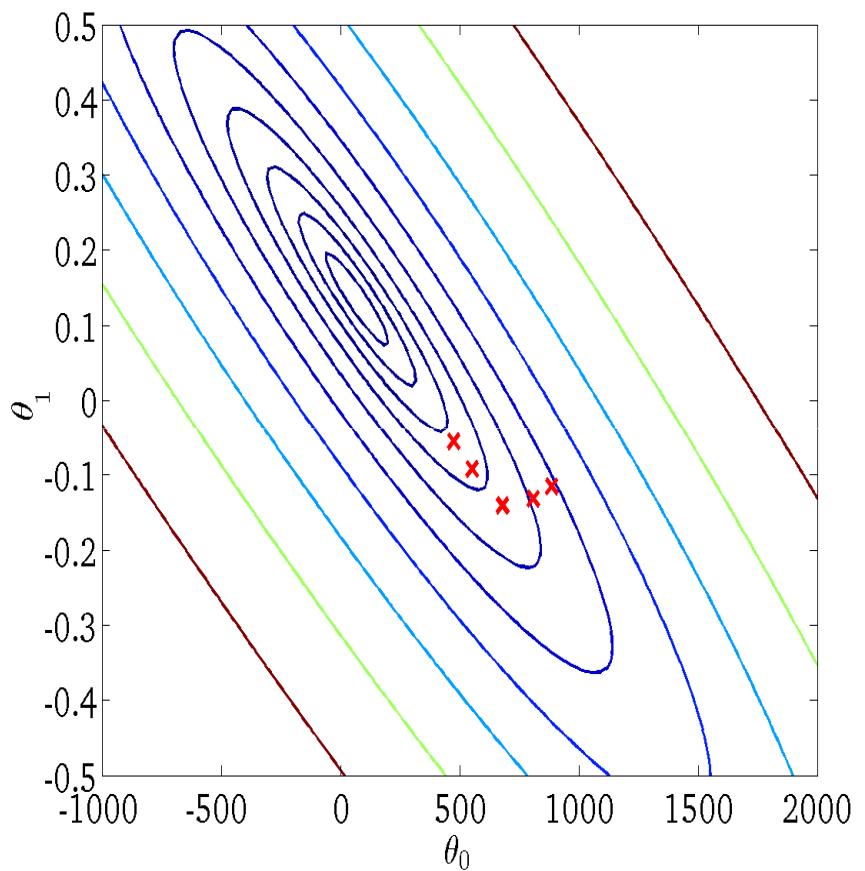
$$h_{\theta}(x)$$

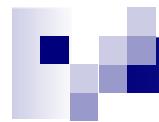
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

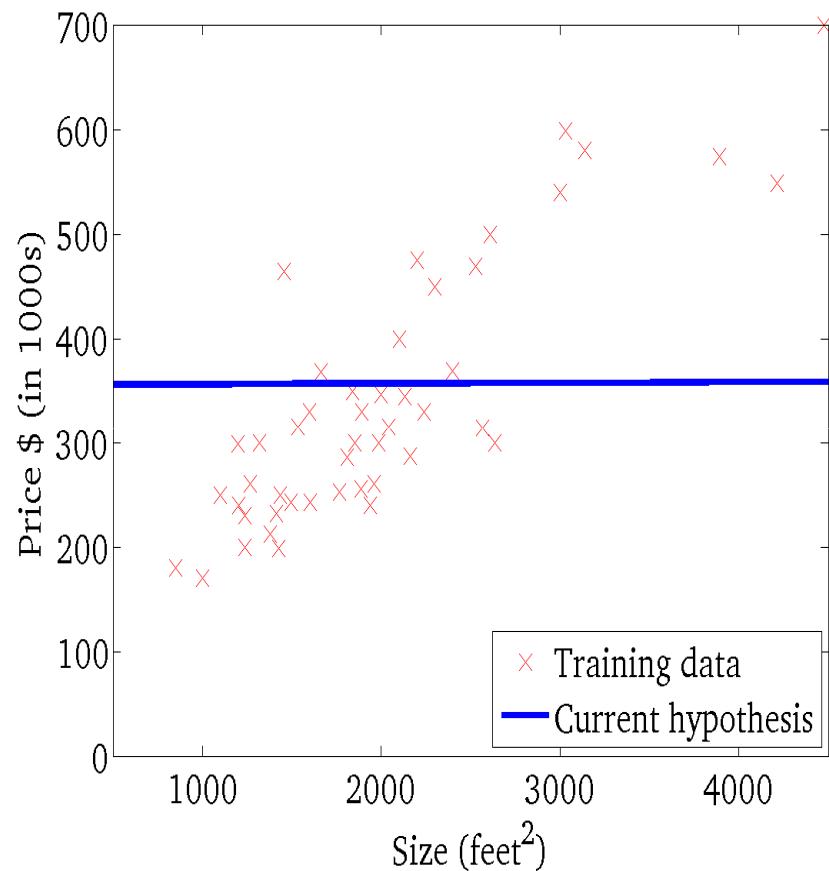
(function of the parameters  $\theta_0, \theta_1$ )





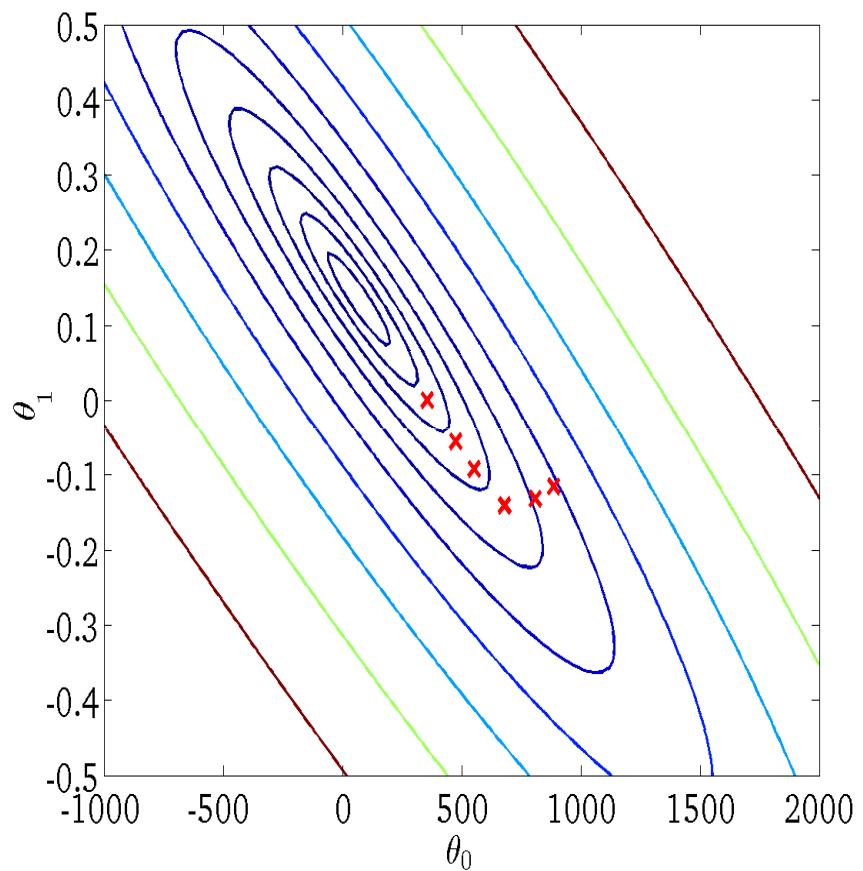
$$h_{\theta}(x)$$

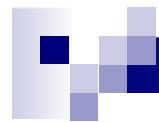
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

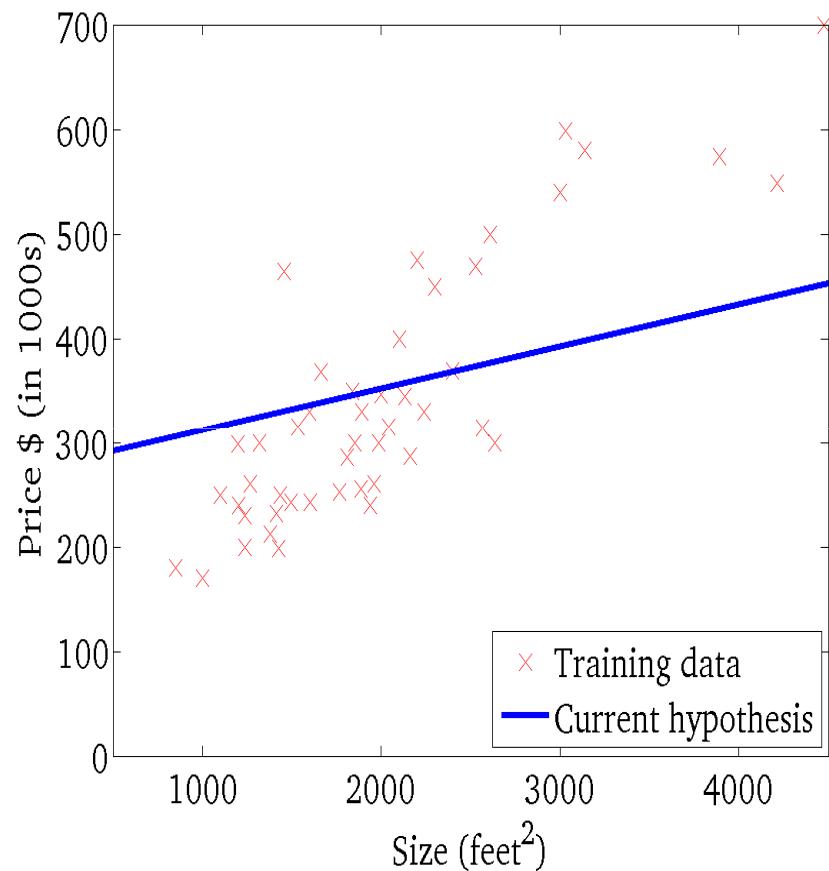
(function of the parameters  $\theta_0, \theta_1$ )





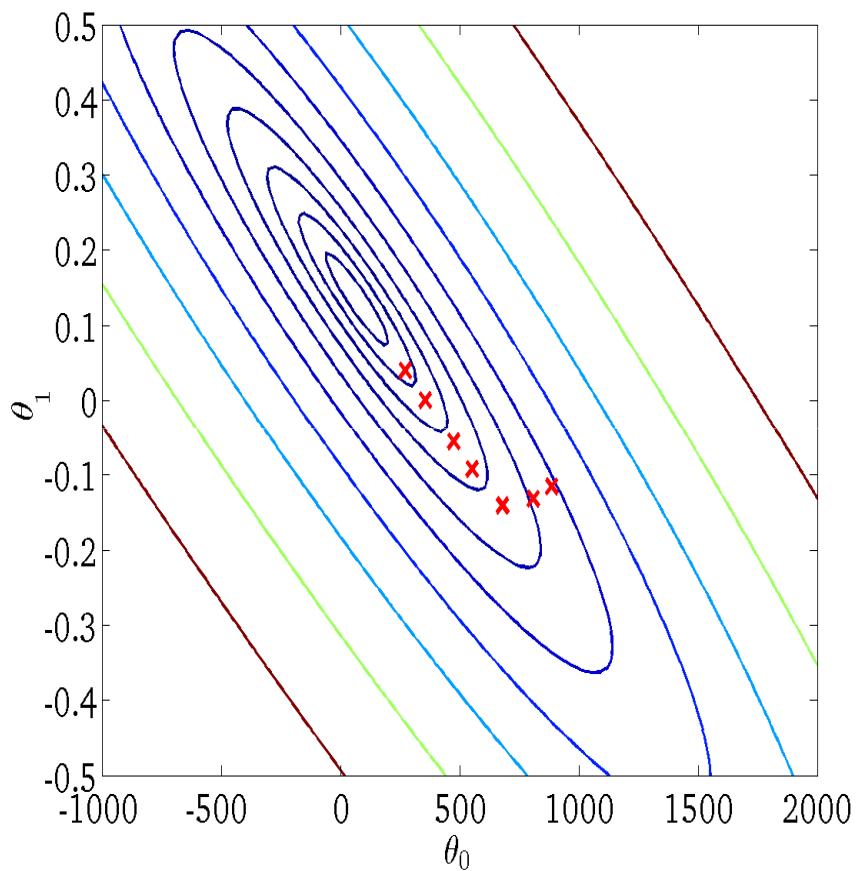
$$h_{\theta}(x)$$

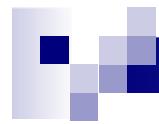
(for fixed  $\theta_0, \theta_1$ , this is a function of  $x$ )



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

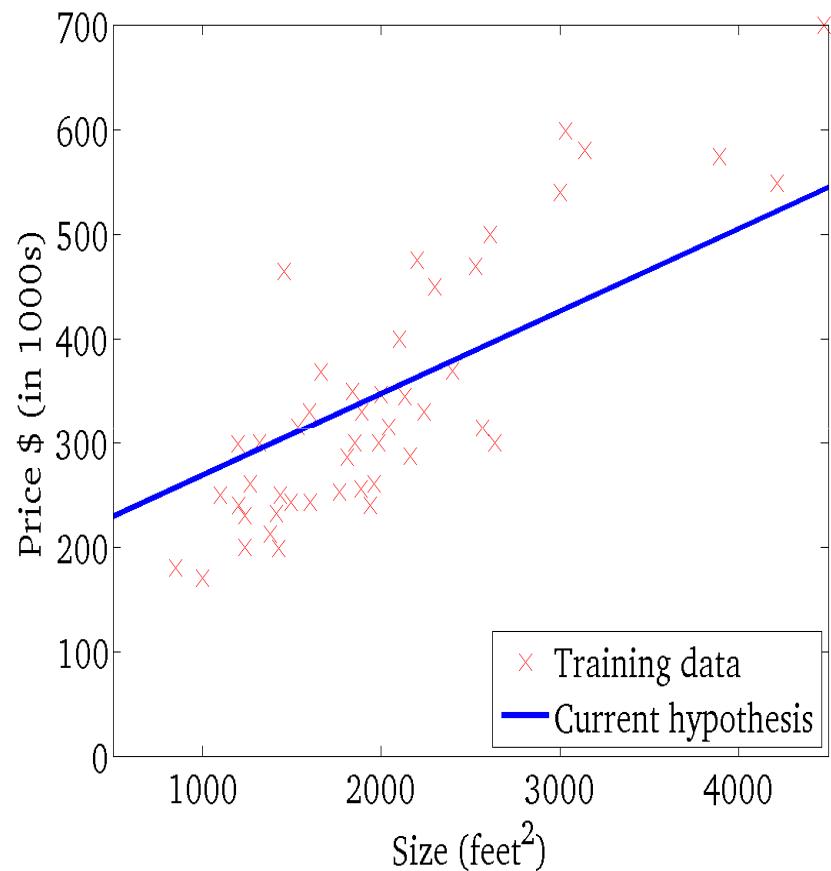
(function of the parameters  $\theta_0, \theta_1$ )





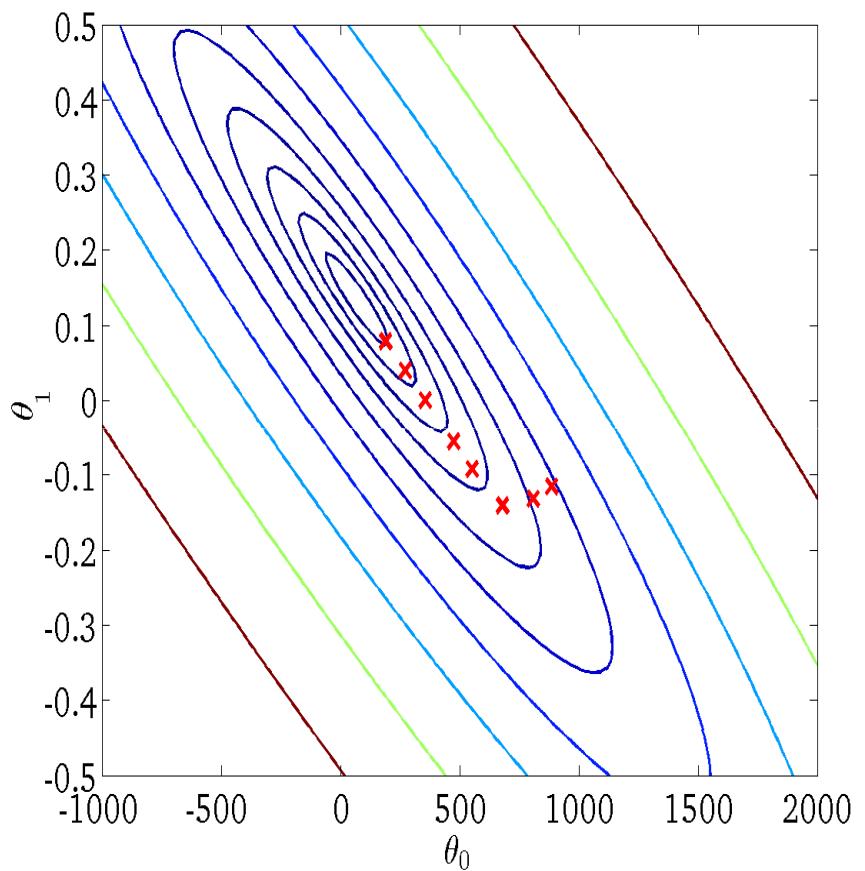
$$h_{\theta}(x)$$

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



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

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

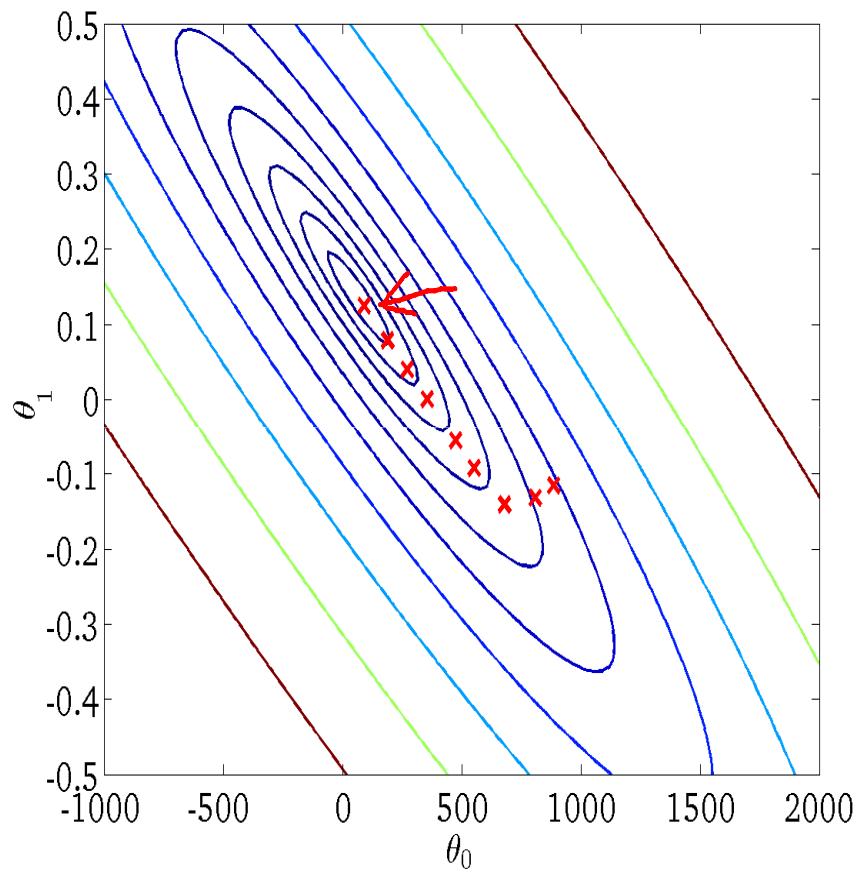
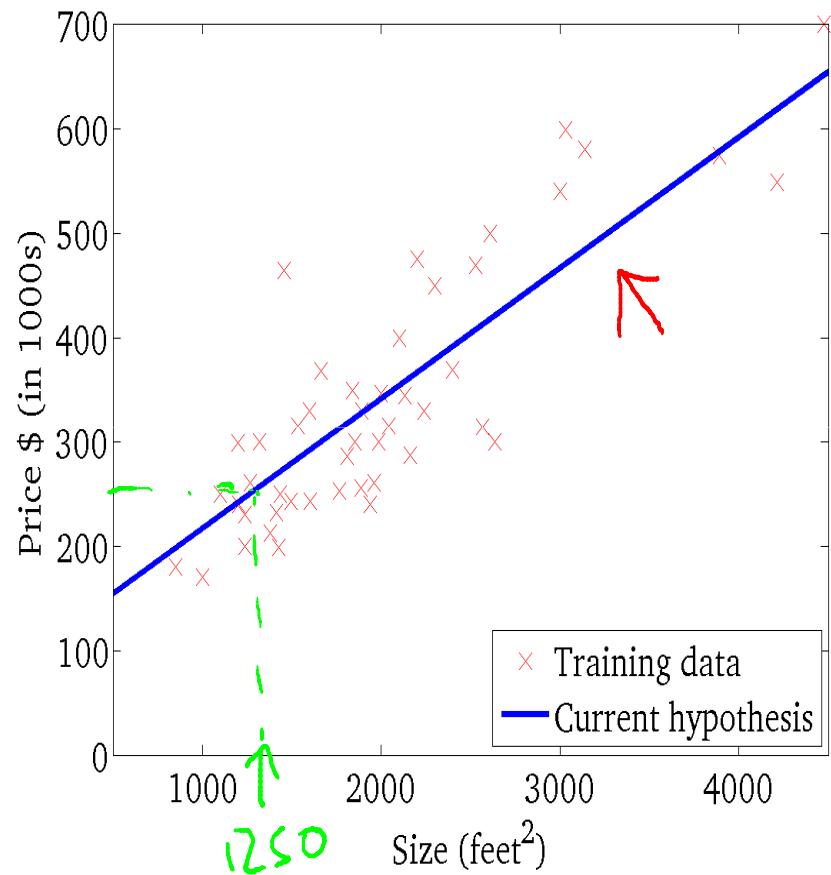


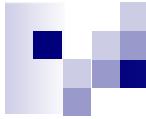
$$h_{\theta}(x)$$

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

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

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





## Introdução à Regressão Linear Múltipla

**Conjunto de  
treinamento  
Para o preço dos  
imóveis  
(Portland, OR)**

**Multiplas Características (variáveis).**

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
...	...	...	...	...

Hipótese:

Anterior (para uma variável):  $h_{\theta}(x) = \theta_0 + \theta_1 x$

→ Agora:  $\underline{h_{\theta}(x)} = \underline{\theta_0} + \underline{\theta_1 x_1} + \underline{\theta_2 x_2} + \cdots + \underline{\theta_n x_n}$

Por conveniência de notação, definimos  $x_0 = 1$  ( $x_0^{(i)} = 1$ )

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

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

$$\begin{aligned} h_{\theta}(x) &= \underline{\theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n} \\ &= \boxed{\Theta^T x} \end{aligned}$$

$$\underbrace{[\theta_0, \theta_1, \dots, \theta_n]}_{\Theta^T} \quad (n+1) \times 1 \text{ matrix}$$

$x$

Regressão Linear Multivariada. ←

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

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

Função de Custo:

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

Gradiente descendente:

Repeat {

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

}

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

## Gradiente Descendente

Anterior (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  $\theta_0, \theta_1$ )

}

Novo algoritmo: ( $n \geq 1$ )

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  
 $j = 0, \dots, n$ )

}

---

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

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

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

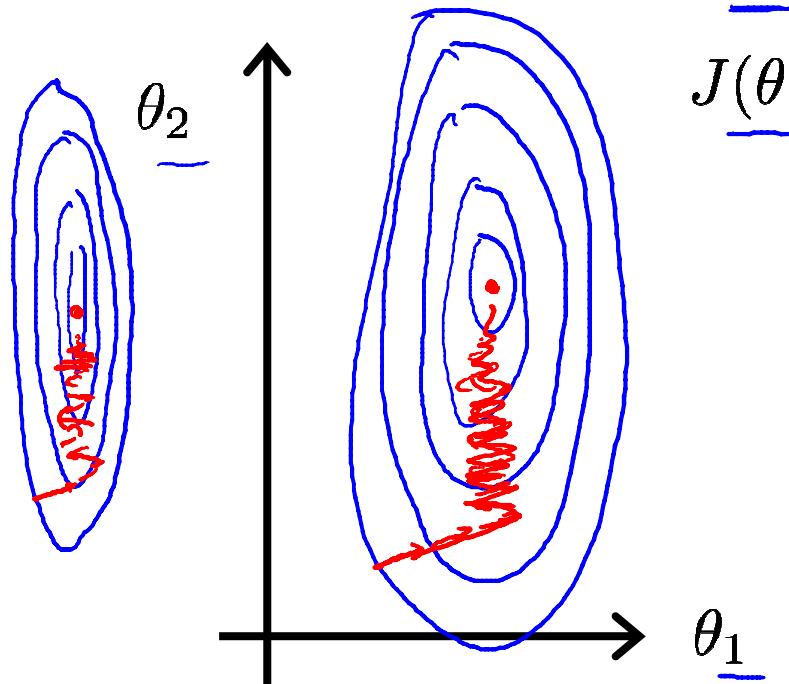
...

## Escalonamento (Normalização) das Variáveis

Ideia: Garantir que as variáveis estejam na mesma escala.

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

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

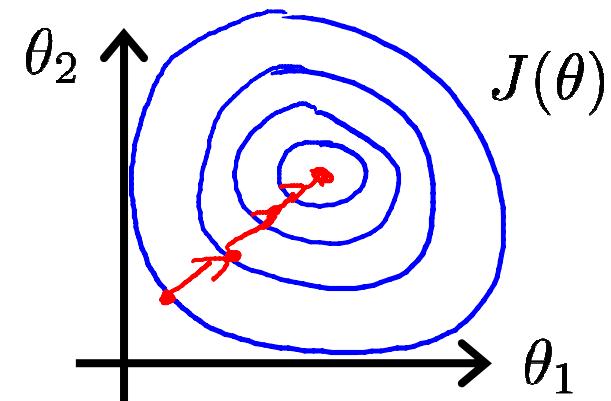


$$\rightarrow x_1 = \frac{\text{size } (\text{feet}^2)}{2000}$$

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

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1$$



## Escalonamento (Normalização) das Variáveis

Normalização linear entre  $y_{min}$  e  $y_{max}$ .  $-1 \leq x_i \leq 1$

$$y = (y_{max} - y_{min}) * (x - x_{min}) / (x_{max} - x_{min}) + y_{min}$$

Normalização para media 0 e desvio padrão 1.

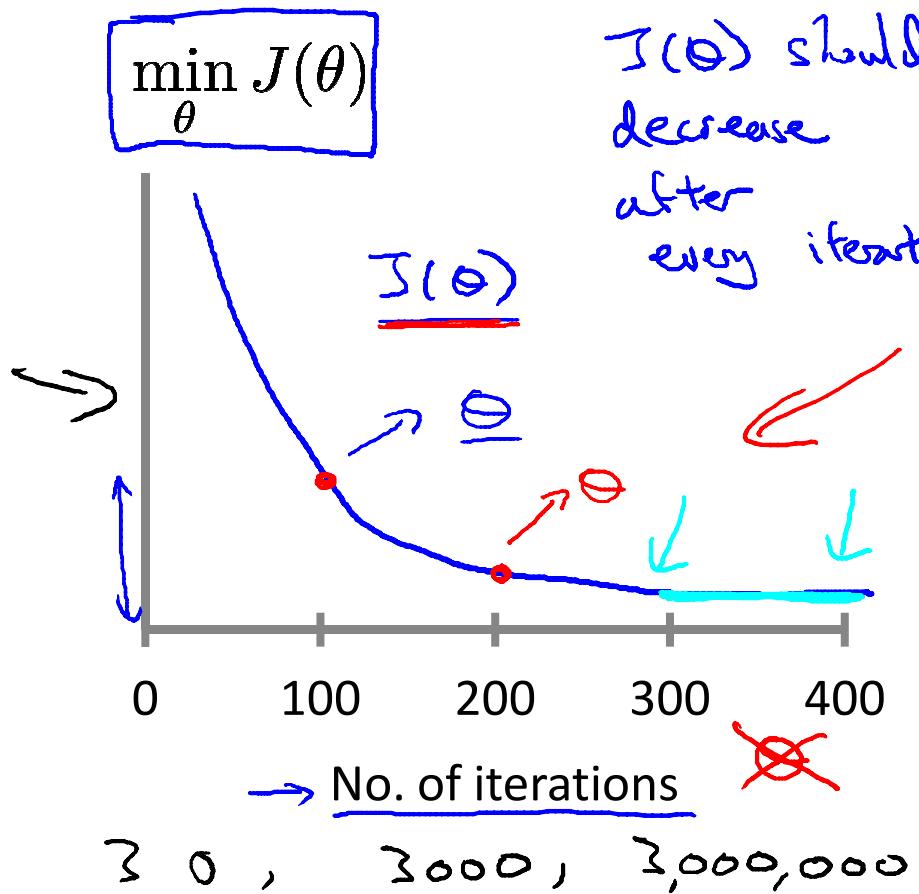
$$y = (x - x_{mean}) * (y_{std}/x_{std}) + y_{mean}$$

## Gradiente descendente

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

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

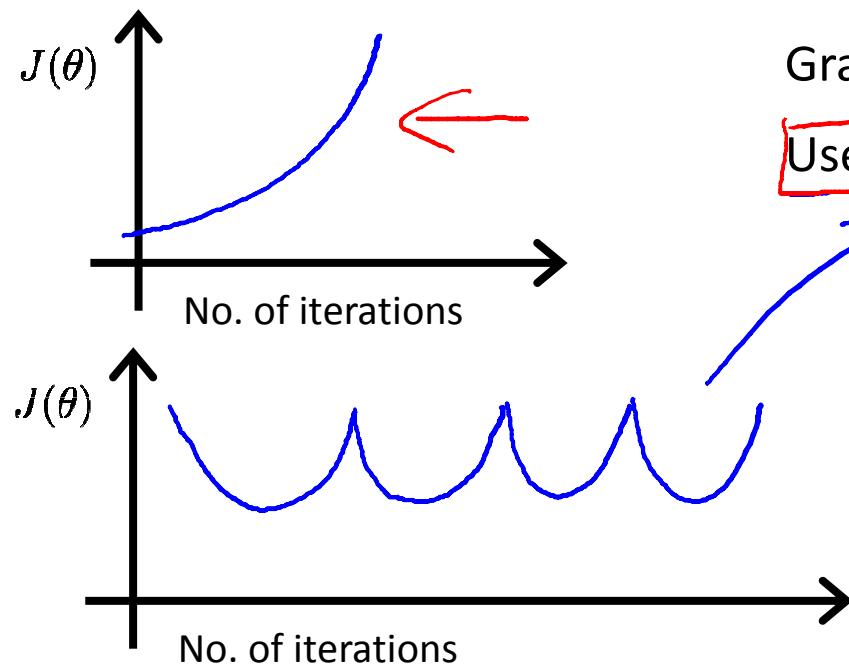
## Making sure gradient descent is working correctly.



→ Example automatic convergence test:

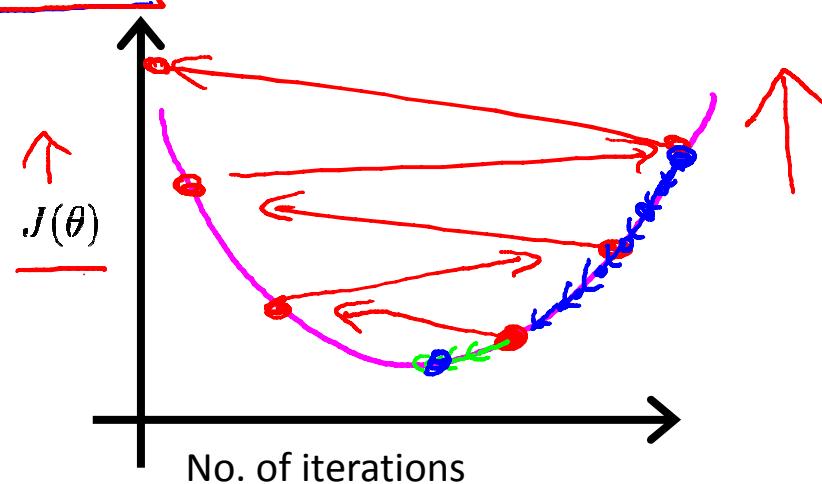
→ Declare convergence if  $\underline{J(\theta)}$  decreases by less than  $\underline{10^{-3}}$  in one iteration.

## 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.

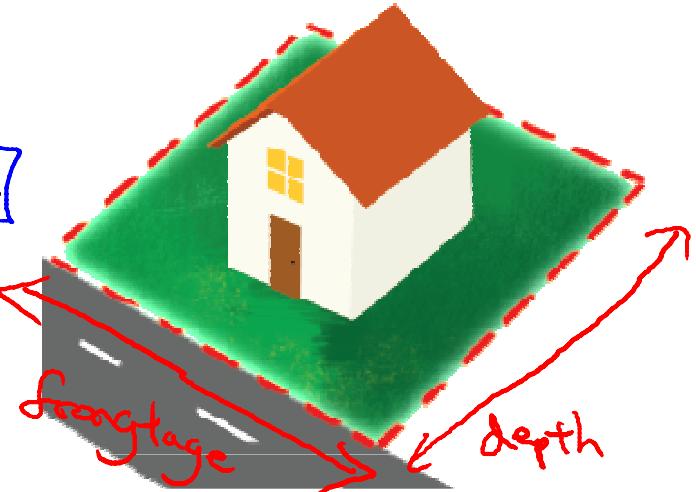
To choose  $\alpha$ , try

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

## Housing prices prediction

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

$x_1$   
 $x_2$



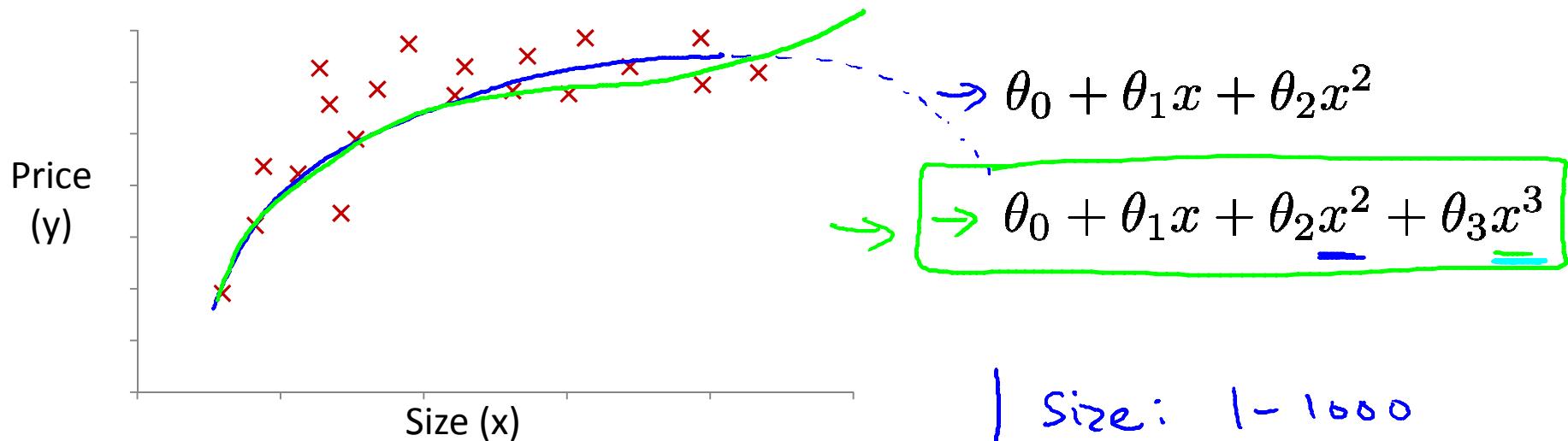
Area

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

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

↑ land area

## Regressão Polinomial



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

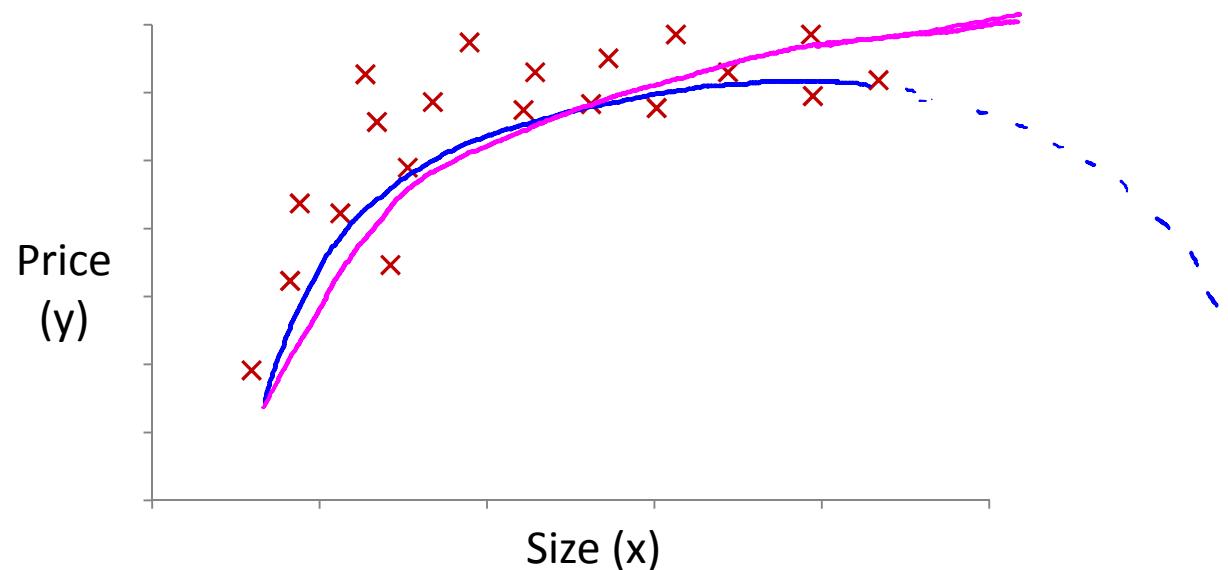
$\rightarrow x_1 = (\text{size})$

$\rightarrow x_2 = (\text{size})^2$

$\rightarrow x_3 = (\text{size})^3$

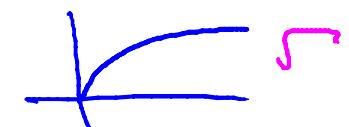
Size: 1 - 1000  
Size<sup>2</sup>: 1 - 1000,000  
Size<sup>3</sup>: 1 - 10<sup>9</sup>

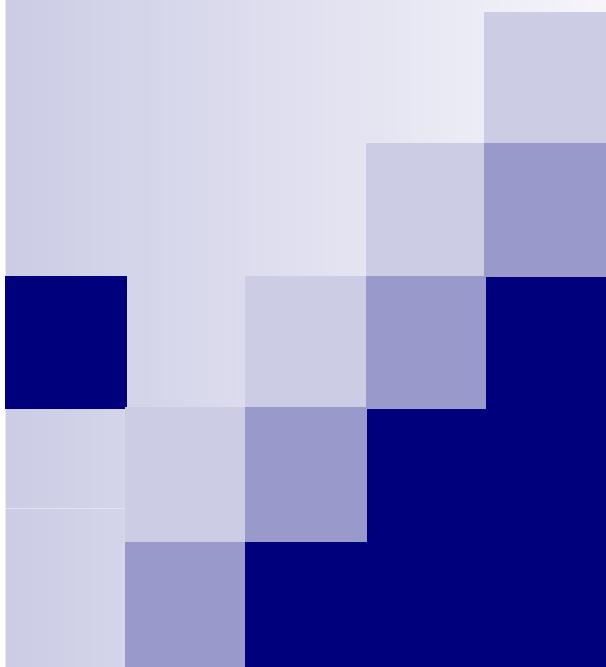
## Outras alternativas para “linearização” das variáveis



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

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





# INTRODUÇÃO AO ESTUDO DAS REDES NEURAIS ARTIFICIAIS

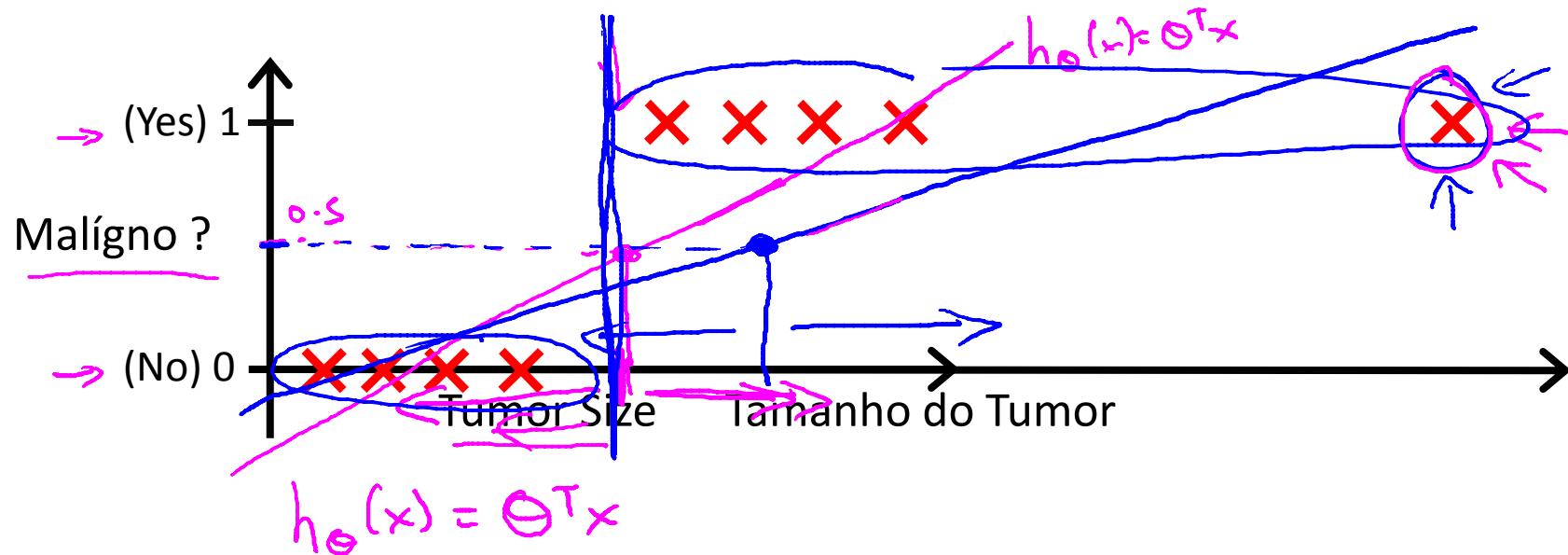
Laboratório de Conexionismo e Ciências

Cognitivas L3C  
Grupo SICRES  
INE - UFSC

## Regressão Logística

## Classification

- Email: Spam / Not Spam?
  - Transações Online: Fraudulenta (Yes / No)?
  - Tumor: Maligno / Benigno ?
- $y \in \{0, 1\}$
- 0: “Negativo” (e.g., tumor benigno)  
1: “Positivo” (e.g., tumor maligno)
- $y \in \{0, 1, 2, 3\}$



→ Saída do classificador  $h_\theta(x)$  para 0.5:

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

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

Classificação:  $y = 0 \text{ or } 1$

$h_\theta(x)$  Pode ser  $> 1$  or  $< 0$

Regressão Logística:  $0 \leq h_\theta(x) \leq 1$

Classification



# Régressão Logística

## Representação da Hipótese

## Modelo de Regressão Logística

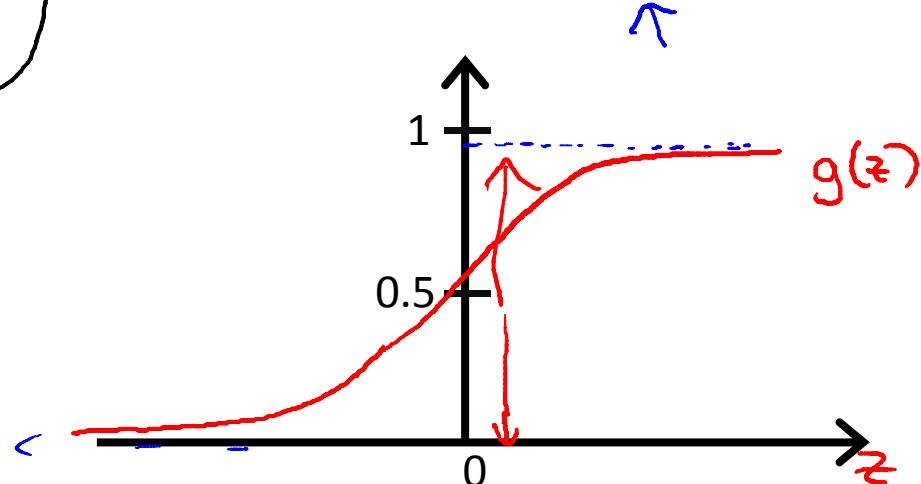
Queremos  $0 \leq h_\theta(x) \leq 1$

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

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

- Função Sísmoide
- Função Logística

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



Parametros  $\theta$ :

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

## Interpretação da Saída da Hipótese

$h_{\theta}(x)$  = probabilidade estimada de  $y = 1$  , entrada  $x$

Exemplo: Se  $x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$

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

70% de chance do tumor ser benígo

“probabilidade de  $y = 1$ , dado  $x$ ,  
parametrizado por  $\theta$ ”

$$\begin{aligned} P(y = 0|x; \theta) + P(y = 1|x; \theta) &= 1 \\ P(y = 0|x; \theta) &= 1 - P(y = 1|x; \theta) \end{aligned}$$



## Regressão Logística

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

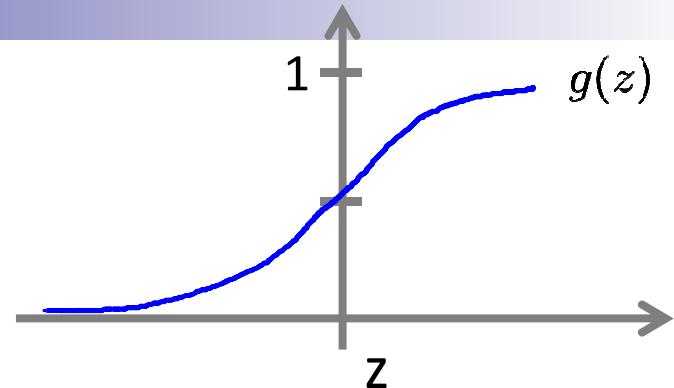
$$g(z) = \frac{1}{1+e^{-z}}$$

Suponha prever “ $y = 1$ ” se  $h_{\theta}(x) \geq 0.5$

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

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

$$\theta^T x < 0$$

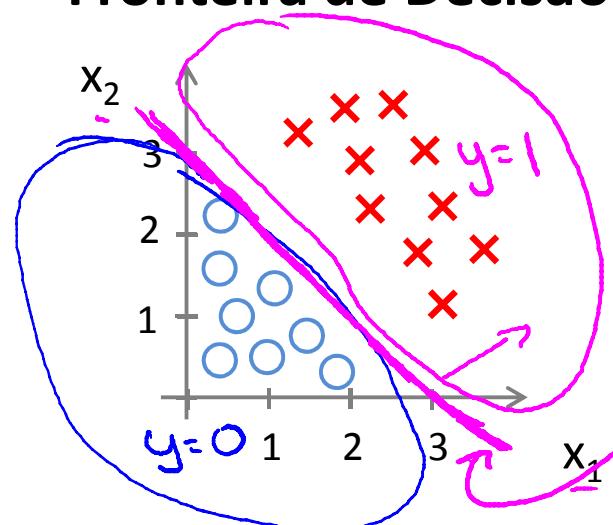


$$g(z) \geq 0.5 \\ \text{when } z \geq 0$$

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

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

## Fronteira de Decisão



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

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

Decision boundary

Prever " $y = 1$ " se  $\underline{\theta^T x} \geq 0$

$$x_1, x_2$$

$$\rightarrow h_{\theta}(x) = 0.5$$

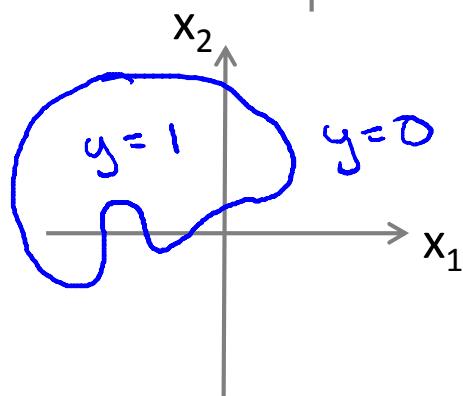
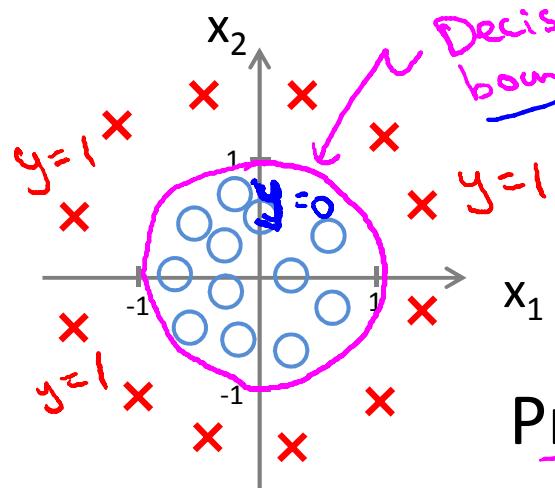
$$\boxed{x_1 + x_2 = 3}$$

$$\underline{x_1 + x_2 \geq 3}$$

$$x_1 + x_2 < 3$$

$$y = 0$$

## Fronteiras de Decisão Não-lineares



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

$\parallel \quad \parallel \quad \parallel$

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

Prever "y = 1" se  $-1 + x_1^2 + x_2^2 \geq 0$

$$x_1^2 + x_2^2 = 1$$

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

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



## Função Custo

Conjunto de Treino:  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

m exemplos

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

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

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

Como escolher o parâmetro  $\underline{\theta}$  ?

# Função Custo

→ Regressão Linear:

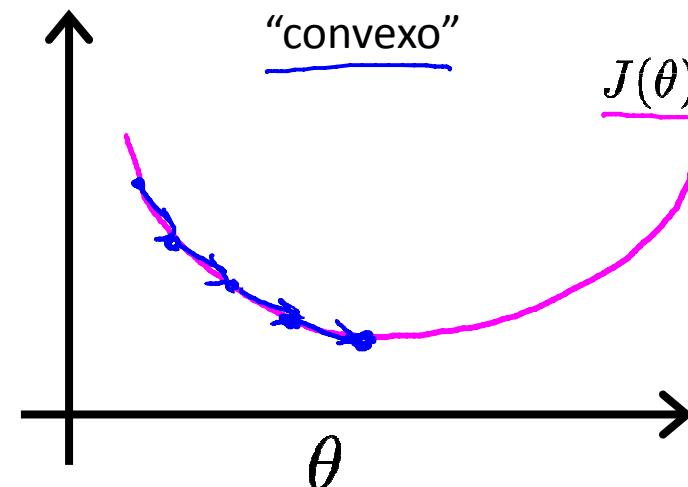
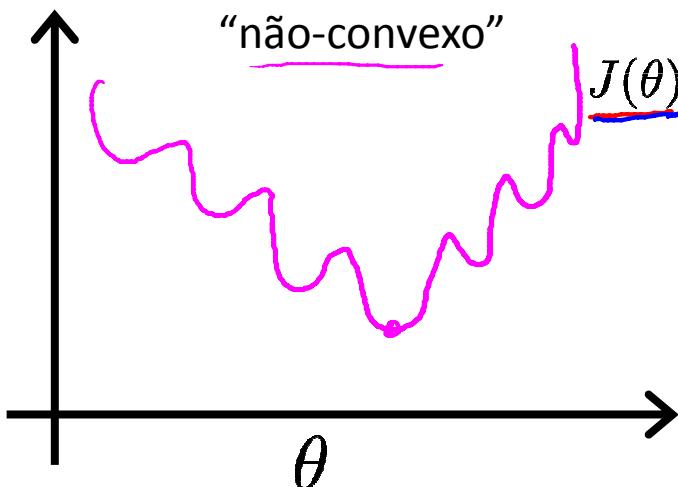
Logistic

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

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

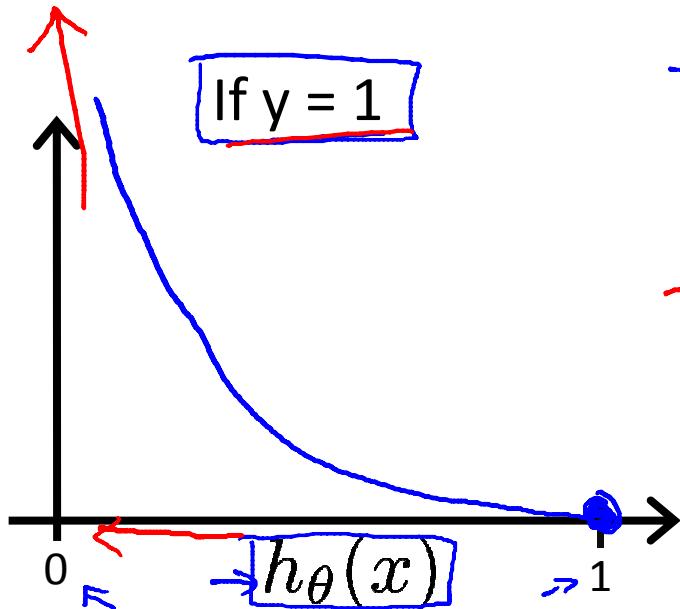
$$\rightarrow \text{Cost}(h_\theta(x^{(i)}), y^{(i)}) = \frac{1}{2} (h_\theta(x^{(i)}) - y^{(i)})^2$$

$$\frac{1}{2} h_\theta(x)^T y$$



## Função Custo da Regressão Logística

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

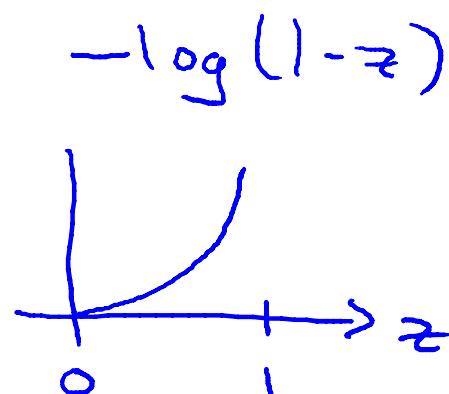
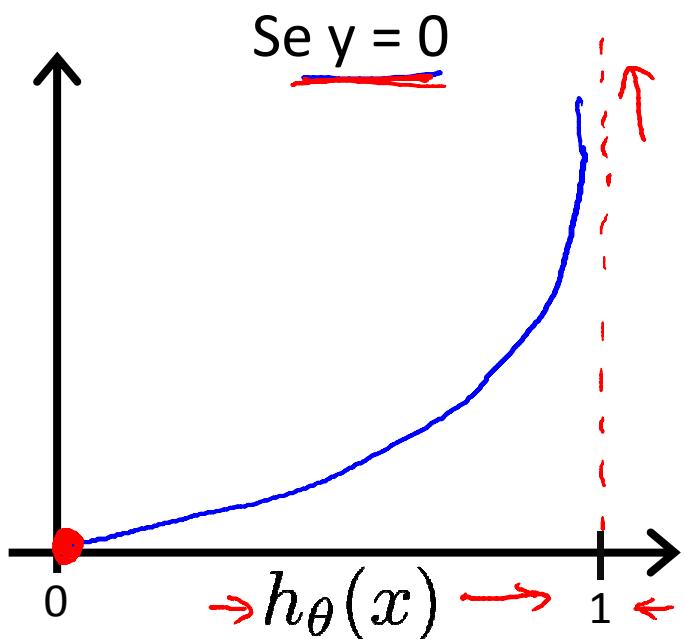


→ Cost = 0 if  $y = 1, h_{\theta}(x) = 1$   
But as  $h_{\theta}(x) \rightarrow 0$   
Cost  $\rightarrow \infty$

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.

## Função Custo da Regressão Logística

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



# Regressão Logística

**Função custo simplificada  
e gradiende descendente**

## Função custo da regressão logística

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

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

Note: y = 0 or 1 always

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

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

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

## Função custo da regressão logística

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

Ajustar o parâmetro  $\theta$  :

$$\boxed{\min_{\theta} J(\theta) \quad \text{Cret } \ominus}$$

Fazer previsão dado um novo  $x$  :

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

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

## Gradiente Descendente

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

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

Repeat {

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

}

(Atualiza simultaneamente  $\theta_j$ )

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

## Gradient Descent

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

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

Repeat {

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

(atualiza simultaneamente  $\theta_j$ )

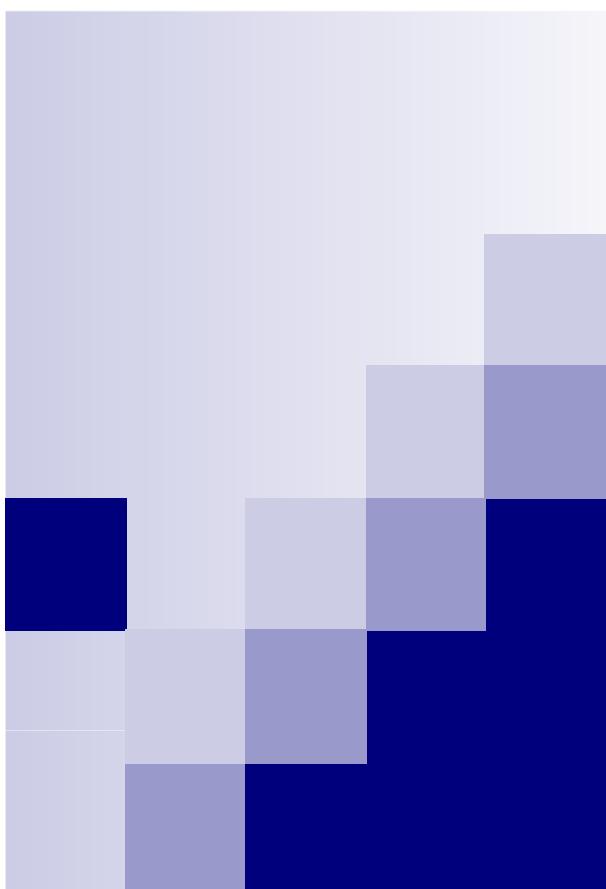
}

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

$$h_\theta(x) = \theta^T x$$

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

Algoritmo parece idêntico à regressão linear!



# Régressão Logística

## Otimização Avançada

## Algoritmo de Otimização

Função custo  $J(\theta)$ . Desejamos  $\min_{\theta} J(\theta)$ .

Dado  $\theta$ , podemos computar

$$\begin{aligned}\rightarrow & - J(\theta) \\ \rightarrow & \frac{\partial}{\partial \theta_j} J(\theta) \quad (\text{for } j = 0, 1, \dots, n)\end{aligned}$$

Gradiente descendente:

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

# Algoritmo de Otimização

Dado  $\theta$ , podemos computar

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

(for  $j = 0, 1, \dots, n$ )

Algoritmos de Otimização:

- - Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

Vantagens:

- Sem atribuição manual de  $\alpha$
- Sempre é mais rápido que o gradient descendente.

Desvantagens:

- Mais complexo

Exemplo:

$$\rightarrow \theta = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix} \quad \min_{\theta} J(\theta)$$

$\theta_1 = 5, \theta_2 = 5$ .

$$\rightarrow J(\theta) = (\theta_1 - 5)^2 + (\theta_2 - 5)^2$$
$$\rightarrow \frac{\partial}{\partial \theta_1} J(\theta) = 2(\theta_1 - 5)$$
$$\rightarrow \frac{\partial}{\partial \theta_2} J(\theta) = 2(\theta_2 - 5)$$

```
→ options = optimset ('GradObj', 'on', 'MaxIter', '100');  
→ initialTheta = zeros(2,1);  
[optTheta, functionVal, exitFlag] ...  
= fminunc(@costFunction, initialTheta, options);  
 $\theta \in \mathbb{R}^d \quad d \geq 2.$ 
```

```
function [jVal, gradient]  
= costFunction(theta)  
jVal = (theta(1)-5)^2 + ...  
      (theta(2)-5)^2;  
gradiente = zeros(2,1);  
gradiente(1) = 2*(theta(1)-5);  
gradiente(2) = 2*(theta(2)-5);
```

$$\underline{\text{theta}} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad \begin{array}{l} \xrightarrow{\hspace{1cm}} \text{theta}(1) \\ \xrightarrow{\hspace{1cm}} \text{theta}(2) \\ \vdots \\ \xrightarrow{\hspace{1cm}} \text{theta}(n+1) \end{array}$$

**function** [jVal, gradient] = costFunction(theta)

jVal = [ código a computar  $J(\theta)$  ] ;

gradiente(1) = [ código a computar  $\frac{\partial}{\partial \theta_0} J(\theta)$  ] ;

gradiente(2) = [ código a computar  $\frac{\partial}{\partial \theta_1} J(\theta)$  ] ;

:

:

gradiente(n+1) = [ código a computar  $\frac{\partial}{\partial \theta_n} J(\theta)$  ] ;

# Regressão Logística

Classificação multiclasse:  
Um-vs-Todos



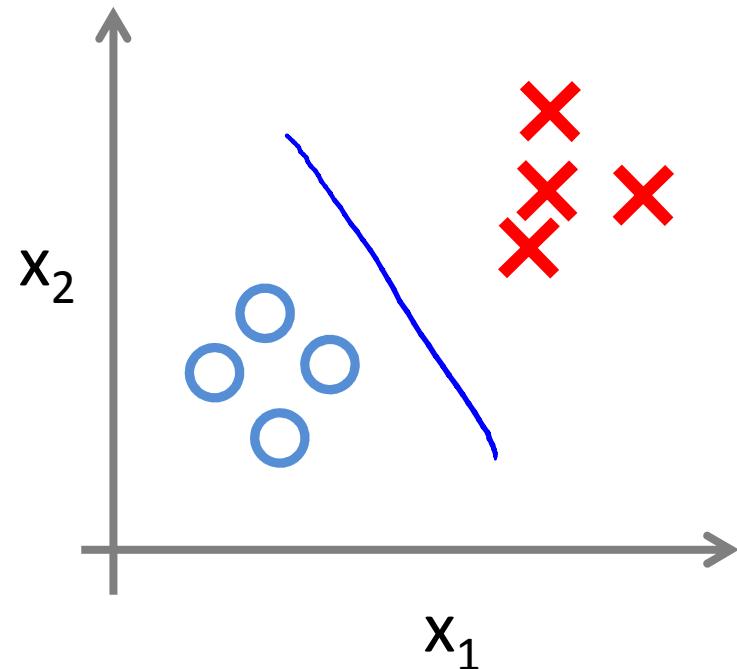
## Classificação Multiclasse

Filtro de email: Trabalho, Amigos, Família

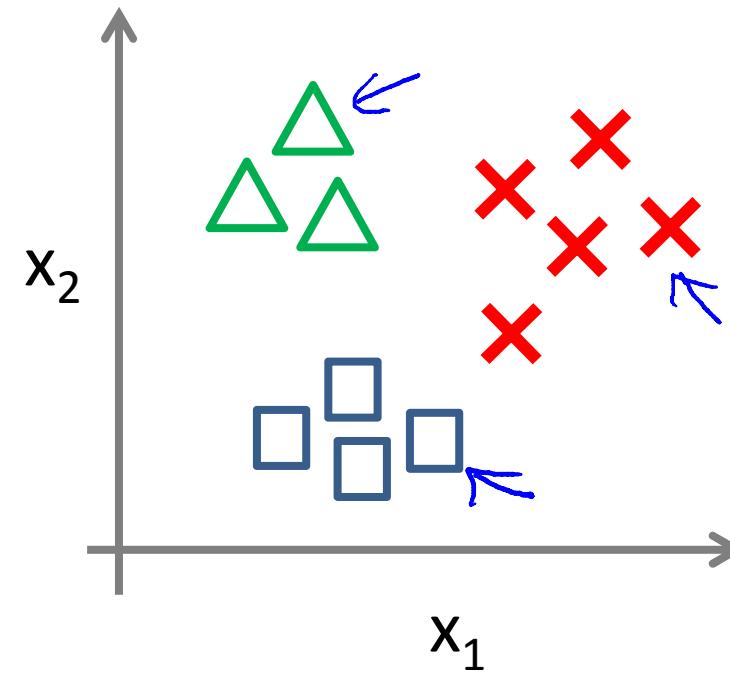
Diagramas médicos: Doente, Resfriado, Gripado

Tempo: Ensolarado, Nublado, Chuvoso, Neve

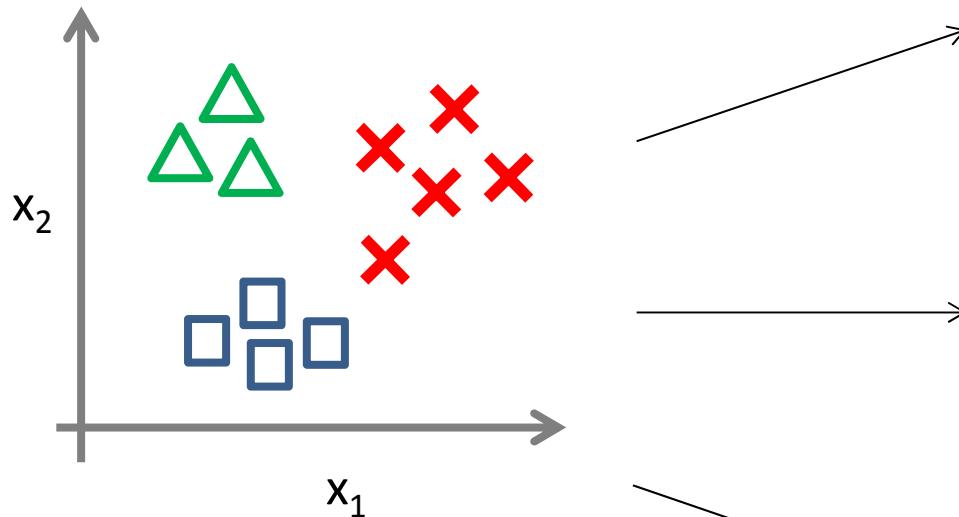
Classificação binária:



Classificação multiclasse



## Um-vs-Todos(um-vs-resto):

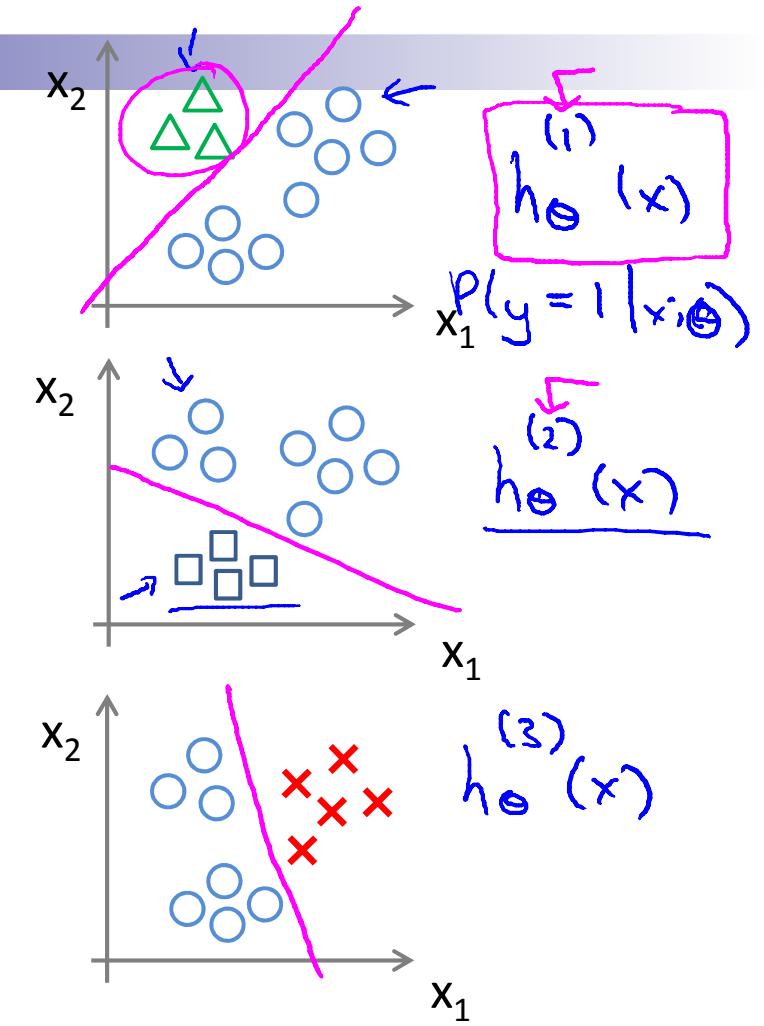


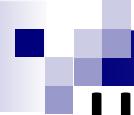
Classe 1:  $\triangle \leftarrow$

Classe 2:  $\square \leftarrow$

Classe 3:  $\times \leftarrow$

$$h_{\theta}^{(i)}(x) = P(y = i|x; \theta) \quad (i = 1, 2, 3)$$





## Um-vs-Todos

Treinar um classificador de regressão logística  $h_{\theta}^{(i)}(x)$  para cada classe  $i$  para prever a probabilidade de  $y = i$

Em uma nova entrada  $x$ , encontrar a classe  $i$  que maximiza

$$\max_i h_{\theta}^{(i)}(x)$$