

# Fundamentos em Redes Neurais

Introdução à Conceitos Avançados de  
Regressão

Profº - Dr. Thales Levi Azevedo Valente  
[thales.l.a.valente@gmail.com.br](mailto:thales.l.a.valente@gmail.com.br)

# Sejam Bem-vindos !



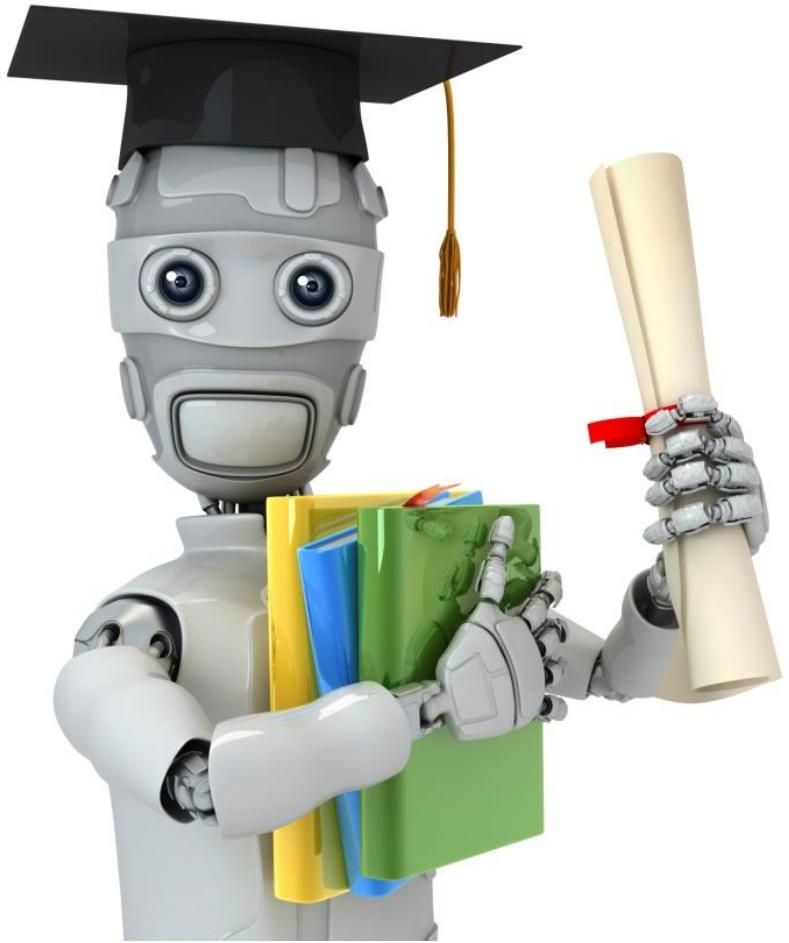
**Os celulares devem  
ficar no silencioso  
ou desligados**

Pode ser utilizado  
apenas em caso  
de emergência



**Boa tarde/noite, por  
favor e com licença  
DEVEM ser usados**

Educação é  
essencial



Machine Learning

Linear Regression with  
multiple variables

---

Multiple features

# Multiple features (variables).

Size (feet <sup>2</sup> )	Price (\$1000)
$\xrightarrow{x}$	$y \xleftarrow{ }$
2104	460
1416	232
1534	315
852	178
...	...

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

---

# Multiple features (variables).

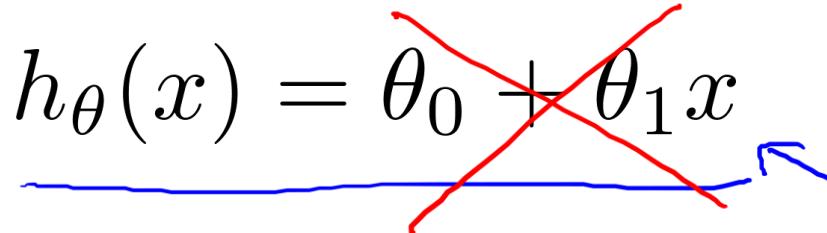
$\rightarrow$ Size (feet <sup>2</sup> )	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$x_1$	$x_2$	$x_3$	$x_4$	$y$
2104	5	1	45	460
$\rightarrow$ 1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...	...	...	...	...
$\uparrow$	$\uparrow$	$\uparrow$	$\uparrow$	
Notation:				
$\rightarrow n =$ number of features				$n = 4$
$\rightarrow x^{(i)}$ input (features) of $i^{\text{th}}$ training example.				
$\rightarrow x_j^{(i)}$ value of feature $j$ in $i^{\text{th}}$ training example.				

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

$x_3^{(2)} = 2$

Hypothesis:

Previously:

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


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

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

↑      ↑      ↑      ↑  
age

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

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

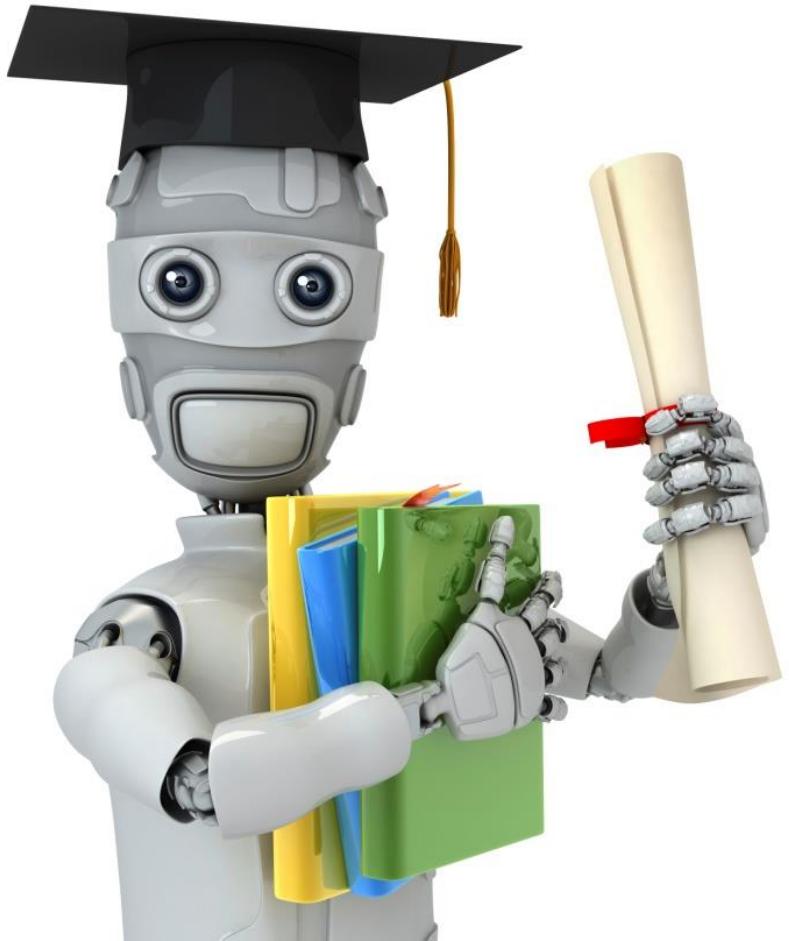
$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

$$h_{\theta}(x) = \underline{\theta_0x_0 + \theta_1x_1 + \cdots + \theta_nx_n} \\ = \boxed{\theta^T x.}$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_n \end{bmatrix} \theta^T \quad (n+1) \times 1 \text{ matrix} \quad \theta^T x$$

Multivariate linear regression.  $\leftarrow$



Machine Learning

# Linear Regression with multiple variables

---

## Gradient descent for multiple variables

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

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

Cost function:

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

Gradient descent:

Repeat {  
     $\rightarrow \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n)$   $\underline{J(\Theta)}$   
    }  
        ↑  
        (simultaneously update for every  $j = 0, \dots, n$ )

# Gradient Descent

Previously (n=1):

Repeat {

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

$\frac{\partial}{\partial \theta_0} J(\theta)$

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

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

}

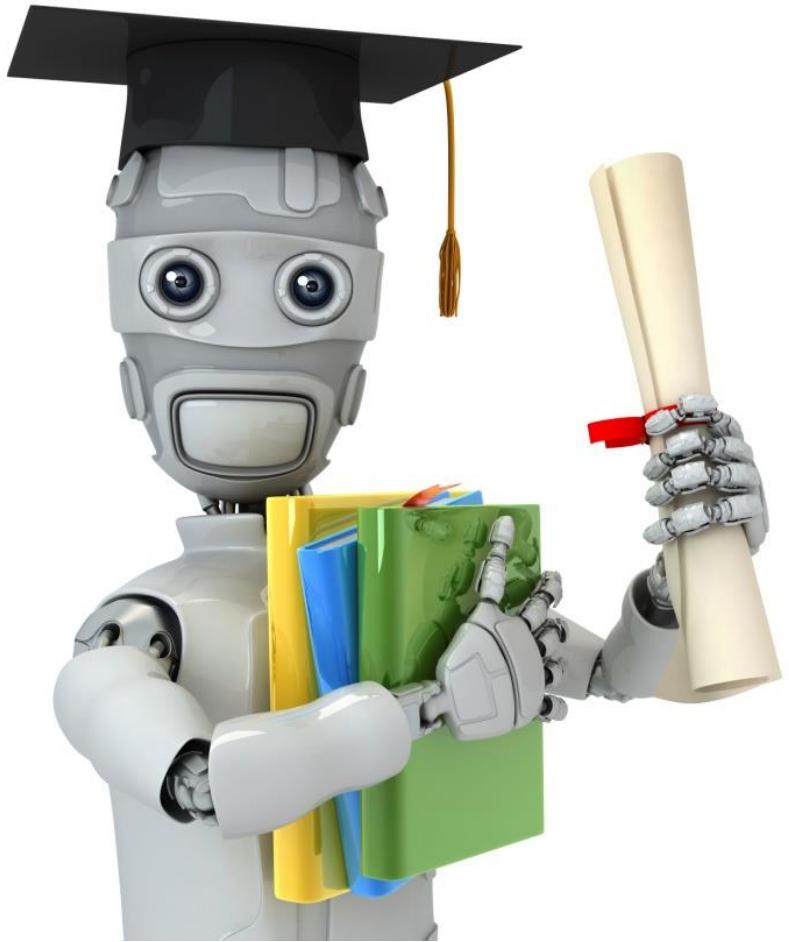
$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

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

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

...

$x_0^{(i)} = 1$



Machine Learning

# Linear Regression with multiple variables

---

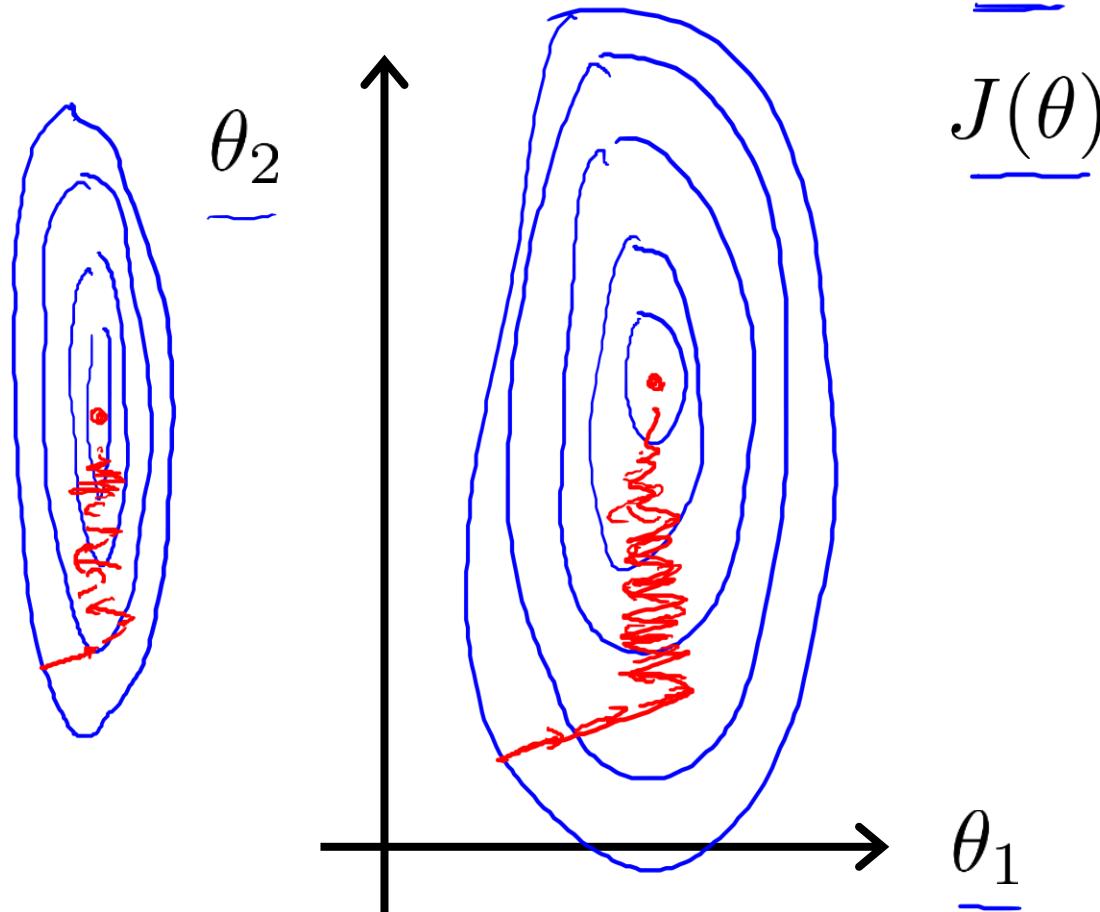
## Gradient descent in practice I: Feature Scaling

# Feature Scaling

Idea: Make sure features are on a similar scale.

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

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

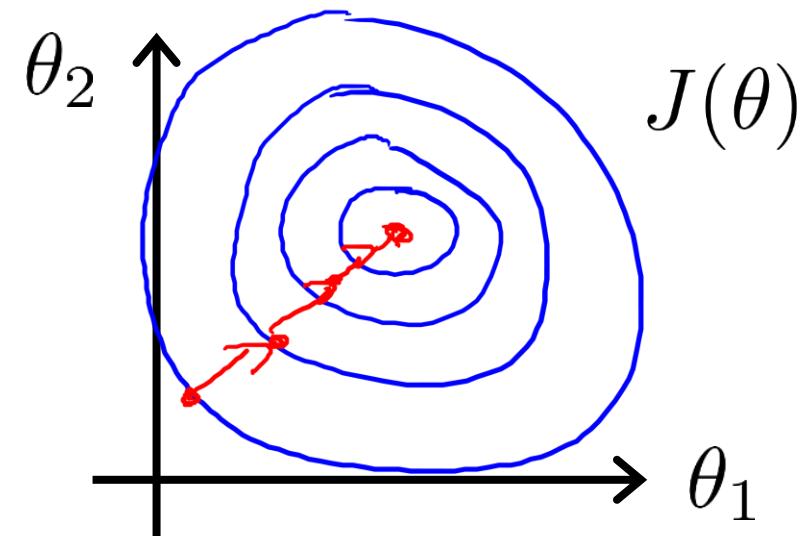


$$\rightarrow x_1 = \frac{\text{size } (\text{feet}^2)}{2000} \quad \swarrow$$

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

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1$$



## Feature Scaling

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

$$x_0 = 1$$

$$6 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-\underline{100} \leq x_3 \leq \boxed{100} \quad \times$$

$$-0.0001 \leq x_4 \leq \boxed{0.0001} \quad \times$$

$$\begin{array}{c} \downarrow \\ -1 \leq x_i \leq 1 \\ \downarrow \end{array}$$

$$-3 \text{ to } 3 \quad \checkmark$$

$$-\frac{1}{2} \text{ to } \frac{1}{2} \quad \checkmark$$

## Mean normalization

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

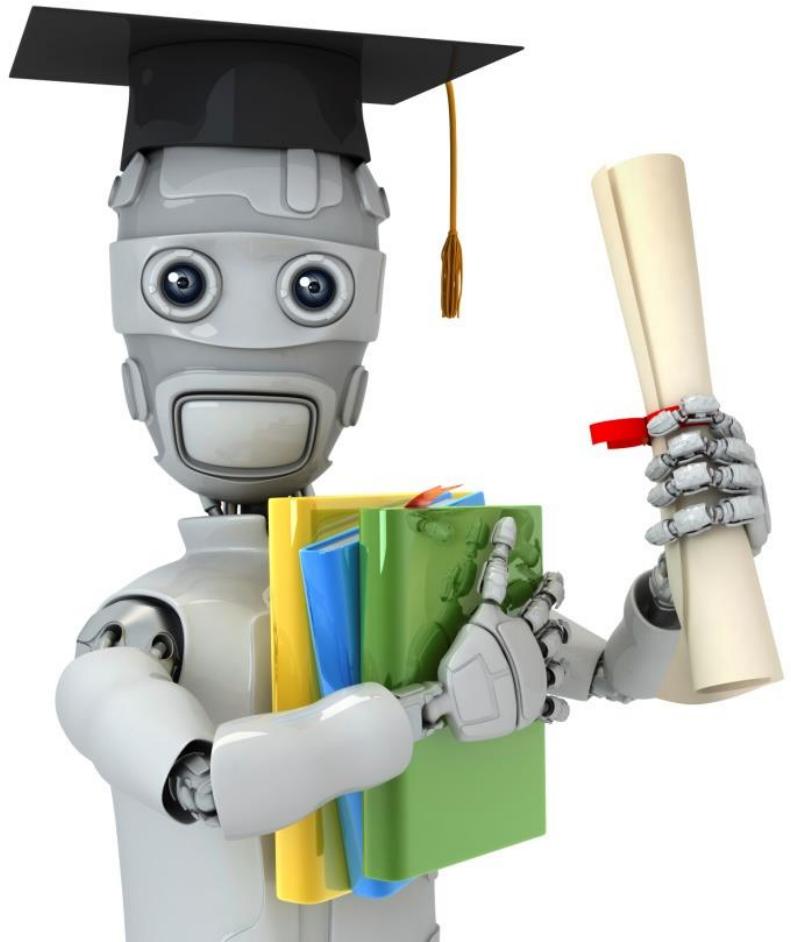
E.g.  $x_1 = \frac{\text{size} - 1000}{2000}$       Average size = 100  
 $x_2 = \frac{\#bedrooms - 2}{5 - 4}$       1-5 bedrooms  
 $\rightarrow [-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5]$

$x_1 \leftarrow \frac{x_1 - \mu_1}{\sigma_1}$        $x_2 \leftarrow \frac{x_2 - \mu_2}{\sigma_2}$

ang value of  $x_1$  in training set

range (max-min) (or standard deviation)





Machine Learning

# Linear Regression with multiple variables

---

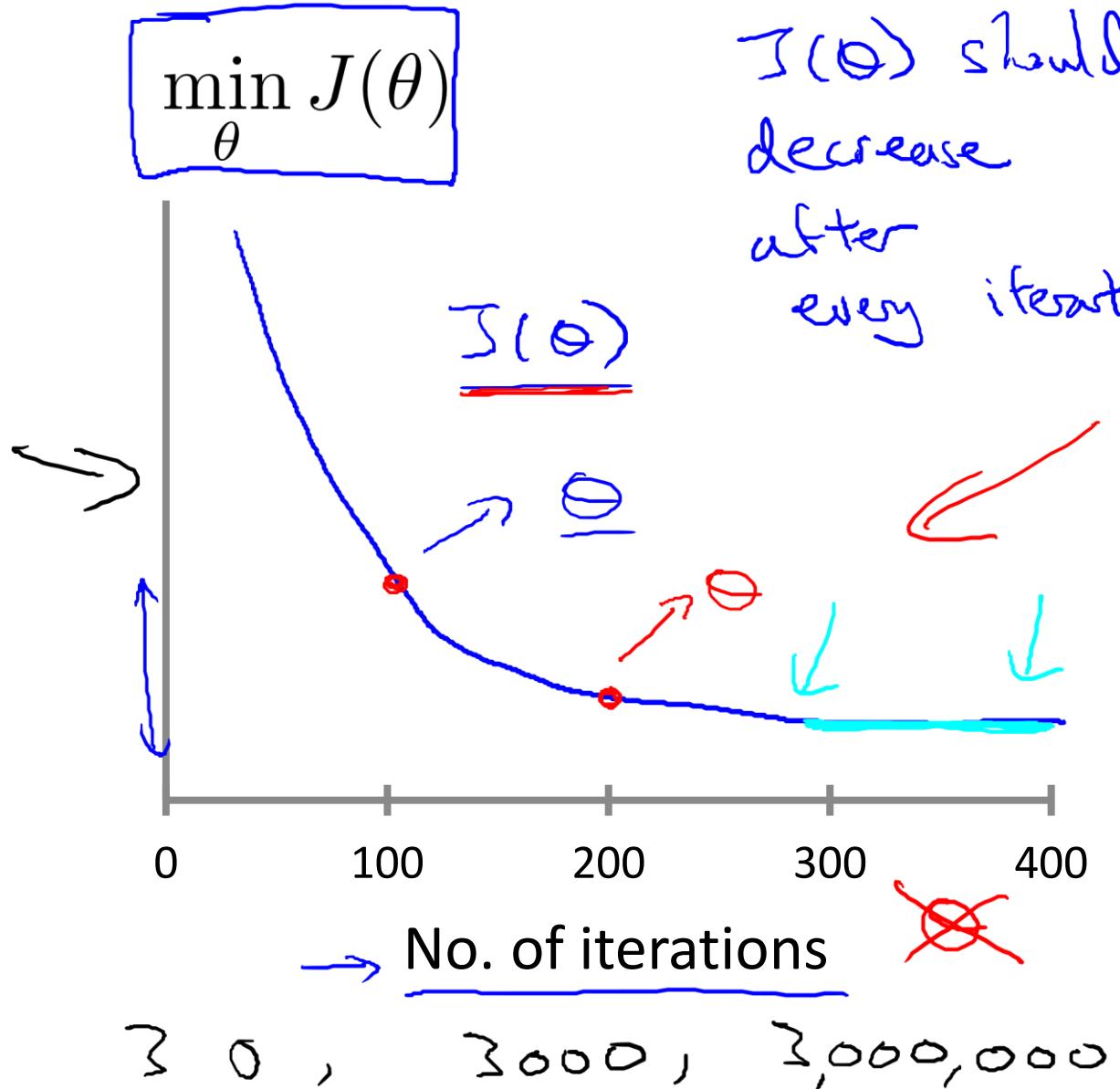
Gradient descent in  
practice II: Learning rate

# Gradient descent

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

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

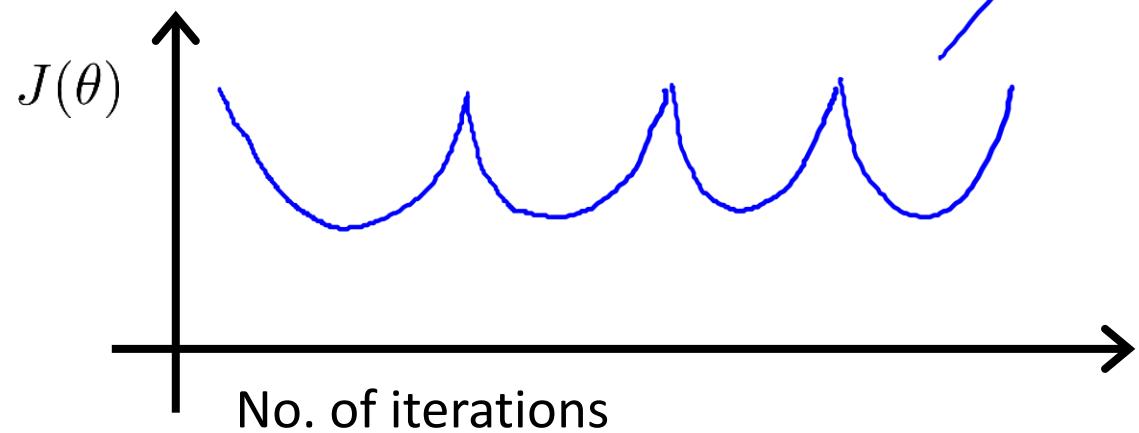
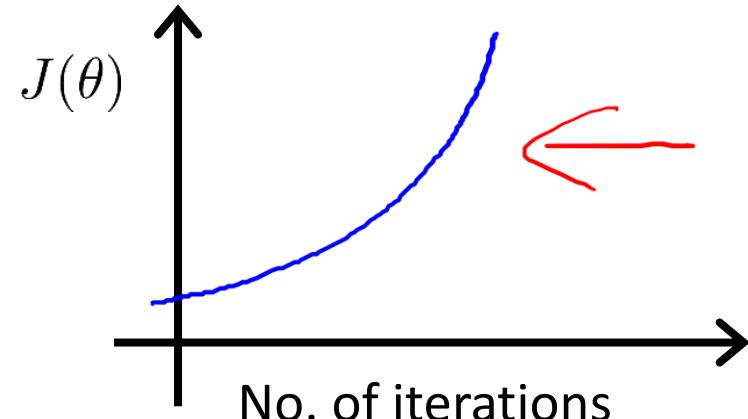
# Making sure gradient descent is working correctly.



→ Example automatic convergence test:

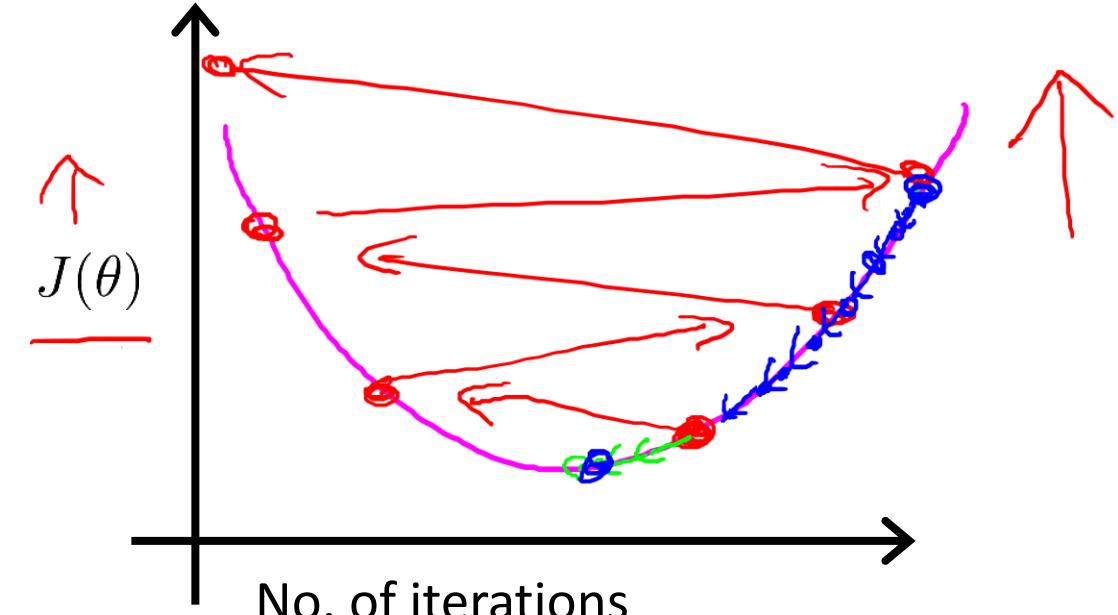
→ Declare convergence if  $\frac{J(\theta)}{\sum}$  decreases by less than  $10^{-3}$  in one iteration.

# Making sure gradient descent is working correctly.



Gradient descent not working.

Use smaller  $\alpha$



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

## Summary:

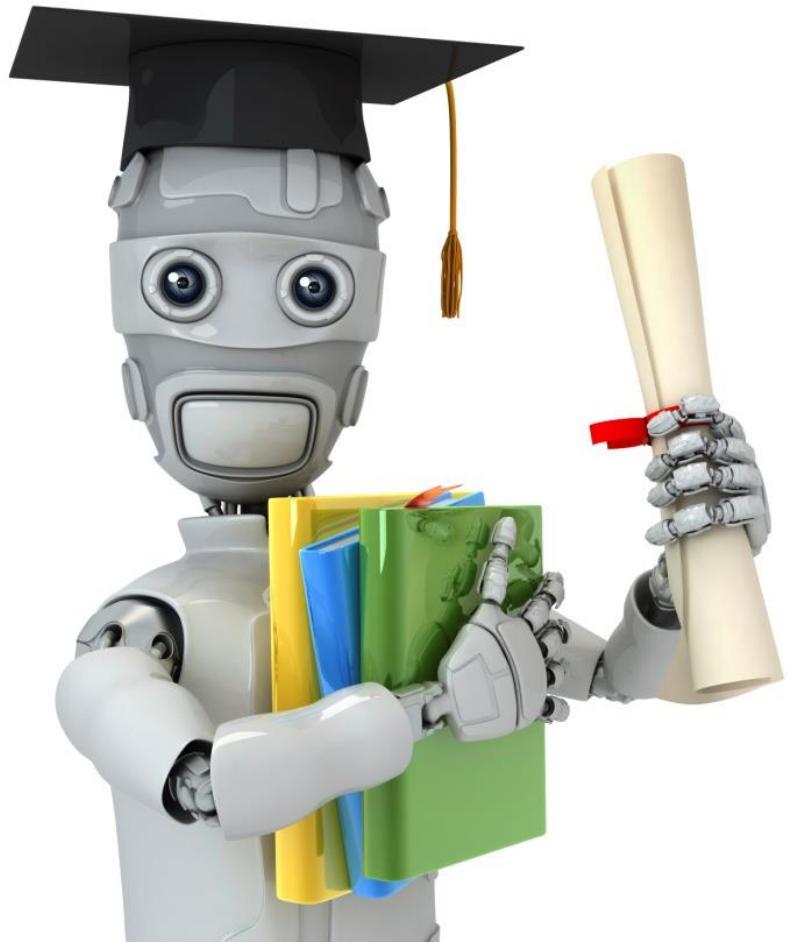
- If  $\alpha$  is too small: slow convergence.
- If  $\alpha$  is too large:  $J(\theta)$  may not decrease on every iteration; may not converge. (Slow converge also possible.)



To choose  $\alpha$ , try

$$\dots, \underbrace{0.001, 0.003}_{\uparrow \curvearrowright}, \underbrace{0.01, 0.03}_{\uparrow \curvearrowright}, \underbrace{0.1, 0.3}_{\uparrow \curvearrowright}, \underbrace{1, \dots}_{\uparrow \curvearrowright}$$

The values are grouped into pairs by blue curly braces. The first pair (0.001, 0.003) has arrows pointing up and right, with "3x" written below the brace. The second pair (0.01, 0.03) has arrows pointing up and right, with " $\approx 3x$ " written below the brace. The third pair (0.1, 0.3) has arrows pointing up and right, with "3x" written below the brace. The fourth value (1) has an arrow pointing up and right, with " $\approx 3x$ " written below the brace.



Machine Learning

# Linear Regression with multiple variables

---

Features and  
polynomial regression

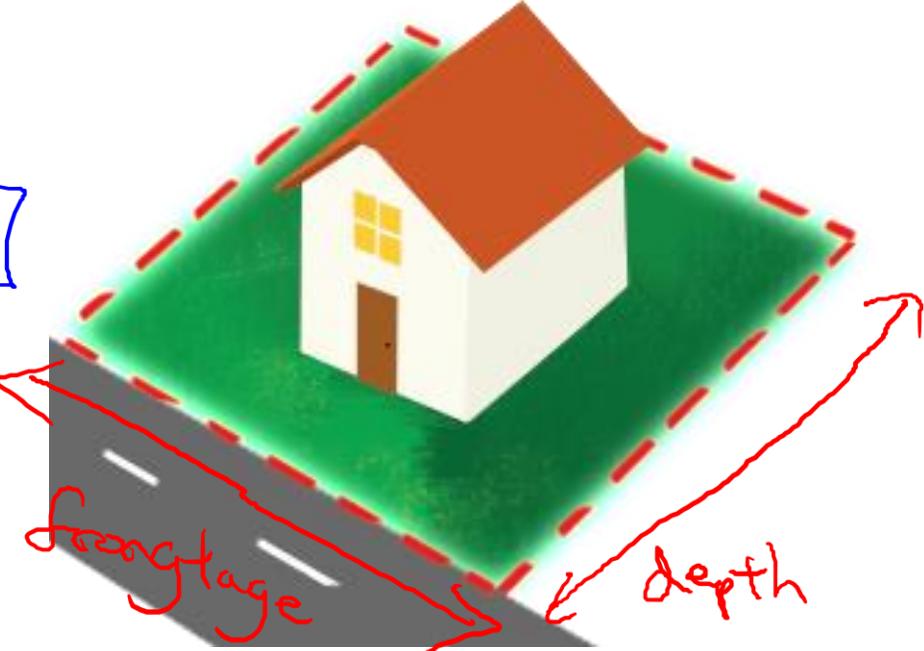
# Housing prices prediction

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

$x_1$   
—  
 $x_2$

Area

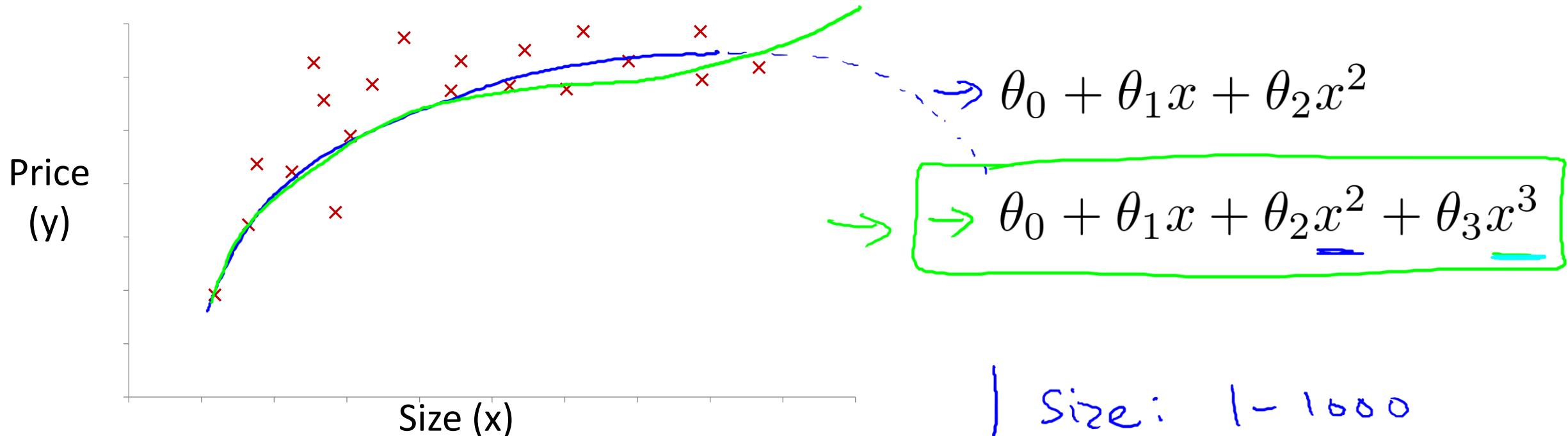
$$x = \underline{\text{frontage} * \text{depth}}$$



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

$\curvearrowleft$  land area

# Polynomial regression



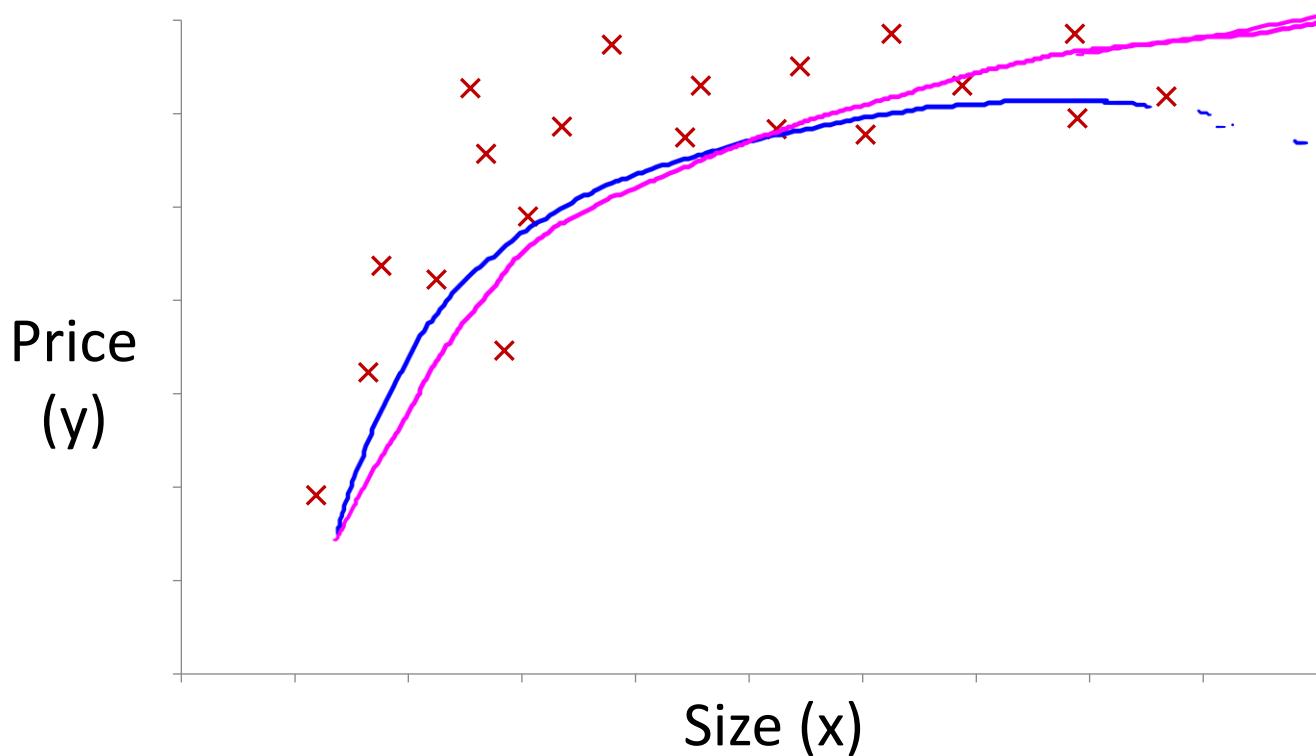
$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3$$

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

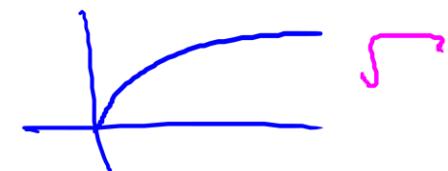
Size: 1 - 1600
Size <sup>2</sup> : 1 - 1000,000
Size <sup>3</sup> : 1 - 10 <sup>9</sup>

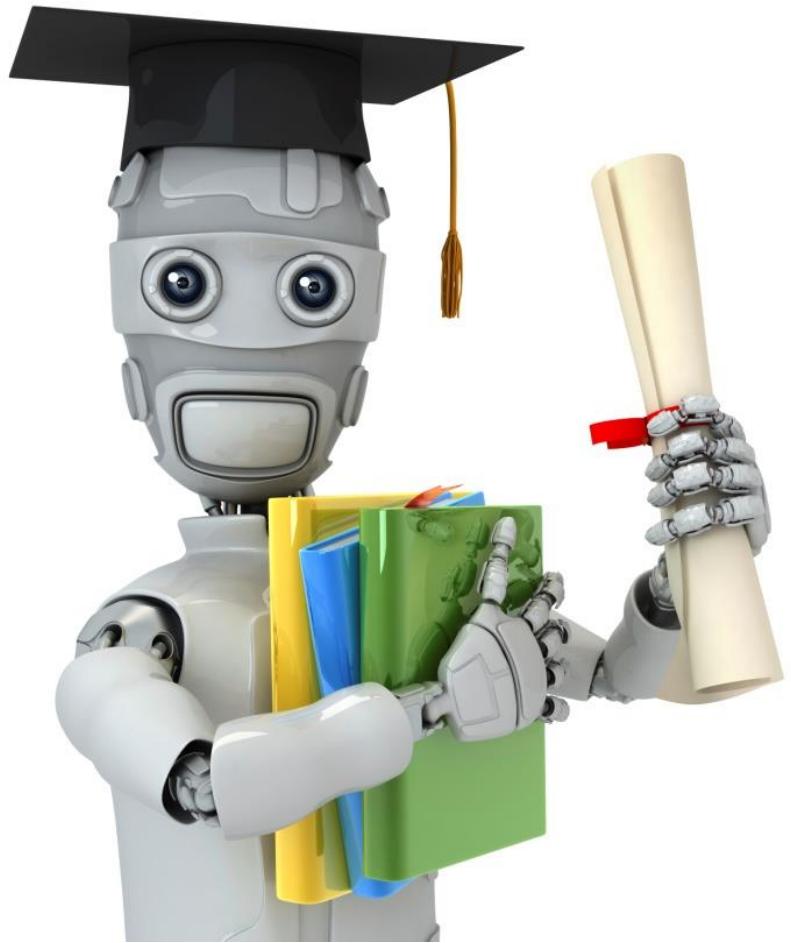
# Choice of features



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

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





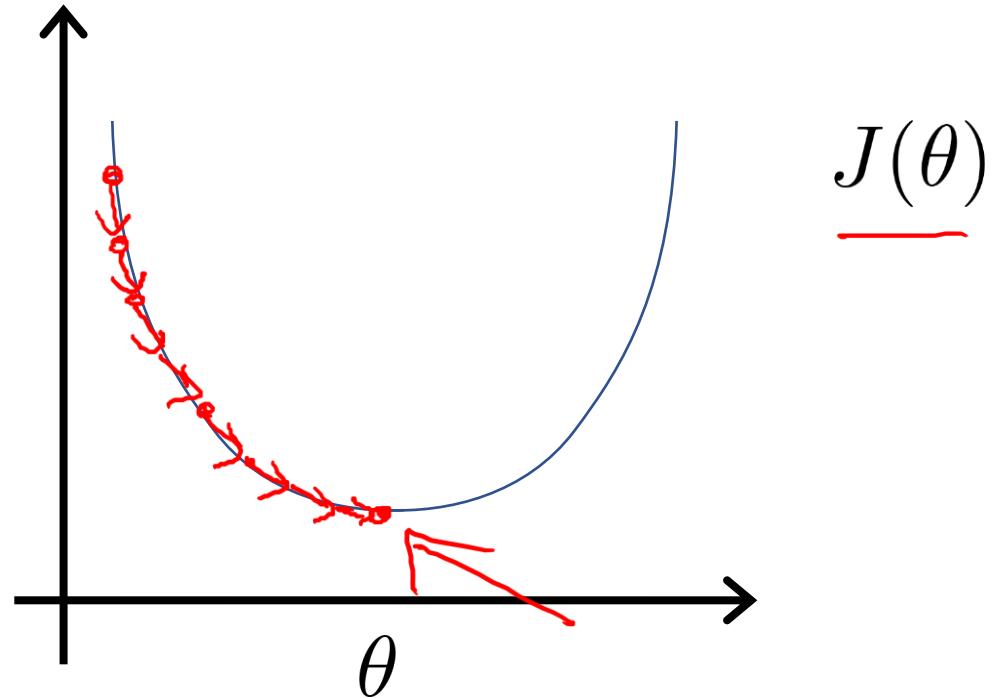
Machine Learning

# Linear Regression with multiple variables

---

Normal equation

# Gradient Descent



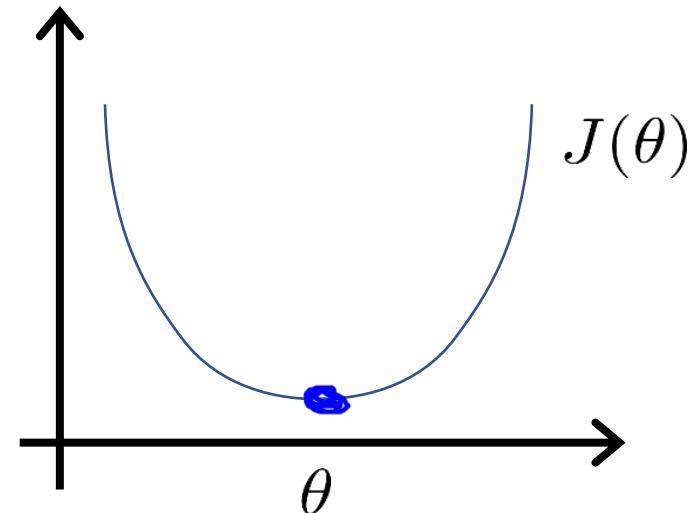
Normal equation: Method to solve for  $\theta$   
analytically.

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

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

$$\frac{\partial}{\partial \theta} J(\theta) = \dots \stackrel{\text{set}}{=} 0$$

Solve for  $\theta$



$$\theta \in \mathbb{R}^{n+1}$$

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

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

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

Examples:  $m = 4$ .

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

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

$m \times (n+1)$

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

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

$m$ -dimensional vector

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

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

$\times$

(design matrix)

$= \begin{bmatrix} \cdots & (x^{(1)})^T & \cdots \\ \cdots & (x^{(2)})^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & (x^{(m)})^T & \cdots \end{bmatrix}$

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

$\times$

$$\underline{\Theta} = (\underline{x}^T \underline{x})^{-1} \underline{x}^T \underline{y}$$

$\underline{x} = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix}$

$\underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix}$

$m \times 2$

$m \times (n+1)$

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



$(X^T X)^{-1}$  is inverse of matrix  $X^T X$ .

Set  $A = X^T X$

$$\boxed{(X^T X)^{-1}} = A^{-1}$$

Octave:  $\text{pinv}(\boxed{X' * X}) * X' * y$

$$\text{pinv}(\boxed{X'^T * X}) * X'^T * y$$

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

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

$X'$        $X^T$

~~Feature Scaling~~

$0 \leq x_1 \leq 1$

$0 \leq x_2 \leq 1000$

$0 \leq x_3 \leq 10^{-5}$  ✓

$m$  training examples,  $n$  features.

## Gradient Descent

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

$$\underline{n = 10^6}$$

← -

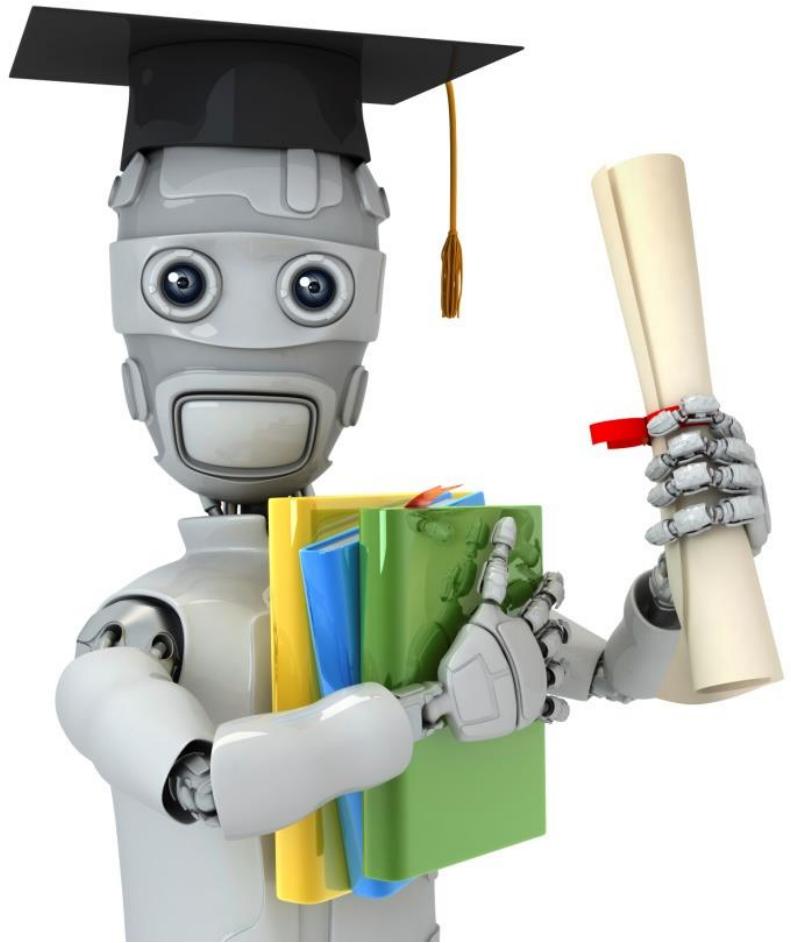
## Normal Equation

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

$$n = 100$$

$$n = 1000$$

$$\dots \quad n = \underline{10000}$$



Machine Learning

# Linear Regression with multiple variables

---

Normal equation and  
non-invertibility  
(optional)

# Normal equation

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

$X^T X$

- What if  $X^T X$  is non-invertible? (singular/  
degenerate)
- Octave: **pinv (X' \* X) \* X' \* y**

pinv (X' \* X) \* X' \* y



|  
|  $\rightarrow$  pinv  
|  $\rightarrow$  inv

What if  $\boxed{X^T X}$  is non-invertible?



- Redundant features (linearly dependent).

E.g.

$$\begin{cases} \underline{x_1} = \text{size in feet}^2 \\ \underline{x_2} = \text{size in m}^2 \\ \underline{x_1} = (3.28)^2 \underline{x_2} \end{cases}$$

$$1m = 3.28 \text{ feet}$$

$$\rightarrow \underline{m = 10} \leftarrow$$

$$\rightarrow \underline{n = 100} \leftarrow$$

$$\Theta \in \mathbb{R}^{101}$$

- Too many features (e.g.  $\underline{m \leq n}$ ).

- Delete some features, or use regularization.

later

↓  
later

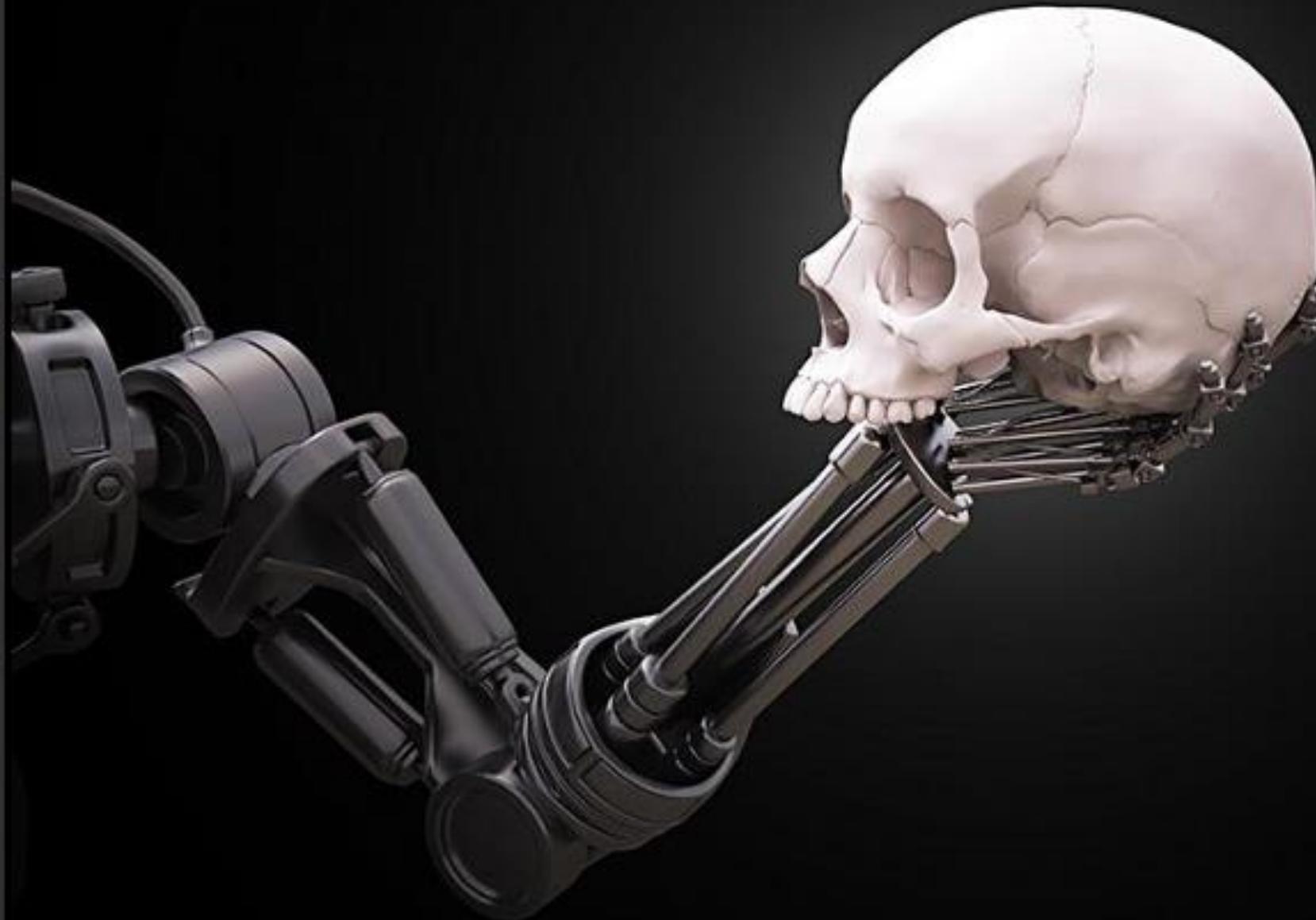
# Referências

NG, Andrew. *Machine Learning*. Stanford University, 2011. Curso oferecido via Coursera. Disponível em:  
<https://www.coursera.org/learn/machine-learning>.



Dúvidas?





Até a próxima...



Apresentador

# Thales Levi Azevedo Valente

E-mail:

[thales.l.a.valente@gmail.com](mailto:thales.l.a.valente@gmail.com)