

# **Aprendizado de Máquina**

## **Regressão**

**Slides adaptados dos Profs. Jesús Mena-Chalco, Ronaldo Prati e David Correa**

# Predição/Regressão X Classificação

- Classificação tenta determinar a qual classe um exemplo pertence, baseado em exemplos cujas classes são conhecidas, gerando um modelo que pode ser aplicado a novos exemplos
- O modelo gerado pode estar representado em diferentes maneiras (regras, árvores, grafos, vetores de suporte).
  - A saída é a classe em que o novo exemplo é predito.
- A classe em classificação é um atributo **nominal**
- E se o valor que eu quero prever é **numérico**, ao invés de um conjunto nominal de valores?
- Então o nosso problema passa a ser de **regressão** ao invés de classificação.

# Predição/Regressão

- Regressão gera uma fórmula para prever um valor numérico a partir dos dados.
- Um caso especial é prever a probabilidade de um exemplo pertencer a uma classe (e.g. regressão logística)
- O caso geral, no entanto, é **estimar/prever valores** de um atributo numérico diretamente **a partir da fórmula** a novos exemplos
- Também existem adaptações de algoritmos de classificação para o problema da regressão, como as Árvores de Regressão

# Predição/Regressão

- Por exemplo, ao invés de prever se o tempo amanhã vai ser “quente”, “ameno”, “frio” ou “gelado”, podemos prever diretamente a temperatura
  - 28,1 graus
  - 7,3 graus
  - mesmo que 28,1 ou 7,3 nunca tenham aparecido no conjunto de treinamento
- Ou o “nível de stress” de uma estrutura em diferentes condições, preço de um imóvel de acordo com suas características, o número de segundos até o próximo ônibus aparecer, ou o número de gols que a seleção irá marcar no próximo jogo, ou ... qualquer outro valor numérico que você queira prever

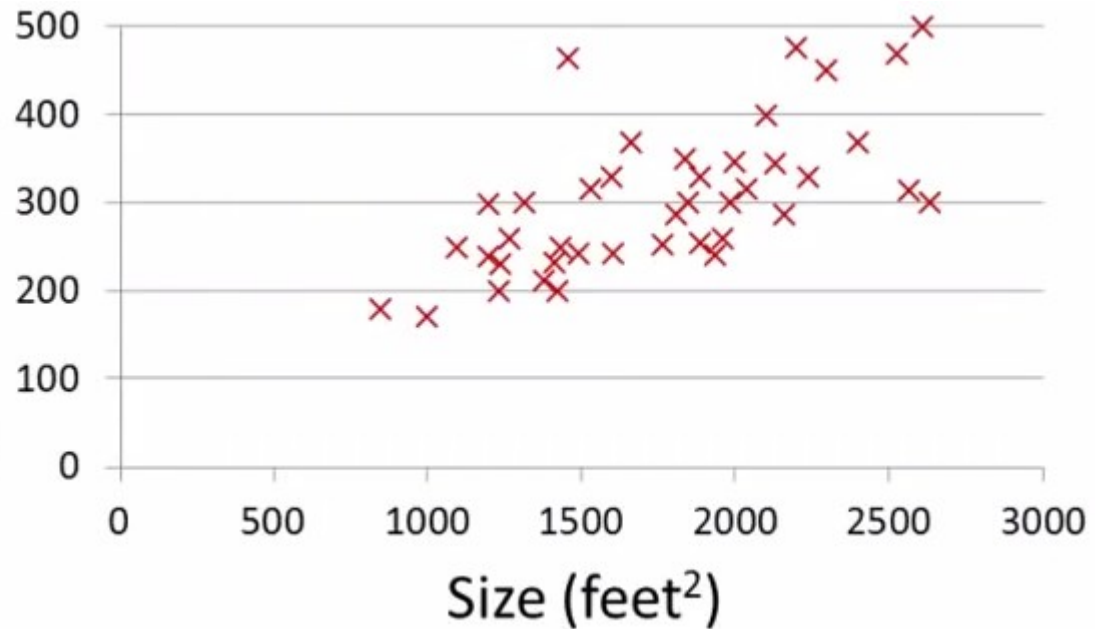
# **Regressão linear com uma variável**

## **(A) Representação do modelo**

# Representação do modelo

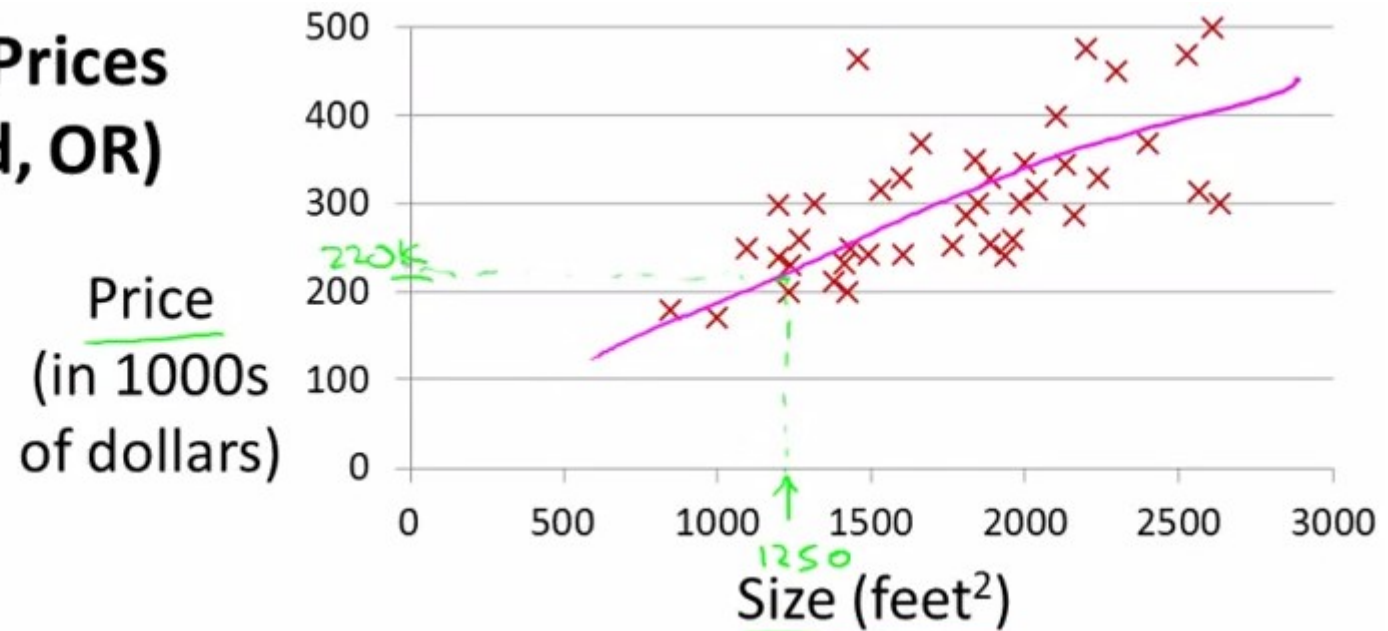
## Housing Prices (Portland, OR)

Price  
(in 1000s  
of dollars)



# Representação do modelo

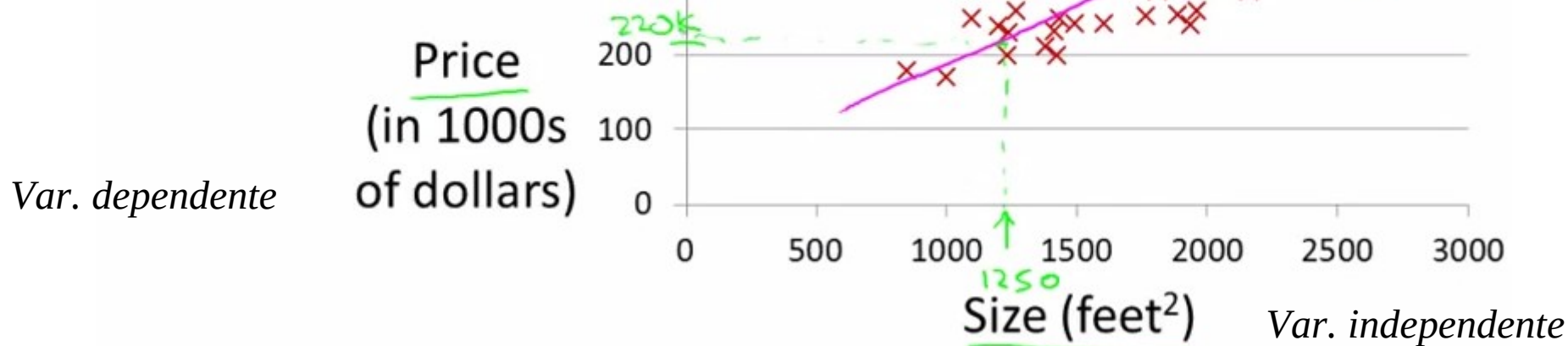
## Housing Prices (Portland, OR)



Pode se ajustar um modelo (uma reta)

# Representação do modelo

## Housing Prices (Portland, OR)



- Este é um exemplo de **Aprendizado Supervisionado**.
- Dando a “resposta correta” para cada um dos dados coletados.

- Este é um exemplo de um **Problema de Regressão** (inferência de relação)
- Usado para Inferir uma relação entre uma var. independente e outra(s) dependente(s).



# Representação do modelo

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
	...	...

## Notação:

- **m**: número de exemplos (amostras) de treinamento.
- **x**'s: variáveis de entrada / atributos / características
- **y**'s: variáveis de saída / variáveis “alvos” (*target*)

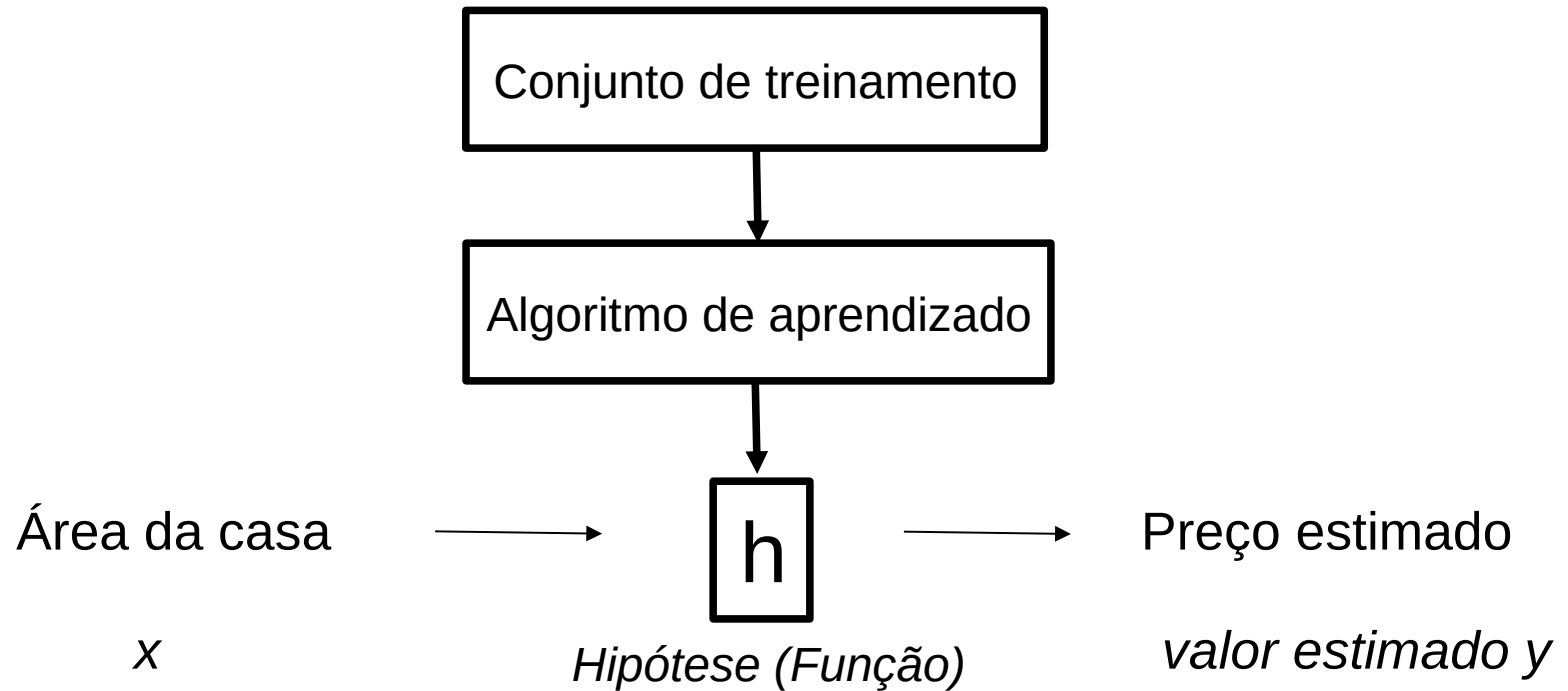
# Representação do modelo

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
	...	...

- $(x, y)$ : denota uma amostra qualquer do conjunto de treinamento.
- $(x^{(i)}, y^{(i)})$ : denota uma amostra específica (i-ésima amostra).  
Onde  $i$  é o índice.

$$y^{(3)} = 315$$

# Representação do modelo



**$h$  é uma função que mapeia  $x$  a  $y$**

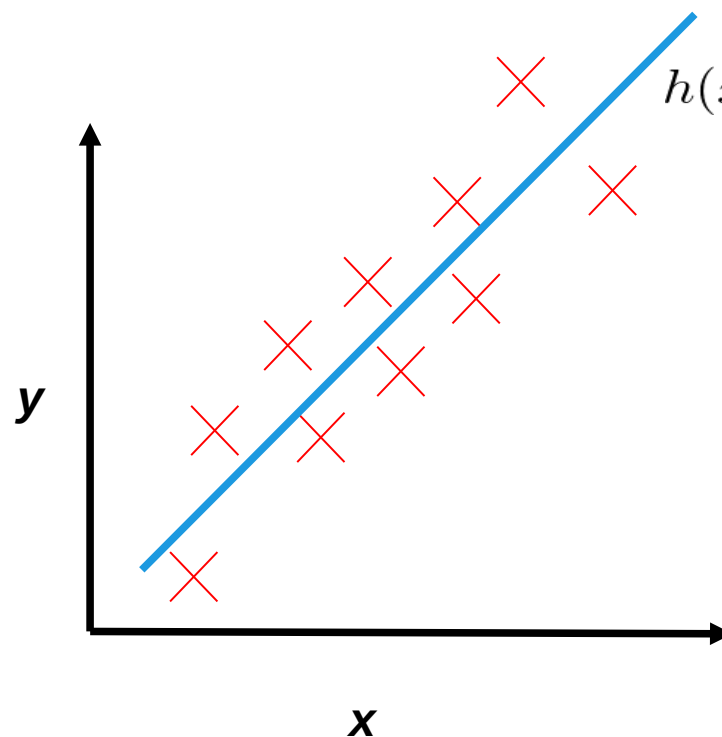
# Representação do modelo

Como representar  $h$ ?

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

$$h(x) = \theta_0 + \theta_1 \cdot x$$

← *forma abreviada*



Este modelo é uma:  
**Regressão linear com  
uma variável**

**= Regressão linear  
univariada.**

# Regressão linear

- Expressa a 'classe' (atributo alvo) como uma combinação linear dos outros atributos com alguns pesos, p.ex:

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

em que  $\theta_i$  é um peso, e  $x_i$  é um atributo

- O valor previsto para um exemplo é calculado com base nos valores dos atributos do exemplo.
  - Então é necessário **aprender os pesos que minimizam o erro** entre o valor real e o valor previsto (**função custo**) nos exemplos de treinamento.

# **Regressão linear com uma variável**

## **(B) Função custo**

# Função custo

- A função custo nos permitirá definir a melhor função utilizada para aproximar a linha reta para o conjunto de treinamento.

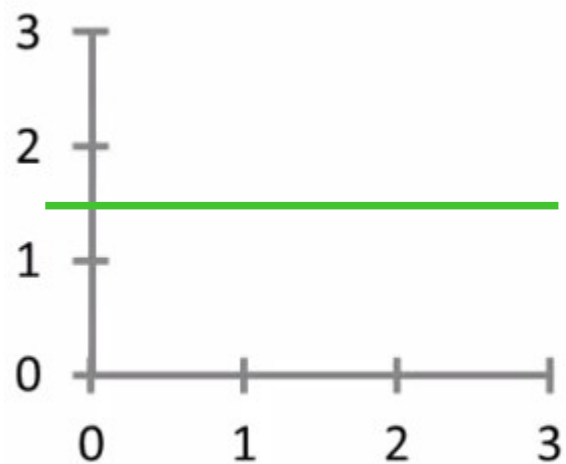
Conjunto de treinamento:	Size in feet <sup>2</sup> (x)	Price (\$) in 1000's (y)
	2104	460
	1416	232
	1534	315
	852	178
	...	...

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

$\theta_i$  são os parâmetros (pesos). Como definir esses parâmetros?

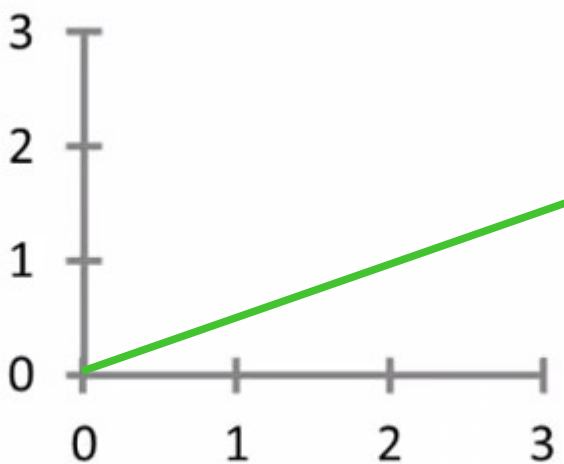
# Função custo

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



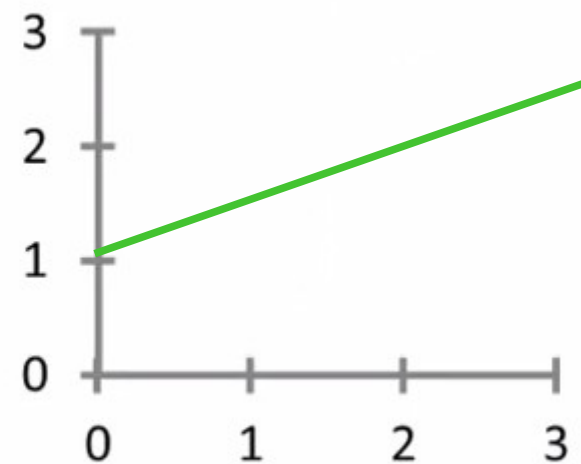
$$\theta_0 = 1.5$$
$$\theta_1 = 0$$

$$h(x) = 1.5$$



$$\theta_0 = 0$$
$$\theta_1 = 0.5$$

$$h(x) = 0.5x$$

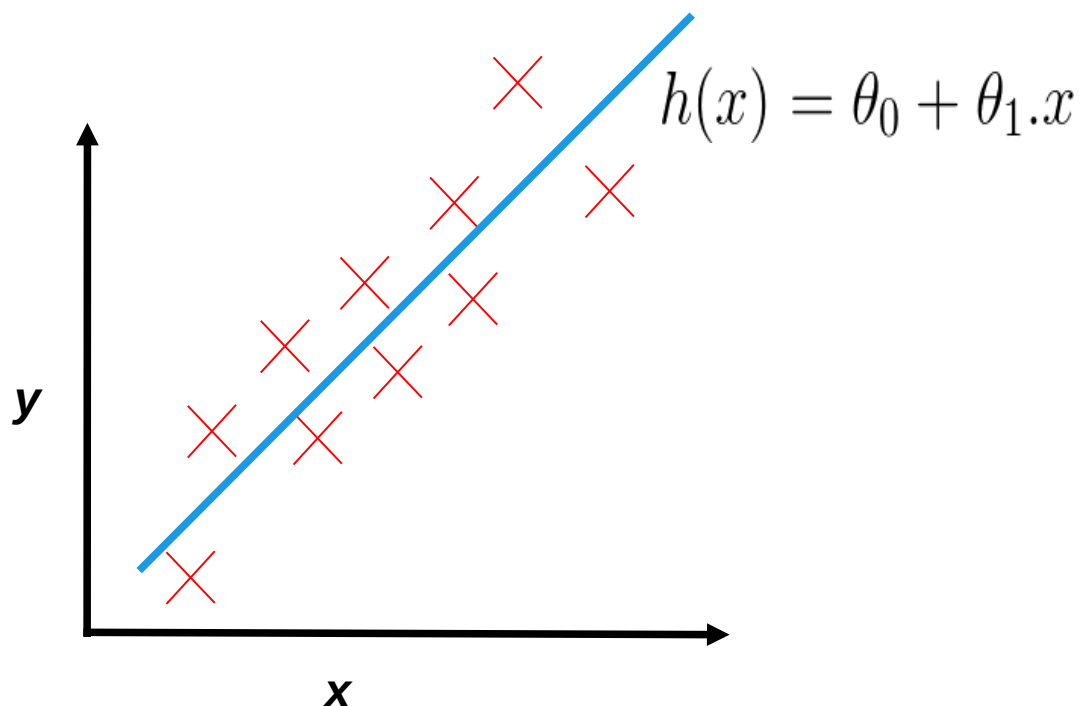


$$\theta_0 = 1$$
$$\theta_1 = 0.5$$

$$h(x) = 1 + 0.5x$$



# Função custo



- Como escolher os parâmetros ( $\theta_0$  e  $\theta_1$ ) que permitam ajustar a melhor reta aos dados do conjunto de treinamento?
  - **IDEIA:** Escolher  $\theta_0$  e  $\theta_1$  tal que  $h(x)$  **seja o mais próximo** dos valores de  $y$  de acordo com o **conjunto de treinamento**  $(x,y)$ .

# Função custo

Minimizar  $\theta_0$  e  $\theta_1$   $h_{\theta}(x^{(i)}) - y^{(i)}$

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

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

# Função custo

Minimizar  $\theta_0$  e  $\theta_1$   $h_{\theta}(x^{(i)}) - y^{(i)}$

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

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

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

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

# Função custo

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

*Objetivo: Descobrir  $\theta_0$  e  $\theta_1$  que minimize  $J(\theta_0, \theta_1)$*

Essa função custo também é conhecida como Erro Quadrático Médio (EQM).

# Função custo – intuição

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

**Parâmetros:**  $\theta_0, \theta_1$

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

**Objetivo:** Descobrir  $\theta_0$  e  $\theta_1$  que minimize  $J(\theta_0, \theta_1)$

# Função custo – intuição

**Hipótese *simplificada*:**  $h_{\theta}(x) = \theta_1 x$   
 $\theta_0 = 0$

**Parâmetros:**  $\theta_1$

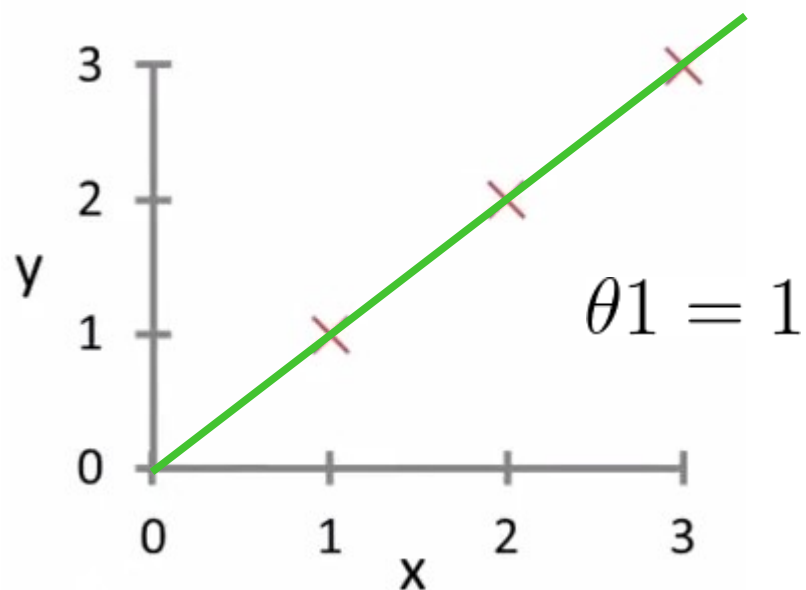
**Função custo:**  $J(\theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$

**Objetivo:** Descobrir  $\theta_1$  que minimize  $J(\theta_1)$

# Função custo – intuição

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

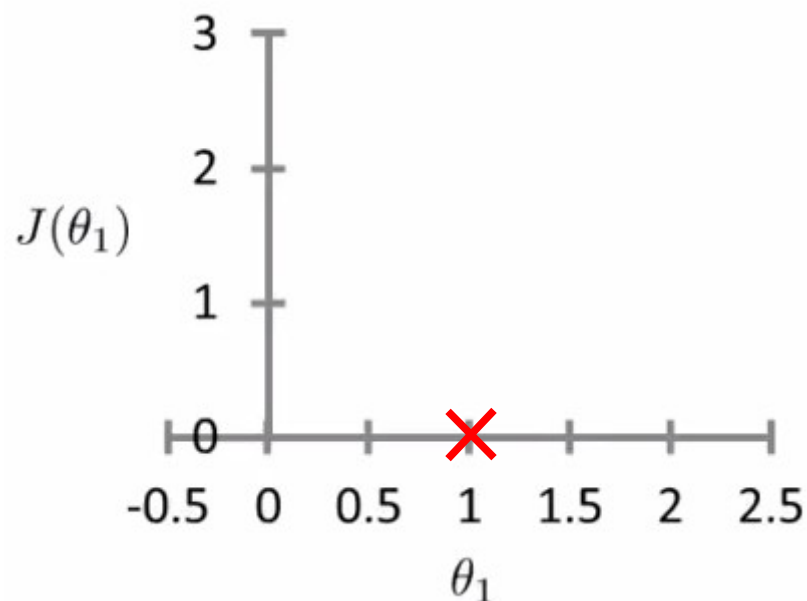
(para um  $\theta_1$ , é uma função de  $x$ )



$$J(\theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$
$$= \frac{1}{m} (0^2 + 0^2 + 0^2)$$

$$J(\theta_1)$$

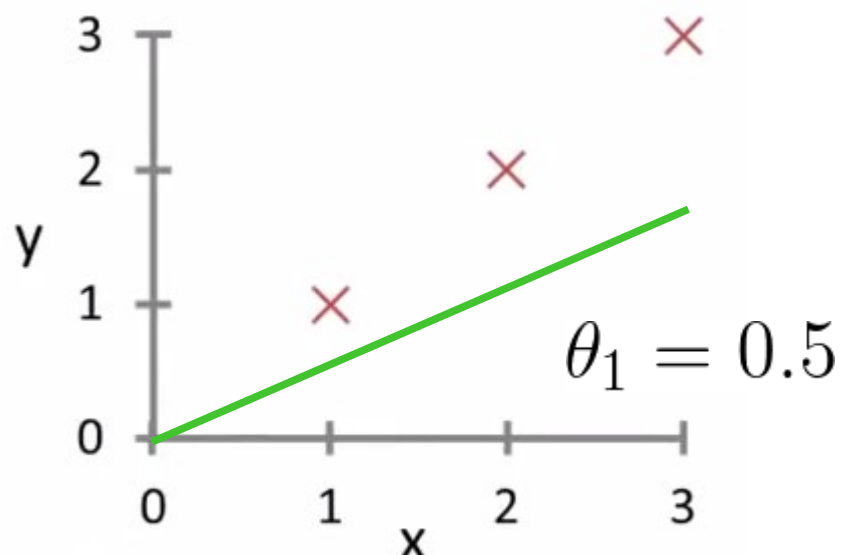
(função com parâmetro  $\theta_1$ )



# Função custo – intuição

$h_{\theta}(x)$

(para um  $\theta_1$ , é uma função de  $x$ )

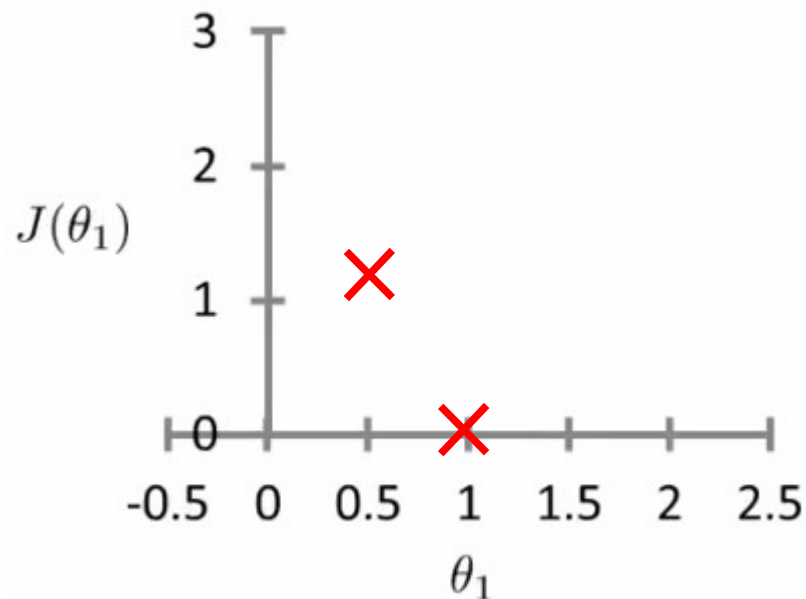


$$J(\theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

$$J(0.5) = \frac{1}{3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \approx 1.17$$

$J(\theta_1)$

(função com parâmetro  $\theta_1$ )

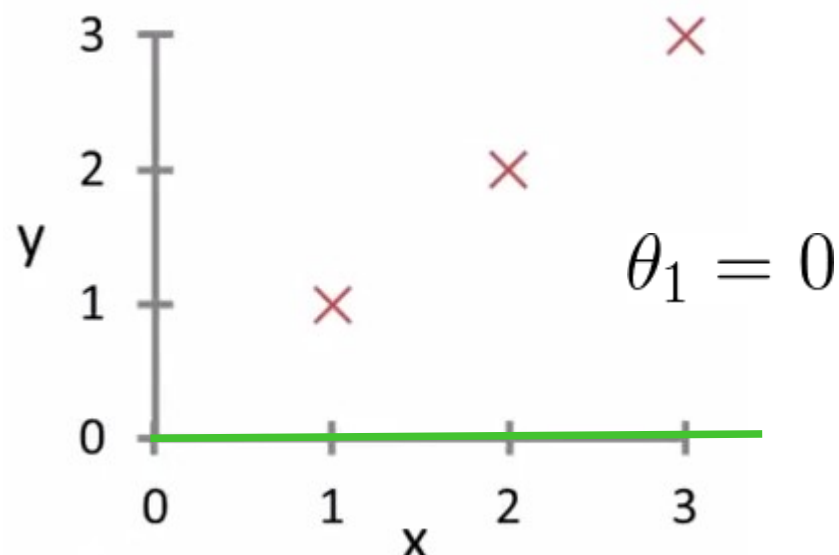




# Função custo – intuição

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

(para um  $\theta_1$ , é uma função de  $x$ )

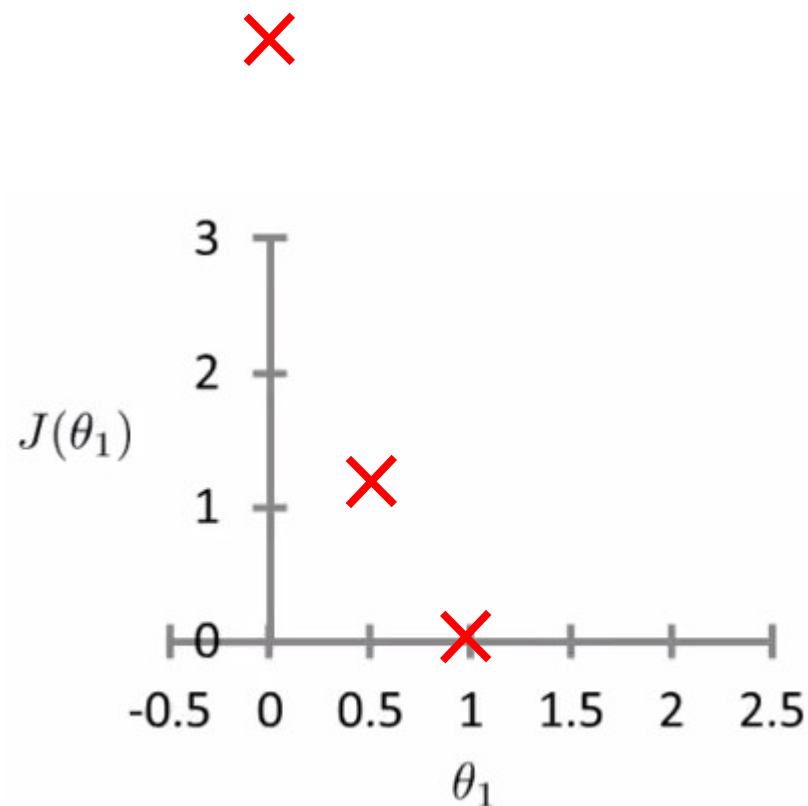


$$J(\theta_1) = \frac{1}{m} \sum_{i=1}^m (\theta_1 x^{(i)} - y^{(i)})^2$$

$$J(0.5) = \frac{1}{3} (1^2 + 2^2 + 3^2) \approx 4.67$$

$$J(\theta_1)$$

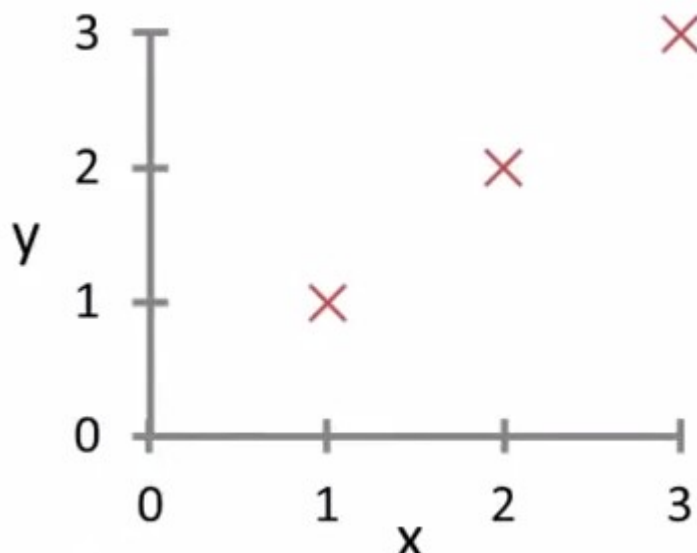
(função com parâmetro  $\theta_1$ )



# Função custo – intuição

$h_{\theta}(x)$

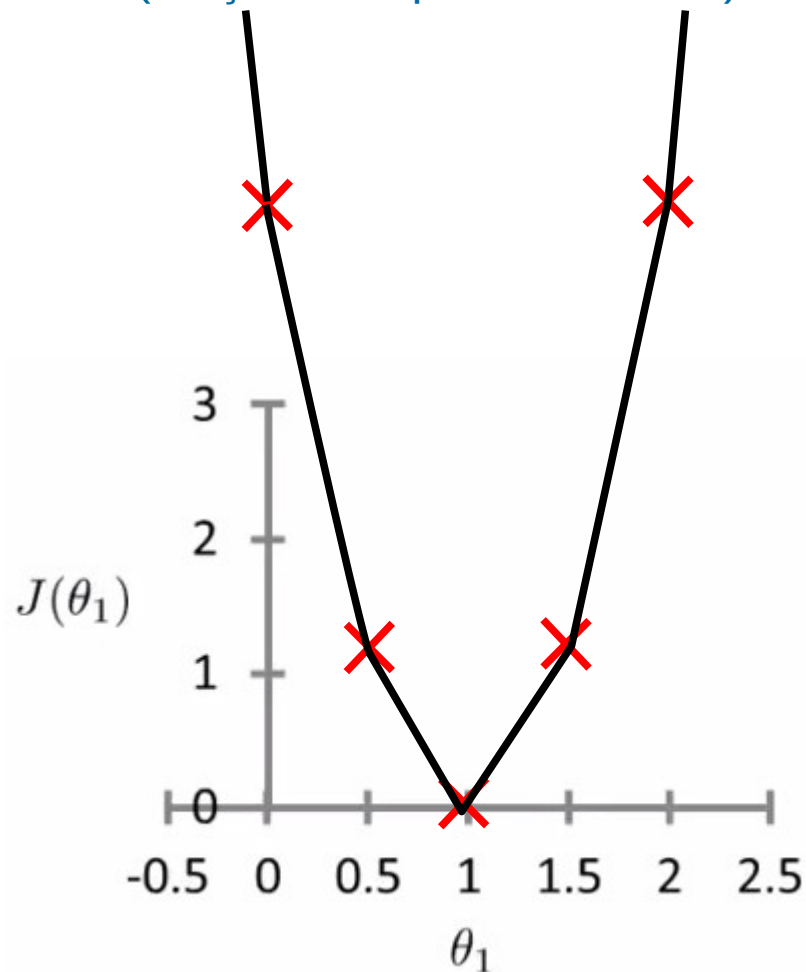
(para um  $\theta_1$ , é uma função de  $x$ )



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

$J(\theta_1)$

(função com parâmetro  $\theta_1$ )



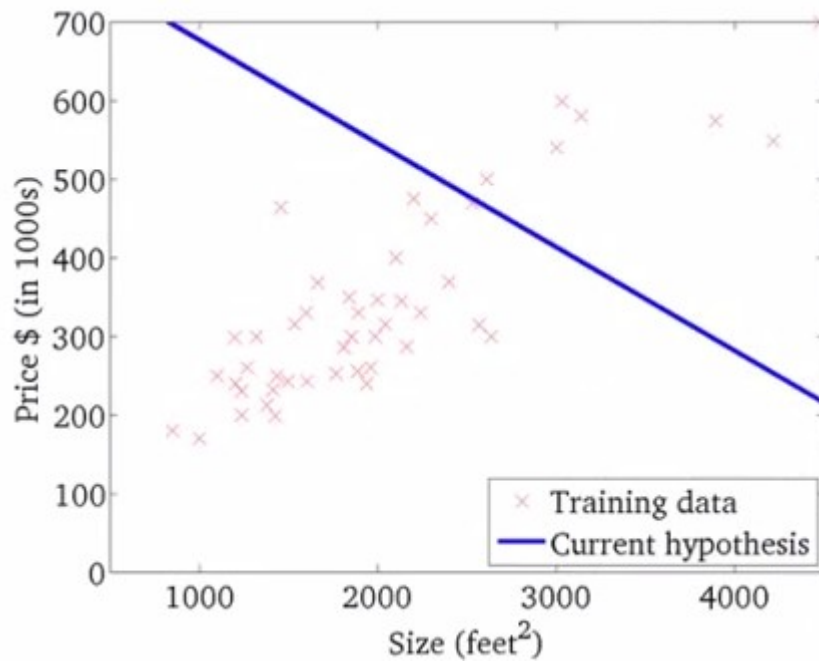
# Como buscar os parâmetros $\theta$ ?

- Como buscar os parâmetros  $\theta$  que minimizam a função custo?
  - **Função Gradiente Descendente**

# Função custo

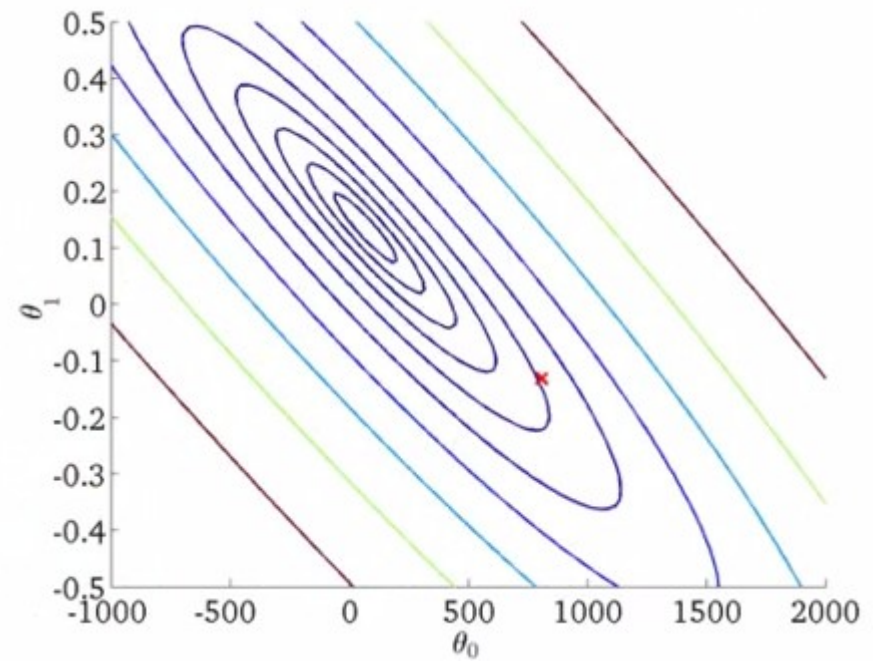
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

(função com parâmetro  $\theta_0, \theta_1$ )

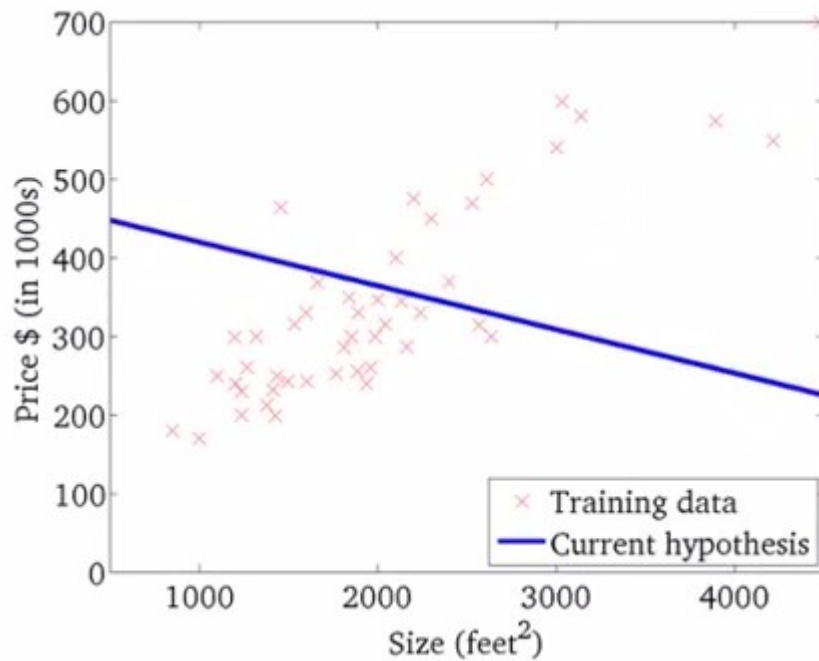


*bacia*

# Função custo

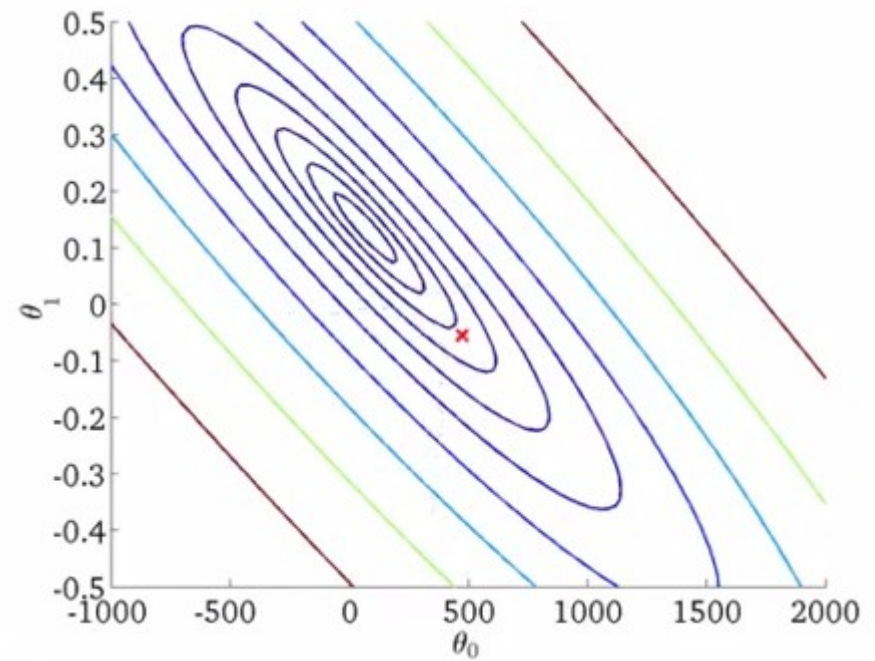
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

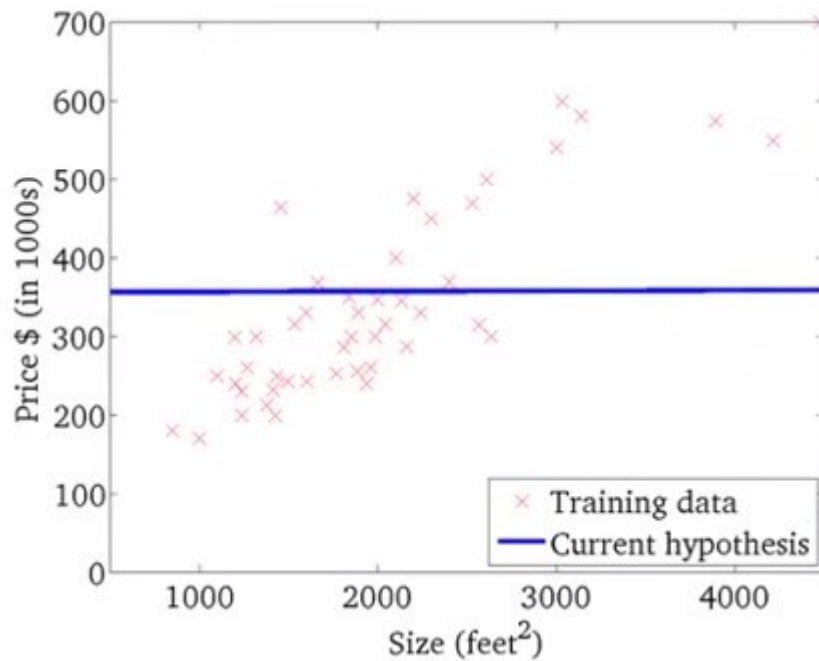
(função com parâmetro  $\theta_0, \theta_1$ )



# Função custo

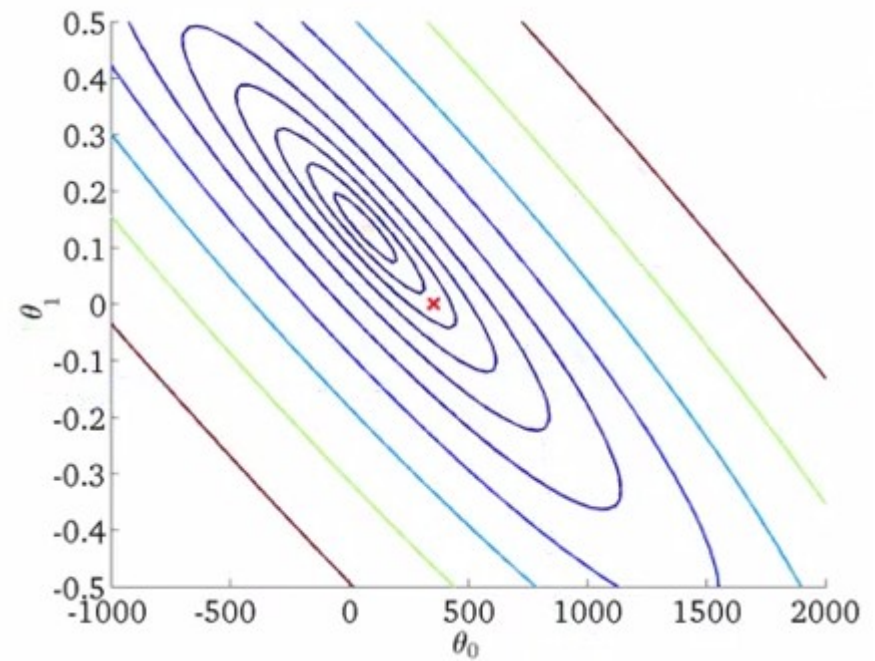
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

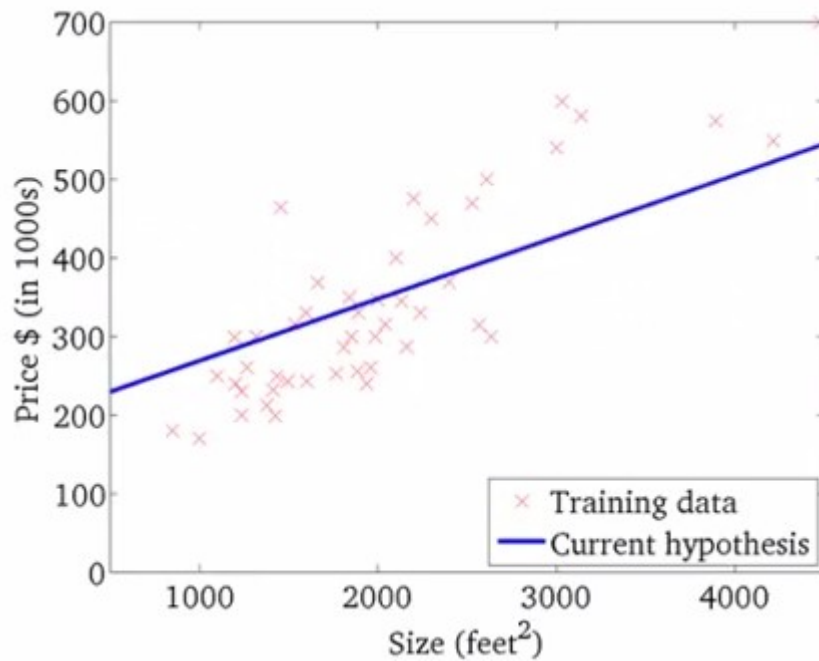
(função com parâmetro  $\theta_0, \theta_1$ )



# Função custo

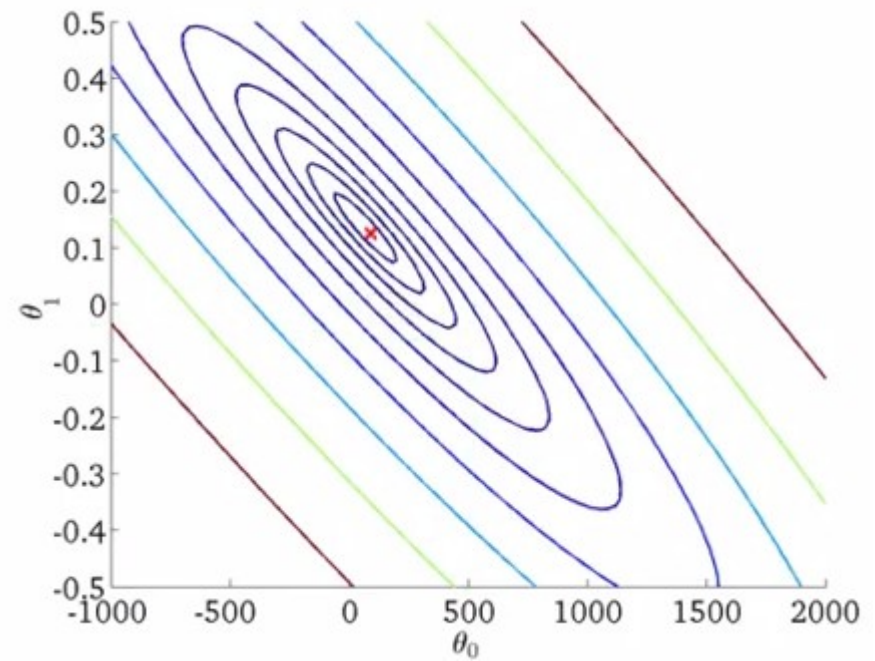
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente

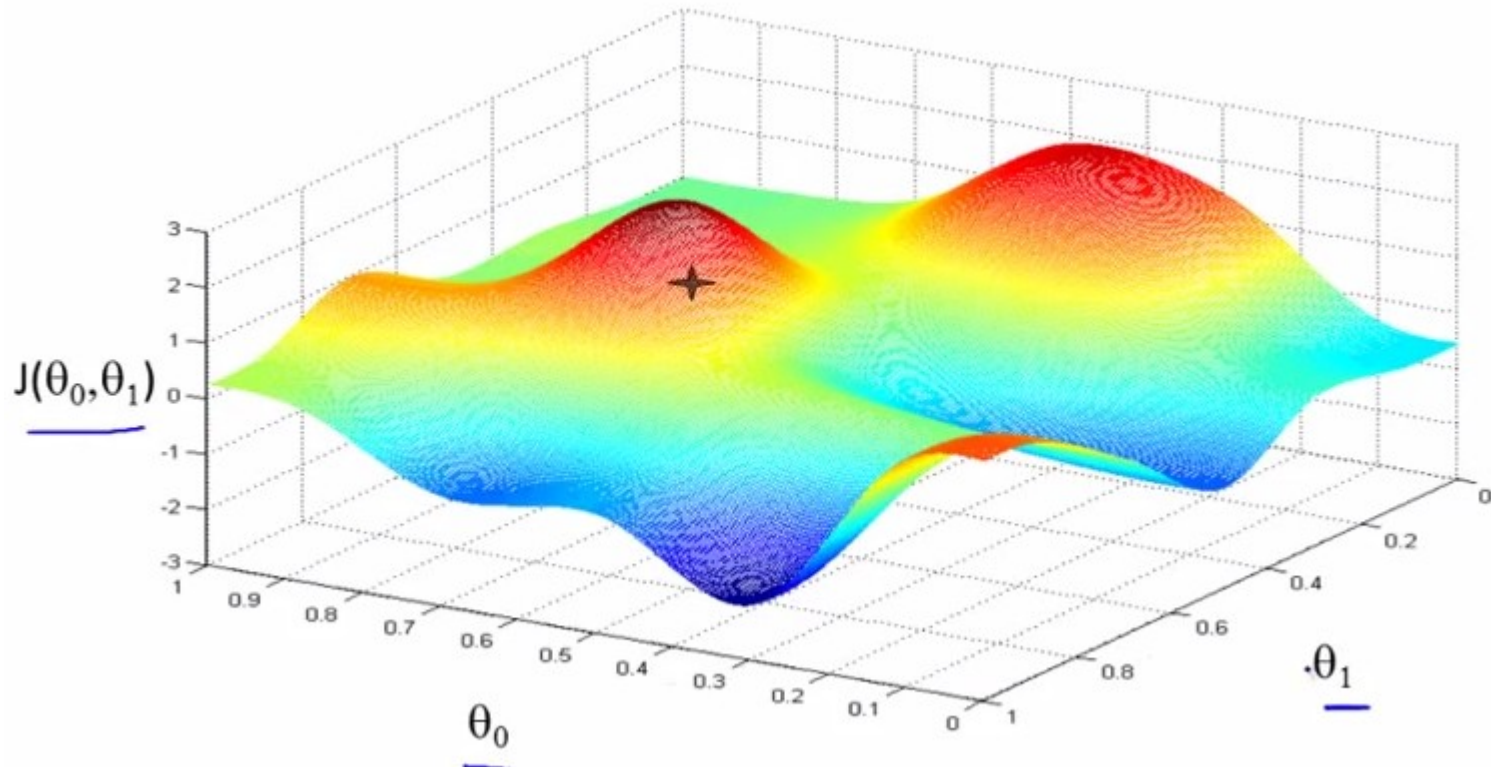
- Temos uma função  $J(\theta_0, \theta_1)$
- Desejamos descobrir  $\theta_0$  e  $\theta_1$  que minimize  $J(\theta_0, \theta_1)$
- **Ideia:**
  - Iniciar com alguns valores para  $\theta_0$  e  $\theta_1$
  - Modifique os valores de  $\theta_0$  e  $\theta_1$  para reduzir  $J(\theta_0, \theta_1)$  até atingir um valor mínimo.



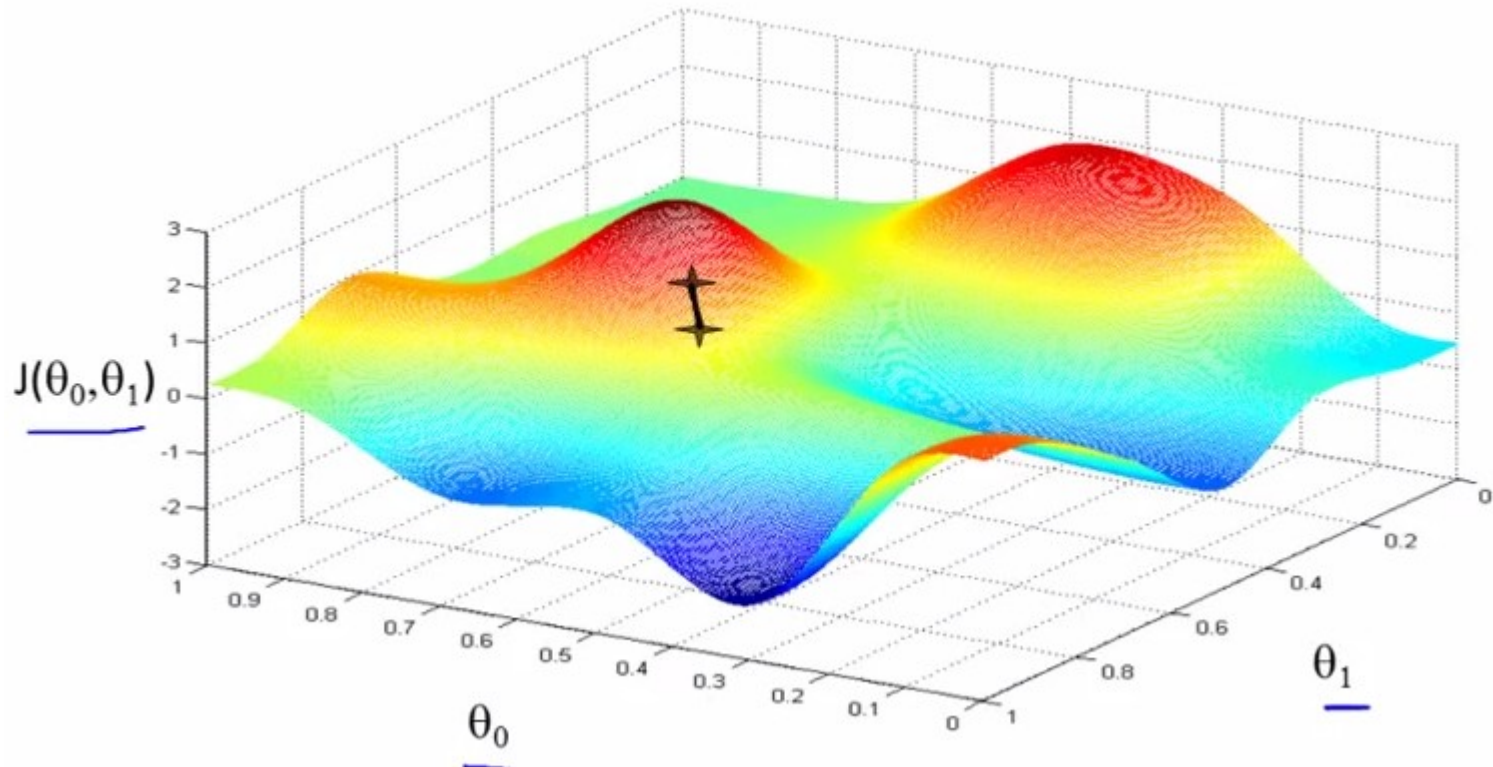
# Gradiente descendente

- Generalização para  $n$  dimensões
- Temos uma função  $J(\theta_0, \theta_1, \dots, \theta_n)$
- Desejamos descobrir  $\theta_0, \theta_1, \dots, \theta_n$  que minimize  $J(\theta_0, \theta_1, \dots, \theta_n)$
- **Ideia:**
  - Iniciar com alguns valores para  $\theta_0, \theta_1, \dots, \theta_n$
  - Modifique os valores de  $\theta_0, \theta_1, \dots, \theta_n$  para reduzir  $J(\theta_0, \theta_1, \dots, \theta_n)$  até atingir um valor mínimo.

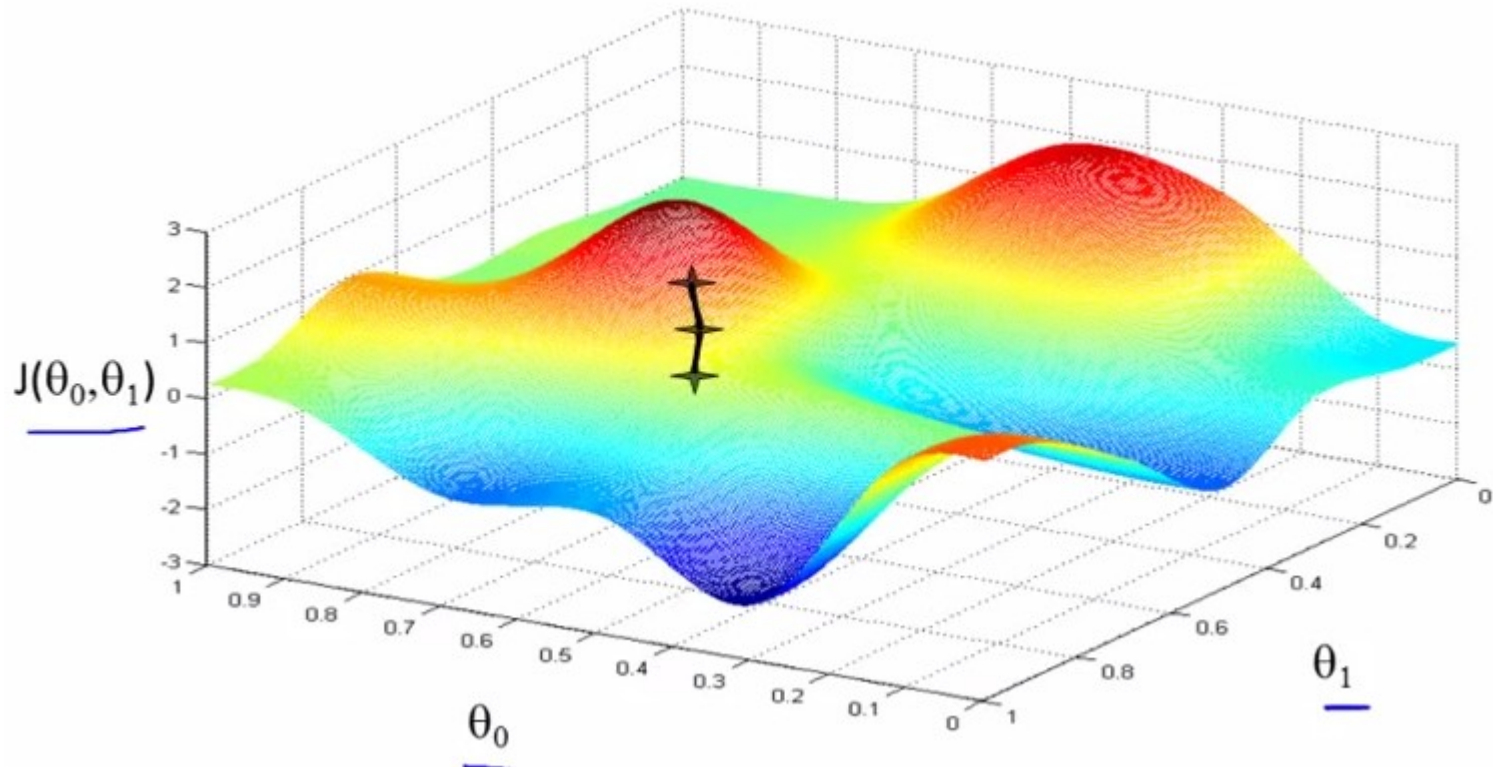
# Gradiente descendente



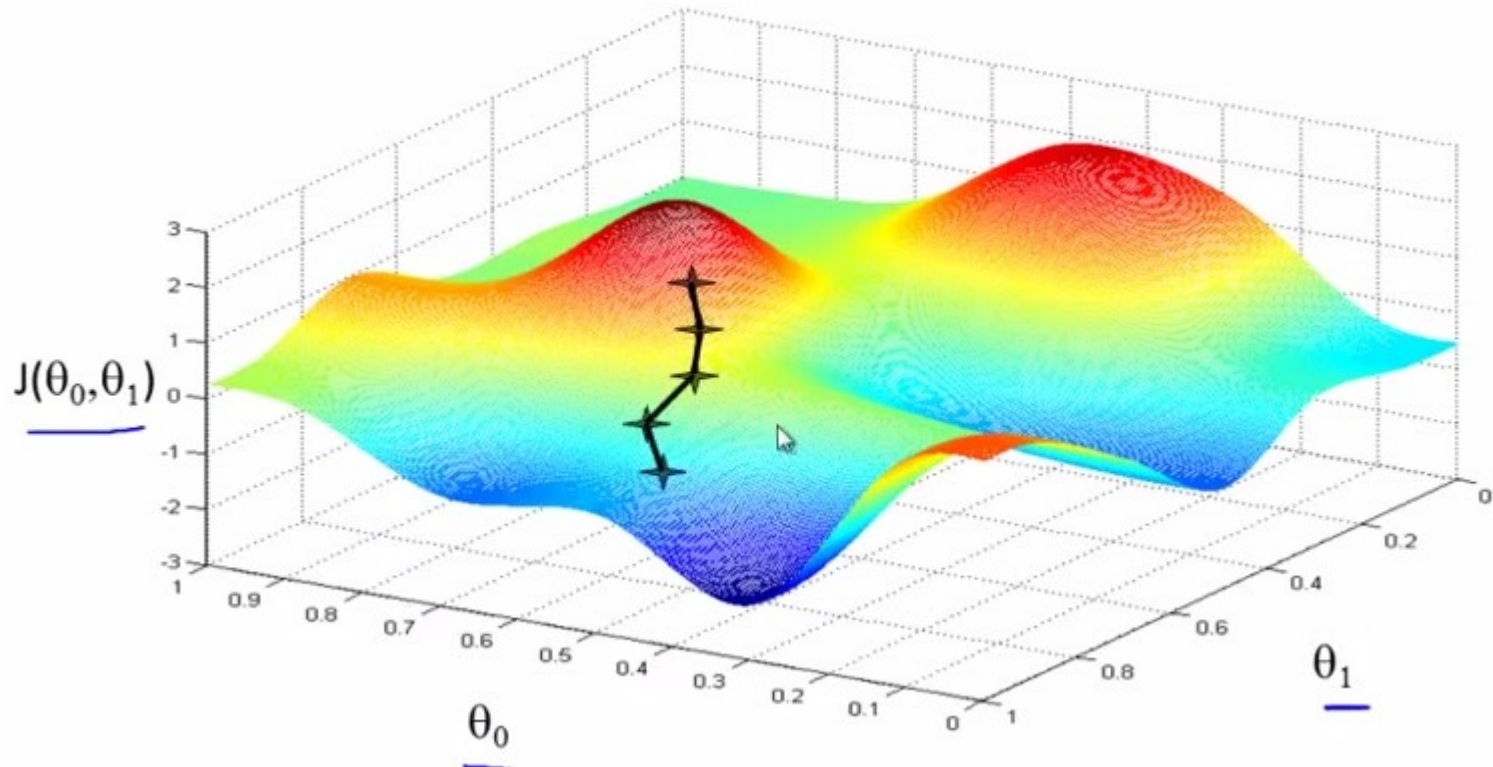
# Gradiente descendente



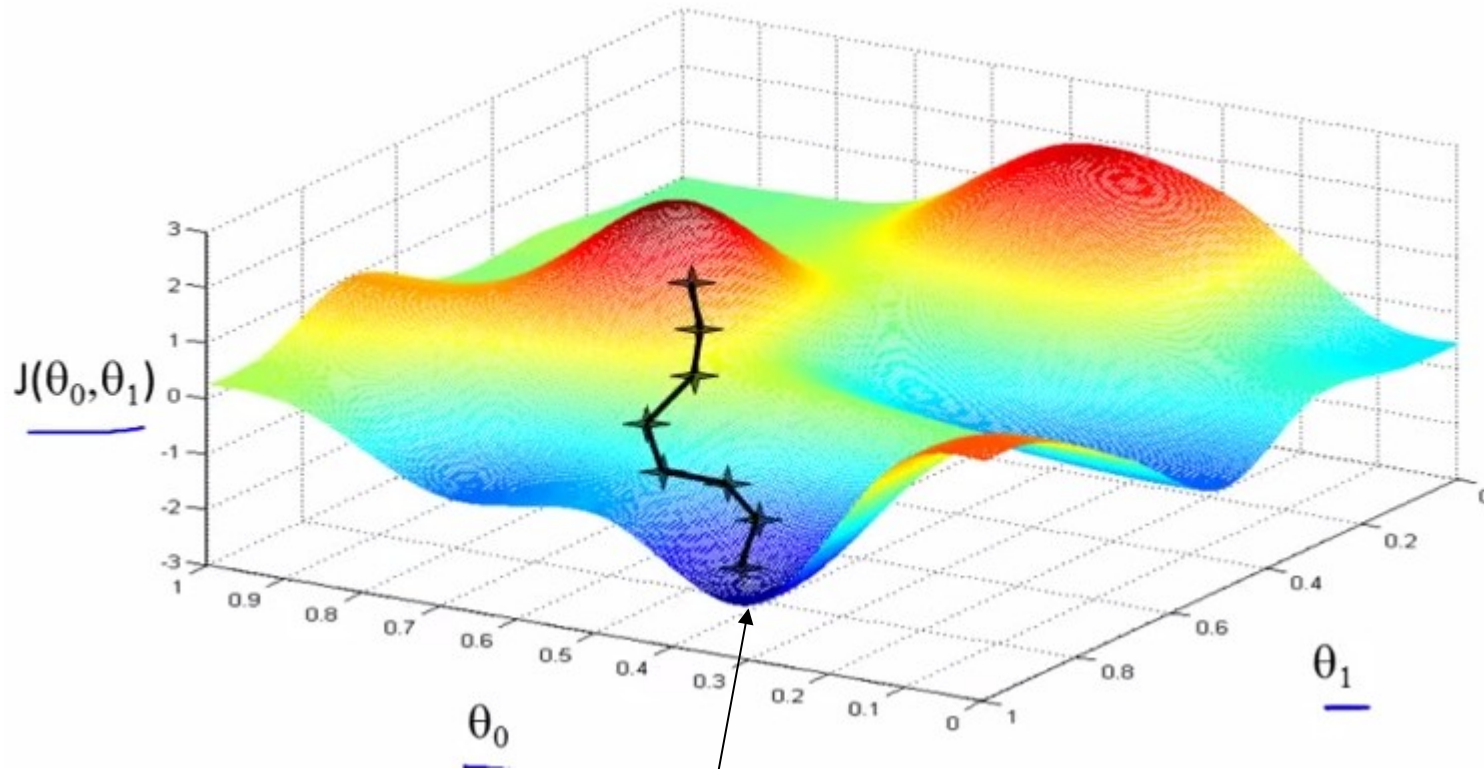
# Gradiente descendente



# Gradiente descendente



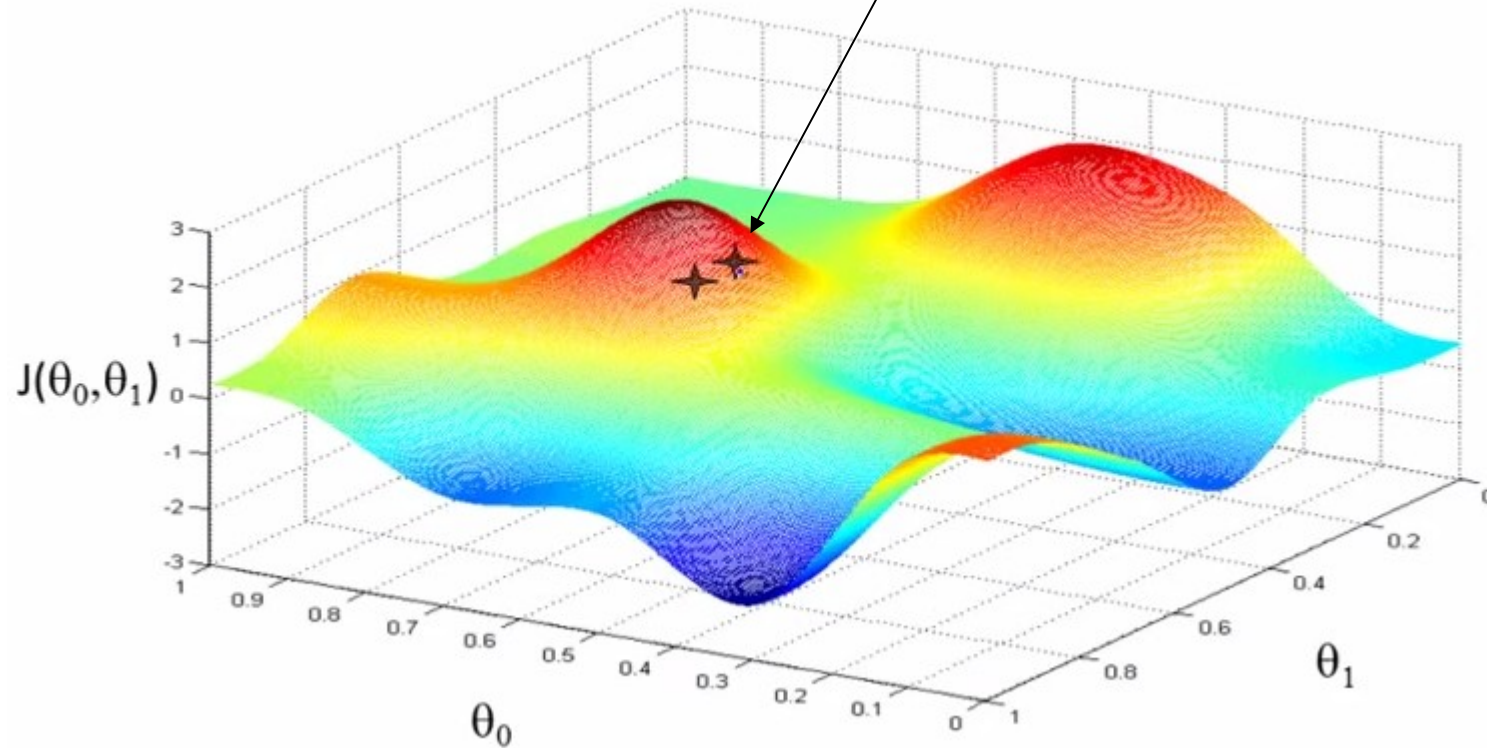
# Gradiente descendente



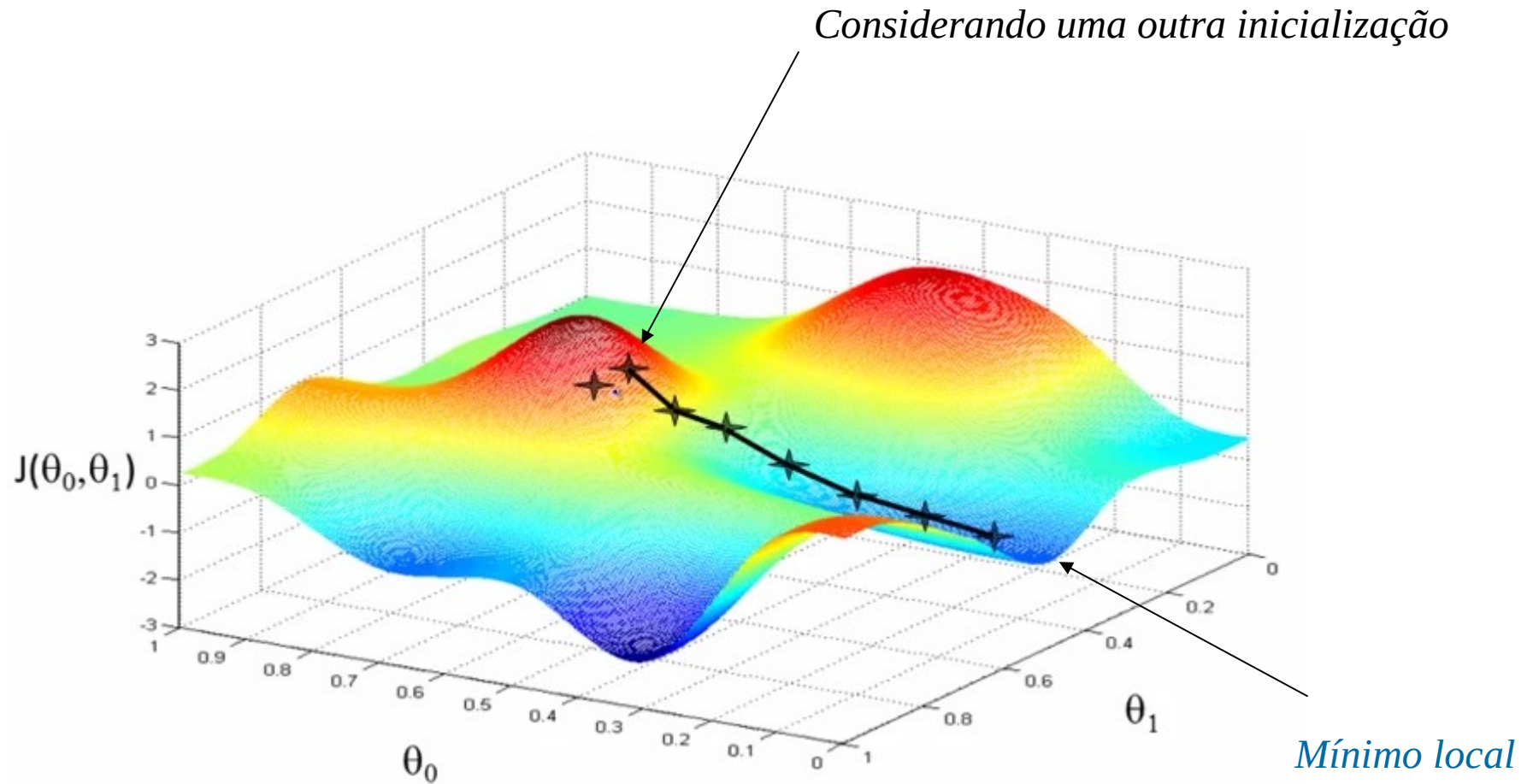
*Mínimo local*

# Gradiente descendente

*Considerando uma outra inicialização*



# Gradiente descendente





# Algoritmo gradiente descendente

- Repetir até a convergência:

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

- onde  $\alpha$  é uma constante que indica a taxa de aprendizagem
  - controla o tamanho do passo (step)
- A implementação correta considera atualização simultânea:

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

*Atribuição correta*

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

*Atribuição incorreta*

# Algoritmo gradiente descendente

Repetir até a convergência:

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

- Nos próximos exemplos apenas consideraremos a versão simplificada da função custo:

Minimizar  $\theta_1$  de  $J(\theta_1)$

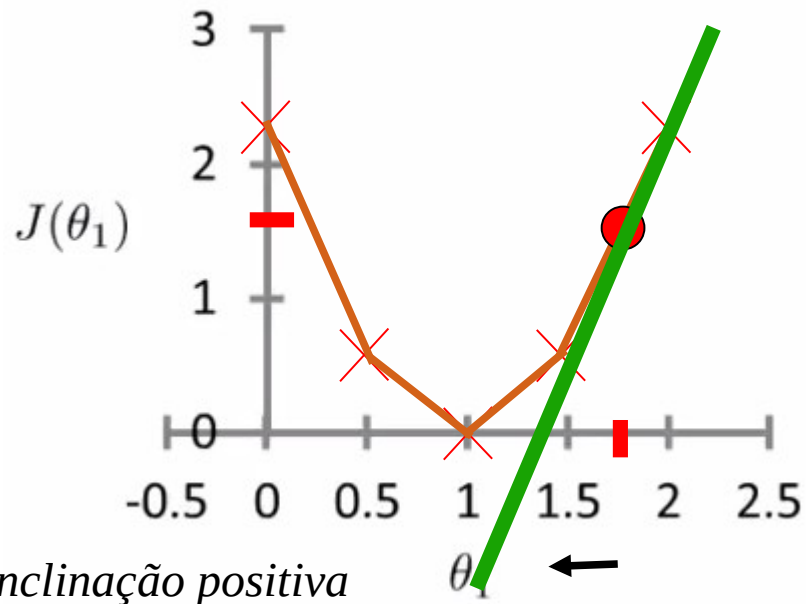
$$\theta_1 \in \mathbb{R}$$

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

# Algoritmo gradiente descendente

$J(\theta_1)$

(função com parâmetro  $\theta_1$ )



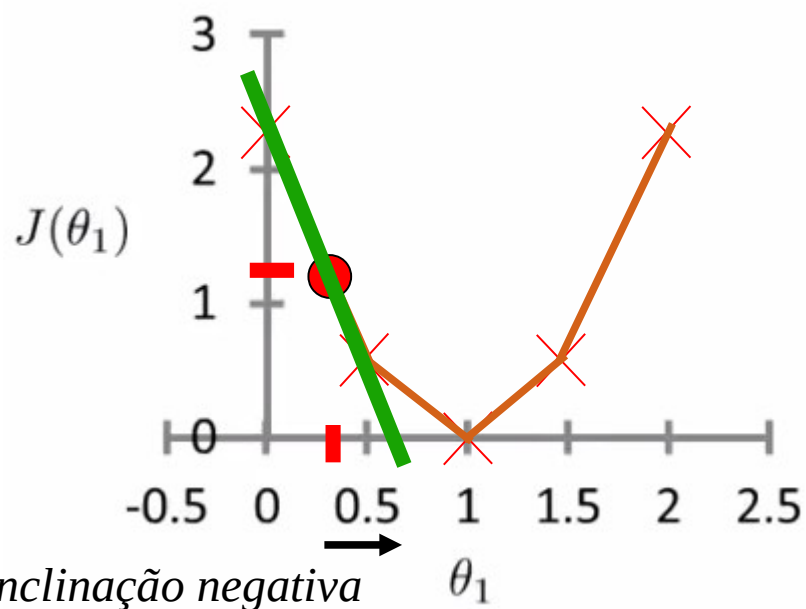
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha(\text{um Valor Positivo})$$

# Algoritmo gradiente descendente

$$J(\theta_1)$$

(função com parâmetro  $\theta_1$ )



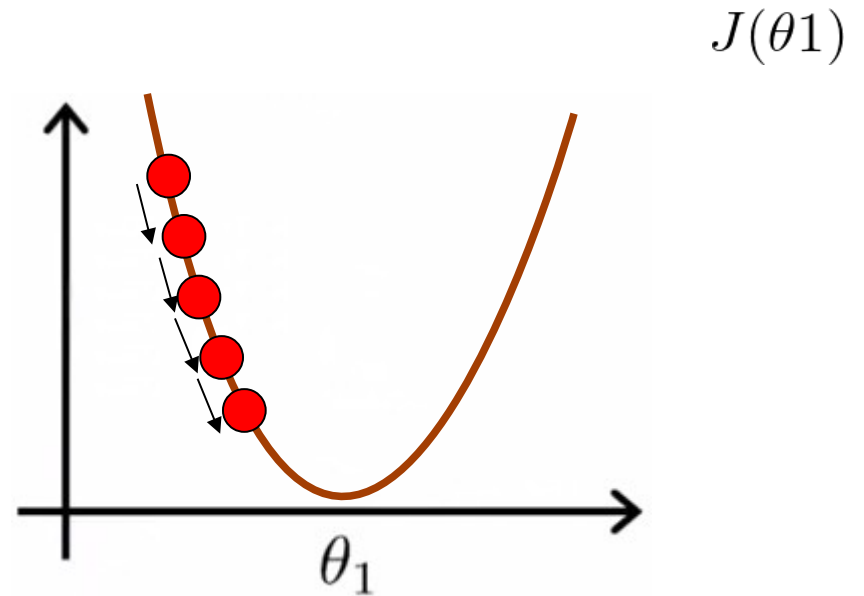
$$\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$$

$$\theta_1 := \theta_1 - \alpha(\text{um Valor Negativo})$$

# Algoritmo gradiente descendente

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

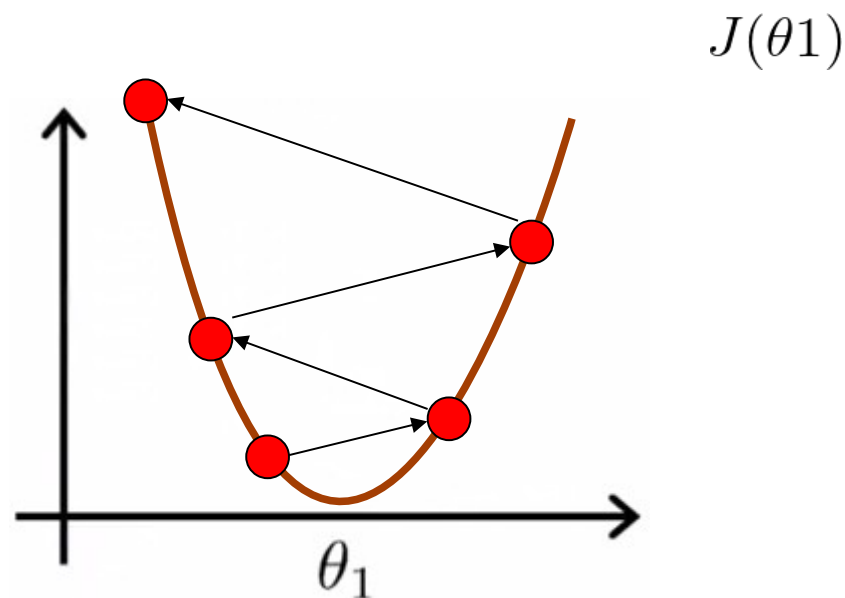
- Se  $\alpha$  for **muito pequeno** o gradiente descendente pode ser **lento**.



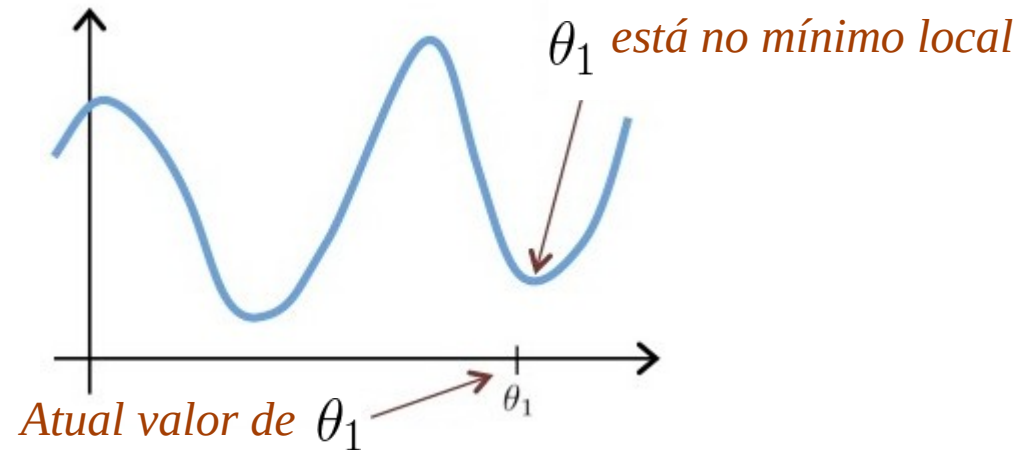
# Algoritmo gradiente descendente

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

- Se  $\alpha$  for **muito grande** o gradiente descendente pode **ultrapassar o mínimo**
- Pode falhar na convergência, ou até pode **divergir**.

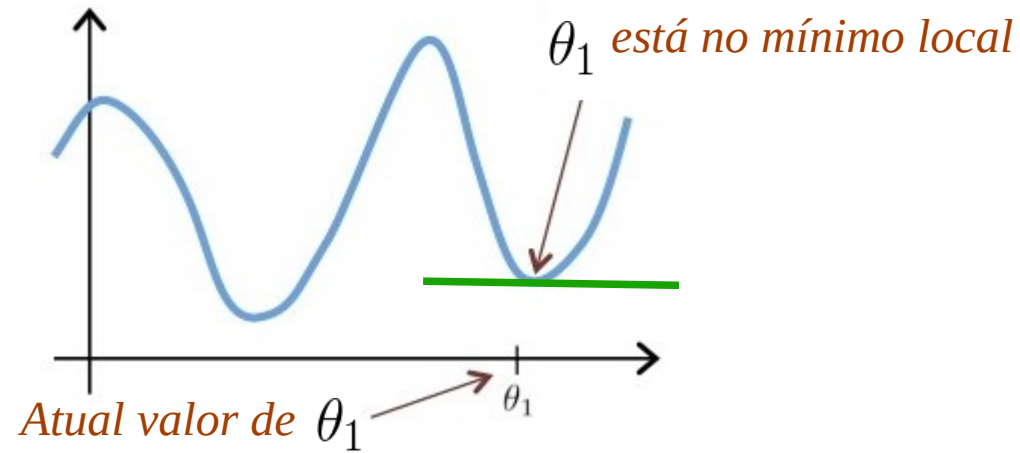


Suponha  $\theta_1$  que é um ótimo local de  $J(\theta_1)$ , como apresentado na figura.



**O seguinte passo do gradiente descendente:**  $\theta_1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_1)$

- (A) Diminuirá o valor de  $\theta_1$ .
- (B) Deixará  $\theta_1$  sem modificação.
- (C) Moverá  $\theta_1$  em uma direção aleatória.
- (D) Moverá  $\theta_1$  na direção do mínimo global de  $J(\theta_1)$ .



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

$$\theta_1 := \theta_1 - \alpha(0)$$

**Resposta Correta (B)**



# Gradiente descendente + Regressão linear

- Algoritmo gradiente descendente

Repetir até a convergência:

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

para  $j=0$  e  $j=1$

- Modelo de regressão linear

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

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

•

# Gradiente descendente + Regressão linear

•

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

$$= \frac{\partial}{\partial \theta_j} \cdot \frac{1}{2m} \sum (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

$$\frac{\partial(x^n)}{\partial x} = n \cdot x^{(n-1)} \frac{\partial x}{\partial x} = n \cdot x^{(n-1)}$$

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

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

# Gradiente descendente + Regressão linear

- Algoritmo gradiente descendente

Repetir até a convergência:

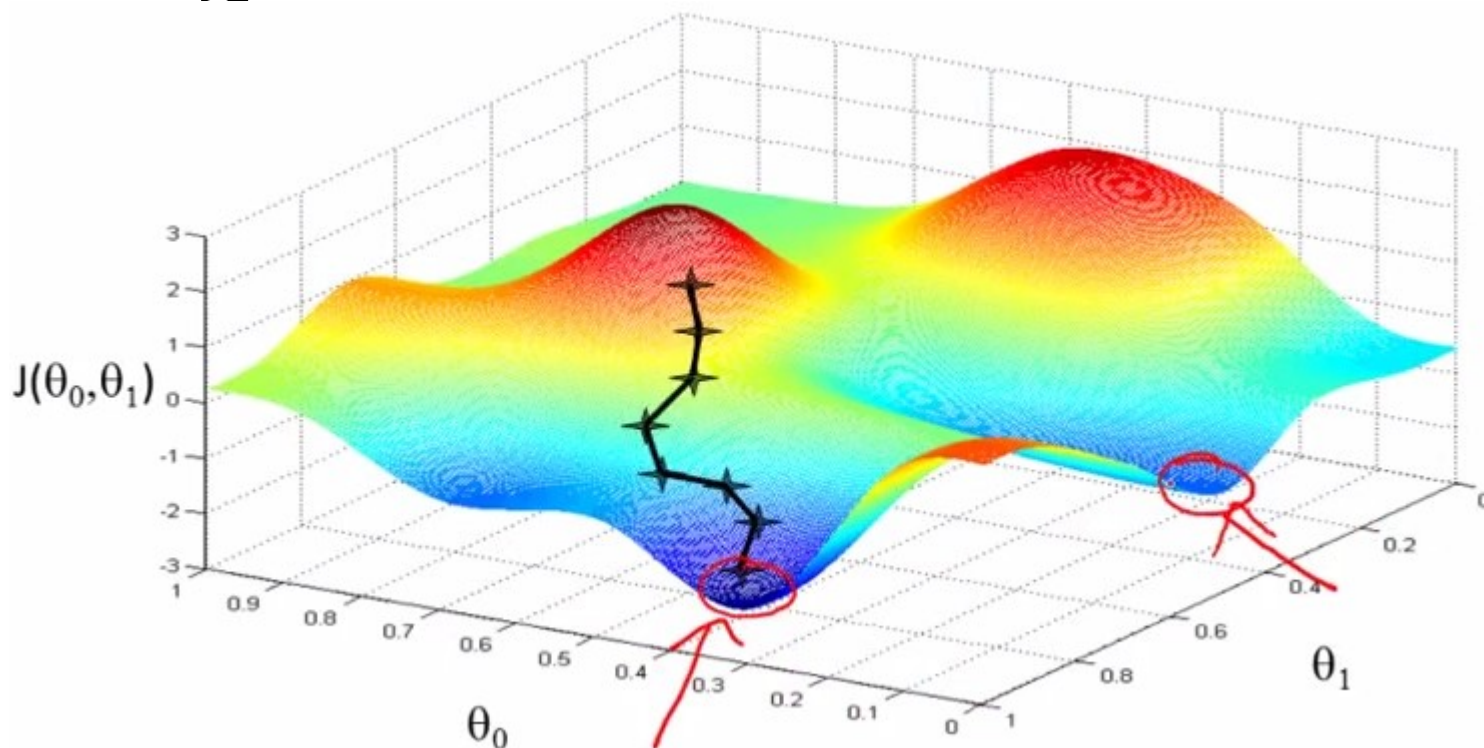
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

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

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

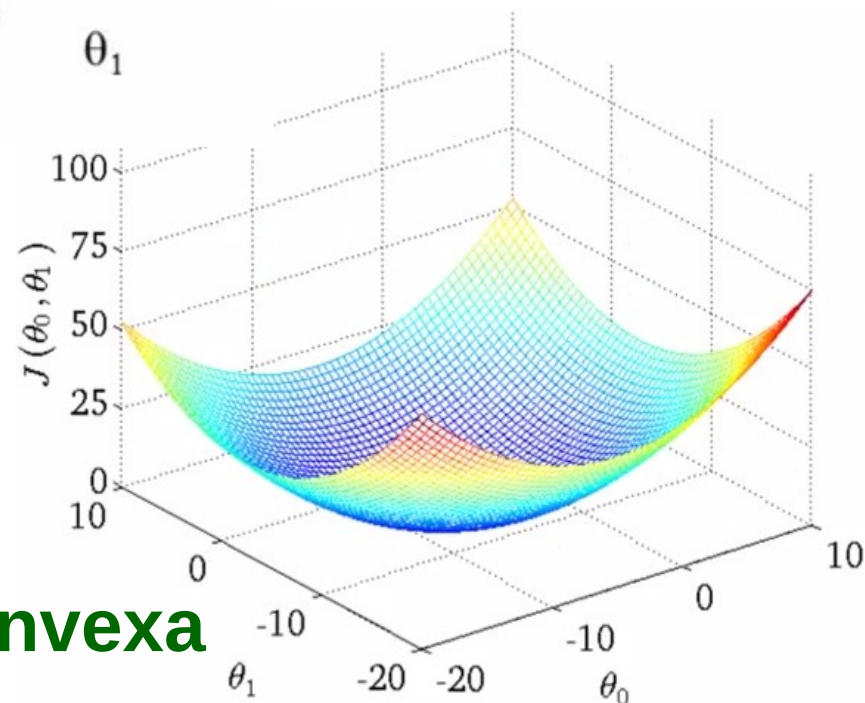
**Repetir  
simultaneamente**

# Gradiente descendente + Regressão linear



Não-convexa

- A função custo para uma **regressão linear**, sempre será uma função na forma de **parabola**.
  - função **convexa**
  - gradiente descendente garante o **ótimo global**

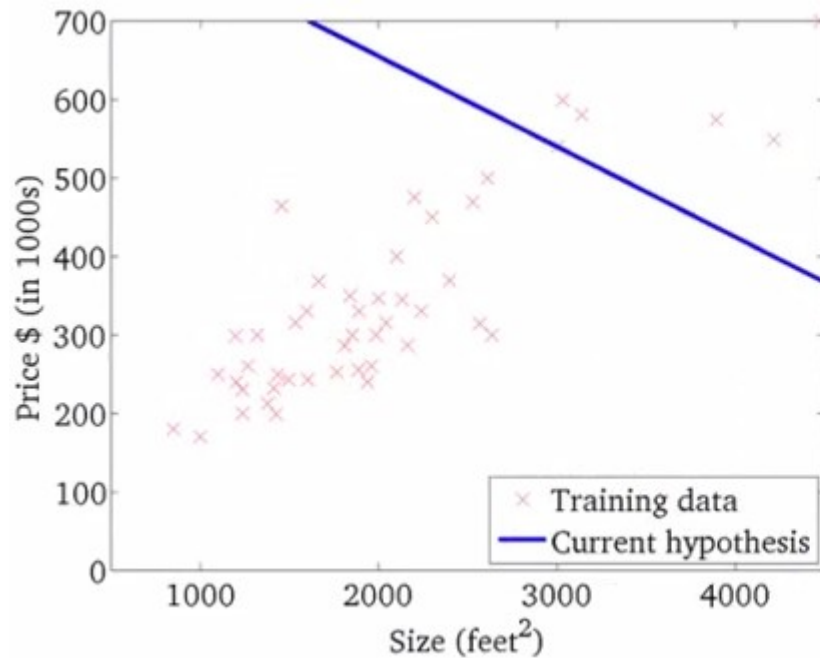


Convexa

# Gradiente descendente + Regressão linear

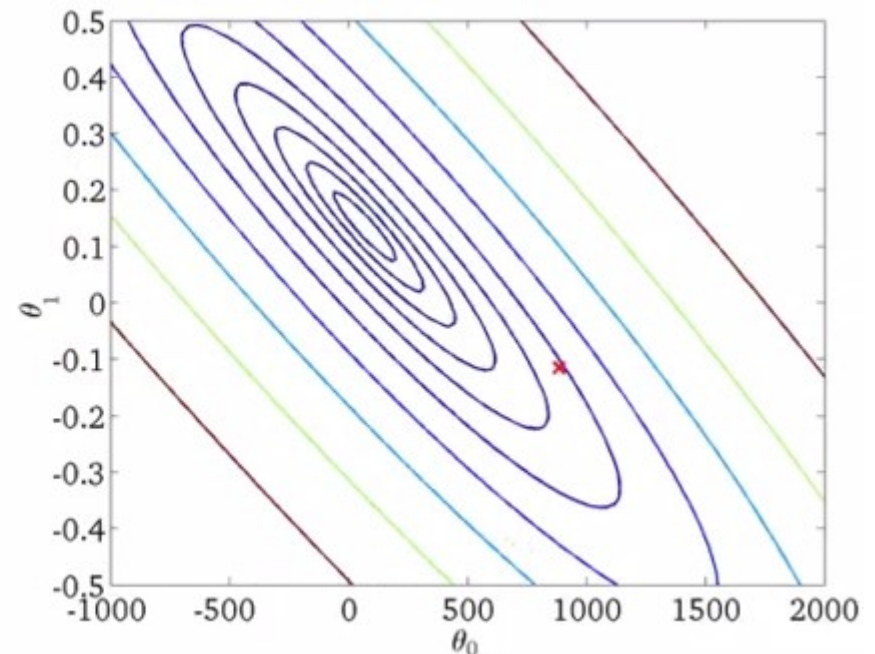
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

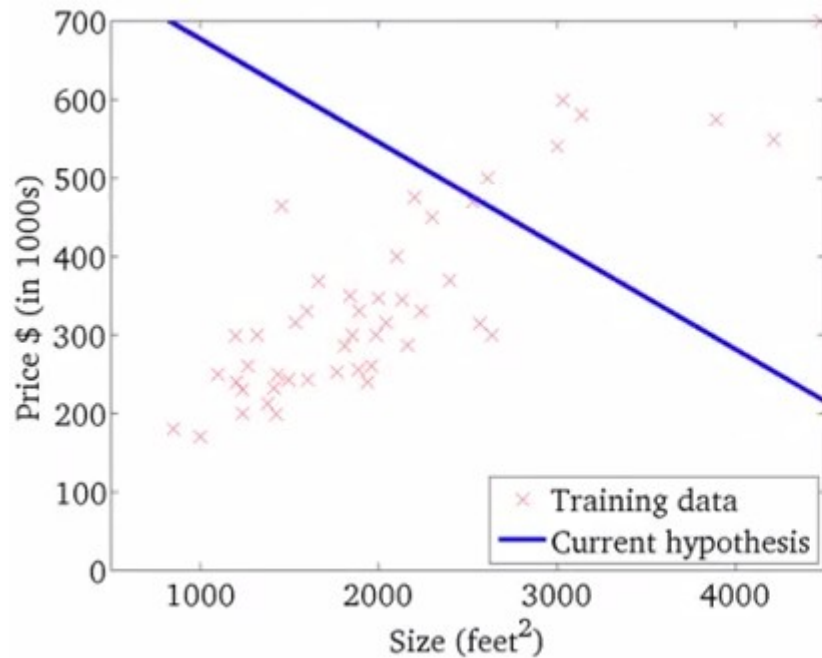
(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente + Regressão linear

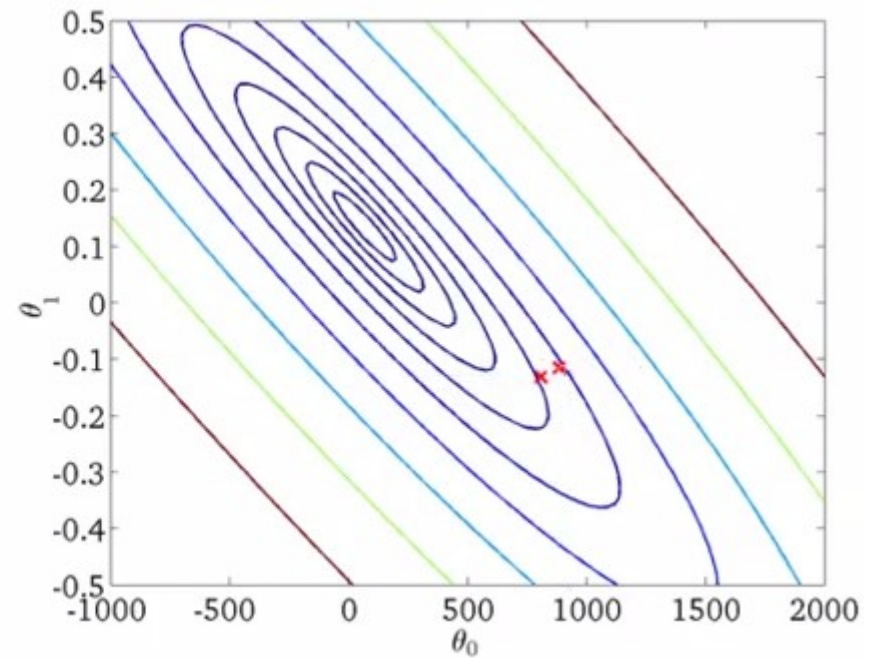
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

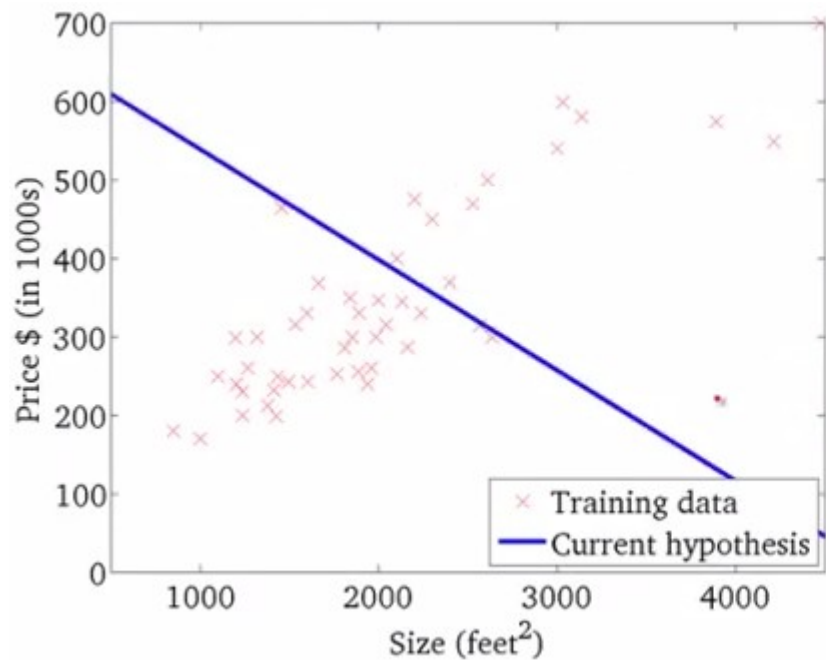
(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente + Regressão linear

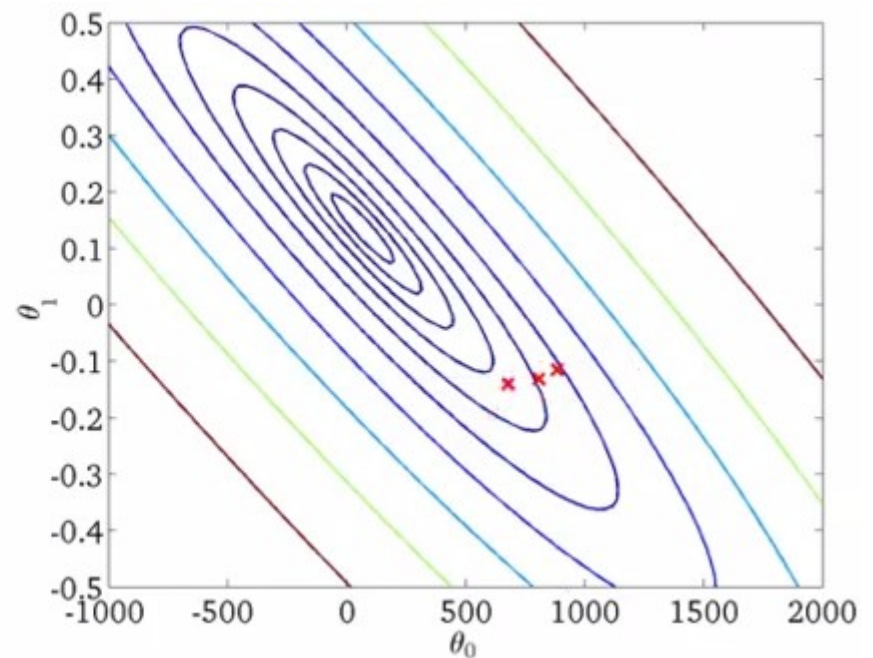
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

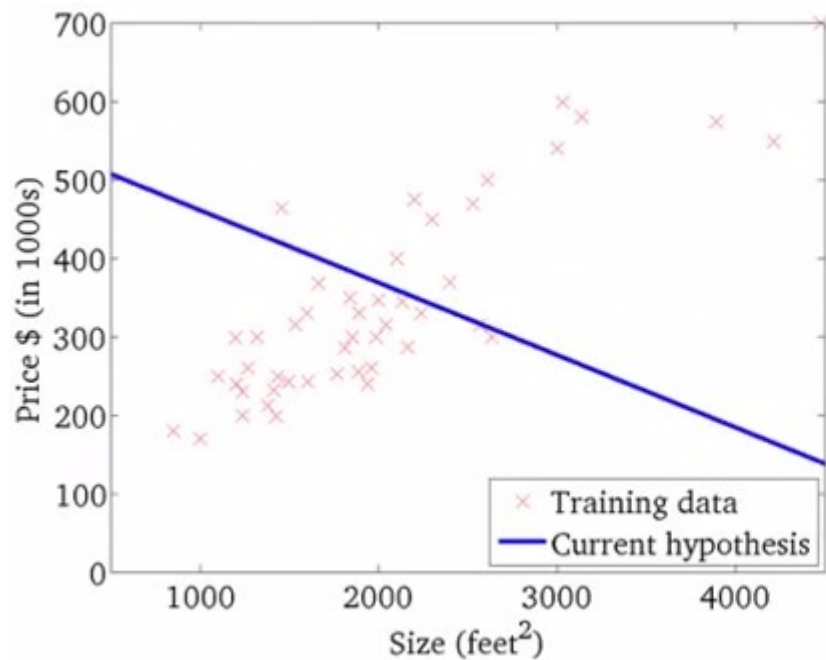
(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente + Regressão linear

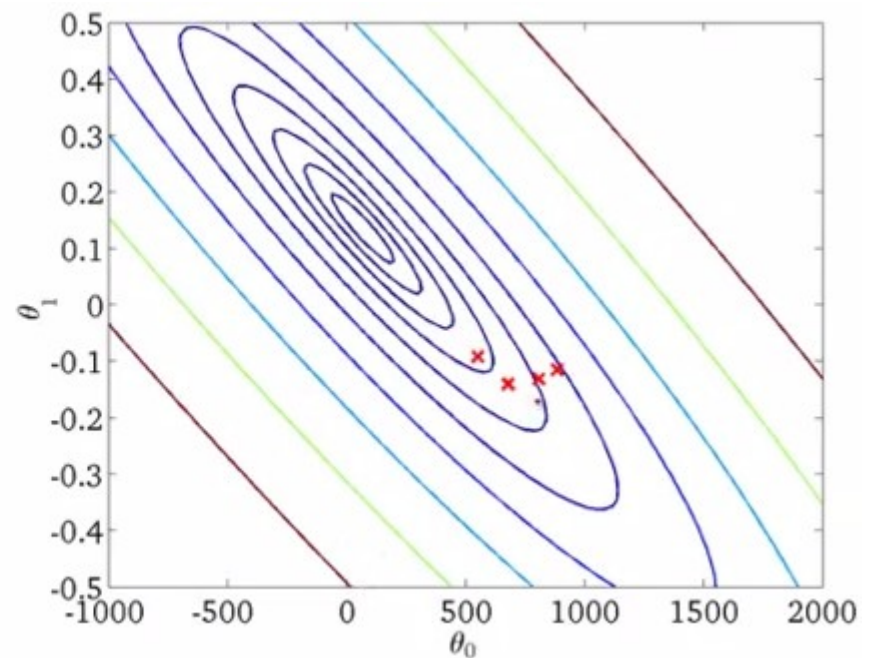
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

(função com parâmetro  $\theta_0, \theta_1$ )

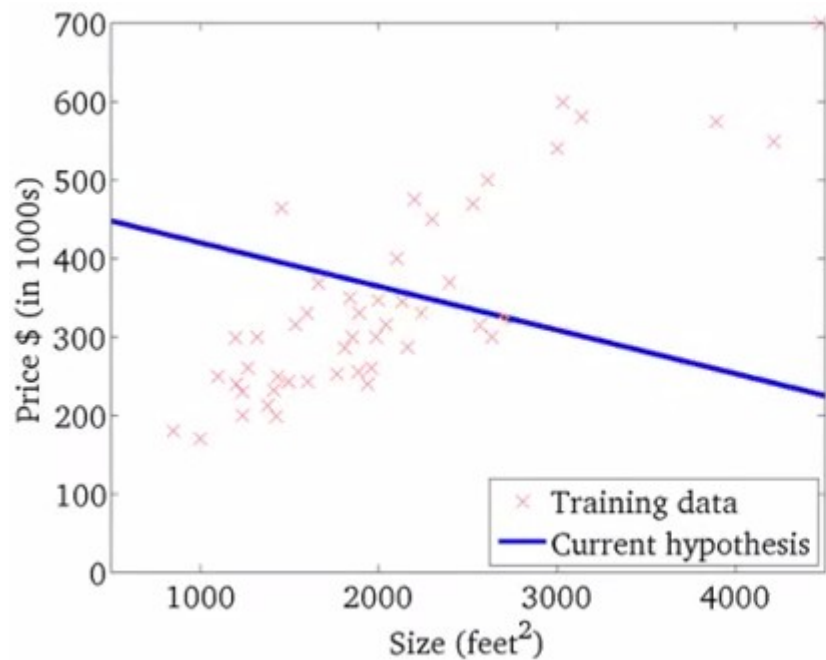




# Gradiente descendente + Regressão linear

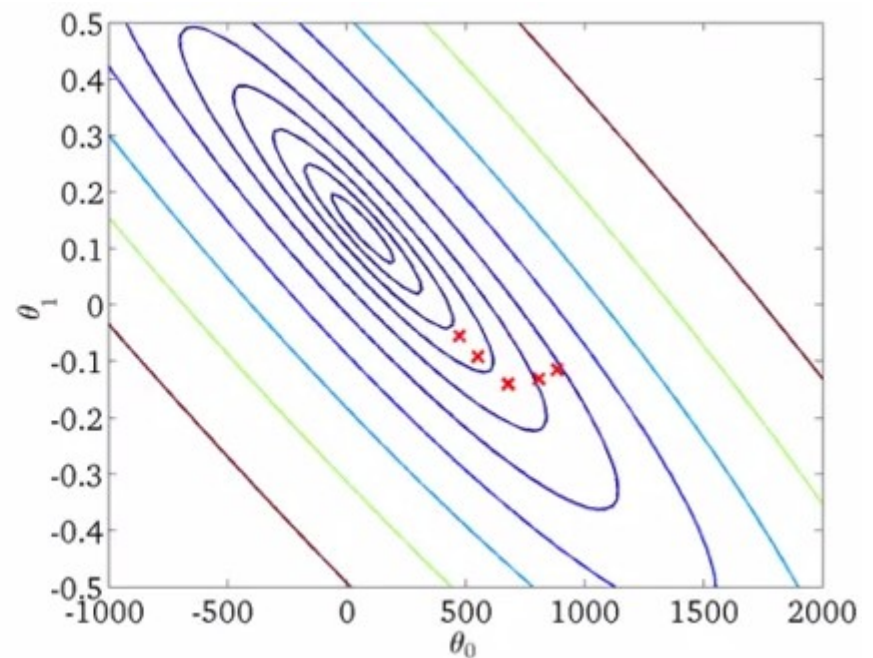
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

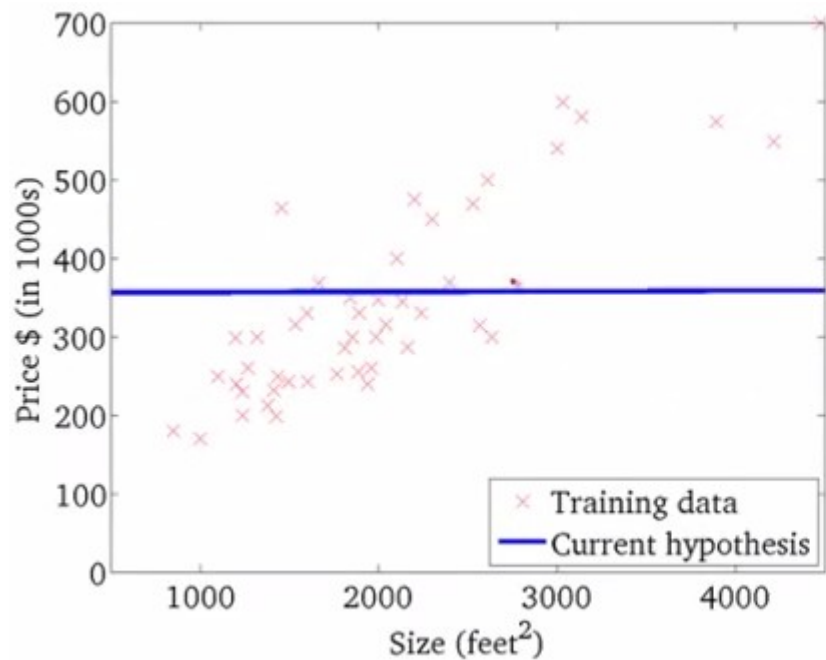
(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente + Regressão linear

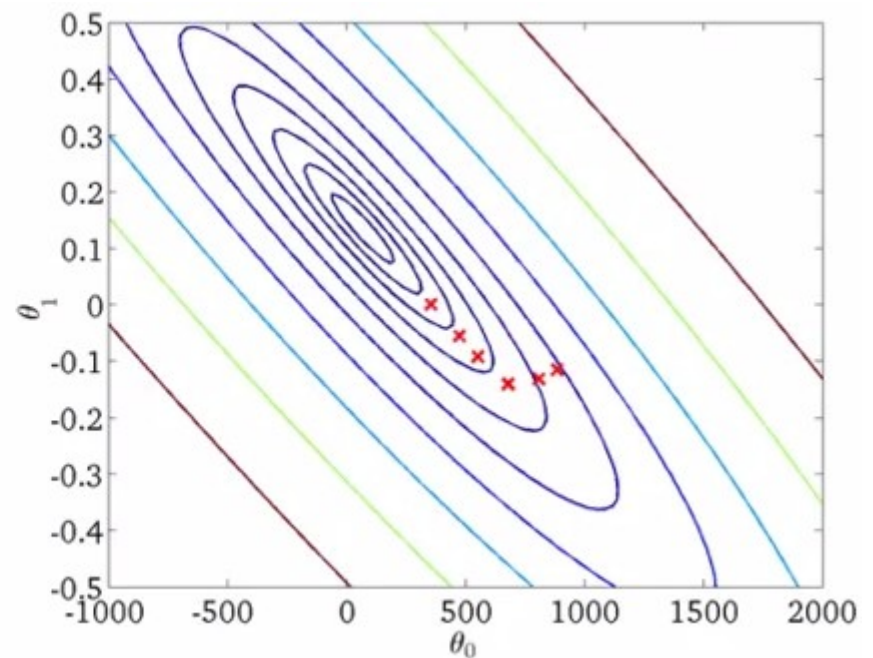
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

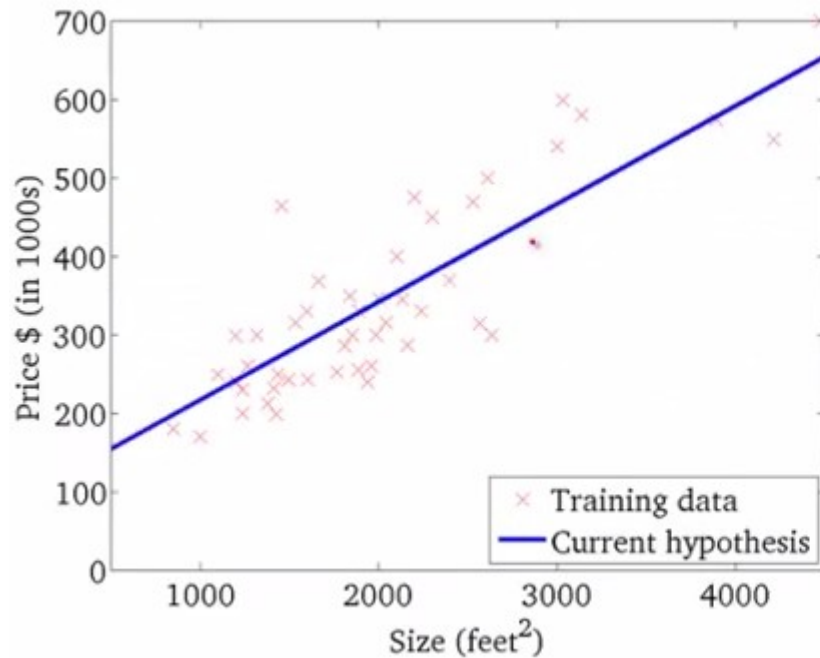
(função com parâmetro  $\theta_0, \theta_1$ )



# Gradiente descendente + Regressão linear

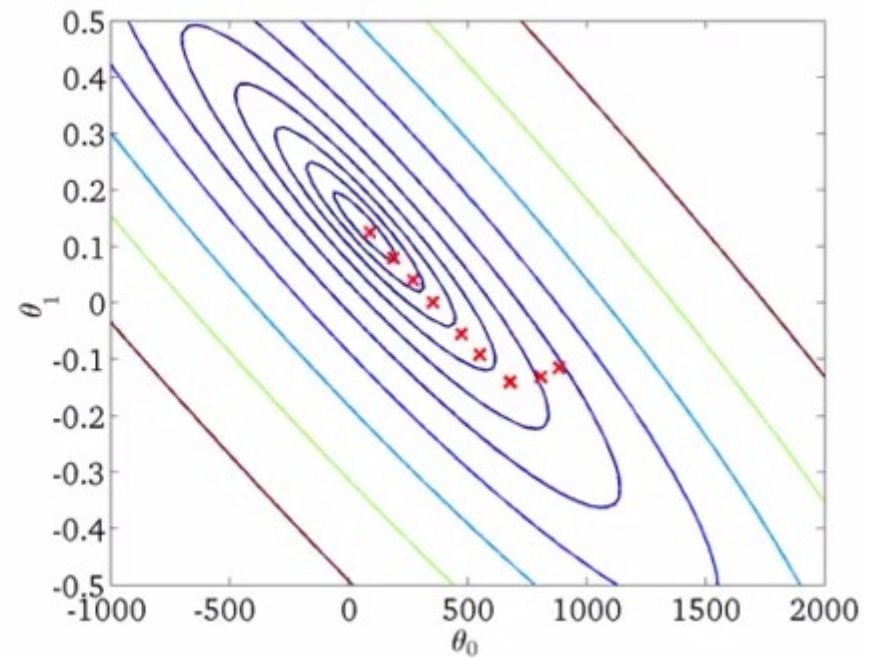
$$h_{\theta}(x)$$

(para  $\theta_0, \theta_1$ , é uma função de  $x$ )



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

(função com parâmetro  $\theta_0, \theta_1$ )



# **Outros tipos de regressão**

# Regressão não-linear

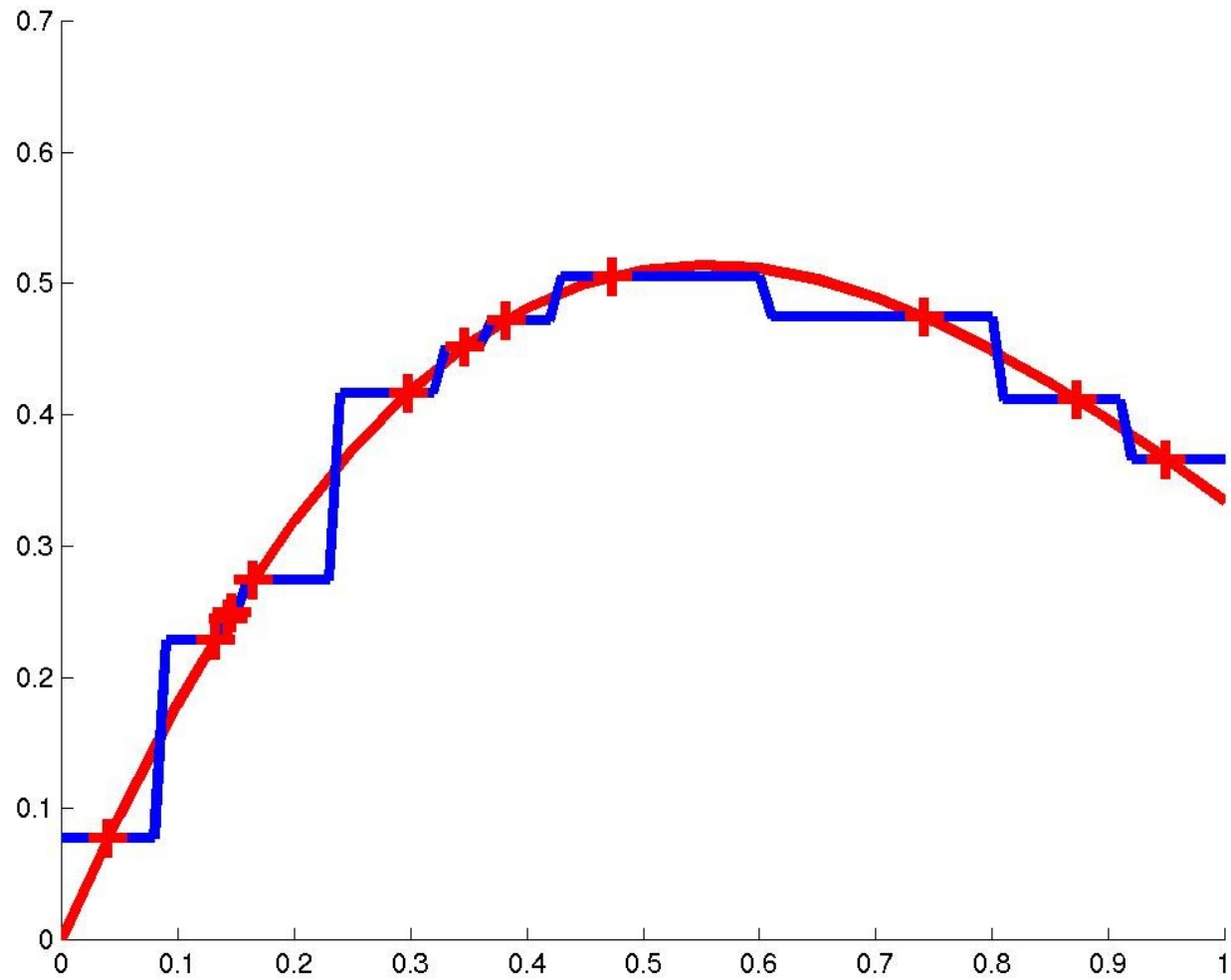
- Podemos aplicar uma função para cada atributo ao invés de multiplicar pelo peso
  - Na verdade, regressão linear é um caso particular da regressão não-linear
    - $y = c + f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$
    - Funções podem ser quadrado, log, raiz, módulo, exponencial, etc...
- Determinar as funções apropriadas é um problema!
- Splines, modelos lineares generalizados etc

# KNN

- Calcula a média dos k-vizinhos mais próximos para estimar o valor
  - Assuma os seguintes exemplos de treinamento  
 $(x^{(1)}, 1); (x^{(2)}, 1); (x^{(3)}, 0)$
  - Desejamos classificar  $x^{(4)}$  usando  $K=2$  vizinhos mais próximos, sabendo que
    - $\text{dist}(x^{(1)}, x^{(4)}) = 5$
    - $\text{dist}(x^{(2)}, x^{(4)}) = 2$
    - $\text{dist}(x^{(3)}, x^{(4)}) = 4$
    - Portanto  $x^{(4)}$  está mais próximo a  $x^{(2)}$  e  $x^{(3)}$
    - Calculando a média, obtém-se  $(1+0)/2 = 0.5$

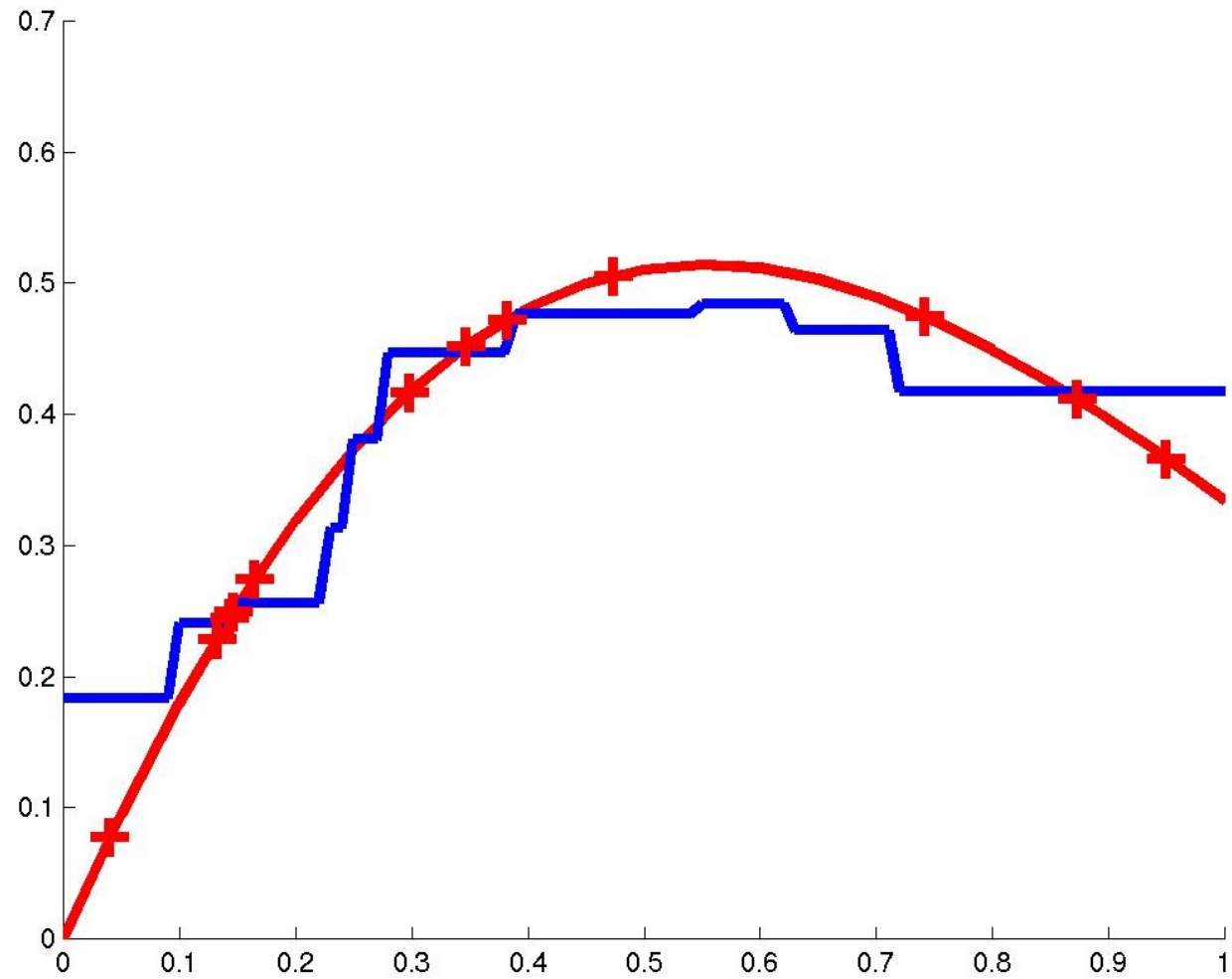
# kNN

- $k = 1$



# kNN

- $k = 3$





# kNN

- $k = 5$

