

Linear Regression

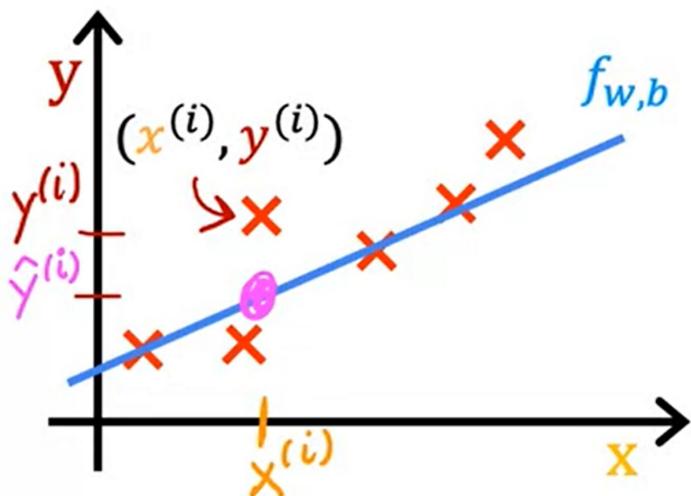
07 September 2021 00:12

Training set

features	targets
size in feet ² (x)	price \$1000's (y)
2104	460
1416	232
1534	315
852	178
...	...

$$\text{Model: } f_{w,b}(x) = wx + b$$

w, b : parameters
coefficients
weights



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

$$f_{w,b}(x^{(i)}) = wx^{(i)} + b$$

Cost function: Squared error cost function

$$\bar{J}(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

error

m = number of training examples

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

model:

$$f_{w,b}(x) = wx + b$$

parameters:

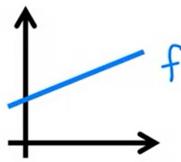
$$\underline{w, b}$$

cost function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

goal:

$$\underset{w, b}{\text{minimize}} J(w, b)$$



simplified

$$f_w(x) = \underline{wx} \quad b = \emptyset$$

$$w$$

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

$$\underset{w}{\text{minimize}} \underline{J(w)}$$

$\nwarrow \underline{wx^{(i)}}$

Have some function $J(w, b)$ for linear regression or any function

Want $\underset{w, b}{\text{minimize}} J(w, b)$ $\min_{w_1, \dots, w_n, b} J(w_1, w_2, \dots, w_n, b)$

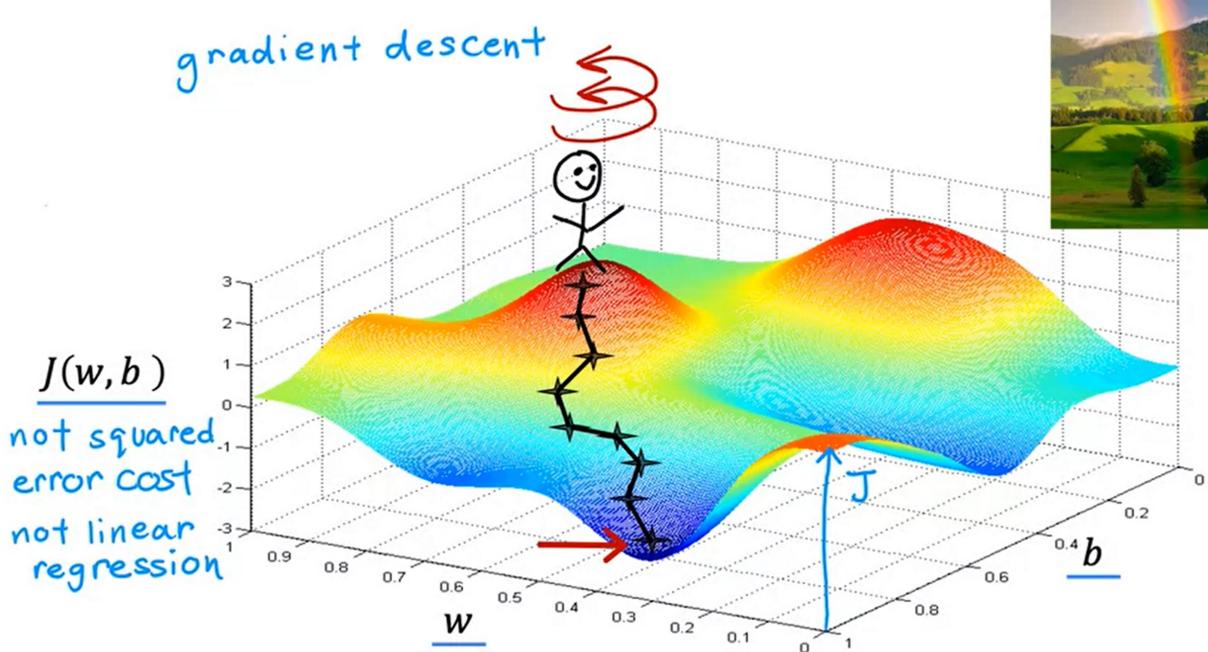
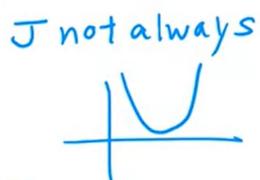
Outline:

Start with some w, b (set $w=0, b=0$)

Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

may have >1 minimum



Gradient descent algorithm

$$w = w - \alpha \frac{\partial}{\partial w} J(w, b)$$

$$b = b - \alpha \frac{\partial}{\partial b} J(w, b)$$

Simultaneously
update w and b

Correct: Simultaneous update

$$\begin{aligned} \text{tmp_w} &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \text{tmp_b} &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w &= \text{tmp_w} \\ b &= \text{tmp_b} \end{aligned} \quad \left. \right\}$$

Linear regression model

$$f_{w,b}(x) = wx + b \quad J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Cost function

Gradient descent algorithm

repeat until convergence {

$$\begin{aligned} w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)} \\ b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \end{aligned}$$

repeat until convergence {

$$\begin{aligned} w &= w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)} \\ b &= b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \end{aligned}$$

}

Multiple features (variables)

Size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price (\$) in \$1000's
x_1	x_2	x_3	x_4	
2104	5	1	45	460
i=2 1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$x_j = j^{\text{th}}$ feature
 $n = \text{number of features}$
 $\vec{x}^{(i)} = \text{features of } i^{\text{th}} \text{ training example}$
 $x_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example}$

$j=1\dots 4$
 $n=4$
 $\vec{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$
 $x_3^{(2)} = 2$

Model:

$$\text{Previously: } f_{w,b}(x) = wx + b$$

$$f_{w,b}(x) = w_1 x_1 + w_2 x_2 + w_3 x_3 + w_4 x_4 + b$$

example

$$f_{w,b}(x) = 0.1 \underset{\substack{\uparrow \\ \text{size}}}{x_1} + 4 \underset{\substack{\uparrow \\ \# \text{bedrooms}}}{x_2} + 10 \underset{\substack{\uparrow \\ \# \text{floors}}}{x_3} - 2 \underset{\substack{\uparrow \\ \text{years}}}{x_4} + 80 \underset{\substack{\uparrow \\ \text{base price}}}{b}$$

$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$ parameters of the model
 b is a number
 vector $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$

$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n + b$$

dot product multiple linear regression
 (not multivariate regression)

	Previous notation	Vector notation
Parameters	w_1, \dots, w_n b	\vec{w} ← vector of length n $w_1 \dots w_n$
Model	$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	b still a number $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(w_1, \dots, w_n, b)$	$J(\vec{w}, b)$ dot product

Gradient descent

```

repeat {
     $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{w_1, \dots, w_n, b})$ 
     $b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n, b})$ 
}

```

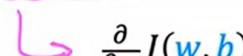
```

repeat {
     $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$ 
     $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$ 
}

```

Gradient descent

repeat { One feature

repeat {
 $\underline{w} = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \underline{x^{(i)}}$

 $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$
 simultaneously update w, b
 }

n features (*n* ≥ 2)

An alternative to gradient descent

→ Normal equation

- Only for linear regression
 - Solve for w , b without iterations

Disadvantages

- Doesn't generalize to other learning algorithms.
 - Slow when number of features is large ($> 10,000$)

What you need to know

- Normal equation method may be used in machine learning libraries that implement linear regression.
 - Gradient descent is the recommended method for finding parameters w, b

Feature and parameter values

$$\widehat{\text{price}} = w_1 x_1 + w_2 x_2 + b$$

size #bedrooms large small
x₁: size (feet²) range: 300 – 2,000 x₂: # bedrooms range: 0 – 5

House: $x_1 = 2000$, $x_2 = 5$, $\text{price} = \$500\text{k}$ one training example

size of the parameters w_1, w_2 ?

$$w_1 = 50, \quad w_2 = 0.1, \quad b = 50$$

$$\widehat{\text{price}} = \frac{50 * 2000}{100,000\text{K}} + \frac{0.1 * 5}{0.5\text{K}} + \frac{50}{50\text{K}}$$

$$\widehat{\text{price}} = \$100,050.5\text{K} = \$100,050,500$$

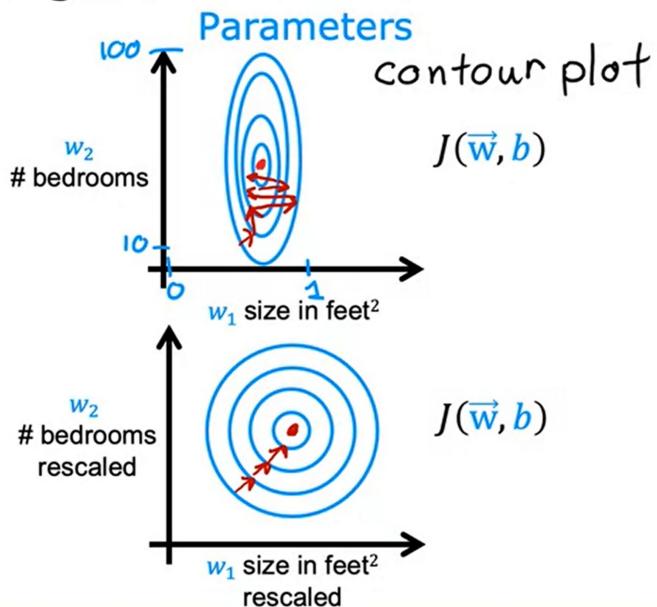
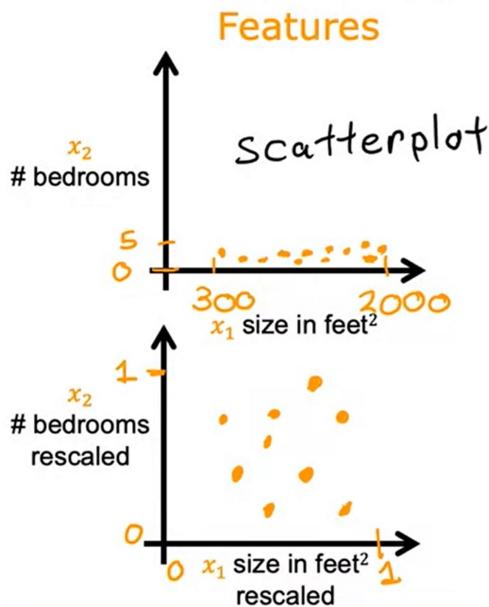
$$w_1 = 0.1, \quad w_2 = 50, \quad b = 50$$

small large

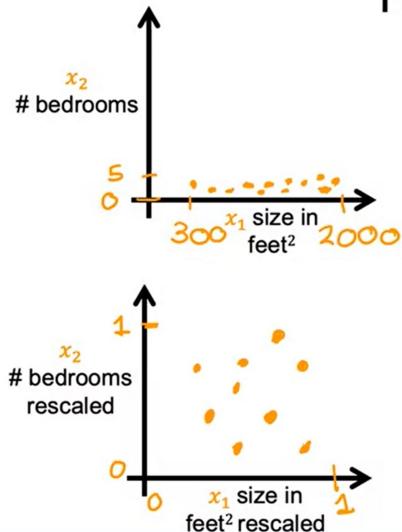
$$\widehat{\text{price}} = \frac{0.1 * 2000\text{k}}{200\text{K}} + \frac{50 * 5}{250\text{K}} + \frac{50}{50\text{K}}$$

$$\widehat{\text{price}} = \$500\text{k} \quad \text{more reasonable}$$

Feature size and gradient descent



Feature scaling



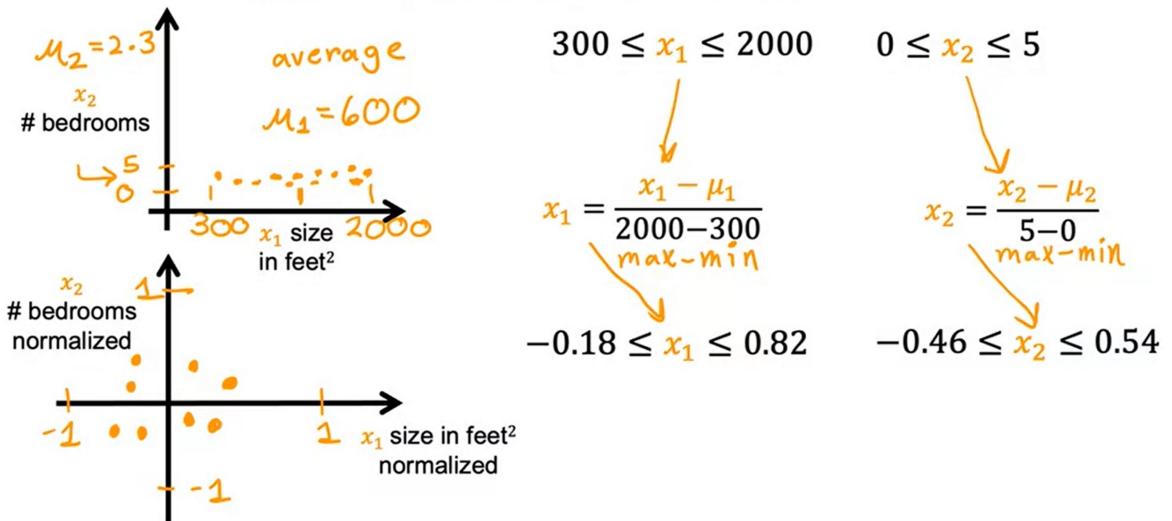
$$300 \leq x_1 \leq 2000 \quad 0 \leq x_2 \leq 5$$

$$x_{1,scaled} = \frac{x_1}{2000 \max}$$

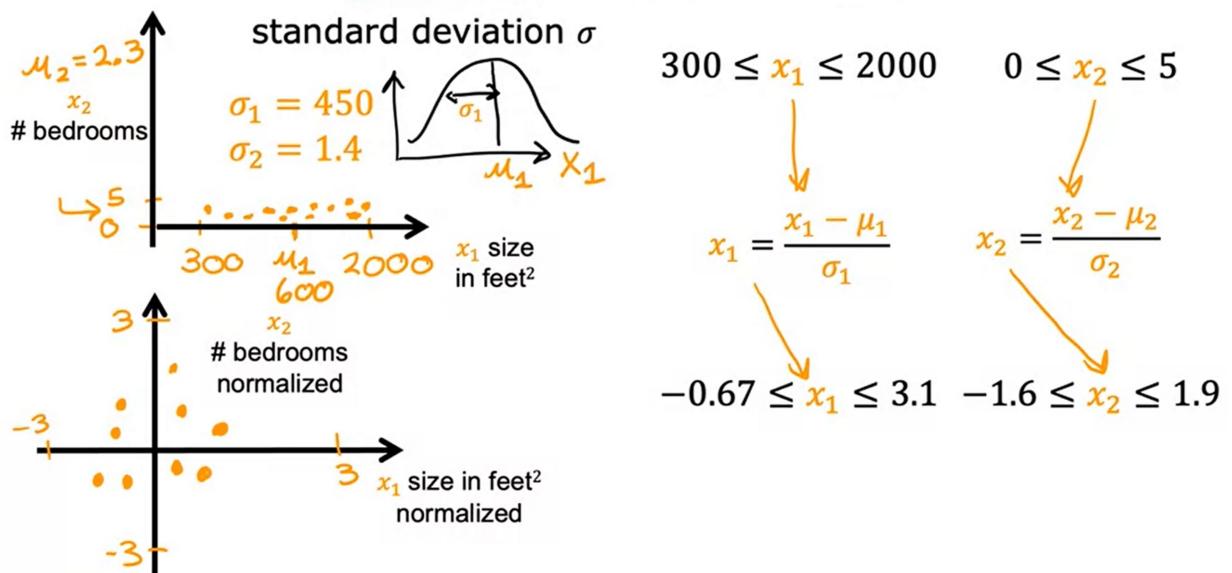
$$x_{2,scaled} = \frac{x_2}{5 \max}$$

$$0.15 \leq x_{1,scaled} \leq 1 \quad 0 \leq x_{2,scaled} \leq 1$$

Mean normalization



Z-score normalization



Feature scaling

aim for about $-1 \leq x_j \leq 1$ for each feature x_j

$-3 \leq x_j \leq 3$	}	acceptable ranges
$-0.3 \leq x_j \leq 0.3$		

$0 \leq x_1 \leq 3$	Okay, no rescaling
$-2 \leq x_2 \leq 0.5$	Okay, no rescaling
$-100 \leq x_3 \leq 100$	too large → rescale
$-0.001 \leq x_4 \leq 0.001$	too small → rescale
$98.6 \leq x_5 \leq 105$	too large → rescale