

## 6. Multiple features (variables)

Previously, we only had one feature. Our model looked something like this.

$$f_{w,b}(x) = w_1 x + b \text{, where } x \text{ is only one feature.}$$

What if we have multiple features as seen in ~~the~~ image.

We call each feature,  $x_1, x_2, x_3$  and  $x_4$ .

~~the~~

Let  $x_j = j^{\text{th}}$  feature.

$n$  = number of features.

$\vec{x}^{(i)}$  = features of  $i^{\text{th}}$  training example. (vector).  
 $\vec{x}^{(2)} = \begin{bmatrix} 1416 & 3 & 2 & 40 \end{bmatrix}$  a row vector. This is our 2nd training example.

$x_j^{(i)}$  = value of feature  $j$  in  $i^{\text{th}}$  training example.

$x_3^{(2)} = 2$ . in our above example.

Model: Previously,  $f_{w,b}(x) = w_1 x + b$ .

Now,  $f_{w,b}(x) = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$ .  
for  $n$  features.

for our new example,  $f_{w,b} = 0.1x_1 + 4x_2 + 10x_3 + (-2)x_4 + 80$ .

Here  $x_1 \rightarrow \text{size}$ ,  $x_2 \rightarrow \text{bedrooms}$ ,  $x_3 \rightarrow \text{floors}$ ,  $x_4 \rightarrow \text{years}$ .

Now, for  $n$  features we will have.

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

let our  $\vec{w} = [w_1, w_2, w_3, \dots, w_n]$  vector } there are  
parameters of  
 $b$  is a number. our model.

$$\vec{x} = [x_1, x_2, x_3, \dots, x_n] \rightarrow \text{input vector.}$$

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

$\vec{w} \cdot \vec{x}$  is the dot product.

We have our model in a compact form.

This model is called multiple Linear Regression model.  
(not multivariate regression)

In our previous example, what will be value of

$$x_1^{(4)} \rightarrow 1^{\text{st}} \text{ feature and } 4^{\text{th}} \text{ row.}$$

$$\rightarrow 852$$

### I. Vectorization.

Parameters and features.

$\vec{w} = [w_1 \ w_2 \ w_3]$   $b$  is a number.

$\vec{x} = [x_1 \ x_2 \ x_3]$   $n=3$ .

In linear algebra: count starts from 1.

In python, indexing starts from 0.

- Without vectorization.

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

In code, see python code.

If value of  $n$  is 100, or 10000.

for

$$f_{\vec{w}, b}(\vec{x}) = \left( \sum_{j=1}^n w_j x_j \right) + b. \quad \begin{cases} \text{we can implement} \\ \text{this using for loop.} \end{cases}$$

- Vectorization.

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b. \quad \begin{cases} \text{we can implement} \\ \text{the same using} \\ \text{numpy.} \end{cases}$$

When  $n$  is large, our code runs very fast.

numpy uses parallel hardware makes it much more faster.  
very useful when  $n$  is very large.

$$\text{numpy} \Rightarrow f = \text{np.dot}(w, x) + b.$$

without vectorization.

Vectorization.

for  $j$  in range(0, 16):

$$f = f + w[j] * x[j].$$

$$\text{to } f + w[0] * x[0]$$

$\text{np.dot}(w, x)$

to. it multiplies the two  
vectors in parallel.

to  $f + w(15) * x(15)$  Add all the 16 values  
at the same time using  
specialized hardware.

so, vectorization is really fast.

efficient on large datasets.

• by gradient Descent.

lets say  $\vec{w} = (w_1, w_2, \dots, w_{16}) - b$  parameters.

$$\vec{d} = (d_1, d_2, \dots, d_{16}),$$

$$w = \text{np.array}([6.5, 1.3, \dots, 3.4])$$

$$d = \text{np.array}([0.3, 0.2, \dots, 0.4])$$

$$\text{compute; } w_j = w_j - 0.1 d_j \quad \text{for } j=1 \dots 16.$$

without vectorization

$$w_1 = w_1 - 0.1 d_1,$$

with vectorization.

$$\vec{w} = \vec{w} - 0.1 \vec{d},$$

$$w_{16} = w_{16} - 0.1 d_{16}.$$

for  $j$  in range(0, 16):

$$w[j] = w[j] - 0.1 * d[j]$$

code:

$$w = w - 0.1 * d. \quad \begin{cases} \text{efficient} \\ \text{arrays.} \end{cases}$$

## II. Gradient Descent for multiple linear regression.

Previous notation

Parameters.  $w_1, w_2, \dots, w_n$   
 $b$

Model

$$f_{w,b}(\vec{x}) = w_1x_1 + \dots + w_nx_n + b$$

Cost function  $J(w_1, \dots, w_n, b)$

$$\vec{w} = [w_1, \dots, w_n] \quad \{ \text{length } n \}$$

$b$  still a number.

$$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$$

↑ dot product.

$$J(\vec{w}, b)$$

Gradient descent.

$$\text{repeat } \{ w_j^0 = w_j^0 - \alpha \frac{d}{dw_j} J(w_1, \dots, w_n, b) \}$$

repeat {

$$w_j^0 = w_j^0 - \alpha \frac{d}{dw_j} J(\vec{w}, b)$$

$$b = b - \alpha \frac{d}{db} J(w_1, w_2, \dots, w_n, b).$$

$$b = b - \alpha \frac{d}{db} J(\vec{w}, b)$$

}. update  $w, b$ .

• Gradient Descent.

$$\text{for one feature. } \{ w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x^{(i)}$$

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

simultaneously update  $w, b$ .

for  $n$  features ( $n > 2$ ).

$$j=1 \quad w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

↳ this is  $\frac{d}{dw_1} J(\vec{w}, b)$ .

$$j=n \quad w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \cdot x_n^{(i)}$$

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

we simultaneously update  
 $w_j$  (for  $j=1, \dots, n$ ) and  $b$ .

### III. An alternative to gradient descent.

Normal Equation: Works only for linear regression.  
 solve for  $w, b$  without iterations.

#### Disadvantages

- Doesn't generalize to other learning algorithms.
- slow when number of features ~~is~~ is large ( $> 1000$ ).

#### - What we 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  $w, b$ .

### IV. Quiz.

1. Given data, what is  $x_{4,3}^{(3)}$ ?  
 4<sup>th</sup> feature, 3<sup>rd</sup> example.

So, 3<sup>rd</sup> row, 4<sup>th</sup> column.

2. Which of the following are potential benefits of vectorization?

- code runs fast
- code shorter.
- code runs more easily on parallel compute hardware
- All  $\rightarrow$  Answer d). all.

3) To make gradient descent converge about twice as fast, a technique that almost works is to double the learning rate alpha  $\alpha$ ?  
 $\rightarrow$  False.