



Multiple Linear Regression

↗ Course	
↗ Topic	
▼ Type	Lecture
☑ Done	☑
📅 Date	@July 17, 2022
🔗 Materials	https://drive.google.com/file/d/1nJT4W5YT7dK_6JI6uvikMsLkWsy_OK6z/view?usp=sharing

▼ Summary

Introduction

Representing a multiple linear regression model

Vectorization

Gradient Descent for multiple linear regression

Alternative to Gradient Descent: Normal Equation

Introduction:

- In Multiple Linear Regression, **the output variable ‘y’ is dependent** on more than 1 input variable/feature.

For example, what if you had the size of the house, the number of bedrooms, the number of floors and the age of the home in years to be able to predict the house price? It will give us a lot of new information to work with.

	Size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price (\$) in \$1000's	$j=1...4$ $n=4$
	X_1	X_2	X_3	X_4		
$i=2$	2104	5	1	45	460	
	1416	3	2	40	232	
	1534	3	2	30	315	
	852	2	1	36	178	
	

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

$\vec{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$
 $x_3^{(2)} = 2$

Figure 1: A simple example of multiple linear regression.

- We're going to use the variables X_1, X_2, X_3 and X_4 , to denote the four features. Thus for any feature j , X_j is used to denote the list of features. For example, X_1 represents size in feet.

Multiple features (variables)

	Size in feet ²	Number of bedrooms	Number of floors	Age of home in years	Price (\$) in \$1000's	$j=1...4$
	X_1	X_2	X_3	X_4		
	2104	5	1	45	460	
	1416	3	2	40	232	
	1534	3	2	30	315	
	852	2	1	36	178	
	

$x_j = j^{th}$ feature

Figure 2: The yellow rectangle(size of feet) represents X_1 feature

- $\vec{x}^{(i)}$ is used to denote all the features of the i^{th} training example. For example, column 2 [1416 3 2 40] is represented by $\vec{x}^{(2)} = [1416 \ 3 \ 2 \ 40]$
- To refer to a specific feature in the i^{th} training example, we will write $\vec{X}_j^{(i)}$, so for example, $\vec{X}_3^{(2)}$ will be the number of floors in the second training example. So, $\vec{X}_3^{(2)} = 2$.

Representing a multiple linear regression model

- Previously, for linear regression in one variable, the predicted value of a input feature was:

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

Now for multiple linear regression,

$$\mathbf{f}_{\mathbf{w},\mathbf{b}}(\mathbf{x}) = \mathbf{w}_1\mathbf{x}_1 + \mathbf{w}_2\mathbf{x}_2 + \mathbf{w}_3\mathbf{x}_3 + \dots + \mathbf{w}_n\mathbf{x}_n + \mathbf{b}$$

- We can represent $w_1, w_2, w_3 \dots \dots \dots w_n$ as a single vector \vec{w} :

$$\vec{w} = [w_1 \quad w_2 \quad w_3 \quad w_4 \quad \dots \dots \dots w_n]$$

- We can represent $x_1, x_2, x_3 \dots \dots \dots x_n$ (the features of model like size of house, no of bathrooms etc) as a single vector $\vec{x} = [x_1 \quad x_2 \quad x_3 \quad x_4 \quad \dots \dots \dots x_n]$
- $w_1x_1 + w_2x_2 + w_3x_3 + \dots \dots \dots + w_nx_n$ can be represented by the dot product of \vec{w} and \vec{x} i.e $\vec{w} \cdot \vec{x}$

Therefore,

$$\mathbf{f}_{\mathbf{w},\mathbf{b}}(\vec{\mathbf{x}}) = \vec{\mathbf{w}} \cdot \vec{\mathbf{x}} + \mathbf{b}$$

$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$
 $\vec{w} = [w_1 \ w_2 \ w_3 \ \dots \ w_n]$ parameters of the model
 b is a number
 vector $\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$
 $f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b$
 dot product multiple linear regression
 (not multivariate regression)

Figure 3: Representing a multiple linear regression model

Vectorization

- Using vectorization will both make our code shorter and also make it run much more efficiently.

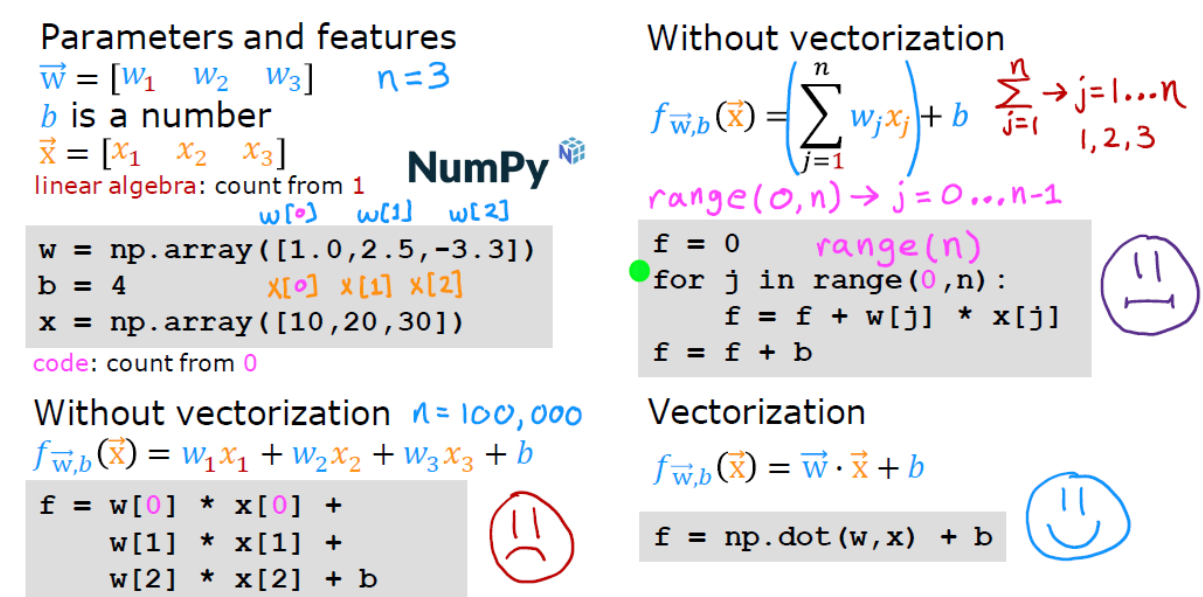


Figure 5: Vectorization is faster than normal calculations

Why is Vectorization faster?

- Without vectorization,
 - If j ranges from 0 to say 15, code performs operations one after another for 15 times.
 - On the first timestamp t_0 , it first operates on the values at index 0.
 - At the next time-step, it calculates values corresponding to index 1 and so on until the 15th step, where it computes that.
 - In other words, it calculates these computations one step at a time, one step after another.
-
- With vectorization, In contrast, `np.dot()` is implemented in the computer hardware with vectorization.
 - The computer can get all values of the vectors w and x , and in a single-step, it multiplies each pair of w and x with each other all at the same time in parallel.
 - Then after that, the computer takes these 16 numbers and uses specialized hardware to add them altogether very efficiently, rather than needing to carry out distinct additions one after another to add up these 16 numbers. This is all done in one step.-
 - This means that codes with vectorization can perform calculations in much less time than codes without vectorization. This matters more when you're running algorithms on large data sets or trying to train large models, which is often the case with machine learning. That's why being able to vectorize implementations of learning algorithms, has been a key step to getting learning algorithms to run efficiently, and therefore scale well to large datasets that many modern machine learning algorithms now have to operate on

Without vectorization

```
for j in range(0,16):
    f = f + w[j] * x[j]
```

t_0

$$f + w[0] * x[0]$$

t_1

$$f + w[1] * x[1]$$

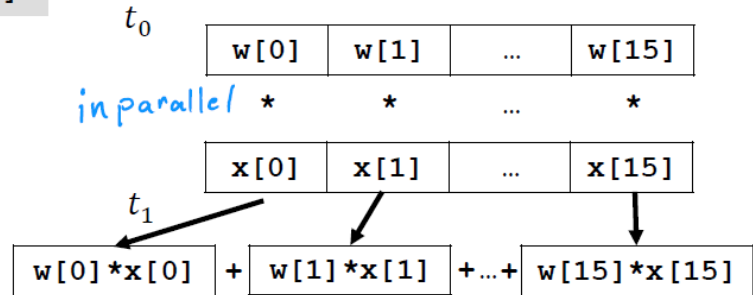
...

t_{15}

$$f + w[15] * x[15]$$

Vectorization

```
np.dot(w,x)
```



efficient → scale to large datasets

Figure 6: Explaining why vectorization is faster

Gradient Descent for multiple linear regression

	Previous notation	Vector notation
Parameters	w_1, \dots, w_n b	\vec{w} ← vector of length n $\vec{w} = [w_1 \dots w_n]$ b still a number
Model	$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(\underbrace{w_1, \dots, w_n}_{\text{vector}})$	$J(\underbrace{\vec{w}}_{\text{dot product}}, b)$
Gradient descent	repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{w_1, \dots, w_n}_{\text{vector}}, b)$ $b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n}_{\text{vector}}, b)$ }	repeat { $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\underbrace{\vec{w}}_{\text{dot product}}, b)$ $b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{\vec{w}}_{\text{dot product}}, b)$ }

Figure 7: How gradient descent changes for multiple linear regression

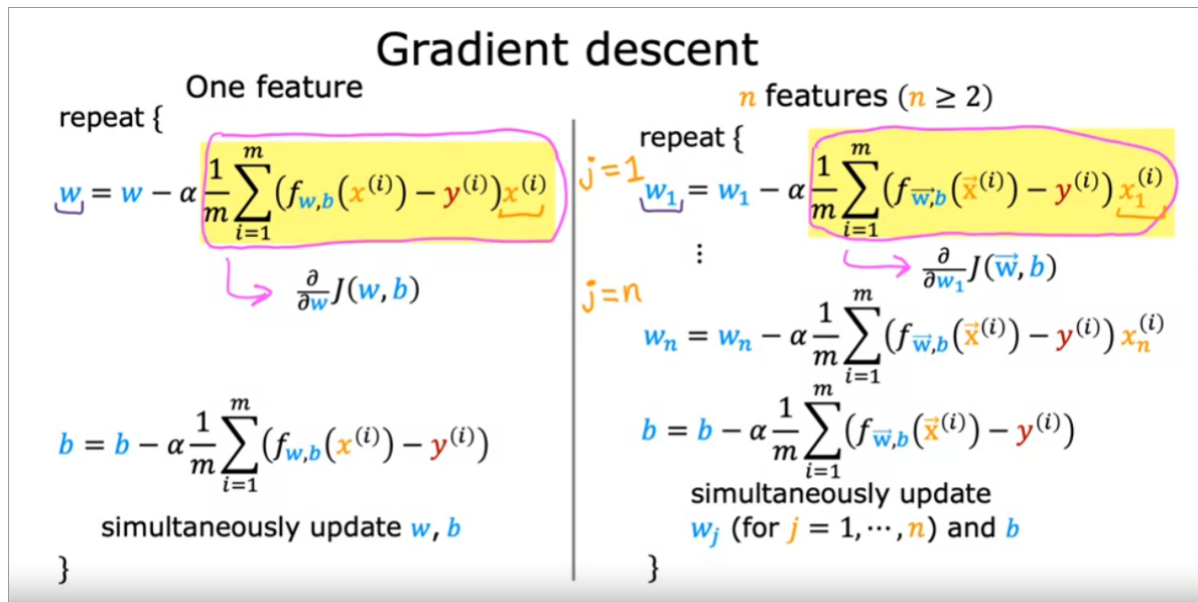


Figure 8: How gradient descent changes for multiple linear regression

Alternative to Gradient Descent

Normal equation

- An alternative way for finding w and b for linear regression is the normal equation
- Only for linear regression
- Solve for w and b all in one go without iterations.
- Normal equation method may be used in machine learning libraries that implement linear regression.

Disadvantages

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



Gradient descent is the recommended method for finding parameters w, b .