

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	
1416	3	2	40	232	
1534	3	2	30	315	
852	2	1	36	178	
...	

$j = 1 \dots 4$
 $n = 4$

$i = 2$

$x_j = j^{th}$ 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

Raw vector(List)

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

X superscript i, subscript j

Model:

Single linear regression model:

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

$$f_{w,b}(X) = w_1X_1 + w_2X_2 + w_3X_3 + w_4X_4 + b$$

example

Multiple linear regression model example

$$f_{w,b}(X) = 0.1X_1 + 4X_2 + 10X_3 + -2X_4 + 80$$

size #bedrooms #floors years base price

$$f_{w,b}(X) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

80 → Base price while assuming that the house has no size, no bedrooms, no floors, no age parameters. (Start price).

0.1 X1 → Means that the price will increase by $(0.1 * 1000 = 100 \$)$ for each additional square feet.

4 X2 → Means that the price will increase by $(4 * 1000 = 4000 \$)$ for each additional bedroom.

10 X3 → Means that the price will increase by $(10 * 1000 = 10000 \$)$ for each additional floor.

-2 X4 → Means that the price will decrease by $(2 * 1000 = 2000 \$)$ for each additional year of age.

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + \dots + w_nx_n + b$$

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

b is a number

parameters
of the model

vector

$$\vec{x} = [x_1 \ x_2 \ x_3 \ \dots \ x_n]$$

This arrow is optional

Note:
Raw vector = List

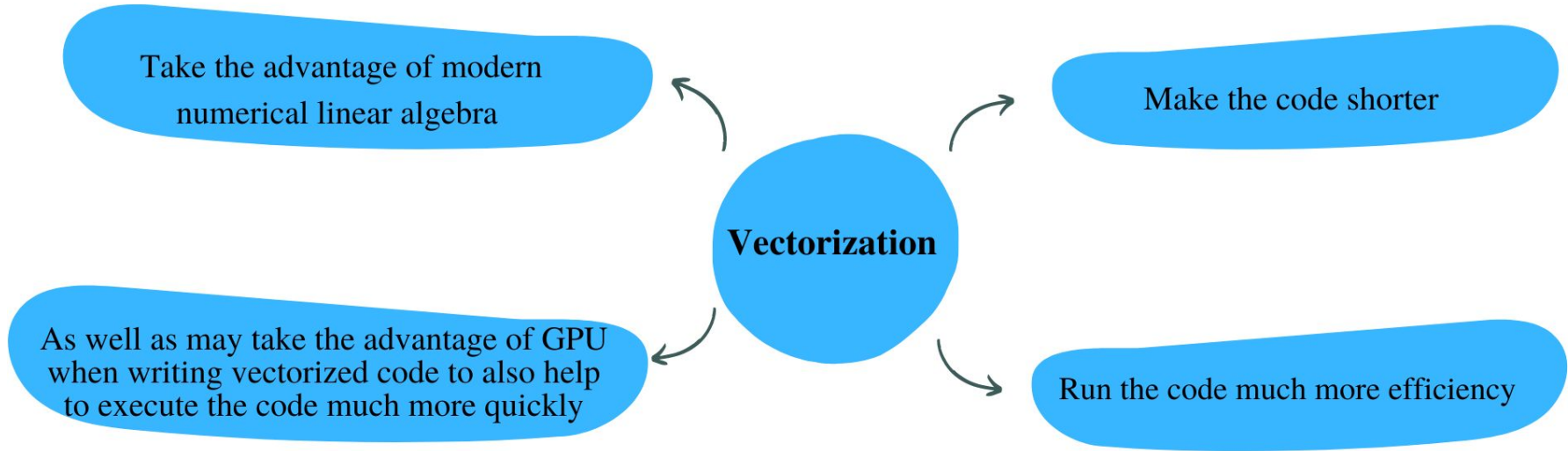
$$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b = w_1x_1 + w_2x_2 + w_3x_3 + \dots + w_nx_n + b$$

dot product

multiple linear regression

(not multivariate regression)

Vectorization



Note:

- GPU (Graphics Processing Unit) is a hardware objectively designed to speed up the computer graphics in our computer.
- GPU is often used to accelerate ml jobs.

Note:

- In **linear algebra**, the **count** starts from **1**, but starts from **0** in **python**.

Parameters and features

$$\vec{w} = [w_1 \ w_2 \ w_3] \quad n=3$$

b is a number

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

linear algebra: count from 1

NumPy 

$w[0] \ w[1] \ w[2]$

```
w = np.array([1.0, 2.5, -3.3])
```

```
b = 4  $x[0] \ x[1] \ x[2]$ 
```

```
x = np.array([10, 20, 30])
```

code: count from 0

Without vectorization $n=100,000$

$$f_{\vec{w},b}(\vec{x}) = w_1x_1 + w_2x_2 + w_3x_3 + b$$

```
f = w[0] * x[0] +  
     w[1] * x[1] +  
     w[2] * x[2] + b
```



Without vectorization

$$f_{\vec{w},b}(\vec{x}) = \left(\sum_{j=1}^n w_j x_j \right) + b \quad \sum_{j=1}^n \rightarrow j=1 \dots n, 1, 2, 3$$

$\text{range}(0, n) \rightarrow j=0 \dots n-1$

```
f = 0  $\text{range}(n)$   
for j in range(0, n):  
    f = f + w[j] * x[j]  
f = f + b
```



Vectorization

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

```
f = np.dot(w, x) + b
```

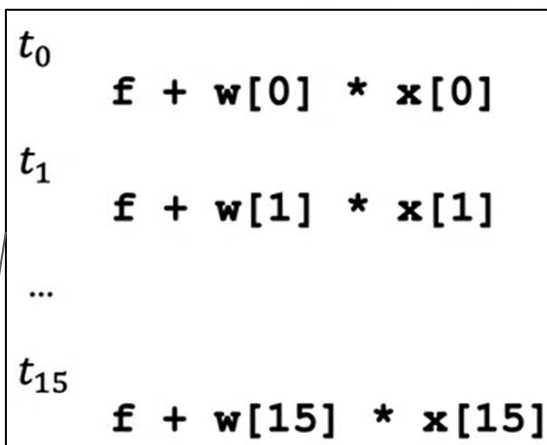


Note that numpy dot function makes the code more efficient and faster as it is able to use parallel hardware in the computer whether we use the CPU(Normal computer) or GPU(Graphics Processing Unit).

Why vectorized code is faster?

Without vectorization

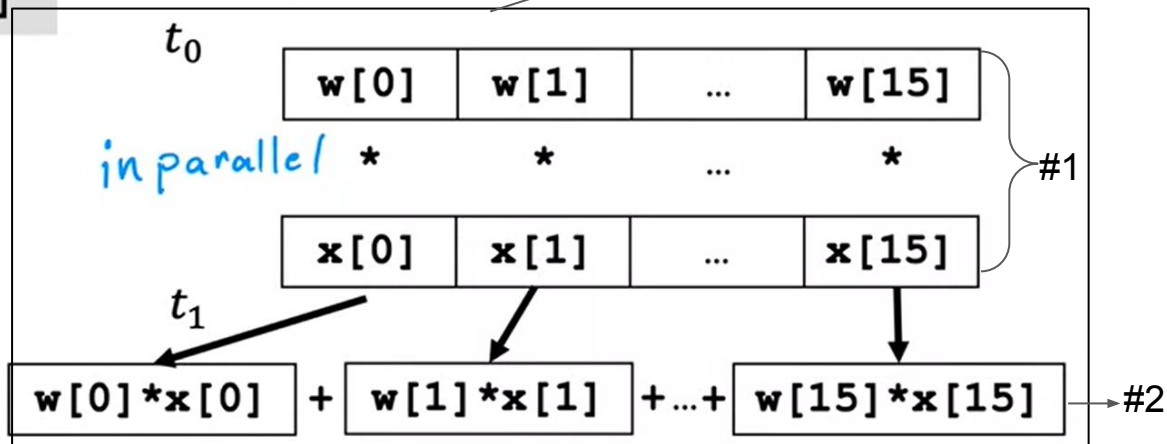
```
for j in range(0,16):  
    f = f + w[j] * x[j]
```



Vectorization

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

2 Steps



efficient \rightarrow scale to large datasets

Gradient descent

$$\vec{w} = (w_1 \quad w_2 \quad \dots \quad w_{16}) \quad \text{parameters}$$

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

```
w = np.array([0.5, 1.3, ... 3.4])
```

```
d = np.array([0.3, 0.2, ... 0.4])
```

compute $w_j = w_j - \underbrace{0.1}_{\text{learning rate } \alpha} d_j$ for $j = 1 \dots 16$

Without vectorization

$$w_1 = w_1 - 0.1d_1$$

$$w_2 = w_2 - 0.1d_2$$

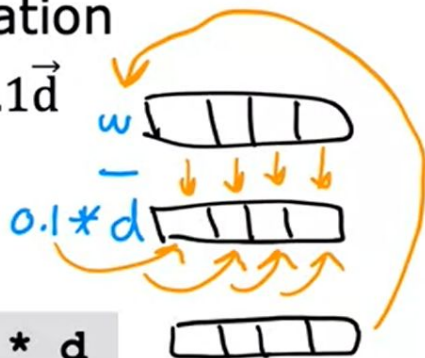
$$\vdots$$

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

```
for j in range(0,16):  
    w[j] = w[j] - 0.1 * d[j]
```

With vectorization

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



```
w = w - 0.1 * d
```

Gradient descent for multiple linear regression with vectorization

Previous notation

Parameters

$$w_1, \dots, w_n$$
$$b$$

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

Cost function $J(\underbrace{w_1, \dots, w_n}_\text{vector}, b)$

Gradient descent

repeat {

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

$$b = b - \alpha \frac{\partial}{\partial b} J(\underbrace{w_1, \dots, w_n}_\text{vector}, b)$$

}

Vector notation

\vec{w} ← vector of length n
 $\vec{w} = [w_1 \ \dots \ w_n]$
 b still a number

$f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
 $J(\underbrace{\vec{w}}_\text{vector}, b)$ dot product

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

One feature

repeat {

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

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

Cost function

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

simultaneously update w, b

}

n features ($n \geq 2$)

repeat {

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

\vdots

$\underline{w_n}$

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

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

simultaneously update

w_j (for $j = 1, \dots, n$) and b

}

An alternative to gradient descent

A complicated method used in the backend of some ml libraries.

→ 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

If we're using mature machine learning library and can't do linear regression, there is a chance of using this method on the backend to solve w and b .

Feature scaling

- When the range of feature's values is large, the (w) parameter will be relatively small.
- And, when the range of feature's values is small, the (w) parameter will be relatively large.

Feature and parameter values

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

x_1 : size (feet²) range: 300 – 2,000 large
 x_2 : # bedrooms range: 0 – 5 small

size #bedrooms

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

size of the parameters w_1, w_2 ?

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

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

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

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

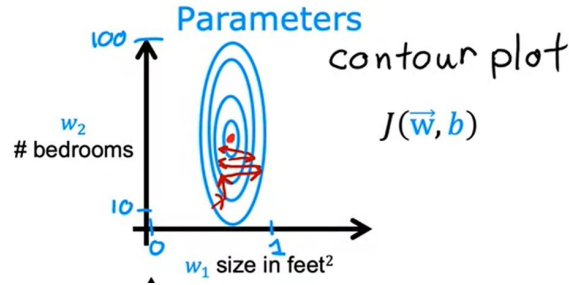
small large

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

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

The problem:

When the dataset features don't have comparable range of values (which means that some features have large scale of values where the other features have small scale of values), the cost function will be tall and skinny like this:



Because the contours are so tall and skinny, gradient descent may end up bouncing back and forth for a long time before it can finally find its way to the global minimum. (Cause gradient descent to run slowly).

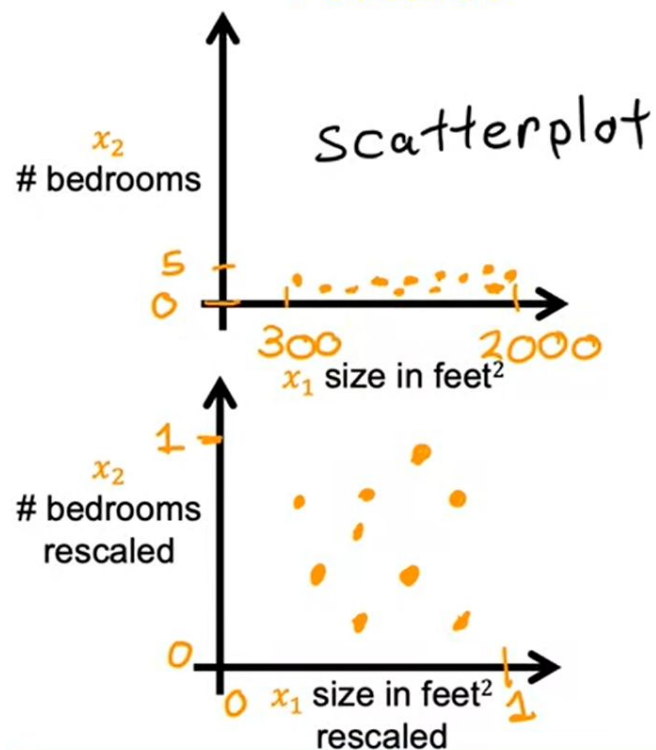


The best action here is features scaling which means performing some transformations of the training data. So, they all take on comparable range of values. Thus, we can speed up gradient descent significantly.

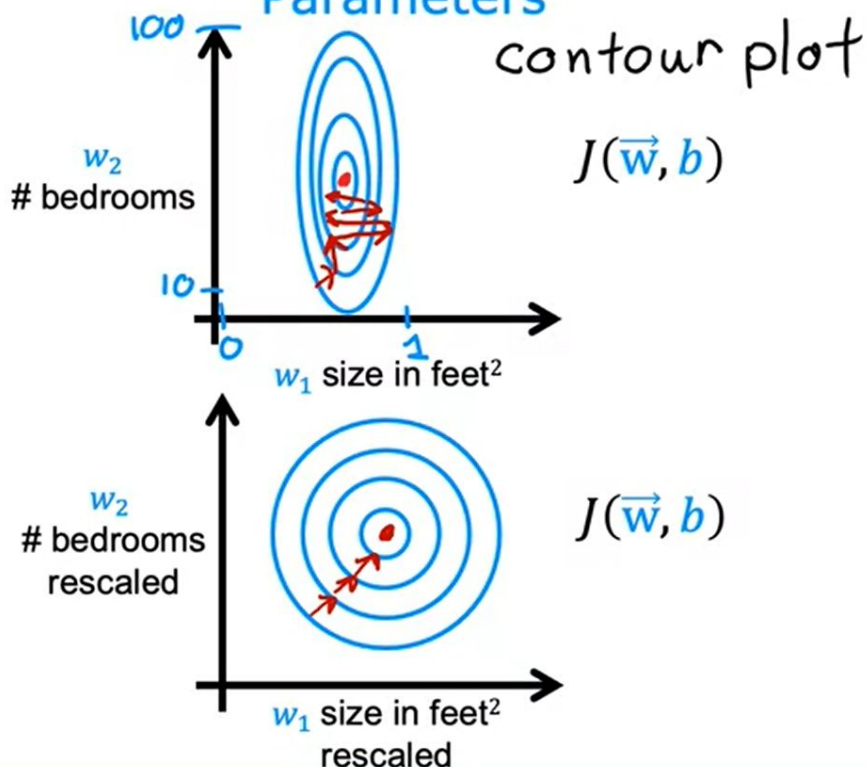


Feature size and gradient descent

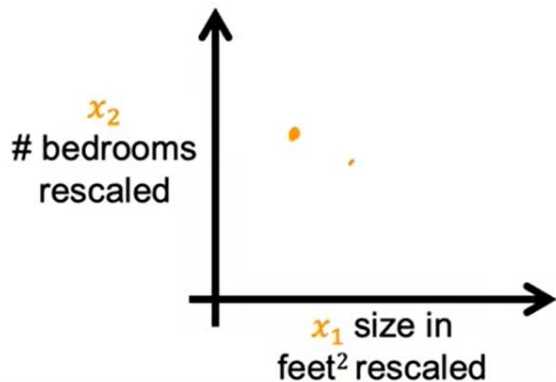
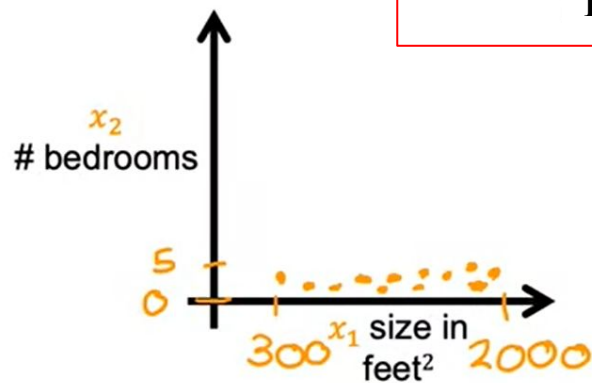
Features



Parameters



Feature scaling methods



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

1

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

max

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

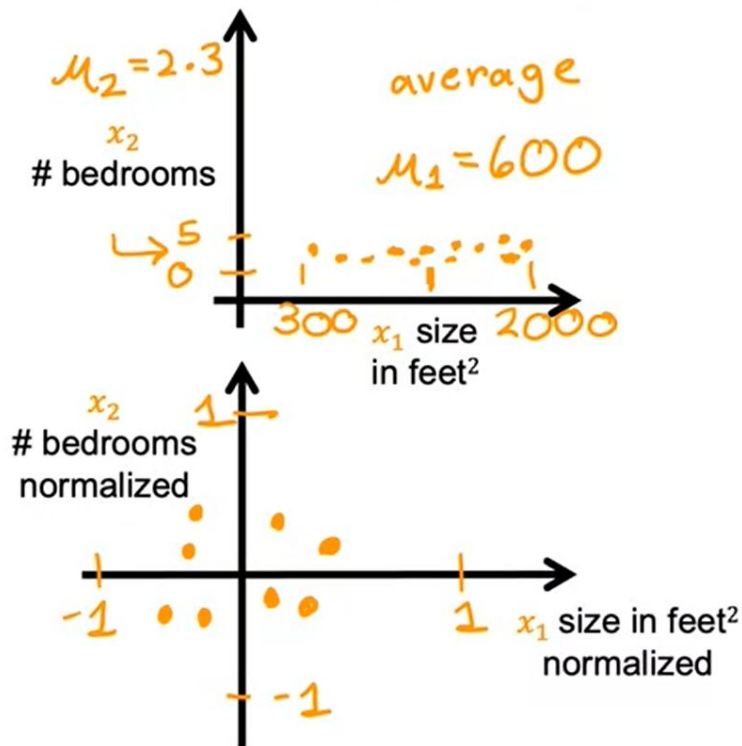
max

$$0.15 \leq x_{1,scaled} \leq 1$$

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

2

Mean normalization



$$300 \leq x_1 \leq 2000$$

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}$$

max-min

$$-0.18 \leq x_1 \leq 0.82$$

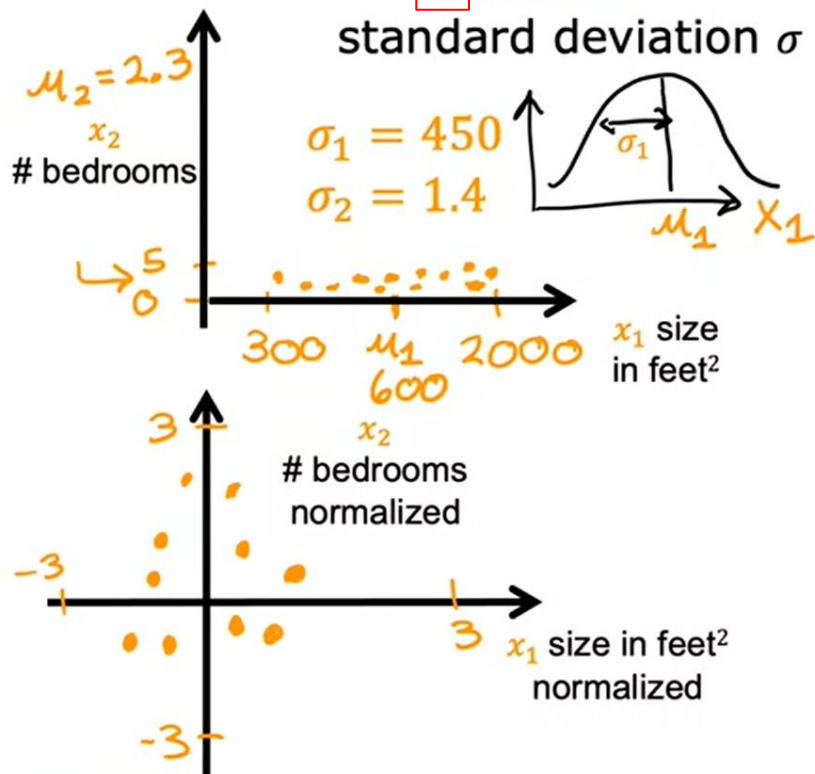
$$0 \leq x_2 \leq 5$$

$$x_2 = \frac{x_2 - \mu_2}{5 - 0}$$

max-min

$$-0.46 \leq x_2 \leq 0.54$$

3 Z-score normalization



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

$$x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

$$-0.67 \leq x_1 \leq 3.1 \quad -1.6 \leq x_2 \leq 1.9$$

Feature scaling

When to do feature scaling?

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

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

$$0 \leq x_1 \leq 3$$

okay, no rescaling

$$-2 \leq x_2 \leq 0.5$$

okay, no rescaling

$$-100 \leq x_3 \leq 100$$

too large \rightarrow rescale

$$-0.001 \leq x_4 \leq 0.001$$

too small \rightarrow rescale

$$98.6 \leq x_5 \leq 105$$

too large \rightarrow rescale

In case of doubt of scaling or not, **scale it**.

Make sure gradient descent is working well:

(The cost function of the training set **MUST** decrease over iterations as the aim is to reach the global minimum)

Learning curve

Automatic convergence test

1. Learning curve

- Plot the cost function (J), which is calculated on the training set at each iteration, iteration means after each simultaneous update of parameters (\mathbf{w} & b).
- $J(\mathbf{w}, b)$ should decrease after each iteration. If it increases after each iteration, that means:

Either Alpha is chosen poorly (It usually means Alpha is too large).

Or a bug in the code.

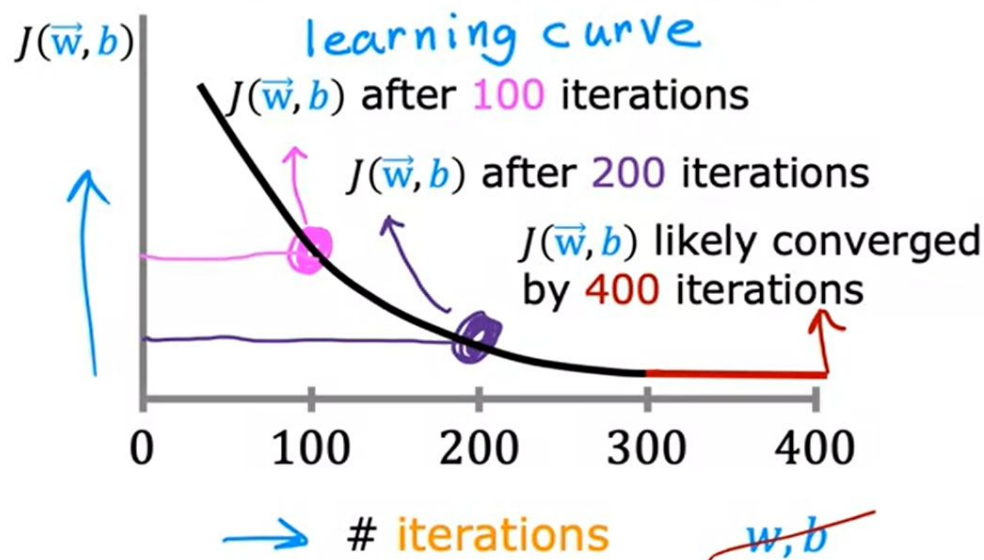
2. Automatic convergence test

If $J(\mathbf{w}, b)$ decreases by $\leq \epsilon$ in one iteration, declare convergence.

The **disadvantage** is that ϵ value is randomly chosen, and it is very difficult to be expected, so learning curve is better to use. In addition, learning curve provides us the ability to track the gradient descent if it is correctly working or not.

Make sure gradient descent is working correctly

objective: $\min_{\vec{w}, b} J(\vec{w}, b)$ $J(\vec{w}, b)$ should **decrease** after every iteration



iterations needed varies 30 1,000 100,000

Automatic convergence test

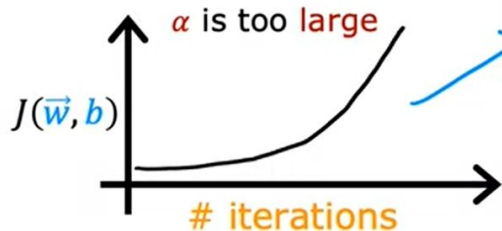
Let ϵ "epsilon" be 10^{-3} .
0.001

If $J(\vec{w}, b)$ decreases by $\leq \epsilon$ in one iteration,
declare convergence.

(found parameters \vec{w}, b to get close to global minimum)

Learning rate (α)

Identify problem with gradient descent



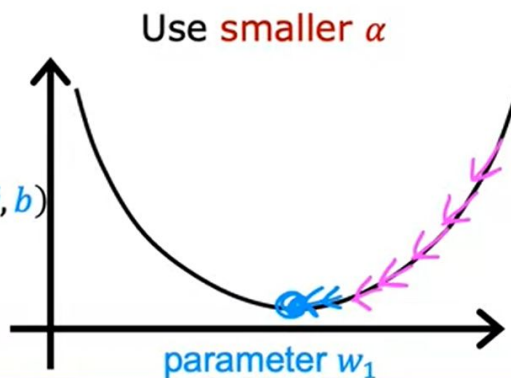
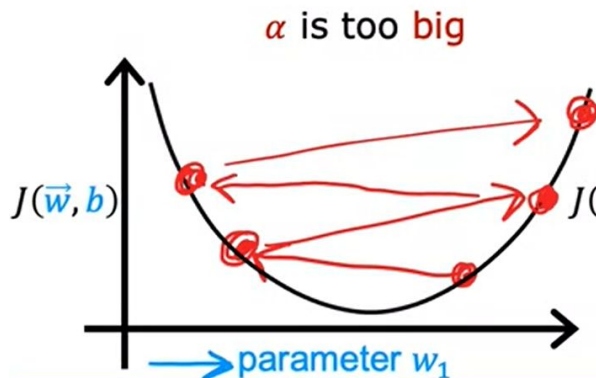
or learning rate is too large

$$w_1 = w_1 + \alpha d_1 \quad \text{!!}$$

use a minus sign

$$w_1 = w_1 - \alpha d_1 \quad \text{!!}$$

Adjust learning rate

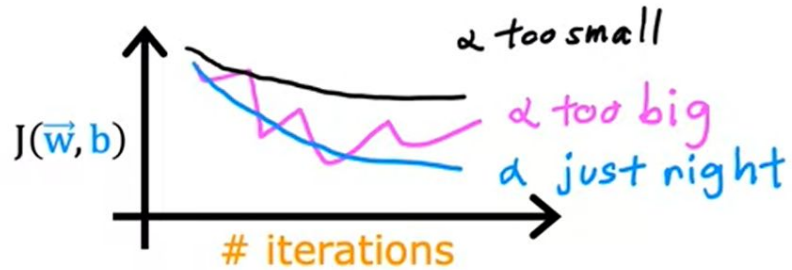
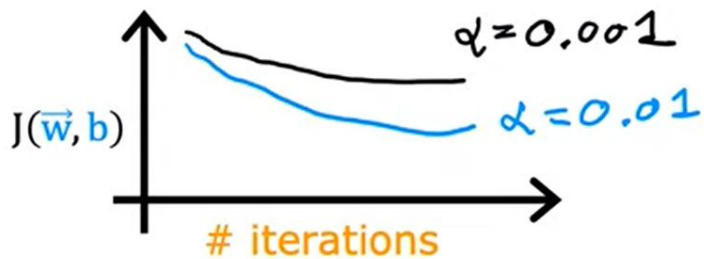


With a small enough α , $J(\vec{w}, b)$ should **decrease** on every iteration

If α is too small, gradient descent takes a lot more iterations to **converge**

Values of α to try:

... 0.001 0.003 0.01 0.03 0.1 0.3 1 ...
 \nearrow \nearrow \nearrow \nearrow \nearrow \nearrow
 $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$



Feature engineering

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1 x_1}_{\text{frontage}} + \underbrace{w_2 x_2}_{\text{depth}} + b$$

$$\text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

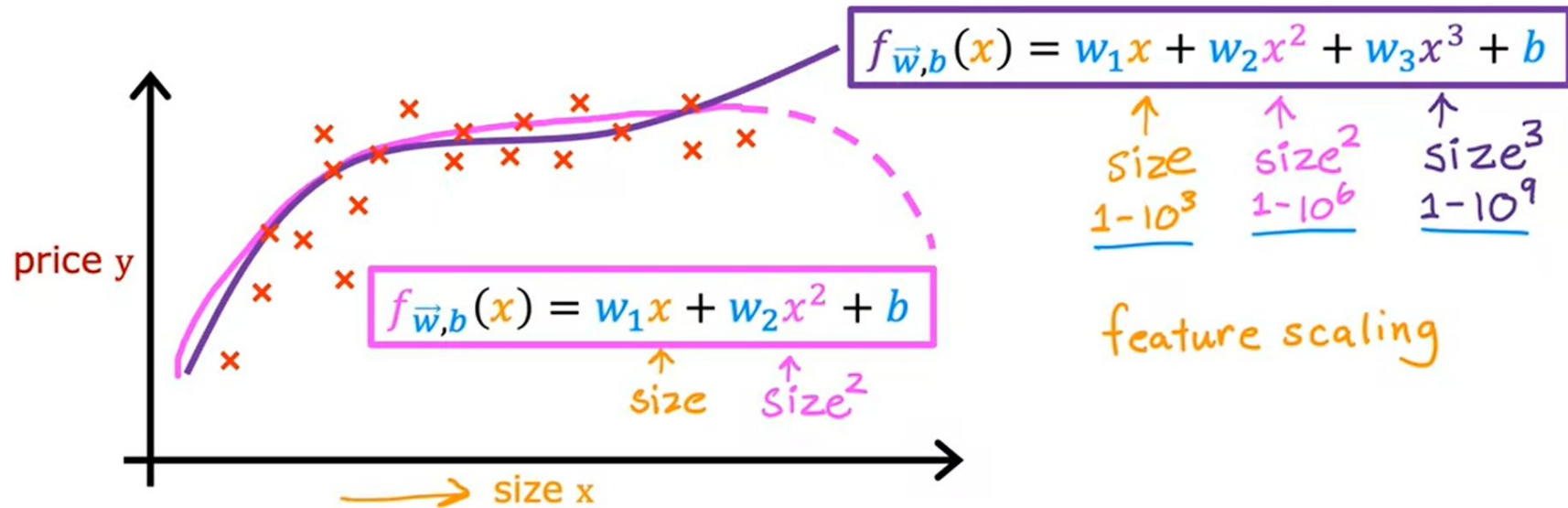
new feature

$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1 x_1}_{\text{frontage}} + \underbrace{w_2 x_2}_{\text{depth}} + \underbrace{w_3 x_3}_{\text{area}} + b$$



Feature engineering:
Using **intuition** to design
new features, by
transforming or combining
original features.

Polynomial regression



Choice of features

