

Week 2

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

$j=1...4$

$n=4$

$x_j = j^{th}$ 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

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

$$x_3^{(2)} = 2$$

emphasize that this is a
vector and not a number.

Model:

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

example $f_{w,b}(x) = w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b$

$$f_{w,b}(x) = 0.1x_1 + 4x_2 + 10x_3 + -2x_4 + 80$$

↑ ↑ ↑ ↑ ↑
size #bedrooms #floors years base price

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

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

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

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)

Question

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
...

$i=2$

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

$x_j = j^{\text{th}}$ 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

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

$$x_3^{(2)} = 2$$

Skip

Continue

Question

In the training set below, what is $x_1^{(4)}$? Please type in the number below (this is an integer such as 123, no decimal points).

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
...

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

$i = 2$

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

n = number of features

$\vec{x}^{(i)}$ = features of i^{th} training example

$x^{(i)}$ = value of feature i in j^{th} training example

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

$$x_3^{(2)} = 2$$

Skip

Continue

Question

$i=2$

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

$x_j = j^{\text{th}}$ 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

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

$$x_3^{(2)} = 2$$

852



Correct

$x_1^{(4)}$ is the first feature (first column in the table) of the fourth training example (fourth row in the table).

Skip

Continue

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

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

NumPy

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

```
b = 4
```

```
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

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

$\sum_{j=1}^n \rightarrow j=1 \dots n$
 $1, 2, 3$

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

```
f = 0  
for j in range(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
```



Question

Which of the following is a vectorized implementation for computing a linear regression model's prediction?

☒ `f = np.dot(w,x) + b`

☐ `f=0`

For `j` in `range(n)`:

`f = f + w[j] * x[j]`

`f=f + b`



Correct

This numpy function uses parallel hardware to efficiently calculate the dot product.

[Skip](#)

[Continue](#)

Without vectorization

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

t_0

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

t_1

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

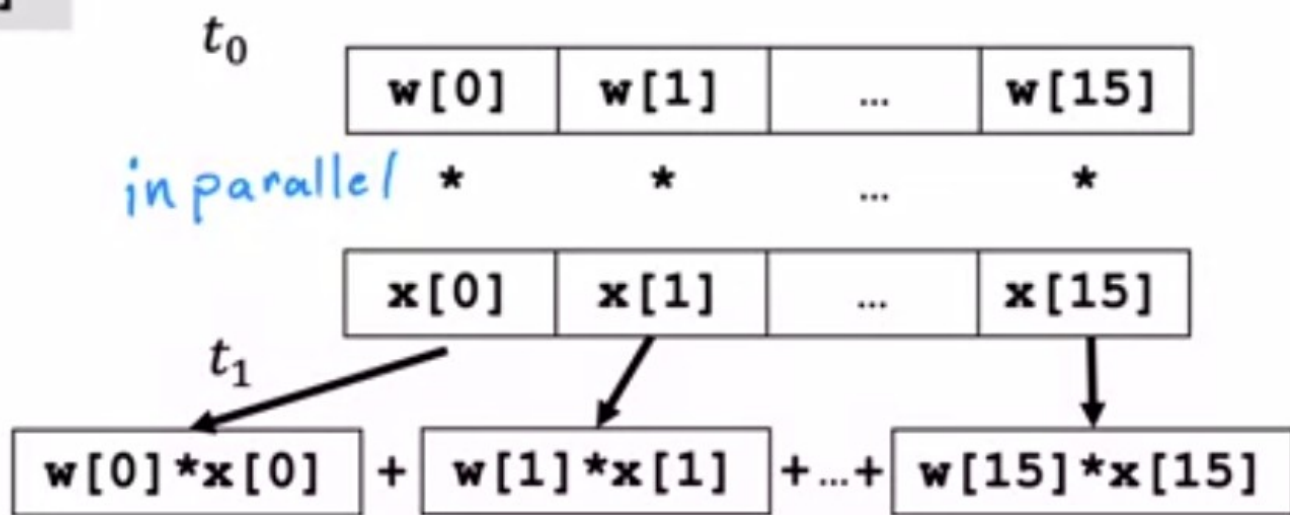
...

t_{15}

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

Vectorization

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



efficient → scale to large datasets

many modern machine
learning algorithms

Gradient descent $\vec{w} = (w_1 \ w_2 \ \dots \ w_{16})$ ~~b~~ parameters
derivatives $\vec{d} = (d_1 \ d_2 \ \dots \ d_{16})$

```
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

$$w_1 = w_1 - 0.1d_1$$

$$w_2 = w_2 - 0.1d_2$$

$$\vdots$$

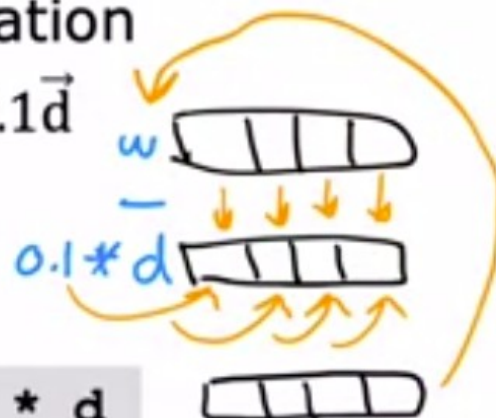
$$w_{16} = w_{16} - 0.1d_{16}$$

```
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
```

versus taking many hours
to do the same thing.

Question

Which of the following is a vectorized implementation for computing a linear regression model's prediction?

☐

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

☐

```
1 f = 0
2 for j in range(n):
3     f = f + w[j] * x[j]
4 f = f + b
```

☒

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

Skip

Continue

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

↖ vector of length n

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

b still a number

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

dot product

$$J(\underbrace{\vec{w}}_\text{vector}, b)$$

repeat {

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

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

}

let's take a look at
the derivative term.

Gradient descent

One feature

repeat {

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

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

simultaneously update w, b

}

n features ($n \geq 2$)

repeat {

$$\begin{aligned} j=1 \quad \underline{w_1} &= w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\bar{w},b}(\bar{x}^{(i)}) - y^{(i)}) \underline{x_1^{(i)}} \\ &\vdots \\ j=n \end{aligned} \quad \leftarrow \frac{\partial}{\partial w_1} J(\bar{w}, b)$$

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

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

simultaneously update

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

}

That's it for gradient descent
for multiple regression.

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

gradient descents offer a
better way to get the job done.

← Back

Practice quiz: Multiple linear regression

Graded Quiz • 15 min

Due Sep 11, 11:59 PM IST

1. In the training set below, what is $x_4^{(3)}$? Please type in the number below (this is an integer such as 123, no decimal points).

1 / 1 point

Size in feet ²	<u>Number of bedrooms</u>	<u>Number of floors</u>	<u>Age of home</u> 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
...

30



Practice quiz: Multiple linear regression

Graded Quiz • 15 min

Due Sep 11, 11:59 PM IST

2.

1 / 1 point

Which of the following are potential benefits of vectorization? Please choose the best option.

- ☐ It makes your code run faster
- ☐ It can make your code shorter
- ☐ It allows your code to run more easily on parallel compute hardware
- ☒ All of the above

✓ Correct

Correct! All of these are benefits of vectorization!

3. True/False? To make gradient descent converge about twice as fast, a technique that almost always works is to double the learning rate *alpha*.

1 / 1 point

- ☒ False
- ☐ True

✓ Correct

Doubling the learning rate may result in a learning rate that is too large, and cause gradient descent to fail to find the optimal values for the parameters w and b .

Feature and parameter values

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

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

size #bedrooms

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

size of the parameters w_1, w_2 ?

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

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

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

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

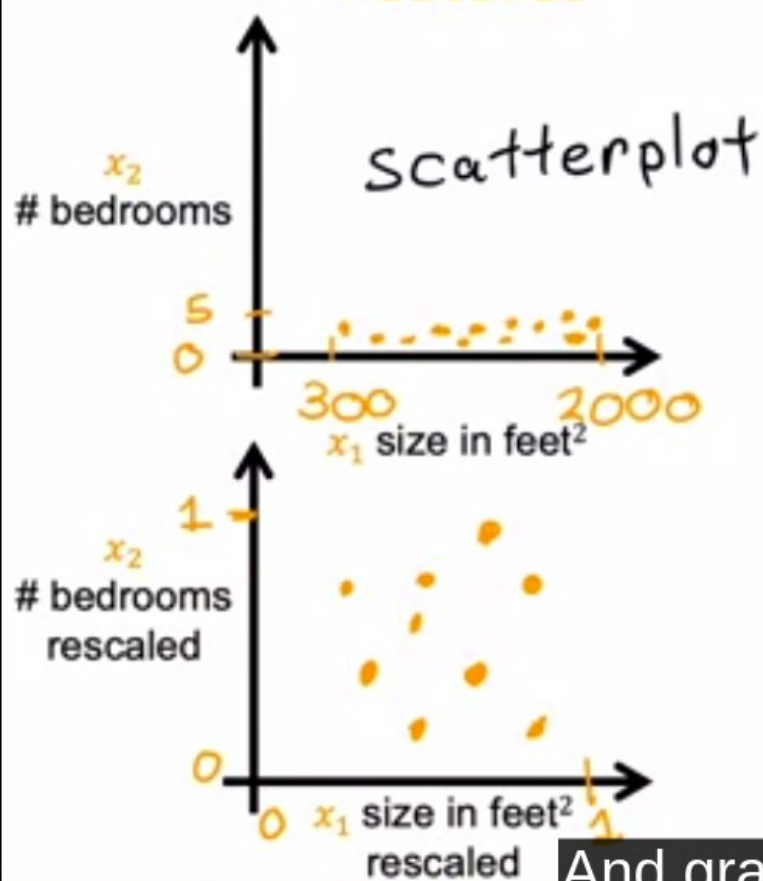
$$\widehat{price} = \underbrace{0.1 * 2000k}_{200K} + \underbrace{50 * 5}_{250K} + \underbrace{50}_{50K}$$

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

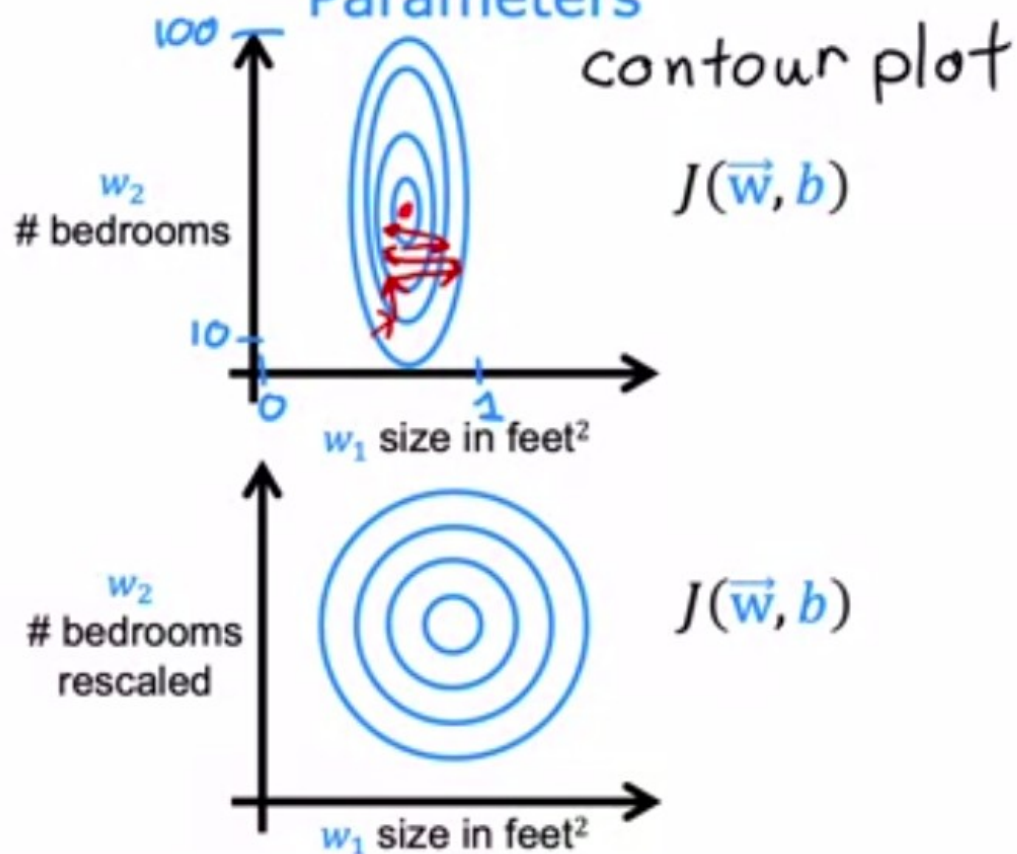
its parameters will be relatively large like 50.

Feature size and gradient descent

Features



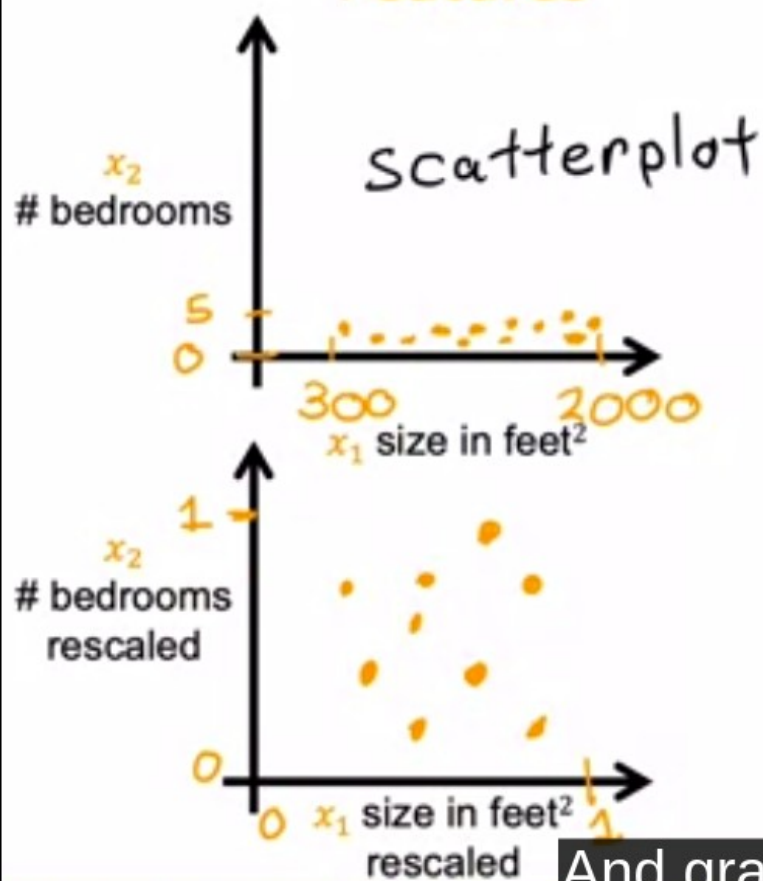
Parameters



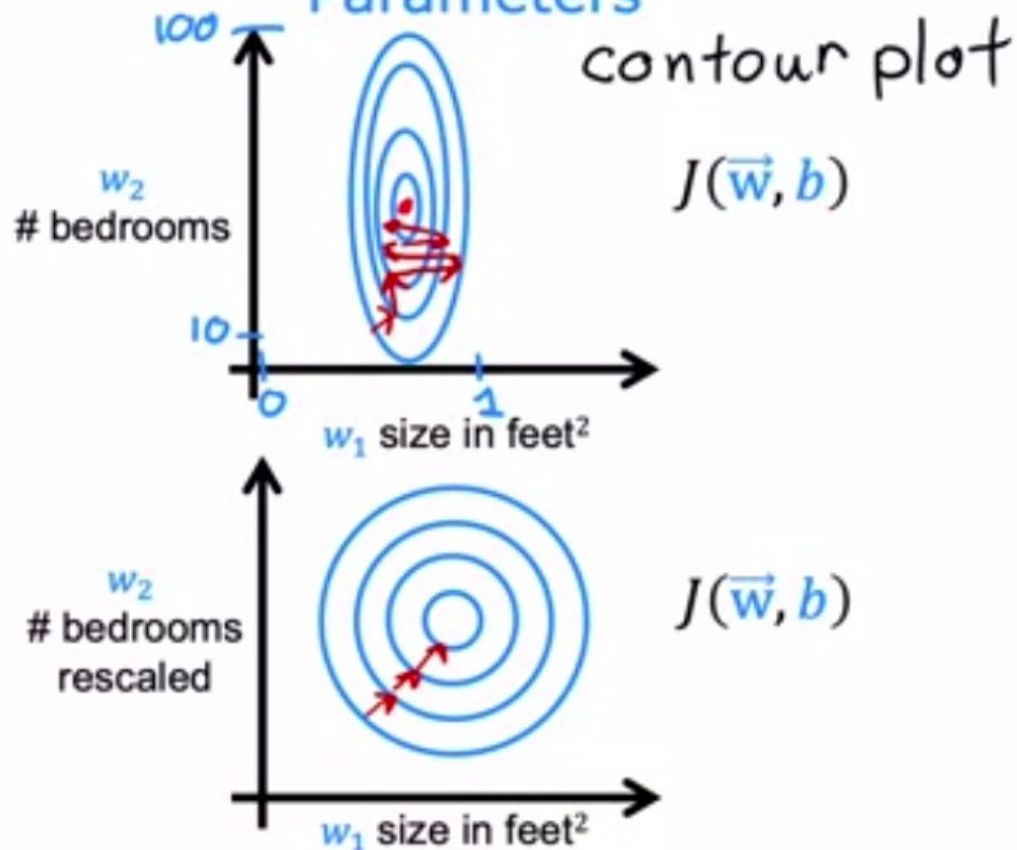
And gradient descent can find a much more direct path to the global minimum.

Feature size and gradient descent

Features

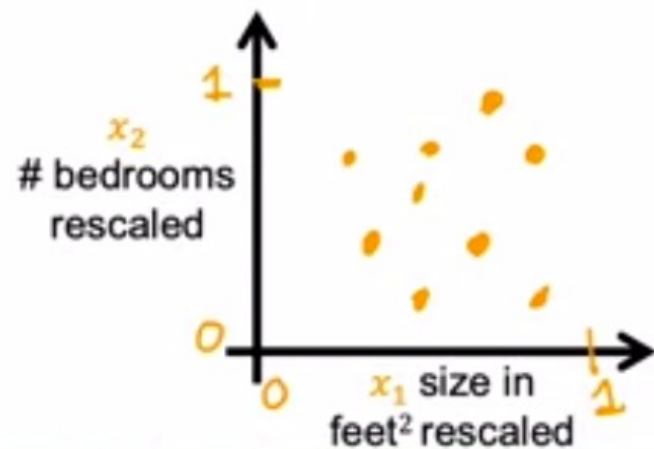
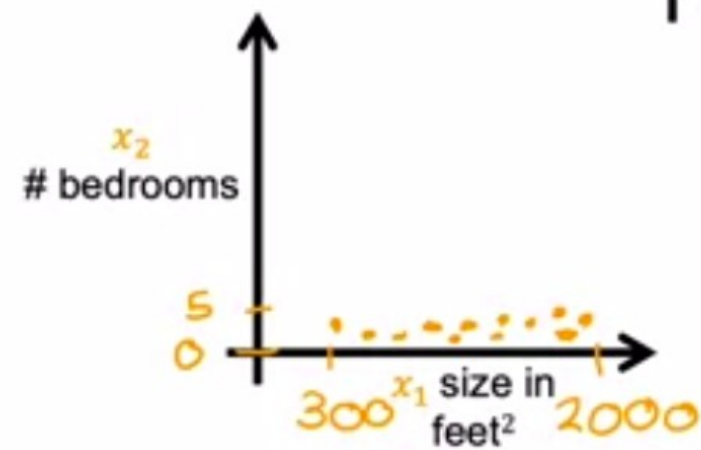


Parameters



And gradient descent can find a much more direct path to the global minimum.

Feature scaling



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

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

max

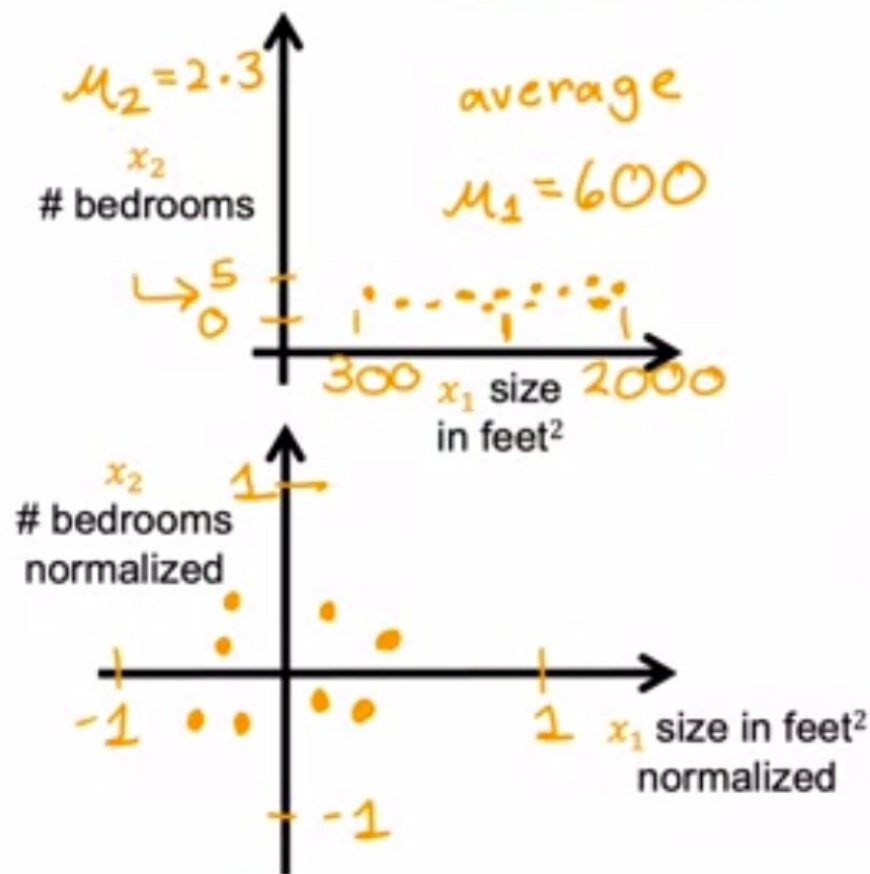
$$x_{2,scaled} = \frac{x_2}{5}$$

max

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

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

Mean normalization



$$300 \leq x_1 \leq 2000$$

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

max-min

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

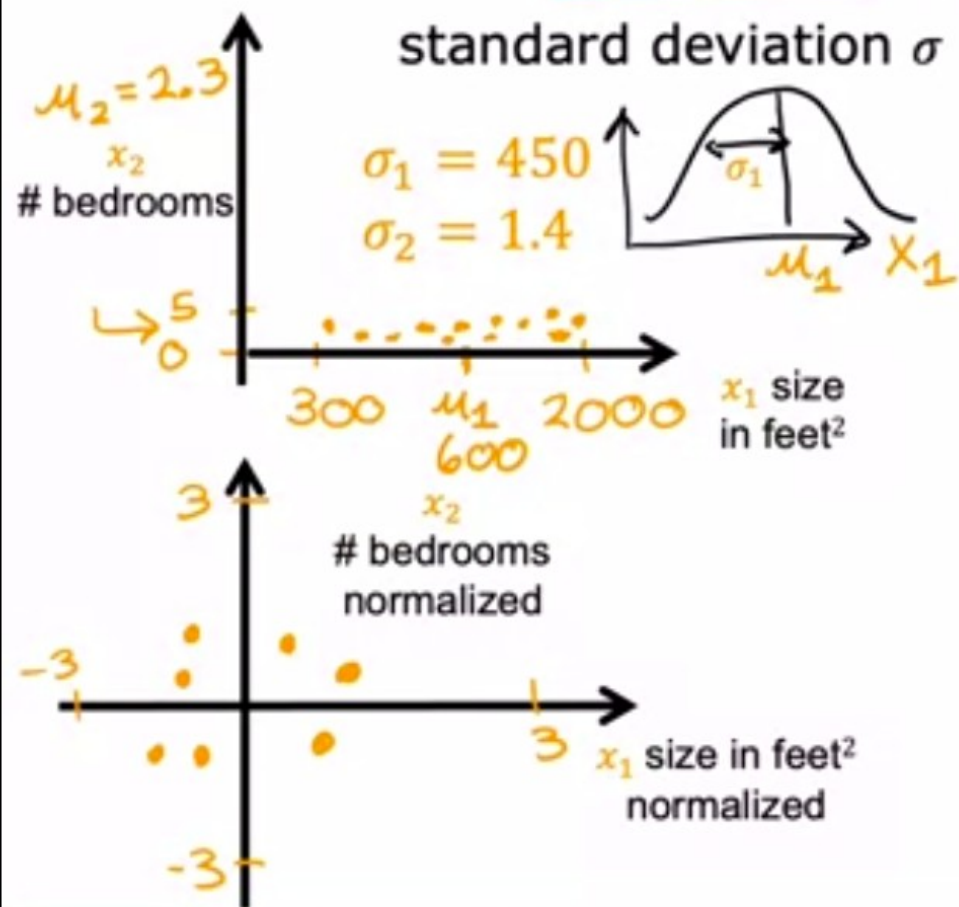
$$0 \leq x_2 \leq 5$$

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

max-min

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

Z-score normalization



$$300 \leq x_1 \leq 2000$$

$$0 \leq x_2 \leq 5$$

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

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

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

Question

x_1 size in
feet² rescaled

Which of the following is a valid step used during feature scaling?

- ☐ Multiply each value by the maximum value for that feature
- ☒ Divide each value by the maximum value for that feature



Correct

By dividing all values by the maximum, the new maximum range of the rescaled features is now 1 (and all other rescaled values are less than 1).

Skip

Continue

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

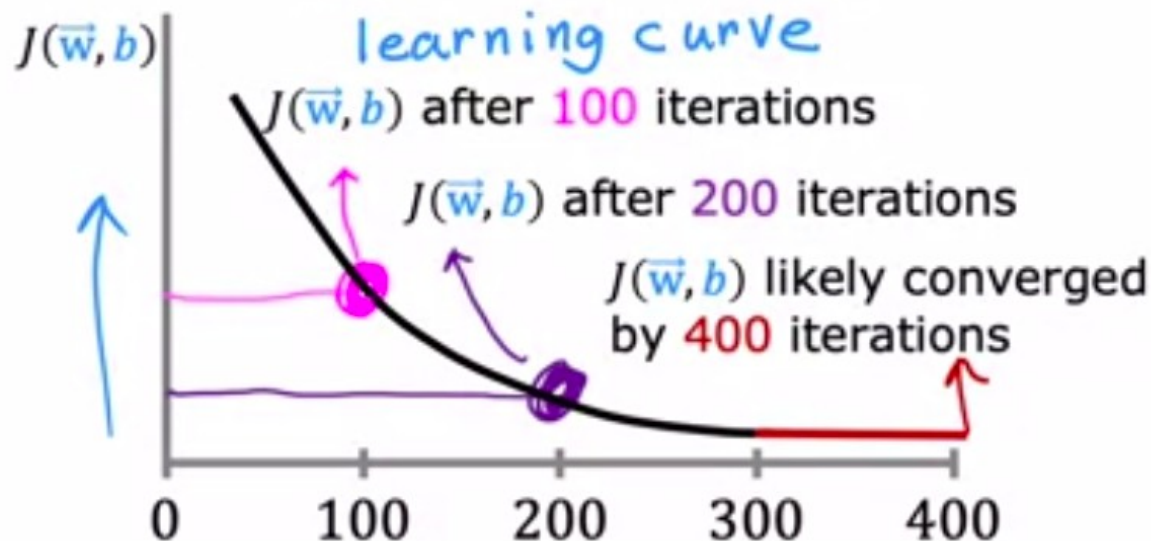
$$98.6 \leq x_5 \leq 105$$

too large \rightarrow rescale

gradient descent to
run much faster.

Make sure gradient descent is working correctly

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



→ # iterations ~~w, b~~

iterations needed varies 30 1,000 100,000

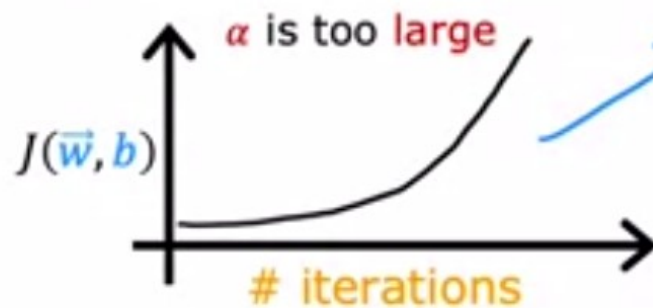
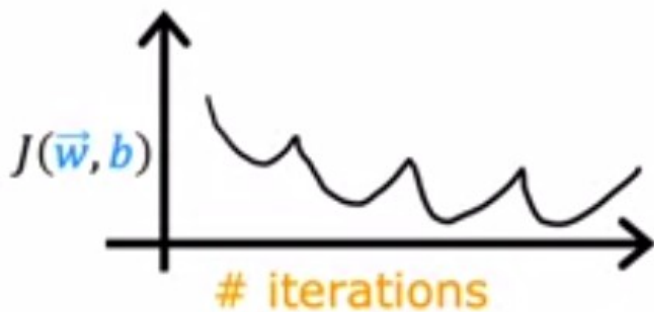
Automatic convergence test

Let ϵ "epsilon" be 10^{-3} .
0.001

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

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

Identify problem with gradient descent



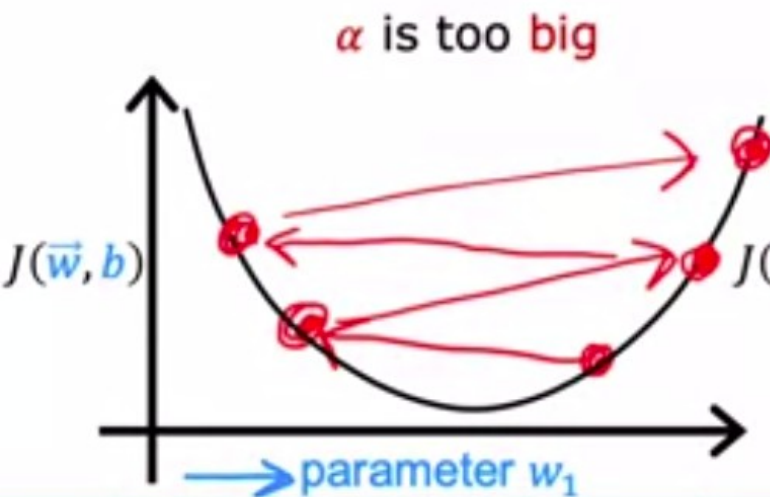
or learning rate is too large

$$w_1 = w_1 + \alpha d_1$$

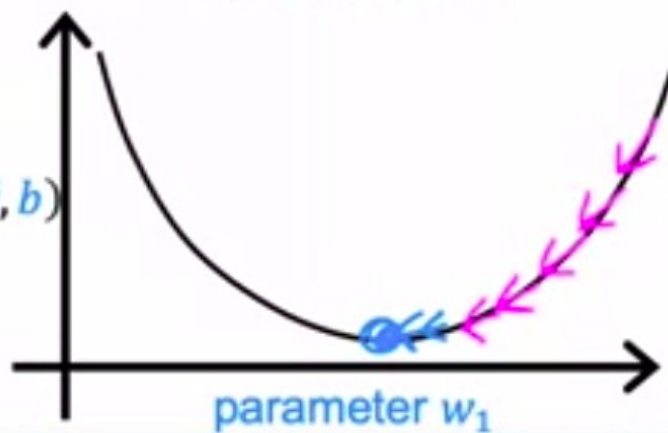
use a minus sign

$$w_1 = w_1 - \alpha d_1$$

Adjust learning rate



Use smaller α

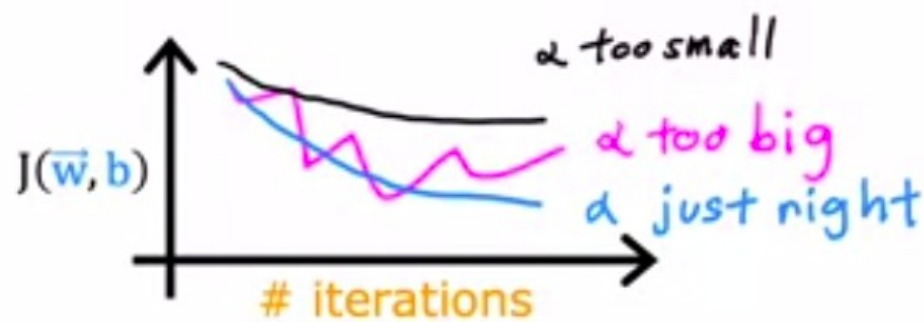
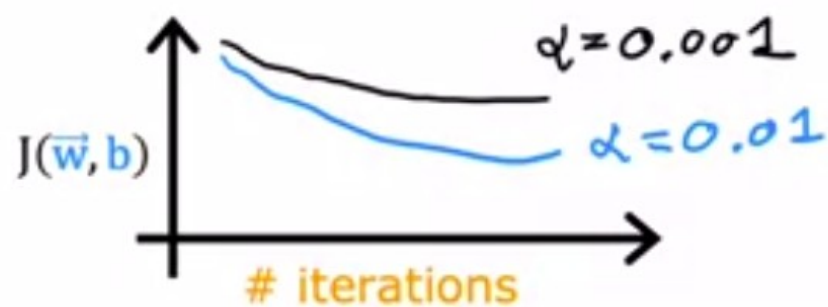


With a small enough α , $J(\bar{w}, b)$ should decrease on every iteration

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
 $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$ $3\times$ $\approx 3\times$



your implementation
of gradient descent.

Question



You run gradient descent for 15 iterations with $\alpha = 0.3$ and compute $J(w)$ after each iteration. You find that the value of $J(w)$ increases over time. How do you think you should adjust the learning rate α ?

- ☐ Try a larger value of α (say $\alpha = 1.0$).
- ☐ Try running it for only 10 iterations so $J(w)$ doesn't increase as much.
- ☐ Keep running it for additional iterations
- ☒ Try a smaller value of α (say $\alpha = 0.1$).



Correct

Since the cost function is increasing, we know that gradient descent is diverging, so we need a lower learning rate.

Skip

Continue

Question

$$x_3 = x_1 x_2$$

new feature

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

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

If you have measurements for the dimensions of a swimming pool (length, width, height), which of the following two would be a more useful engineered feature?

- ☒ $length \times width \times height$
- ☐ $length + width + height$



Correct

The volume of the swimming pool could be a useful feature to use. This is the more useful engineered feature of the two.

Skip

Continue

Feature engineering

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

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

$$x_3 = x_1 x_2$$

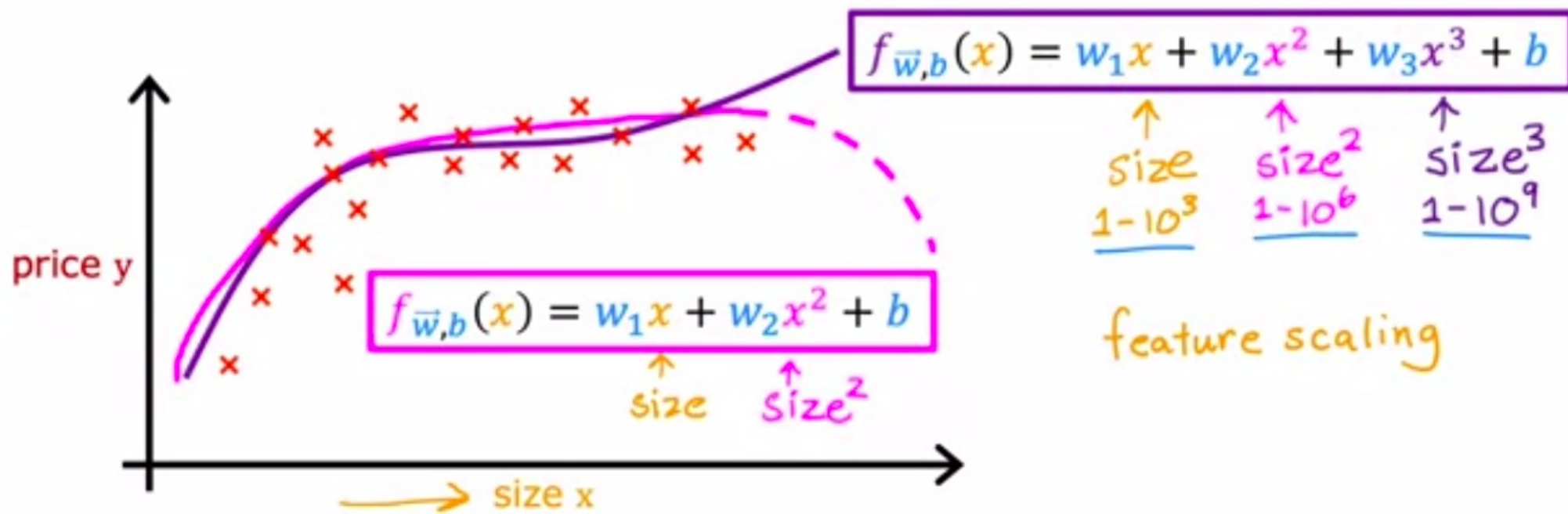
new feature

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



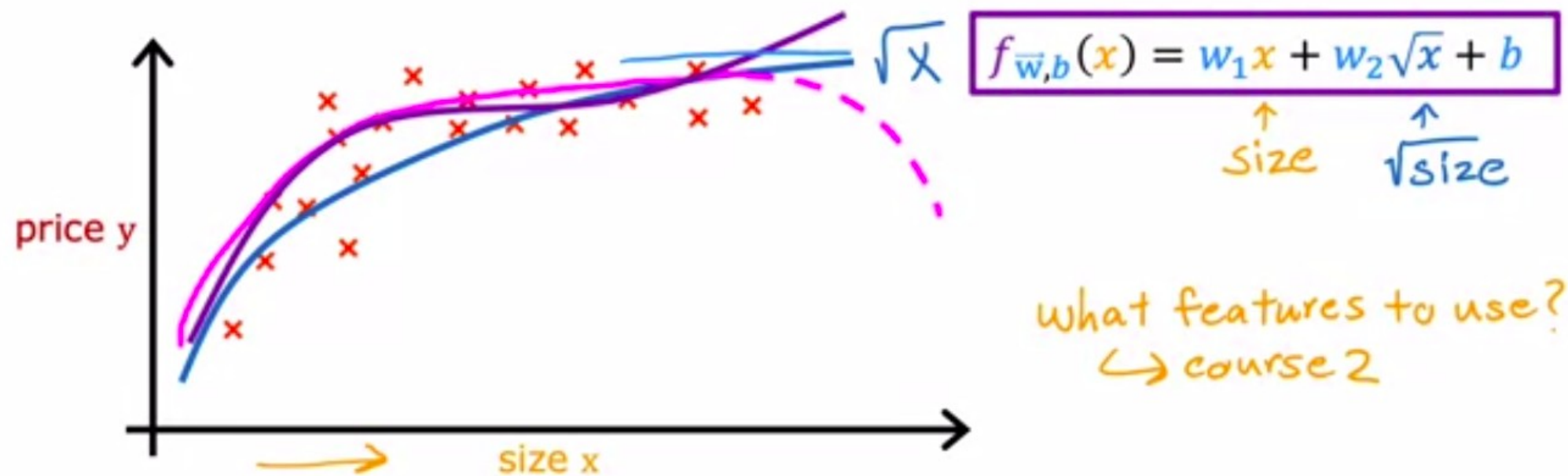
Feature engineering:
Using **intuition** to design
new features, by
transforming or combining
original features.

Polynomial regression



your features into
comparable ranges of values.

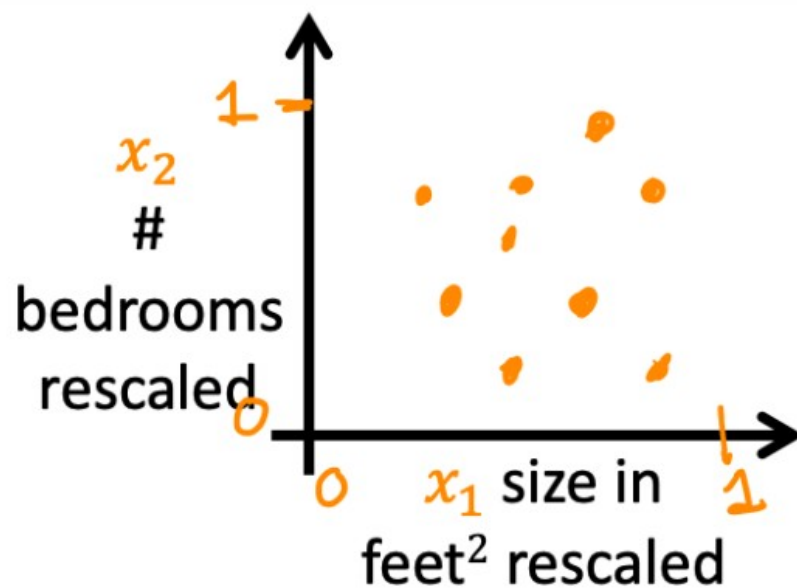
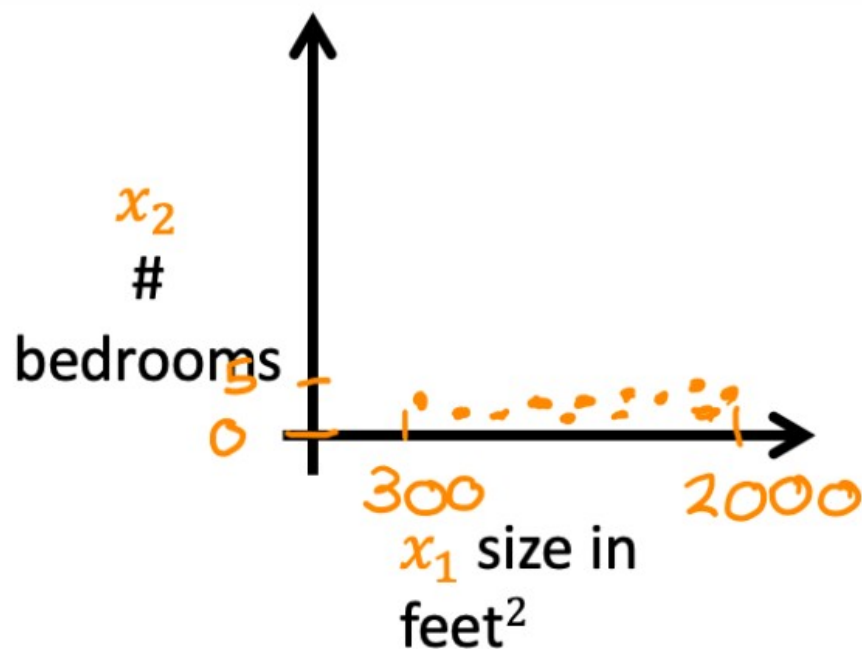
Choice of features



a much better model
for your data.

1.

1 / 1 point



Which of the following is a valid step used during feature scaling?

- ☒ Subtract the mean (average) from each value and then divide by the (max - min).
- ☐ Add the mean (average) from each value and then divide by the (max - min).

✓ Correct

This is called mean normalization.

[Back](#)

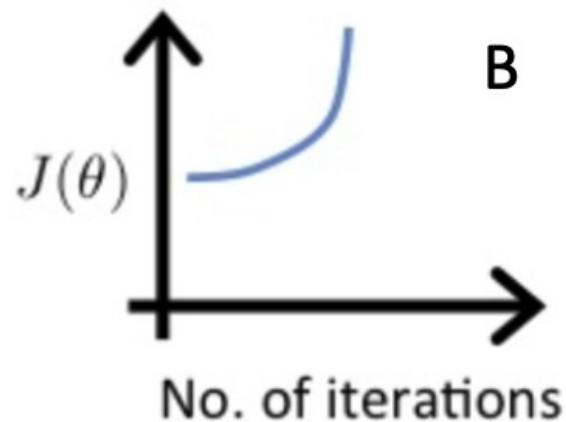
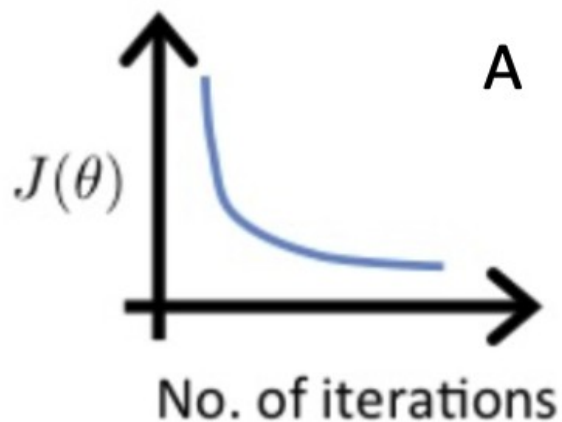
Practice quiz: Gradient descent in practice

Graded Quiz • 30 min

Due Sep 11, 11:59 PM IST

2. Suppose a friend ran gradient descent three separate times with three choices of the learning rate α and plotted the learning curves for each (cost J for each iteration).

1 / 1 point



For which case, A or B, was the learning rate α likely too large?

- ☒ case B only
- ☐ Both Cases A and B
- ☐ case A only
- ☐ Neither Case A nor B

✓ Correct



✓ Correct

The cost is increasing as training continues, which likely indicates that the learning rate α is too large.

3. Of the circumstances below, for which one is feature scaling particularly helpful?

1 / 1 point

- ☐ Feature scaling is helpful when all the features in the original data (before scaling is applied) range from 0 to 1.
- ☒ Feature scaling is helpful when one feature is much larger (or smaller) than another feature.

✓ Correct

For example, the “house size” in square feet may be as high as 2,000, which is much larger than the feature “number of bedrooms” having a value between 1 and 5 for most houses in the modern era.

4.

1 / 1 point

You are helping a grocery store predict its revenue, and have data on its items sold per week, and price per item. What could be a useful engineered feature?

- ☒ For each product, calculate the number of items sold times price per item.
- ☐ For each product, calculate the number of items sold divided by the price per item.

✓ Correct

This feature can be interpreted as the revenue generated for each product.

[Back](#)

Practice quiz: Gradient descent in practice

Graded Quiz • 30 min

Due Sep 11, 11:59 PM IST



Correct

For example, the “house size” in square feet may be as high as 2,000, which is much larger than the feature “number of bedrooms” having a value between 1 and 5 for most houses in the modern era.

4.

1 / 1 point

You are helping a grocery store predict its revenue, and have data on its items sold per week, and price per item. What could be a useful engineered feature?

- ☒ For each product, calculate the number of items sold times price per item.
- ☐ For each product, calculate the number of items sold divided by the price per item.



Correct

This feature can be interpreted as the revenue generated for each product.

5. True/False? With polynomial regression, the predicted values $f_{w,b}(x)$ does not necessarily have to be a straight line (or linear) function of the input feature x .

1 / 1 point

- ☐ False
- ☒ True



Correct

A polynomial function can be non-linear. This can potentially help the model to fit the training data better.