1. Multiple features

# Multiple features (variables)

| Size in feet² | 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

$\vec{x}^{(2)} = \begin{bmatrix} 1416 & 3 & ② & 40 \end{bmatrix}$

$x_3^{(2)} = 2$

## Model:

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

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

size  #bedrooms  #floors  years  base price

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

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

$\vec{w} = \begin{bmatrix} w_1 & w_2 & w_3 & \dots & w_n \end{bmatrix}$   parameters of the model

$b$ is a number

vector $\vec{x} = \begin{bmatrix} x_1 & x_2 & x_3 & \dots & x_n \end{bmatrix}$

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

dot product

multiple linear regression

(not multivariate regression)

2. Vectorization part 1 & 2

## Parameters and features

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

$b$ is a number

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

linear algebra: count from 1

**NumPy**

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

```
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

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

$\sum_{j=1}^{n} \rightarrow j=1...n$

$\quad 1,2,3$

$range(0,n) \rightarrow j = 0...n-1$

```
f = 0          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
```

## Without vectorization

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

$t_0$
$\quad f + w[0] * x[0]$

$t_1$
$\quad f + w[1] * x[1]$

...

$t_{15}$
$\quad f + w[15] * x[15]$

## Vectorization

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

$t_0$

| w[0] | w[1] | ... | w[15] |
|------|------|-----|-------|

in parallel  $*$ $\quad$ $*$ $\quad$ ... $\quad$ $*$

| x[0] | x[1] | ... | x[15] |
|------|------|-----|-------|

$t_1$

| w[0]*x[0] | + | w[1]*x[1] | +...+ | w[15]*x[15] |
|-----------|---|-----------|-------|-------------|

efficient → scale to large datasets

## Gradient descent

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

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

```
w = np.array([0.5, 1.3, … 3.4])
d = np.array([0.3, 0.2, … 0.4])
```
learning rate $\alpha$

compute $w_j = w_j - 0.1d_j$ for $j = 1...16$

| Without vectorization | With vectorization |
|---|---|
| $w_1 = w_1 - 0.1d_1$ <br> $w_2 = w_2 - 0.1d_2$ <br> $\vdots$ <br> $w_{16} = w_{16} - 0.1d_{16}$ | $\vec{w} = \vec{w} - 0.1\vec{d}$ |
| `for j in range(0,16):` <br> `    w[j] = w[j] - 0.1 * d[j]` | `w = w - 0.1 * d` |

$w$

$0.1 * d$

3. Gradient descent for multiple linear regression

|  | Previous notation | Vector notation |
|---|---|---|
| Parameters | $w_1, \cdots, w_n$ <br> $b$ | $\vec{w} = [w_1 \ \cdots \ w_n]$ ↙ vector of length n <br> $b$ still a number |
| Model | $f_{\vec{w},b}(\vec{x}) = w_1 x_1 + \cdots + w_n x_n + b$ | $f_{\vec{w},b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$ <br> ↑ dot product |
| Cost function | $J(w_1, \cdots, w_n, b)$ | $J(\vec{w}, b)$ |

Gradient descent

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

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

### One feature

repeat {
$$w = w - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{w,b}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$
$$\hookrightarrow \frac{\partial}{\partial w} J(w, b)$$

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

simultaneously update $w, b$
}

### $n$ features ($n \geq 2$)

repeat {
j=1
$$w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right) x_1^{(i)}$$
$$\vdots \qquad \hookrightarrow \frac{\partial}{\partial w_1} J(\vec{w}, b)$$
j=n
$$w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^{m} \left( f_{\vec{w},b}(\vec{x}^{(i)}) - y^{(i)} \right) x_n^{(i)}$$

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

simultaneously update
$w_j$ (for $j = 1, \cdots, n$) and $b$
}

## 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

4. Practice quiz 1

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).

| Size in feet² | 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 |
| ... | ... | ... | ... | ... |

> 30

✓ **Correct**
Yes! $x_4^{(3)}$ is the 4th feature (4th column in the table) of the 3rd training example (3rd row in the table).

2.

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 $alpha$.

○ True

● False

✓ **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 $b$.

5. Feature scaling part 1

# Feature and parameter values

$$\widehat{price} = w_1 x_1 + w_2 x_2 + b$$
(size)   (#bedrooms)

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

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

size of the parameters $w_1, w_2$?

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

$$\widehat{price} = \underbrace{50 * 2000}_{100,000K} + \underbrace{0.1 * 5}_{0.5k} + \underbrace{50}_{50K}$$   |   $$\widehat{price} = \underbrace{0.1 * 2000k}_{200K} + \underbrace{50 * 5}_{250K} + \underbrace{50}_{50K}$$

$\widehat{price} = \$100,050.5\underline{k} = \$100,050,500$  |  $\widehat{price} = \$500k$   more reasonable
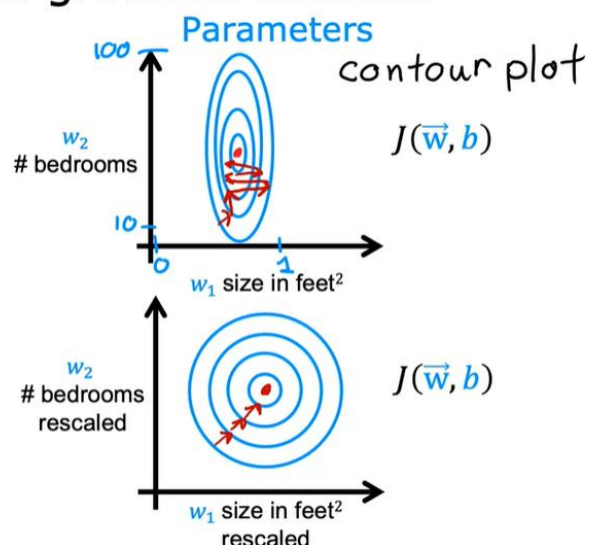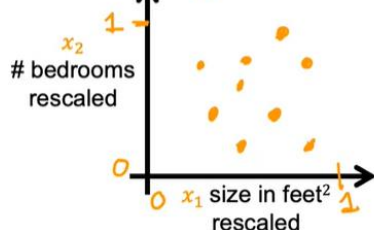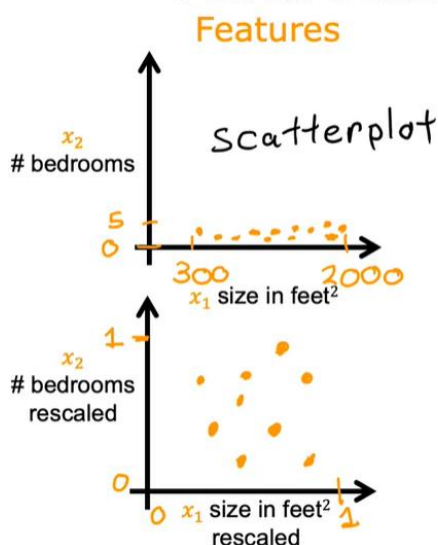
# Feature size and parameter size

|            | size of feature $x_j$ | size of parameter $w_j$ |
|------------|:---------------------:|:-----------------------:|
| size in feet² | ⟷ (large) | ⟷ (small) |
| #bedrooms  | ⟷ (small) | ⟷ (large) |



Features — Scatterplot
$x_2$ # bedrooms (5, 0)  vs  $x_1$ size in feet² (300, 2000)

Parameters — $J(\vec{w}, b)$ contour plot
$w_2$ # bedrooms (100, 10)  vs  $w_1$ size in feet² (1)

# Feature size and gradient descent



Features — Scatterplot
$x_2$ # bedrooms (5, 0)  vs  $x_1$ size in feet² (300, 2000)

Parameters — contour plot $J(\vec{w}, b)$
$w_2$ # bedrooms (100, 10)  vs  $w_1$ size in feet² (1)

$x_2$ # bedrooms rescaled (1, 0)  vs  $x_1$ size in feet² rescaled (0, 1)

Parameters $J(\vec{w}, b)$
$w_2$ # bedrooms rescaled  vs  $w_1$ size in feet² rescaled

6. Feature scaling part 2

# Feature scaling



$$300 \leq x_1 \leq 2000 \qquad 0 \leq x_2 \leq 5$$
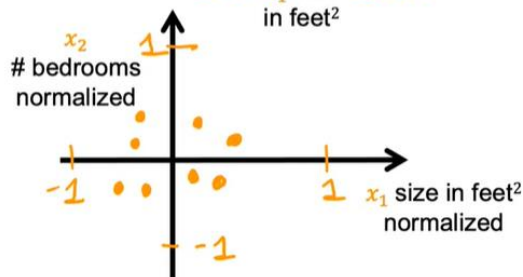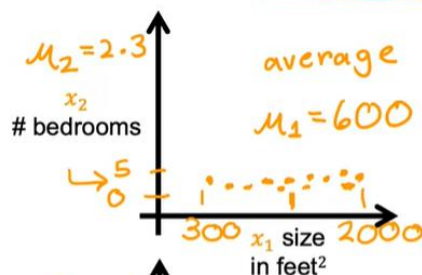
$$x_{1,scaled} = \frac{x_1}{2000}_{max} \qquad x_{2,scaled} = \frac{x_2}{5}_{max}$$

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

# Mean normalization

$\mu_2 = 2.3$

average

$\mu_1 = 600$

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

$$x_1 = \frac{x_1 - \mu_1}{2000 - 300}_{max-min} \qquad x_2 = \frac{x_2 - \mu_2}{5 - 0}_{max-min}$$
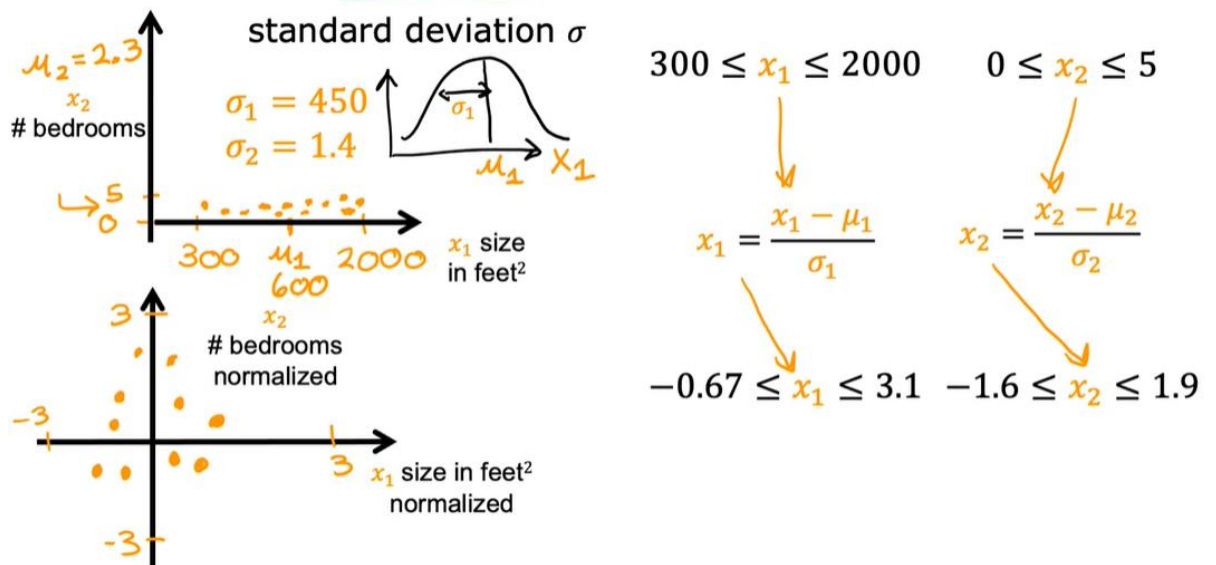
$$-0.18 \leq x_1 \leq 0.82 \qquad -0.46 \leq x_2 \leq 0.54$$

# Z-score normalization

## standard deviation $\sigma$

$\mu_2 = 2.3$

$x_2$
# bedrooms

$\sigma_1 = 450$
$\sigma_2 = 1.4$



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

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1} \qquad x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

$-0.67 \leq x_1 \leq 3.1 \quad -1.6 \leq x_2 \leq 1.9$
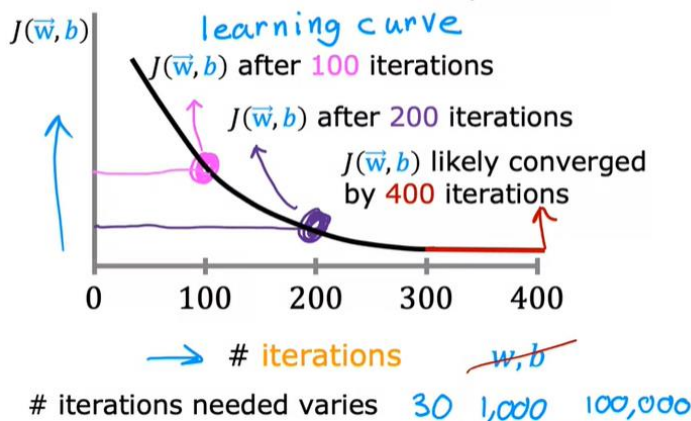
# Feature scaling

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 → rescale

$-0.001 \leq x_4 \leq 0.001$    too small → rescale

$98.6 \leq x_5 \leq 105$    too large → rescale

7.  Checking gradient descent for convergence

# Make sure gradient descent is working correctly

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

**learning curve**

$J(\vec{w}, b)$ after 100 iterations

$J(\vec{w}, b)$ after 200 iterations

$J(\vec{w}, b)$ likely converged by 400 iterations

0    100    200    300    400

→ # iterations    $w, b$

# iterations needed varies    30  1,000  100,000

Automatic convergence test
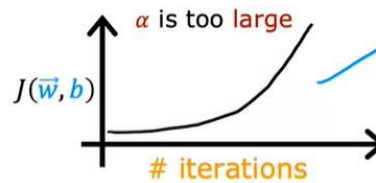Let $\varepsilon$ "epsilon" be $10^{-3}$.
                              0.001

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

(found parameters $\vec{w}, b$ to get close to global minimum)
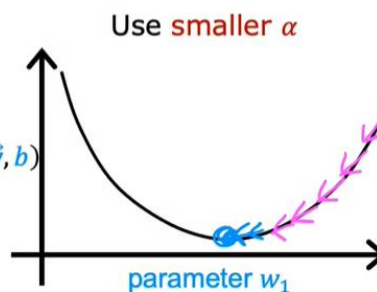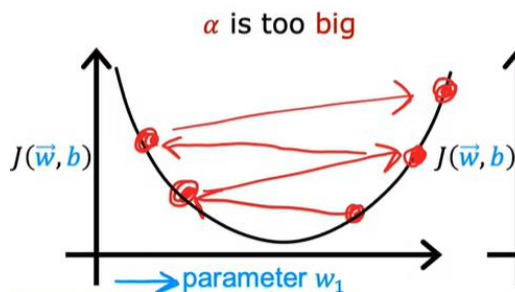
8. Choosing the learning rate

## Identify problem with gradient descent

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

$J(\vec{w}, b)$    $\alpha$ is too large    # iterations

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

$\alpha$ is too big

$J(\vec{w}, b)$    parameter $w_1$

Use smaller $\alpha$

$J(\vec{w}, b)$    parameter $w_1$

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

If $\alpha$ is too small, gradient descent takes a lot more iterations to converge

9. Feature engineering

# Feature engineering

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

$area = frontage \times depth$

$x_3 = x_1 x_2$
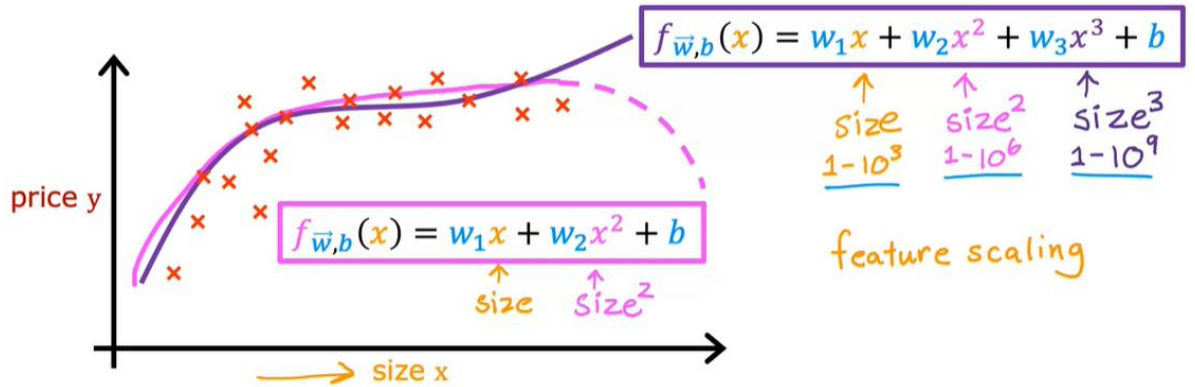new feature

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

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

10. Polynomial regression
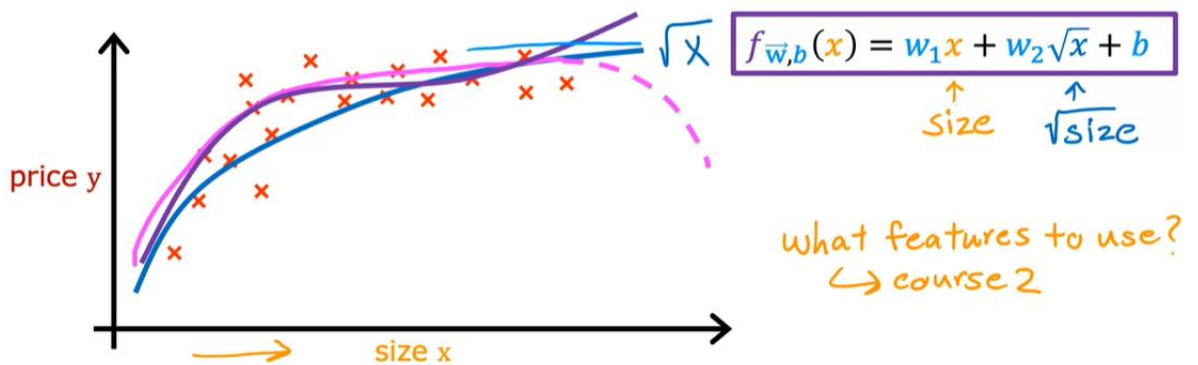
# Polynomial regression



$$f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + b$$

size  size$^2$  size$^3$
1-10$^3$  1-10$^6$  1-10$^9$

feature scaling

$$f_{\vec{w},b}(x) = w_1 x + w_2 x^2 + b$$

size  size$^2$

price y

size x

## Choice of features



$$\sqrt{x} \quad f_{\vec{w},b}(x) = w_1 x + w_2 \sqrt{x} + b$$

size  $\sqrt{size}$

what features to use?
↳ course 2

price y

size x

11. Practice quiz

1.

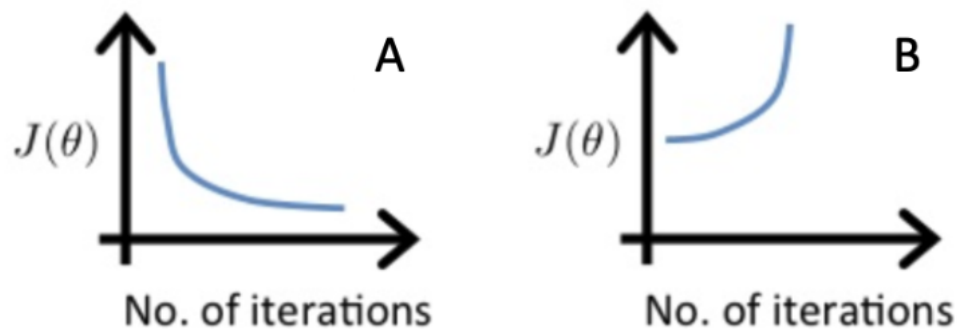

Which of the following is a valid step used during feature scaling?

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

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

✓ **Correct**
   This is called mean normalization.

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



For which case, A or B, was the learning rate $\alpha$ likely too large?

- ● case B only
- ○ Neither Case A nor B
- ○ case A only
- ○ Both Cases A and B

> ✓ **Correct**
> The cost is increasing as training continues, which likely indicates that the learning rate alpha is too large.

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

- ○ 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.

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

○ False

◉ True

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