# Copyright Notice

# Multiple features (variables)

one
feature

| Size in feet$^2$ ($x$) | Price ($) in 1000's ($y$) |
| --- | --- |
| 2104 | 400 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

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

# 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 |
| ... | ... | ... | ... | ... |

$i=2$

$j=1...4$

$n=4$

$x_j = j^{th}$ 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

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

$x_3^{(2)} = 2$

# 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

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

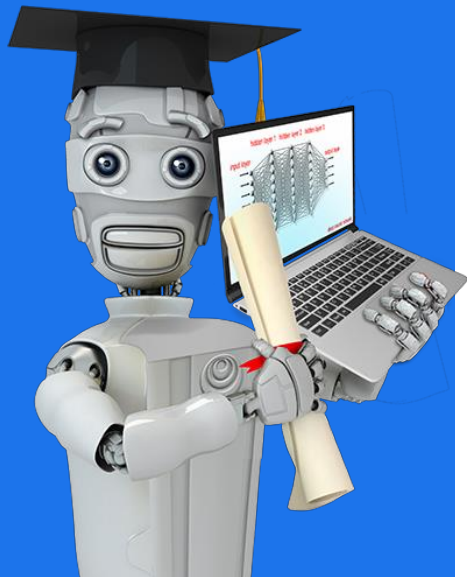$$f_{\vec{w},b}(\vec{x}) = w_1 x_1 + w_2 x_2 + \cdots + w_n x_n + b$$

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

parameters of the model

b is a number

vector $\vec{x} = [x_1 \; x_2 \; x_3 \ldots x_n]$

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

dot product

multiple linear regression

(not multivariate regression)

# 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

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

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

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

$$1, 2, 3$$

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

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

# Without vectorization  $n = 100,000$

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

```
f = w[0] * x[0] +
    w[1] * x[1] +
    w[2] * x[2] + 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]
```
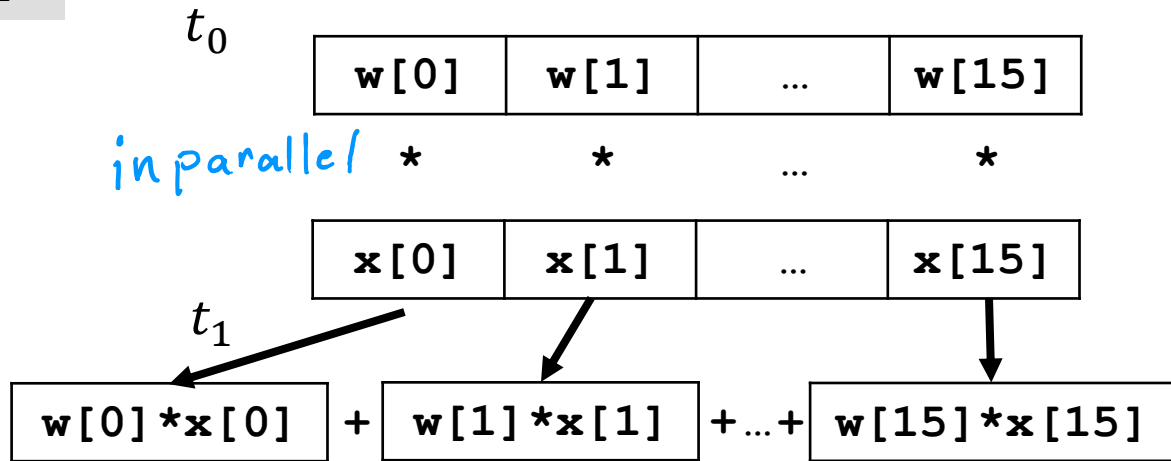
$t_0$

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

$t_1$

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

...

$t_{15}$

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

# Vectorization

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

$t_0$

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

in parallel    *      *     ...     *

| 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

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

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

```
w = np.array([0.5, 1.3, … 3.4])
d = np.array([0.3, 0.2, … 0.4])
```

learning rate $\alpha$

compute $w_j = w_j - 0.1 d_j$ for $j = 1 \ldots 16$

## Without vectorization

$$w_1 = w_1 - 0.1 d_1$$
$$w_2 = w_2 - 0.1 d_2$$
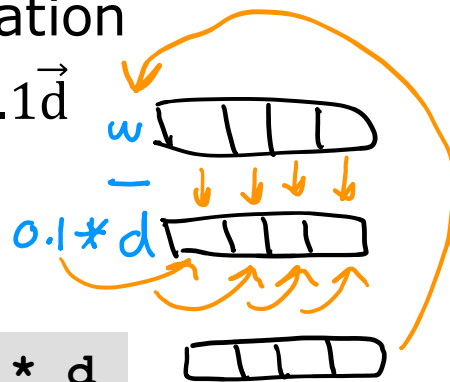$$\vdots$$
$$w_{16} = w_{16} - 0.1 d_{16}$$

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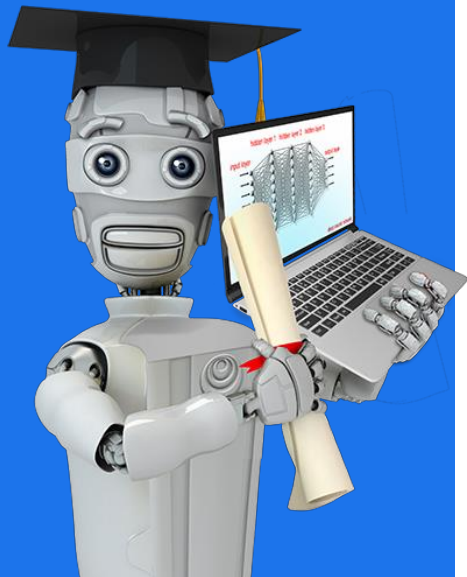
## With vectorization

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



$w$
$0.1 * d$

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

|                | Previous notation | Vector notation |
|----------------|-------------------|-----------------|

**Parameters**

$$w_1, \cdots, w_n$$
$$b$$

$\vec{\mathbf{w}} = [w_1 \quad \cdots \quad w_n]$ ← vector of length n

$b$ still a number

**Model**

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

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

dot product

**Cost function**

$$J(w_1, \cdots, w_n, b)$$

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

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

# Gradient descent

$$y = u^n \longrightarrow y' = n.u'.u^{(n-1)}$$

## One feature

repeat {

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

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

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

simultaneously update $w$, $b$

}

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

repeat {

$j=1$

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

$$\vdots$$

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

$j=n$

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

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

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)
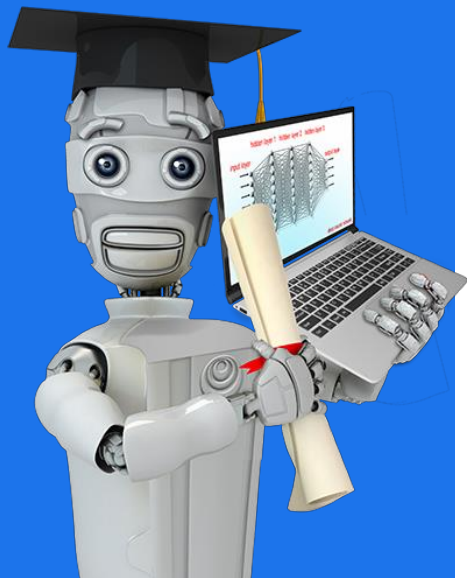
What you need to know
- Normal equation method may be used in machine learning libraries that implement linear regression.
- Gradient descent is the recommended method for finding parameters w,b

https://www.geeksforgeeks.org/ml-normal-equation-in-linear-regression/
https://prutor.ai/normal-equation-in-linear-regression/

# Feature and parameter values

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

size   #bedrooms

$x_1$: size (feet²)
range: $300 - 2,000$   large   small

$x_2$: # bedrooms
range: $0 - 5$

House: $x_1 = 2000$, $x_2 = 5$, $price = \$500k$   one training example
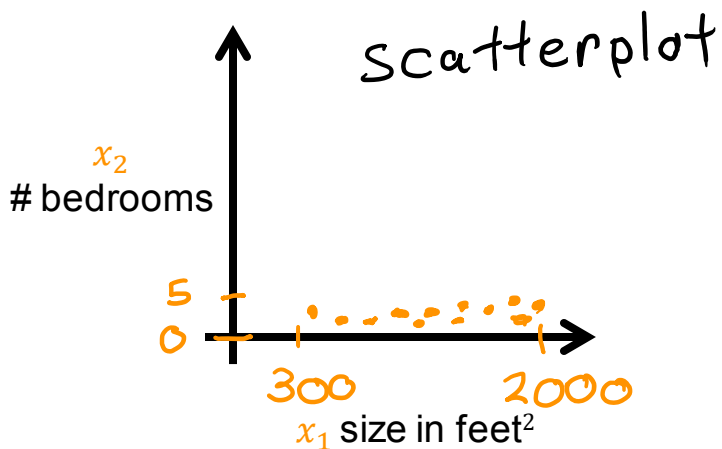
size of the parameters $w_1, w_2$?

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

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

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

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

small   large

$\widehat{price} = \underbrace{0.1 * 2000k}_{200K} + \underbrace{50 * 5}_{250K} + \underbrace{50}_{50k}$
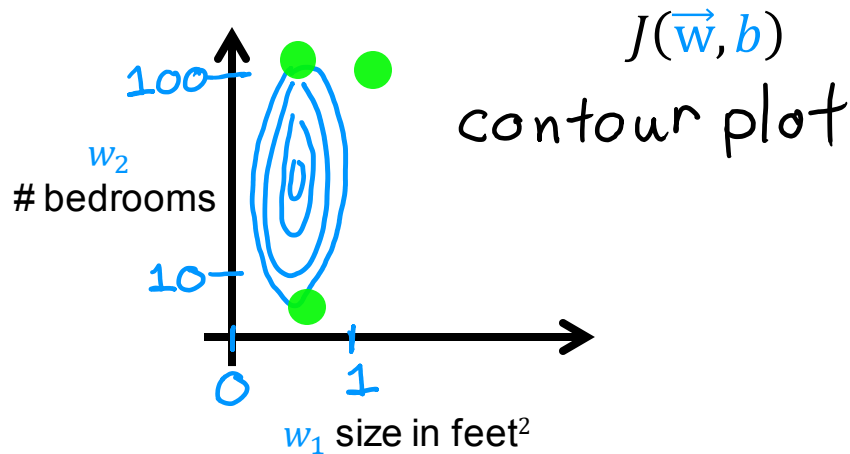
$\widehat{price} = \$500k$   more reasonable

# Feature size and parameter size

|  | size of feature $x_j$ | size of parameter $w_j$ |
|---|---|---|
| size in feet$^2$ | ⟷ | ⟷ |
| #bedrooms | ⟷ | ⟷ |

## Features

scatterplot

$x_2$
# bedrooms

5
0

300    2000

$x_1$ size in feet$^2$

## Parameters

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

contour plot

100

$w_2$
# bedrooms

10

0    1

$w_1$ size in feet$^2$

# Feature size and gradient descent



Features

Parameters

contour plot

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

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

# Feature scaling



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

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

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

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

# Mean normalization



$\mu_2 = 2.3$

average

$\mu_1 = 600$

$x_2$ # bedrooms

$x_1$ size in feet$^2$

$x_2$ # bedrooms normalized

$x_1$ size in feet$^2$ normalized

$300 \leq x_1 \leq 2000$

$0 \leq x_2 \leq 5$

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

max - min

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

max - min

$-0.18 \leq x_1 \leq 0.82$

$-0.46 \leq x_2 \leq 0.54$

# Z-score normalization

standard deviation $\sigma$

$\mu_2 = 2.3$

$x_2$
# bedrooms

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



$\mu_1$ $X_1$

$x_1$ size
in feet$^2$

5
0

300  $\mu_1$  2000
      600

$x_2$
# bedrooms
normalized

3

$-3$

$-3$

3

$x_1$ size in feet$^2$
normalized

$-3$

$300 \leq x_1 \leq 2000$

$0 \leq x_2 \leq 5$

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

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

$-0.67 \leq x_1 \leq 3.1$

$-1.6 \leq x_2 \leq 1.9$

# Feature scaling

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

$$-3 \leq x_j \leq 3$$
$$-0.3 \leq x_j \leq 0.3$$

$\Big\}$ 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

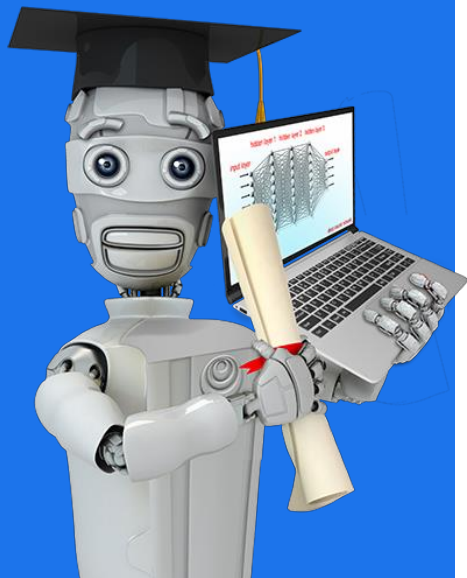$98.6 \leq x_5 \leq 105$     too large → rescale

# Gradient descent

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

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

# Make sure gradient descent is working correctly

objective: $\min_{\vec{w},b} J(\vec{w}, b)$

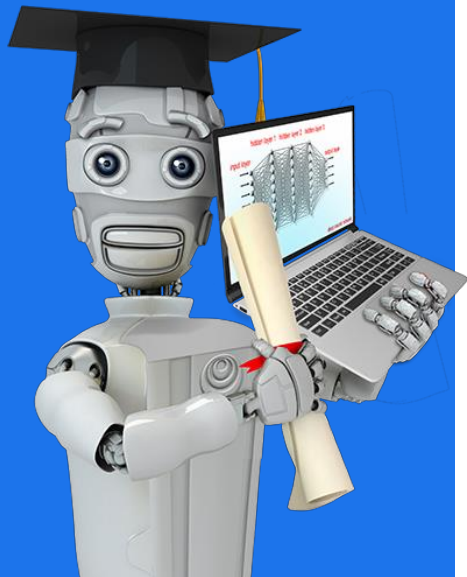$J(\vec{w}, b)$ should **decrease** after every iteration

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

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

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)

# Identify problem with gradient descent
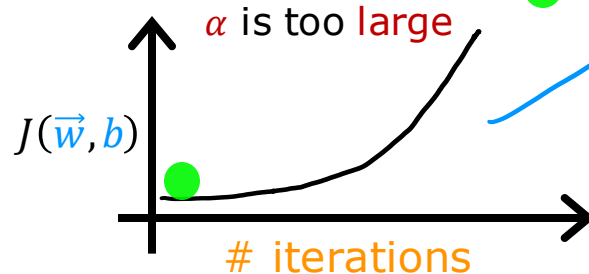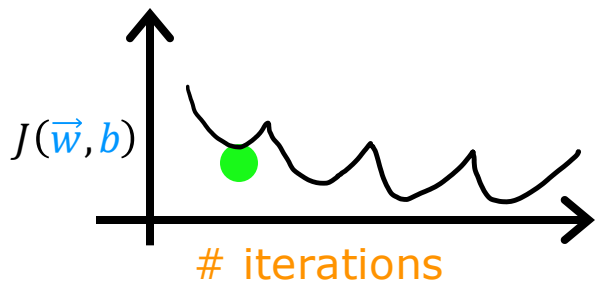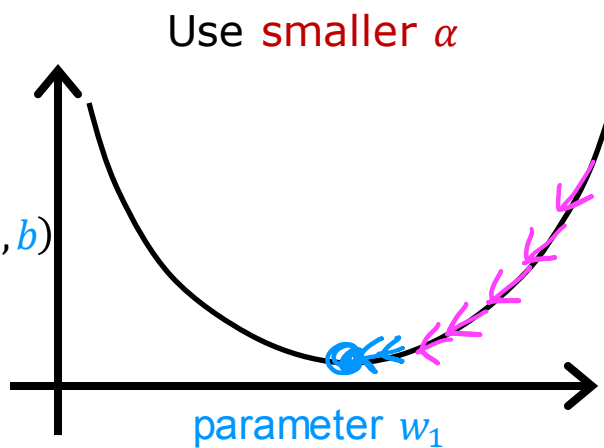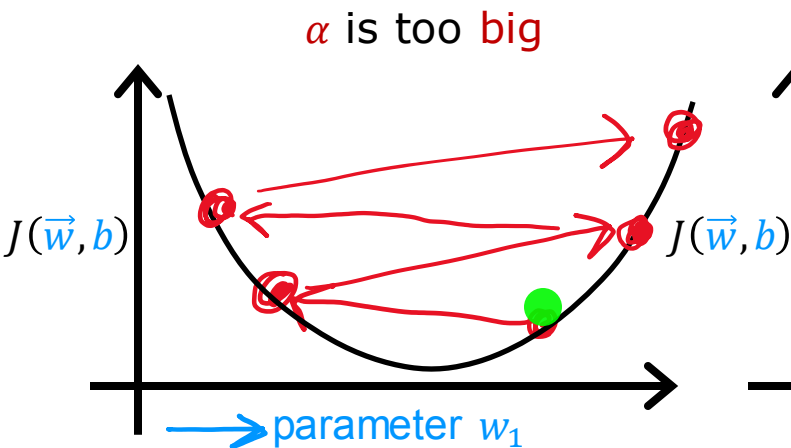


$J(\vec{w}, b)$

# iterations

$\alpha$ is too large

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

# 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 $w_1$
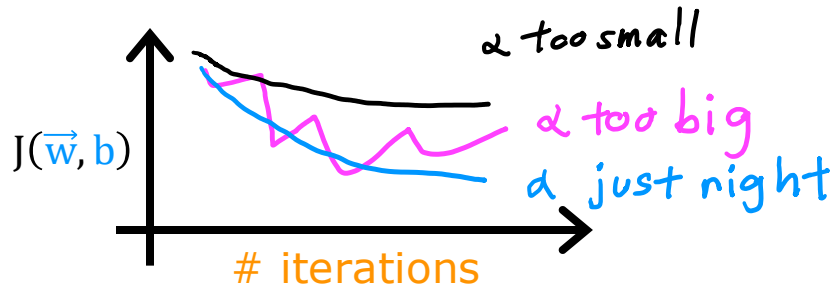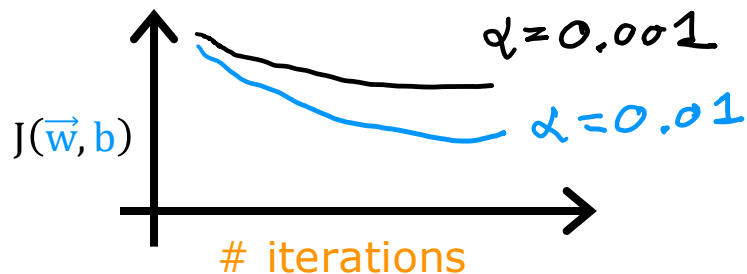
Use smaller $\alpha$

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

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
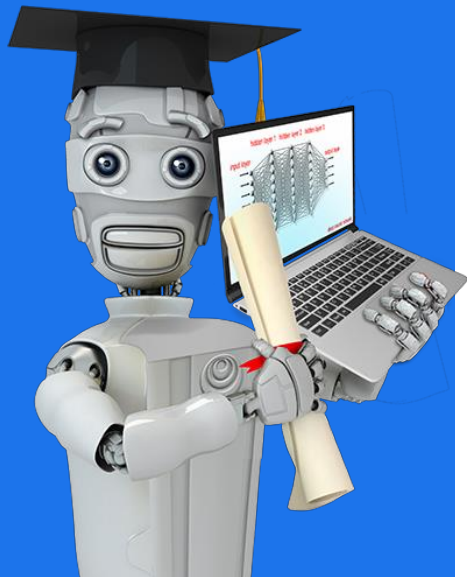
Values of $\alpha$ to try:

... 0.001 0.003   0.01 0.03   0.1   0.3   1 ...

3X   ≈3X   3X   ≈3X   3X   ≈3X



$J(\vec{w}, b)$                                    $\alpha = 0.001$

$\alpha = 0.01$

# iterations

$J(\vec{w}, b)$          $\alpha$ too small

$\alpha$ too big

$\alpha$ just right

# iterations

# Feature engineering

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

$$area = frontage \times depth$$

$$x_3 = x_1 x_2$$
new feature

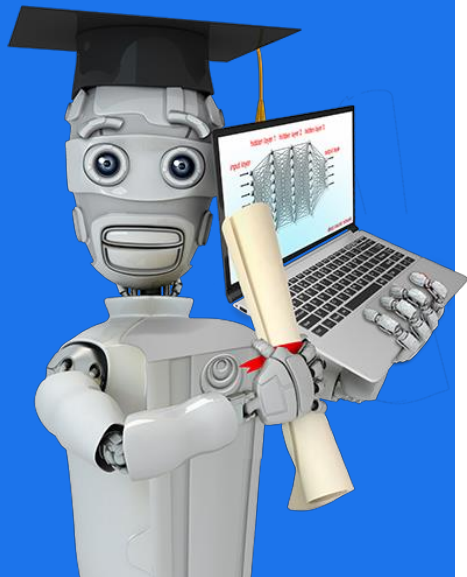$$f_{\vec{w},b}(\vec{x}) = \underbrace{w_1} x_1 + \underbrace{w_2} x_2 + \underbrace{w_3} x_3 + b$$



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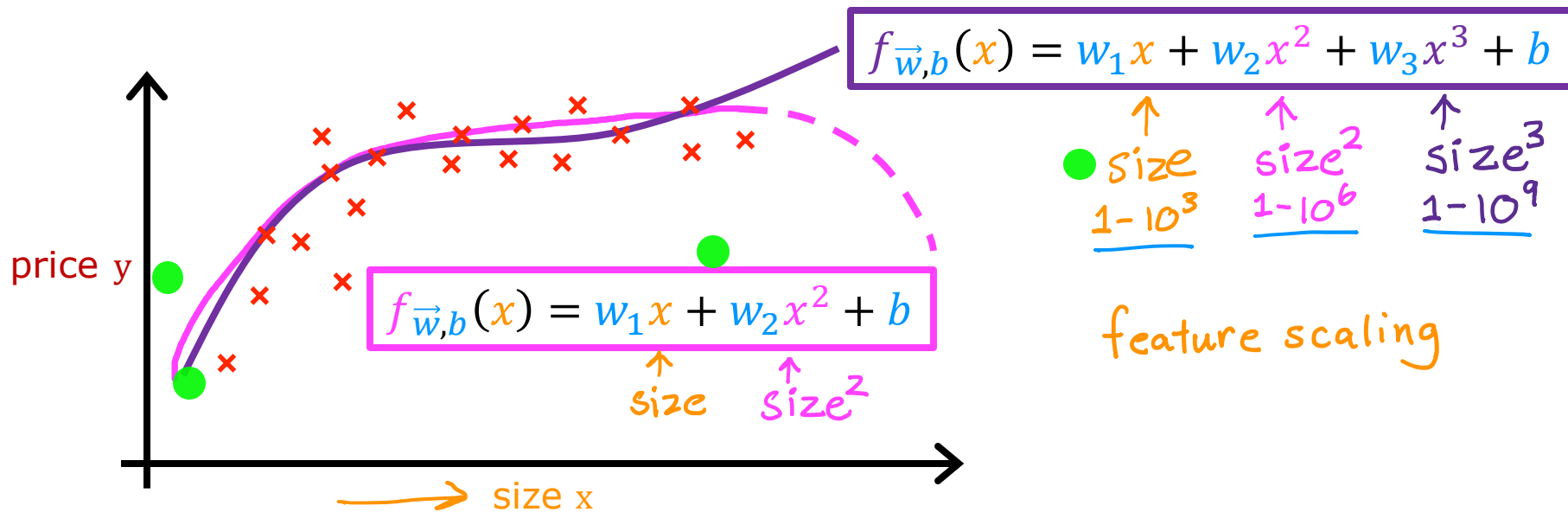Using intuition to design new features, by transforming or combining original features.

# Polynomial regression



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

● size        size²        size³
1−10³       1−10⁶       1−10⁹

feature scaling

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

size      size²

price y

size x

# Choice of features



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

size    $\sqrt{\text{size}}$

price y

size x

$\sqrt{x}$

What features to use?
↳ course 2