



Stanford  
ONLINE

DeepLearning.AI



# Machine Learning Overview

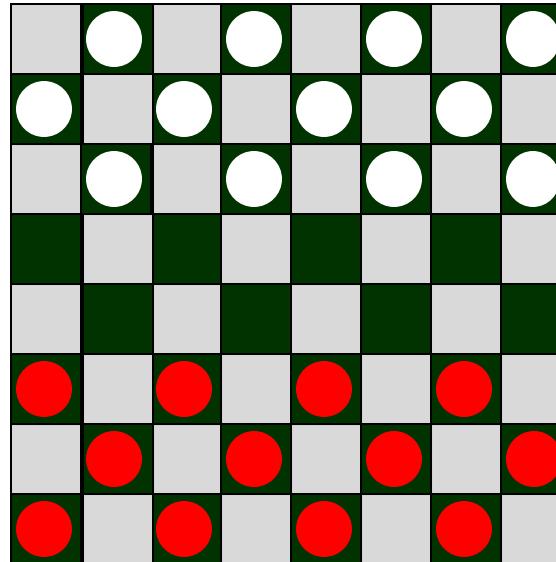
---

What is  
Machine Learning?

# Machine learning

“Field of study that gives computers the ability to learn without being explicitly programmed.”

Arthur Samuel (1959)



# Question

---

If the checkers program had been allowed to play only ten games (instead of tens of thousands) against itself, a much smaller number of games, how would this have affected its performance?

- Would have made it better
  -   Would have made it worse
-

# Machine learning algorithms

rapid advancements

used most in real-world applications

- Supervised learning ← course 1, 2
- Unsupervised learning ←
- Recommender systems
- Reinforcement learning

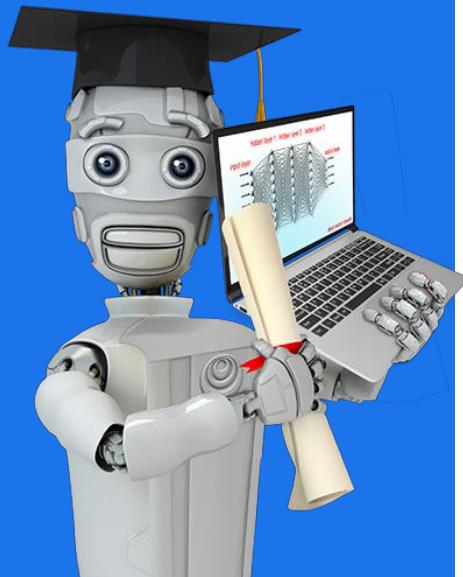
course 3

Practical advice for applying learning algorithms



Stanford  
ONLINE

DeepLearning.AI

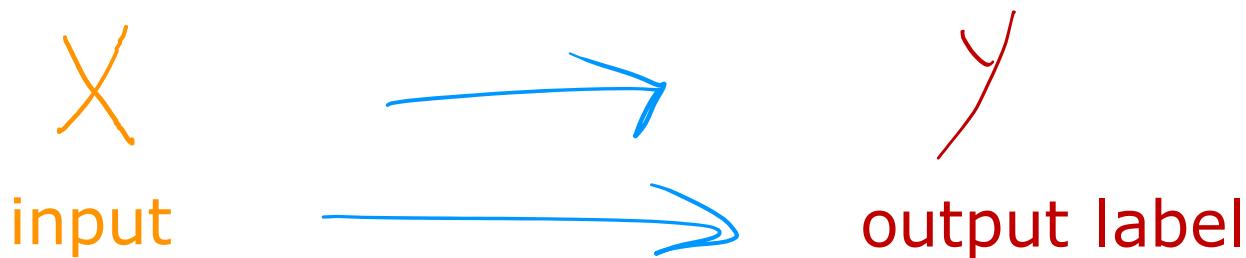


# Machine Learning Overview

---

## Supervised Learning Part 1

# Supervised learning

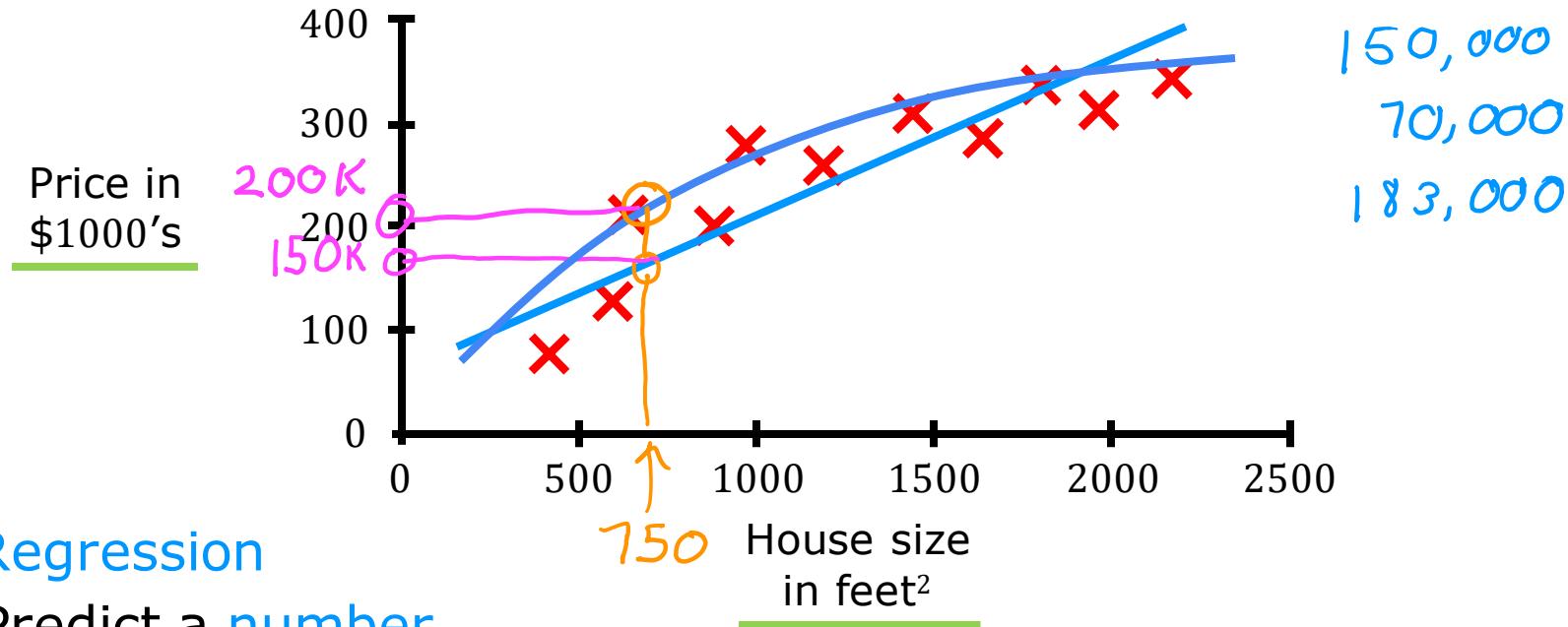


Learns from being given “right answers”

(correct label  $y$  for input  $x$ )

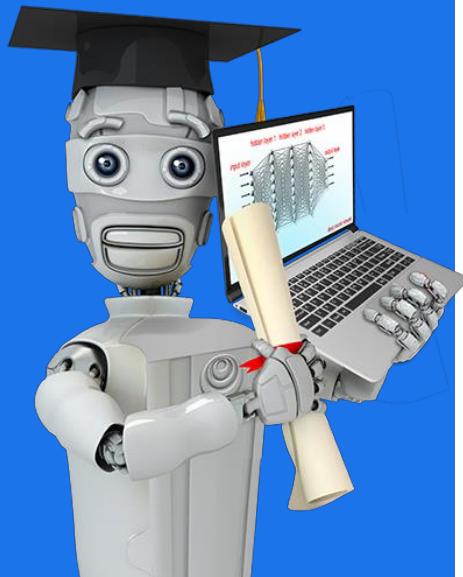
Input (X)	Output (Y)	Application
email	spam? (0/1)	spam filtering
audio	text transcripts	speech recognition
English	Spanish	machine translation
ad, user info	click? (0/1)	online advertising
image, radar info	position of other cars	self-driving car
image of phone	defect? (0/1)	visual inspection

# Regression: Housing price prediction



Stanford  
ONLINE

DeepLearning.AI



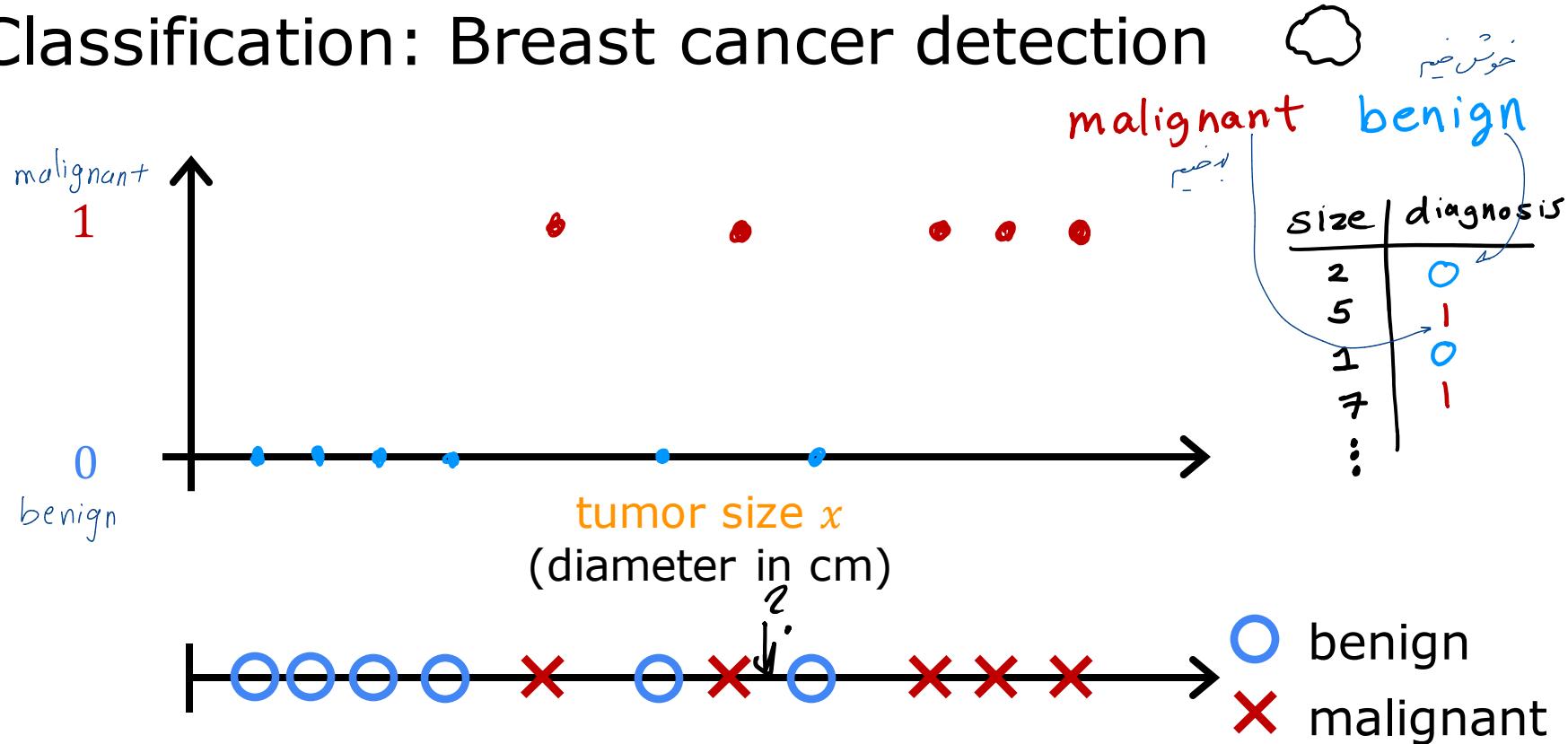
# Machine Learning Overview

---

→ predict  $X$  to  $Y$  mappings

## Supervised Learning Part 2

# Classification: Breast cancer detection



# Classification: Breast cancer detection

○ benign

✗ malignant type 1

△ malignant type 2

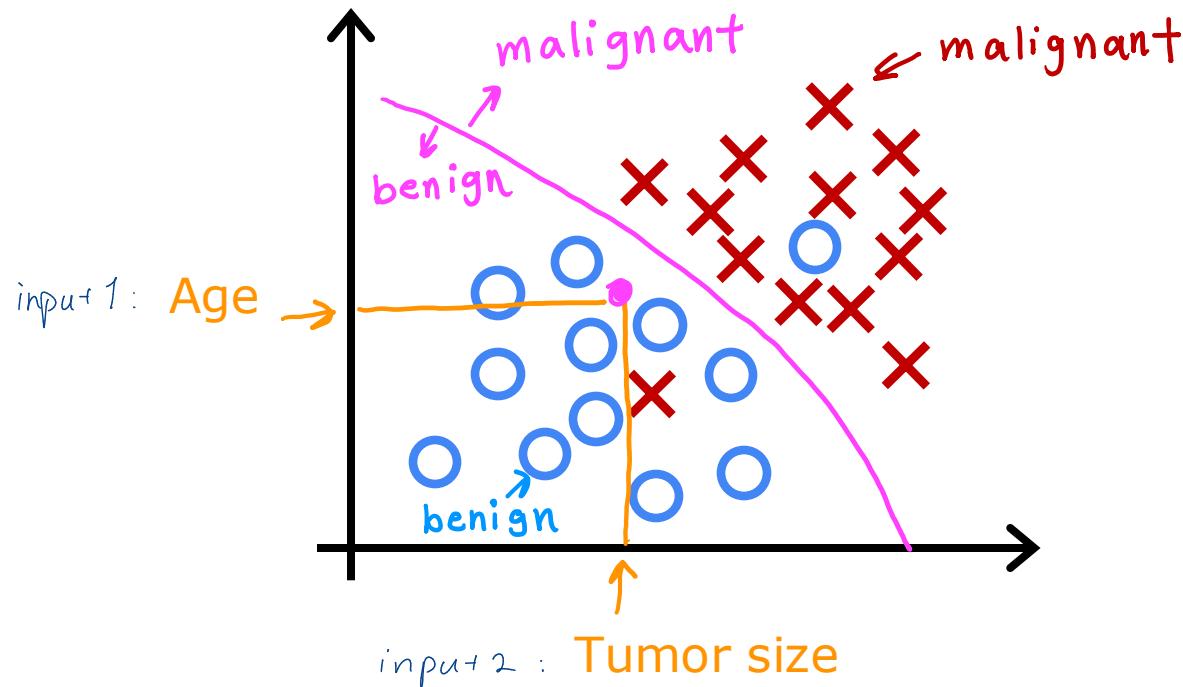


Classification  
predict categories      class / category

cat dog / benign malignant / 0, 1, 2

small number of possible outputs

# Two or more inputs



# Supervised learning

Learns from being given “right answers”

1

Regression

Predict a number

infinitely many possible outputs

2

Classification

predict categories

small number of possible outputs

f<sub>inite</sub>

Stanford  
ONLINE

DeepLearning.AI

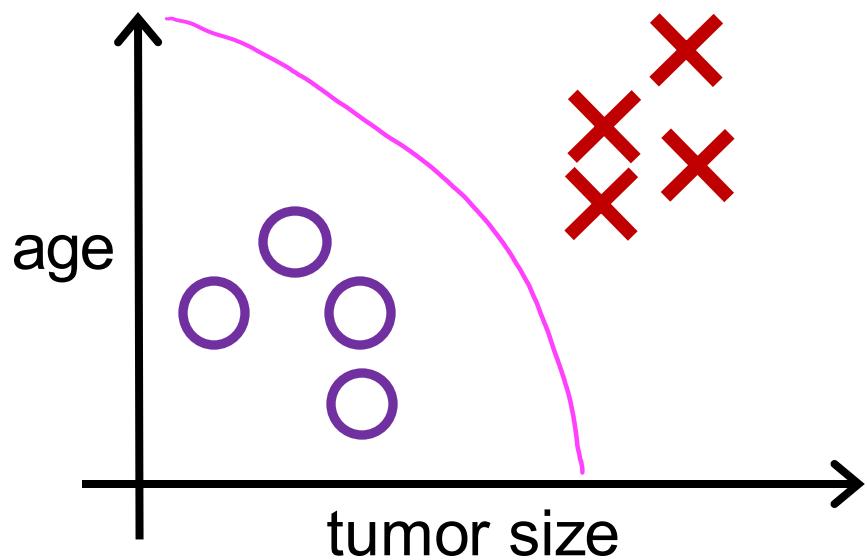


# Machine Learning Overview

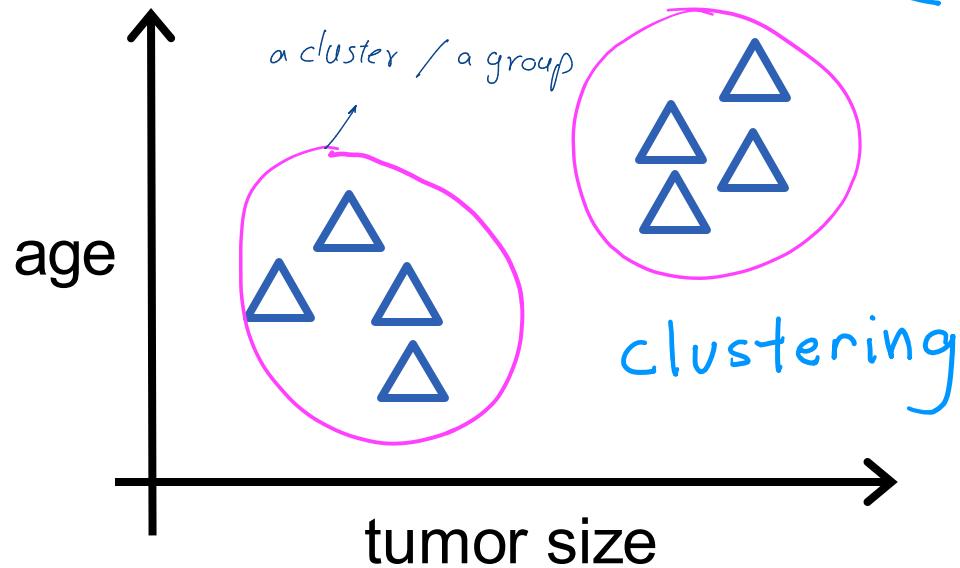
---

## Unsupervised Learning Part 1

Supervised learning  
Learn from data labeled  
with the “right answers”



Unsupervised learning  
Find something interesting  
in unlabeled data.



# Clustering: Google news



Giant **panda** gives birth to rare **twin** cubs at Japan's oldest **zoo**

USA TODAY · 6 hours ago



- Giant **panda** gives birth to **twin** cubs at Japan's oldest **zoo**

CBS News · 7 hours ago

- Giant **panda** gives birth to **twin** cubs at Tokyo's Ueno **Zoo**

WHBL News · 16 hours ago

- A Joyful Surprise at Japan's Oldest **Zoo**: The Birth of **Twin Pandas**

The New York Times · 1 hour ago

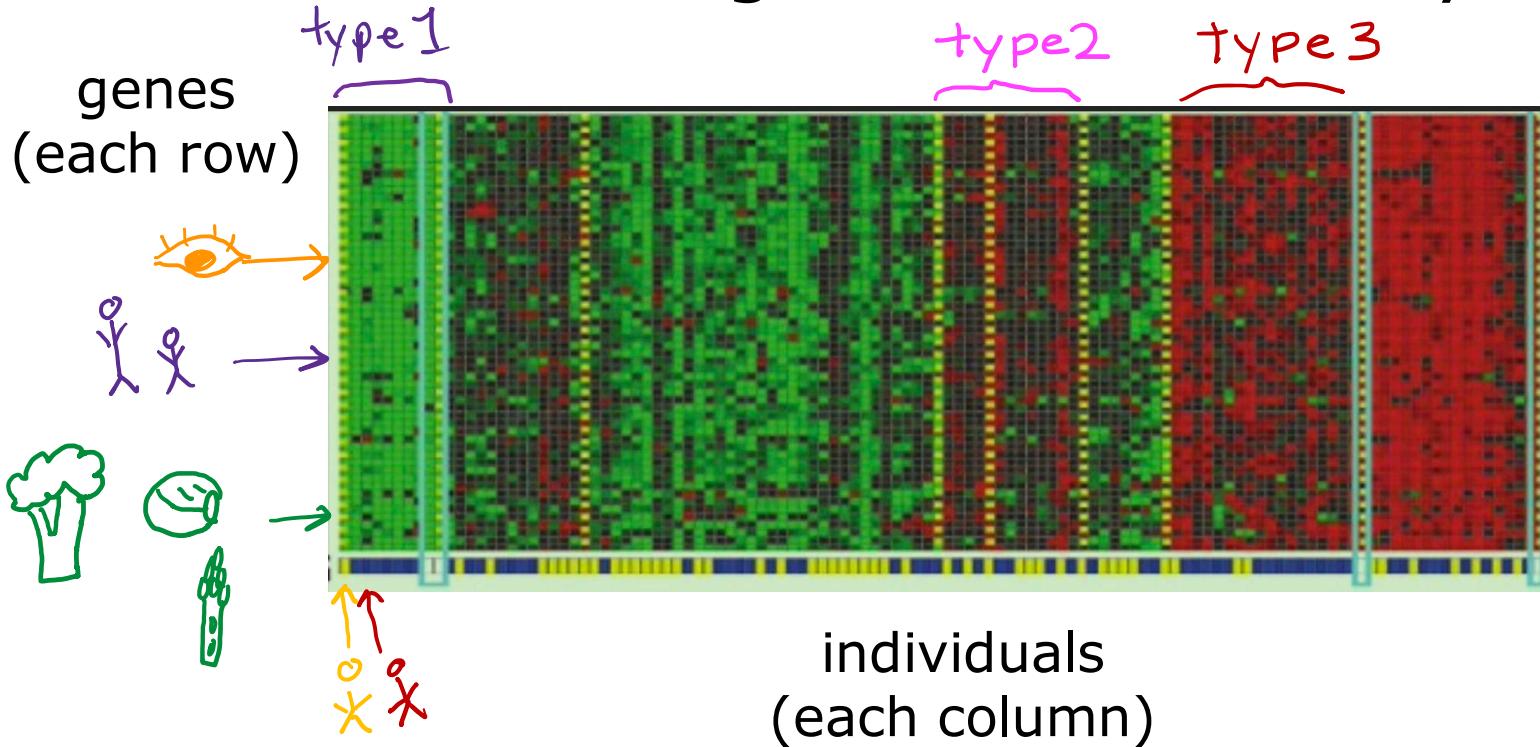
- **Twin** Panda **Cubs** Born at Tokyo's Ueno **Zoo**

PEOPLE · 6 hours ago

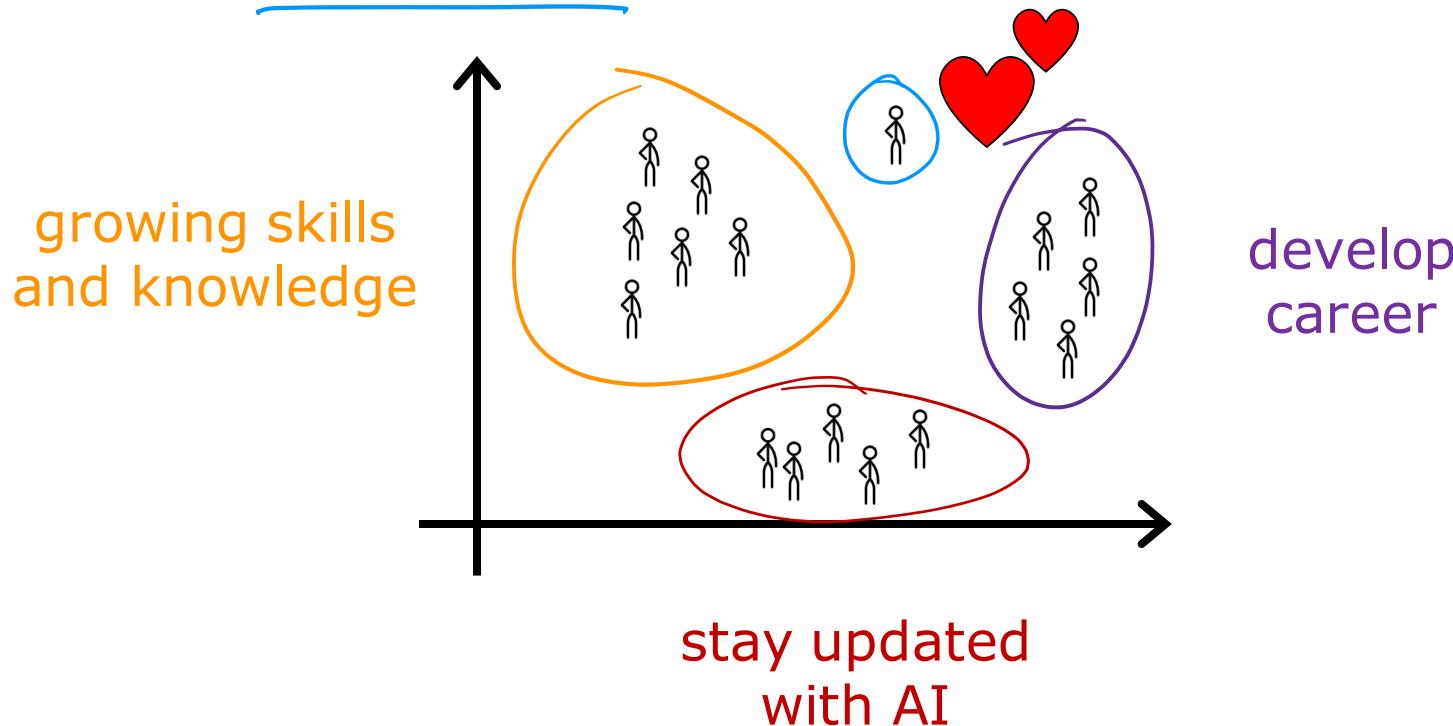
View Full Coverage



# Clustering: DNA microarray

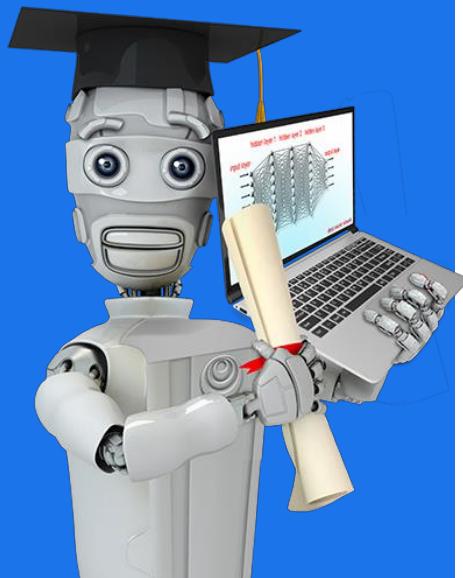


# Clustering: Grouping customers



Stanford  
ONLINE

DeepLearning.AI



# Machine Learning Overview

---

## Unsupervised Learning Part 2

# Unsupervised learning

Data only comes with inputs  $x$ , but not output labels  $y$ .  
Algorithm has to find **structure** in the data.

*or  
pattern*

## Clustering

Group similar data points together.

## Dimensionality reduction

Compress data using fewer numbers.

## Anomaly detection

Find unusual data points.

*or  
events*

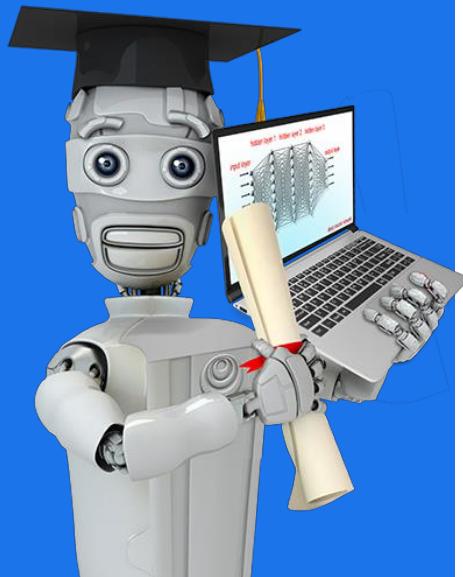
# Question

Of the following examples, which would you address using an **unsupervised** learning algorithm?

- Given email labeled as spam/not spam, learn a spam filter.  
*supervised*
- Given a set of news articles found on the web, group them into sets of articles about the same story.  
*unsupervised*
- Given a database of customer data, automatically discover market segments and group customers into different market segments.
- Given a dataset of patients diagnosed as either having diabetes or not, learn to classify new patients as having diabetes or not

Stanford  
ONLINE

DeepLearning.AI

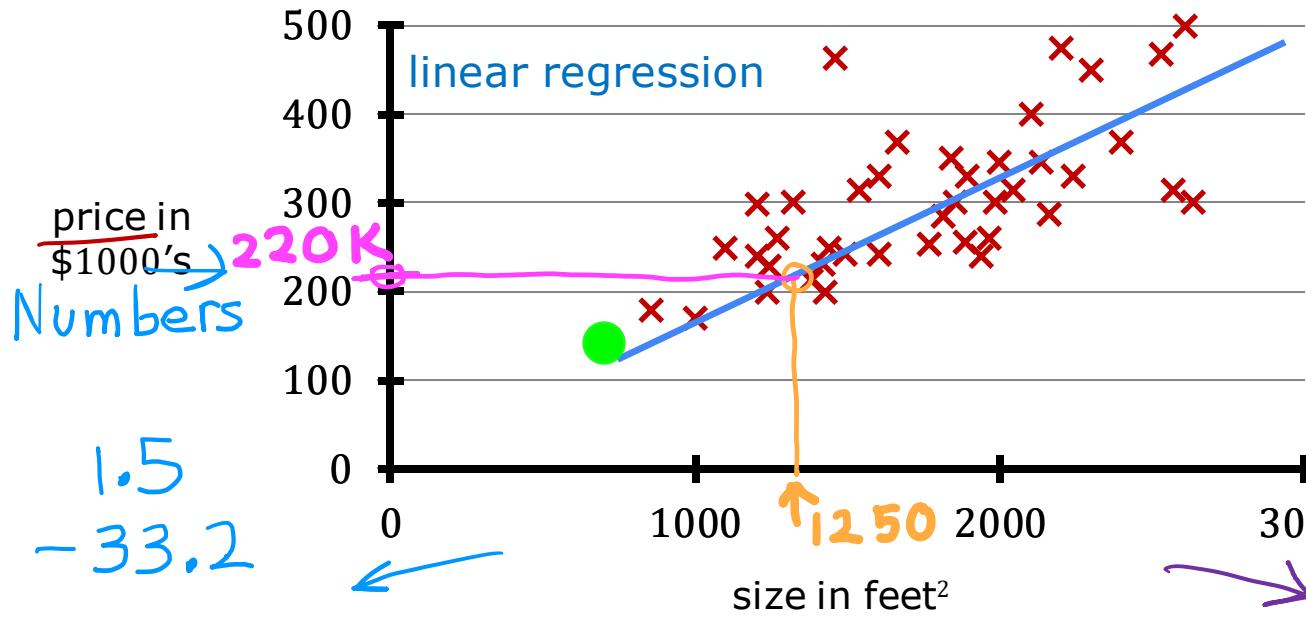


# Linear Regression with One Variable

---

## Linear Regression Model Part 1

# House sizes and prices



Regression model  
Predicts numbers  
Infinitely many possible outputs

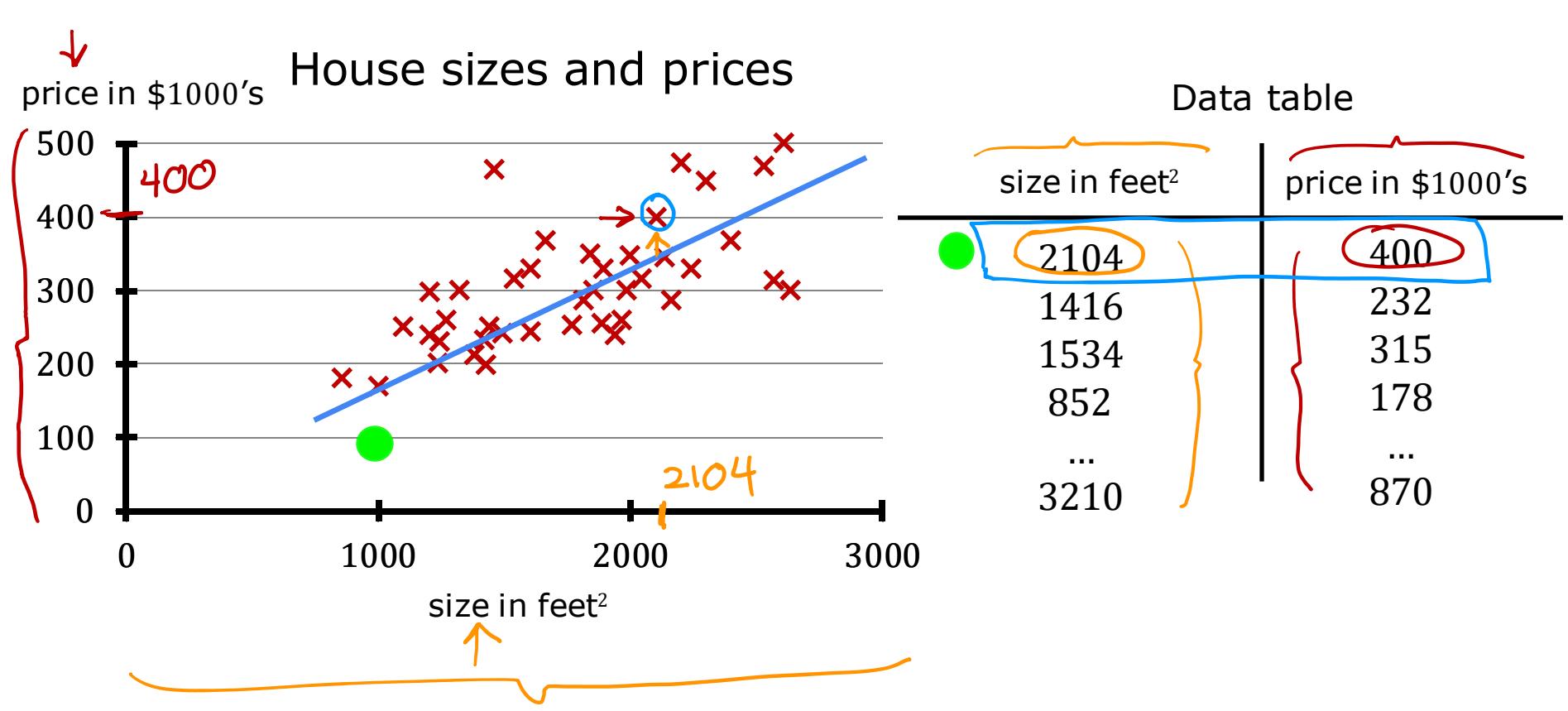
Supervised learning model  
Data has "right answers"

Classification model  
Predicts categories  
Small number of possible outputs

categories  
cat } 2  
dog }

disease  10

(in contrast with regression)



# Terminology

Training set: Data used to train the model

size in feet<sup>2</sup>

	$x$	$y$
(1)	2104	400
(2)	1416	232
(3)	1534	315
(4)	852	178
...	...	...
(47)	3210	870

$$m = 47$$

$$x^{(1)} = 2104$$

$$(x^{(1)}, y^{(1)}) = (2104, 400)$$

$$x^{(2)} = 1416$$

$$X^{(2)} \neq X^2 \text{ not exponent}$$

Notation:

$x$  = "input" variable / input feature

$y$  = "output" variable  
"target" variable

$m$  = number of training examples

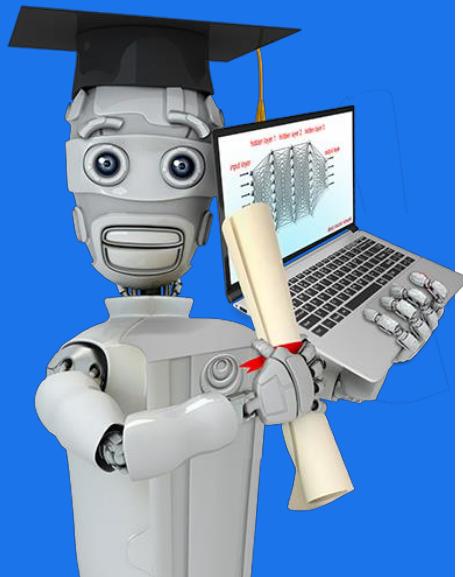
$(x, y)$  = single training example

$(x^{(i)}, y^{(i)})$   $\xrightarrow{\text{specific row}}$

$(x^{(i)}, y^{(i)})$  =  $i^{\text{th}}$  training example  
index  $(1^{\text{st}}, 2^{\text{nd}}, 3^{\text{rd}} \dots)$

Stanford  
ONLINE

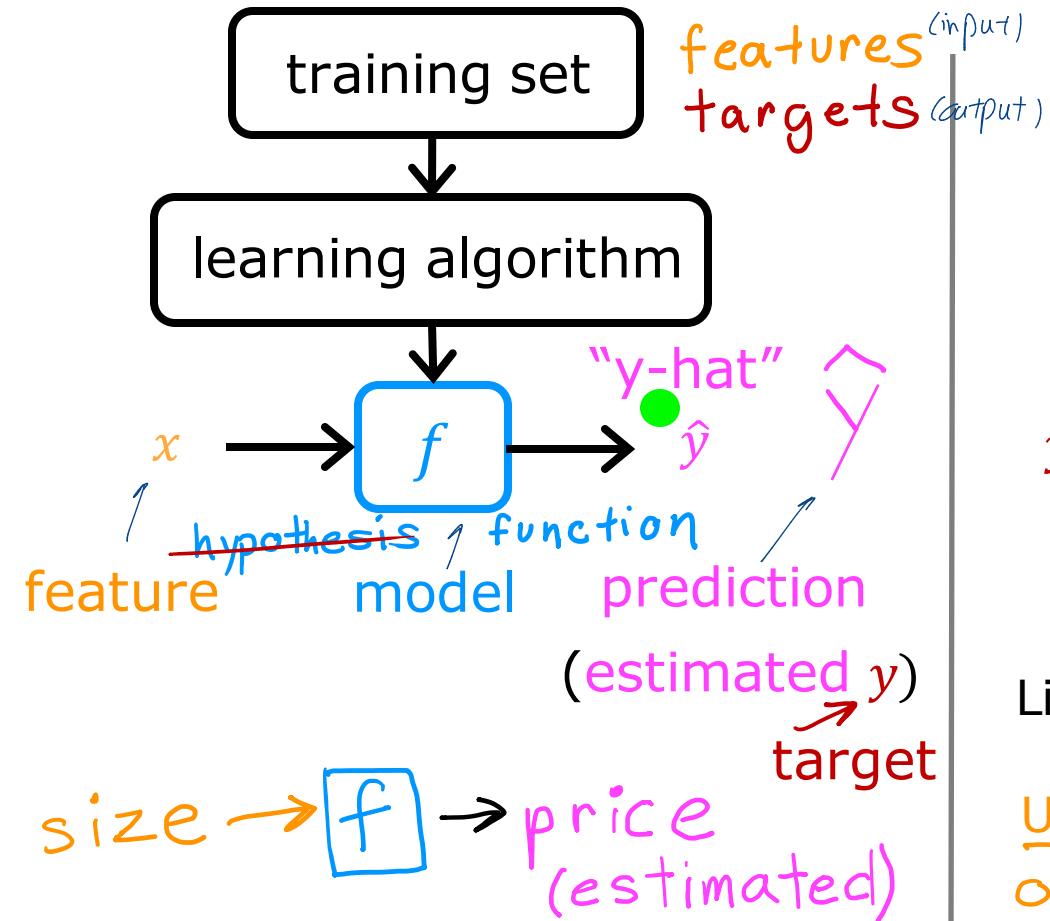
DeepLearning.AI



# Linear Regression with One Variable

---

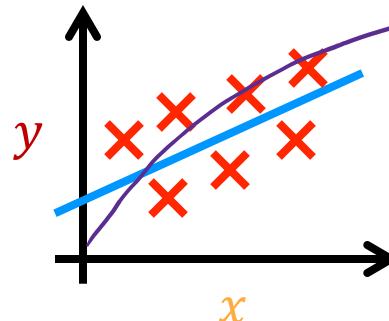
## Linear Regression Model Part 2



# How to represent $f$ ?

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

$$f(x)$$



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

$$f(x) = wx + b$$

linear

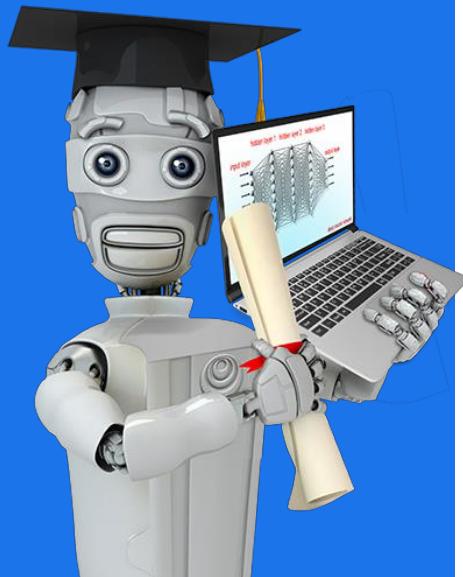
single feature  $x$

Linear regression with one variable.  
size

Univariate linear regression.  
one variable

Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with One Variable

---

## Cost Function

# Training set

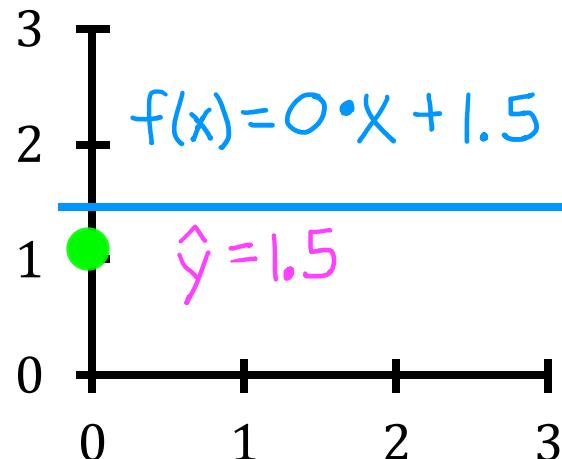
features	targets
size in feet <sup>2</sup> ( $x$ )	price \$1000's ( $y$ )
2104	460
1416	232
1534	315
852	178
...	...

$$\text{Model: } f_{w,b}(x) = wx + b$$

$w, b$ : parameters  
coefficients  
weights

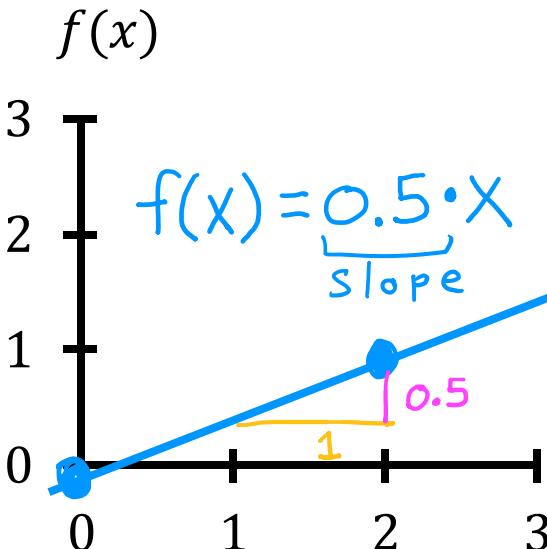
What do  $w, b$  do?

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

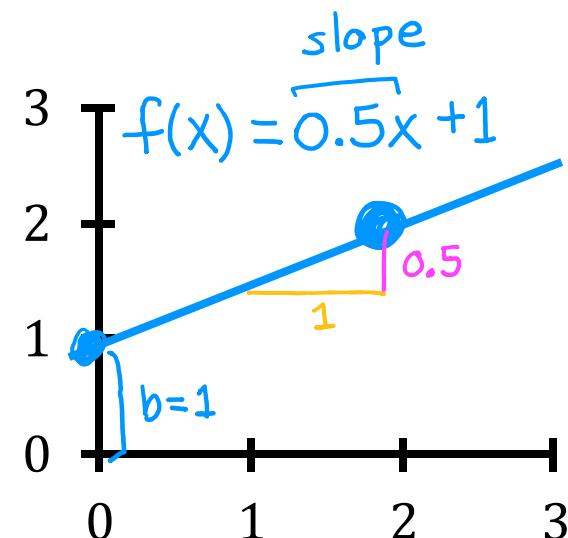


$$\rightarrow w = 0$$
$$\rightarrow b = 1.5$$

*(y-intercept)*

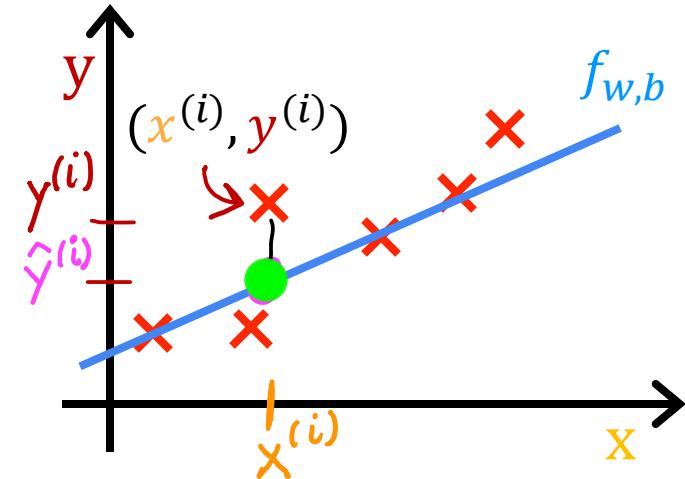


$$\rightarrow w = 0.5$$
$$\rightarrow b = 0$$



$$\rightarrow w = 0.5$$
$$\rightarrow b = 1$$

## Cost function: Squared error cost function



$$\hat{y}^{(i)} = f_{w,b}(x^{(i)})$$

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

$$\bar{J}(w, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2$$

$m$  = number of training examples

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

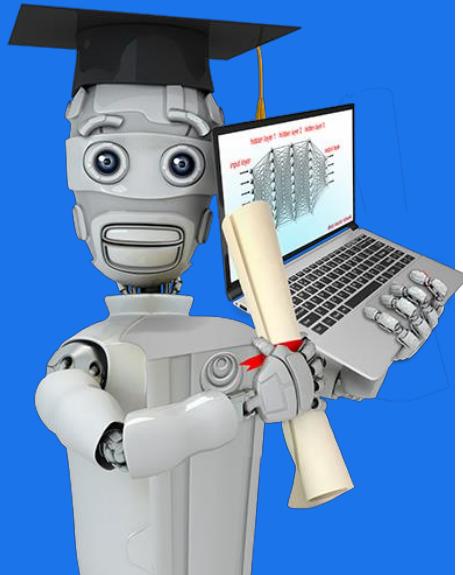
intuition (next!)

Find  $w, b$ :

$\hat{y}^{(i)}$  is close to  $y^{(i)}$  for all  $(x^{(i)}, y^{(i)})$ .

Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with One Variable

---

Cost Function  
Intuition

view / previous / next /

model:

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

parameters:

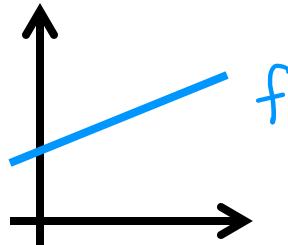
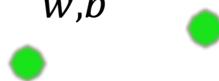
$$\underline{w, b}$$

cost function:

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

goal:

$$\underset{w,b}{\text{minimize}} J(w, b)$$

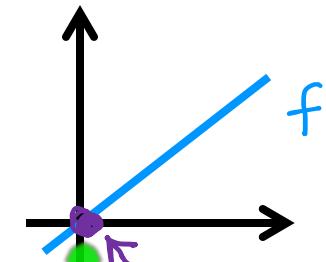


simplified

$$f_w(x) = \underline{wx}$$

$$b = \emptyset^{\text{zero}}$$

$$w$$



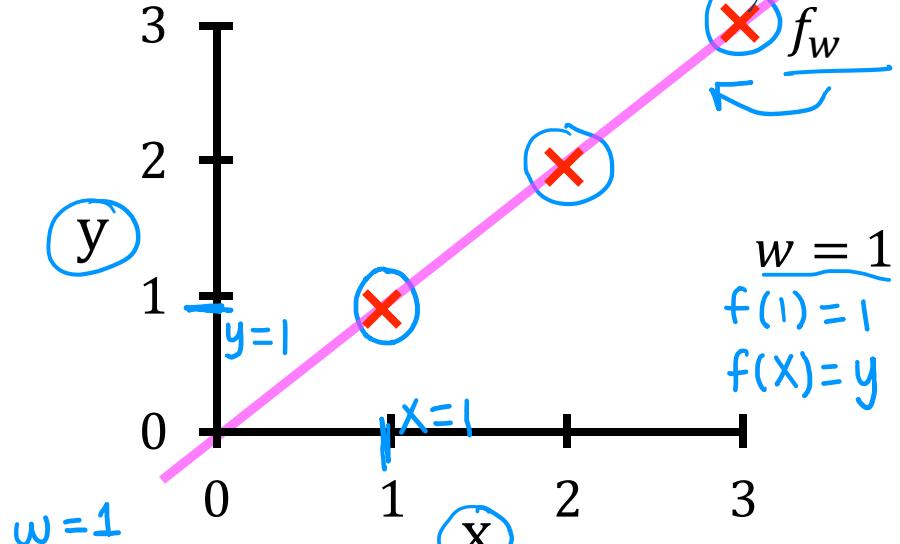
$$\underline{J(w)} = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

$$\underset{w}{\text{minimize}} \underline{J(w)}$$

$$\omega x^{(i)}$$

$\rightarrow f_w(x)$

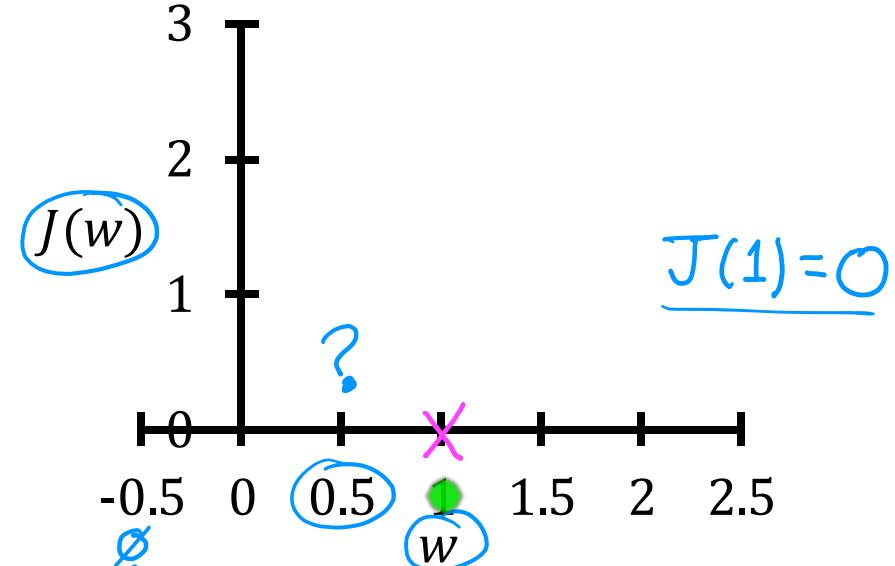
(for fixed  $w$ , function of  $x$ )  
input



$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} - y^{(i)})^2$$

$J(w)$

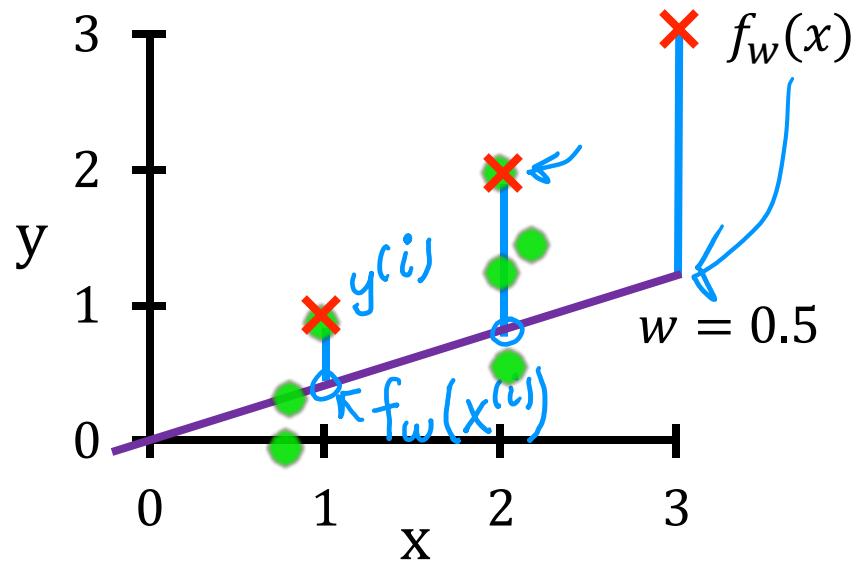
(function of  $w$ )  
parameter



$$= \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$

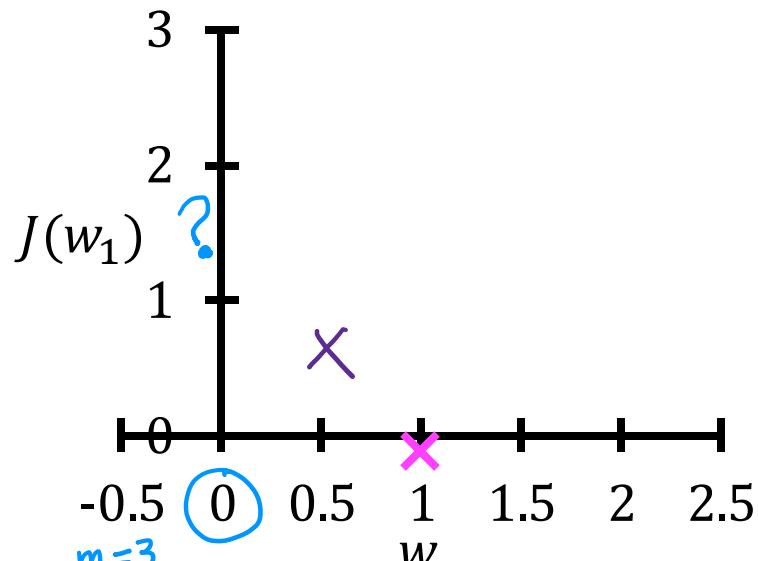
$f_w(x)$

(function of  $x$ )

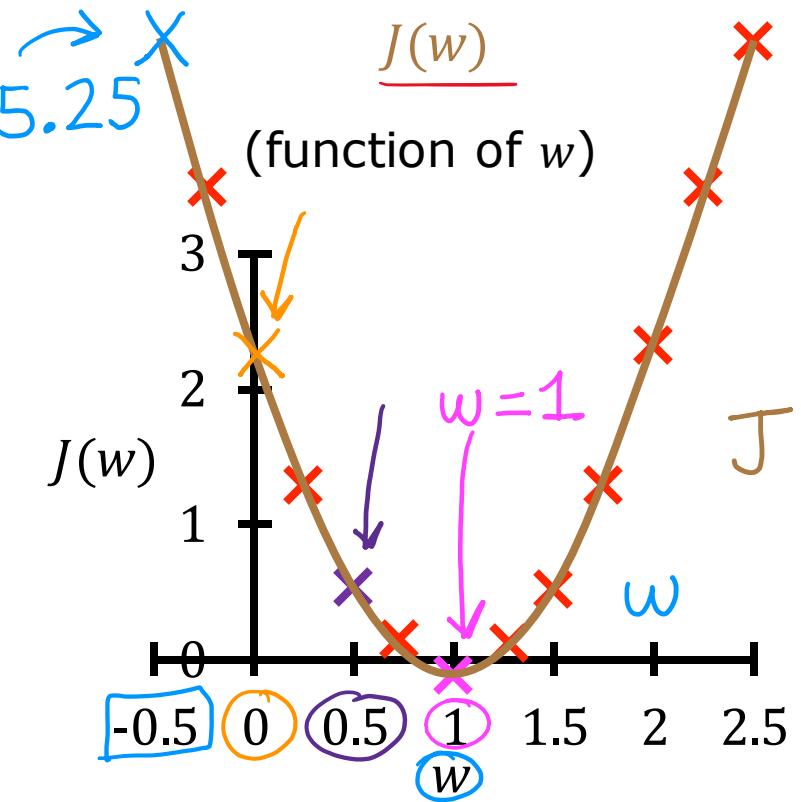
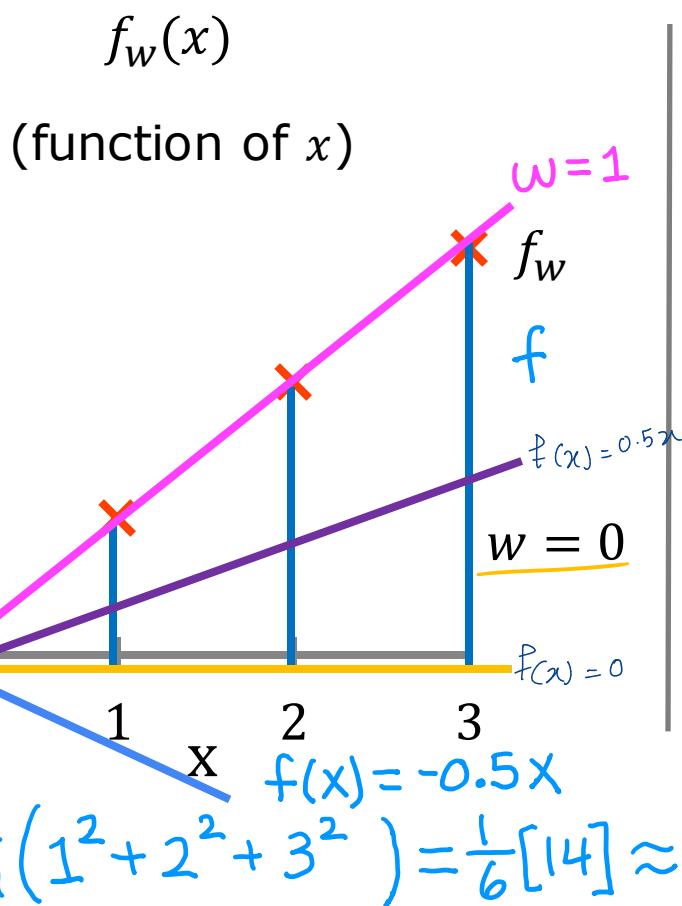


$J(w)$

(function of  $w$ )



$$J(0.5) = \frac{1}{2m} \left[ (0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right] = \frac{1}{2 \times 3} [3.5] = \frac{3.5}{6} \approx 0.58$$

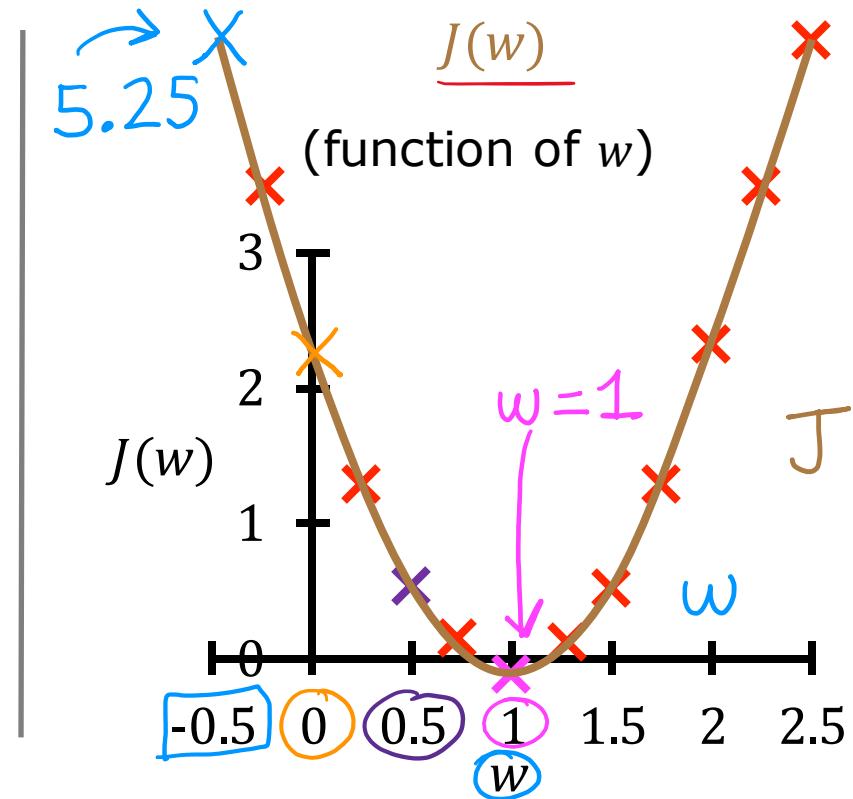


goal of linear regression:

$$\underset{w}{\text{minimize}} J(w)$$

general case:

$$\underset{w,b}{\text{minimize}} J(w, b)$$



choose  $w$  to minimize  $J(w)$

Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with One Variable

---

Visualizing  
the Cost Function

Model

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

Parameters

$w, b$

~~before:  $b=0$~~

Cost Function

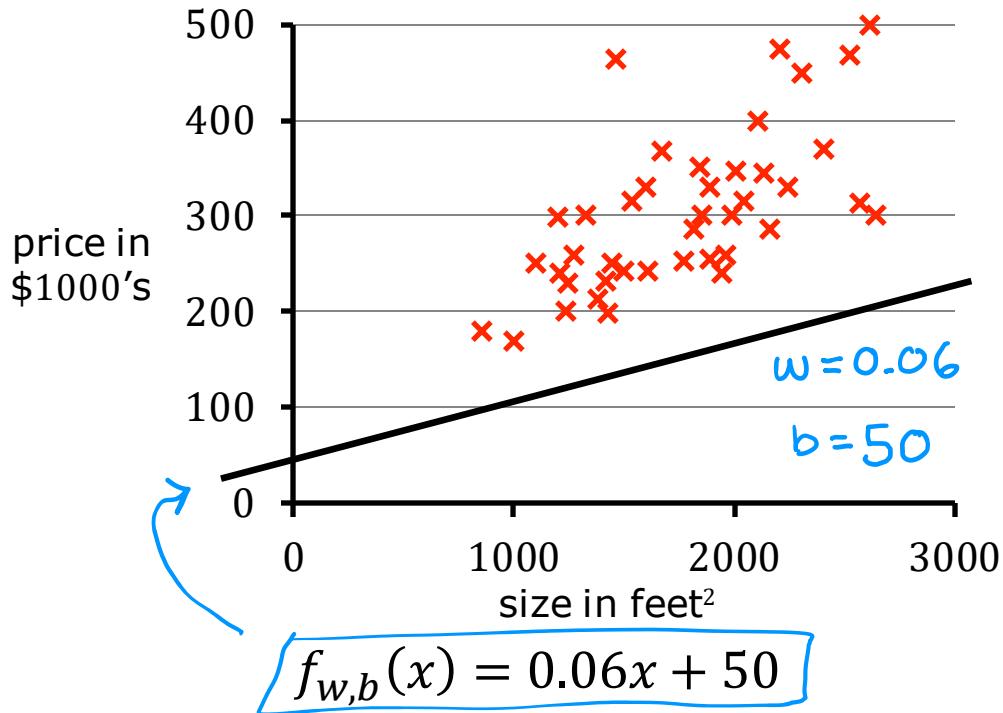
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

Objective

$$\underset{w,b}{\text{minimize}} J(w, b)$$

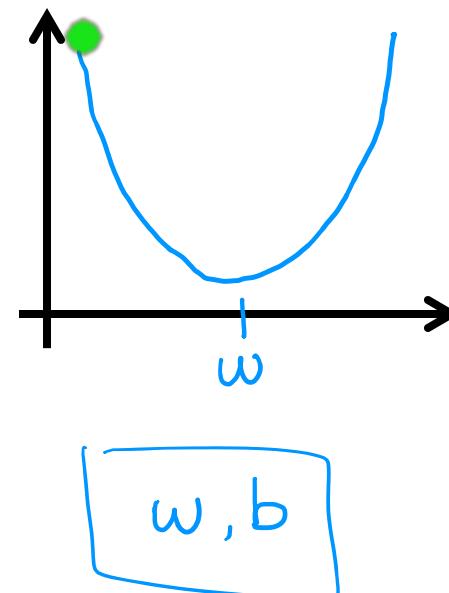
$$\underline{f_{w,b}}$$

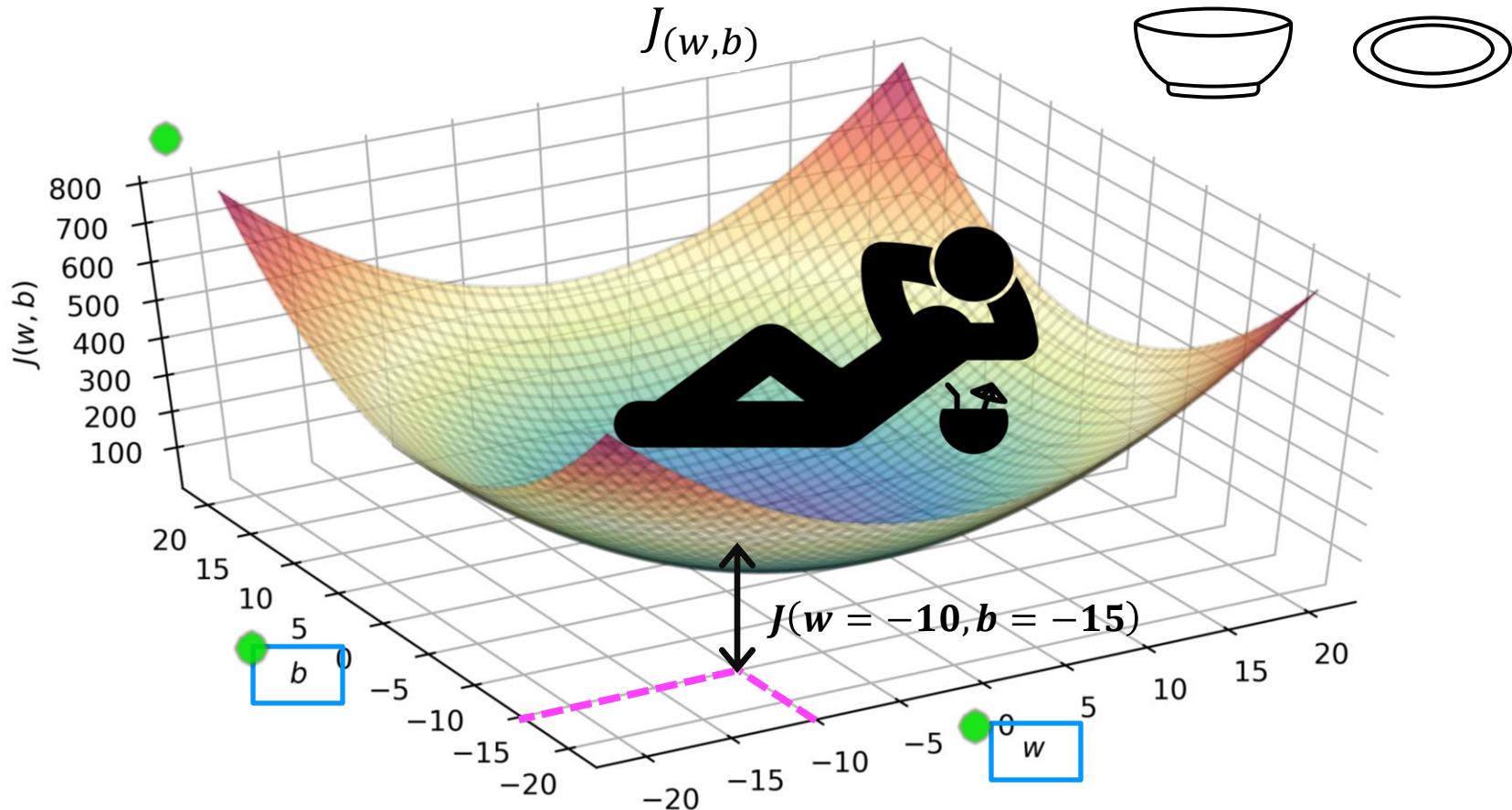
(function of  $x$ )



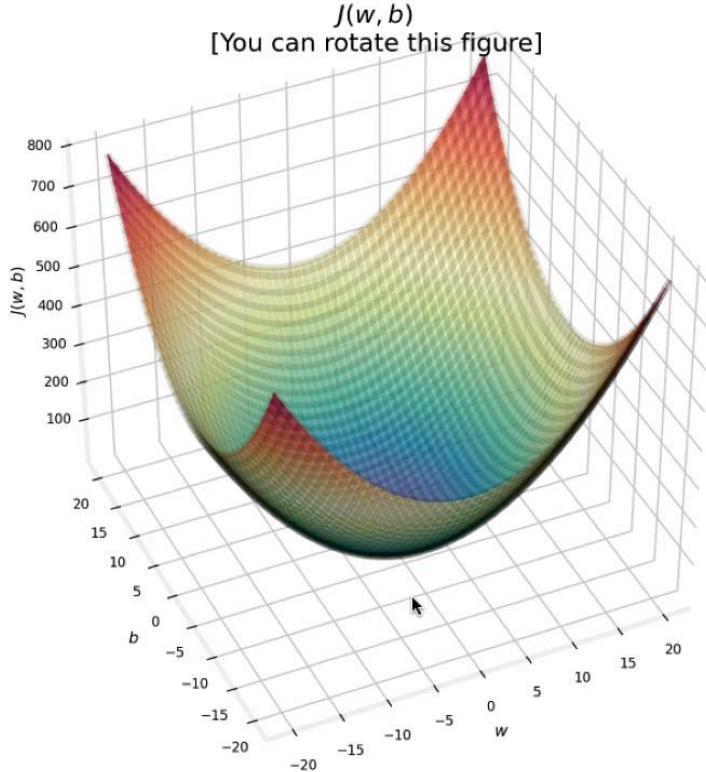
$$\underline{J}$$

(function of  $w, b$ )





# 3D surface plot



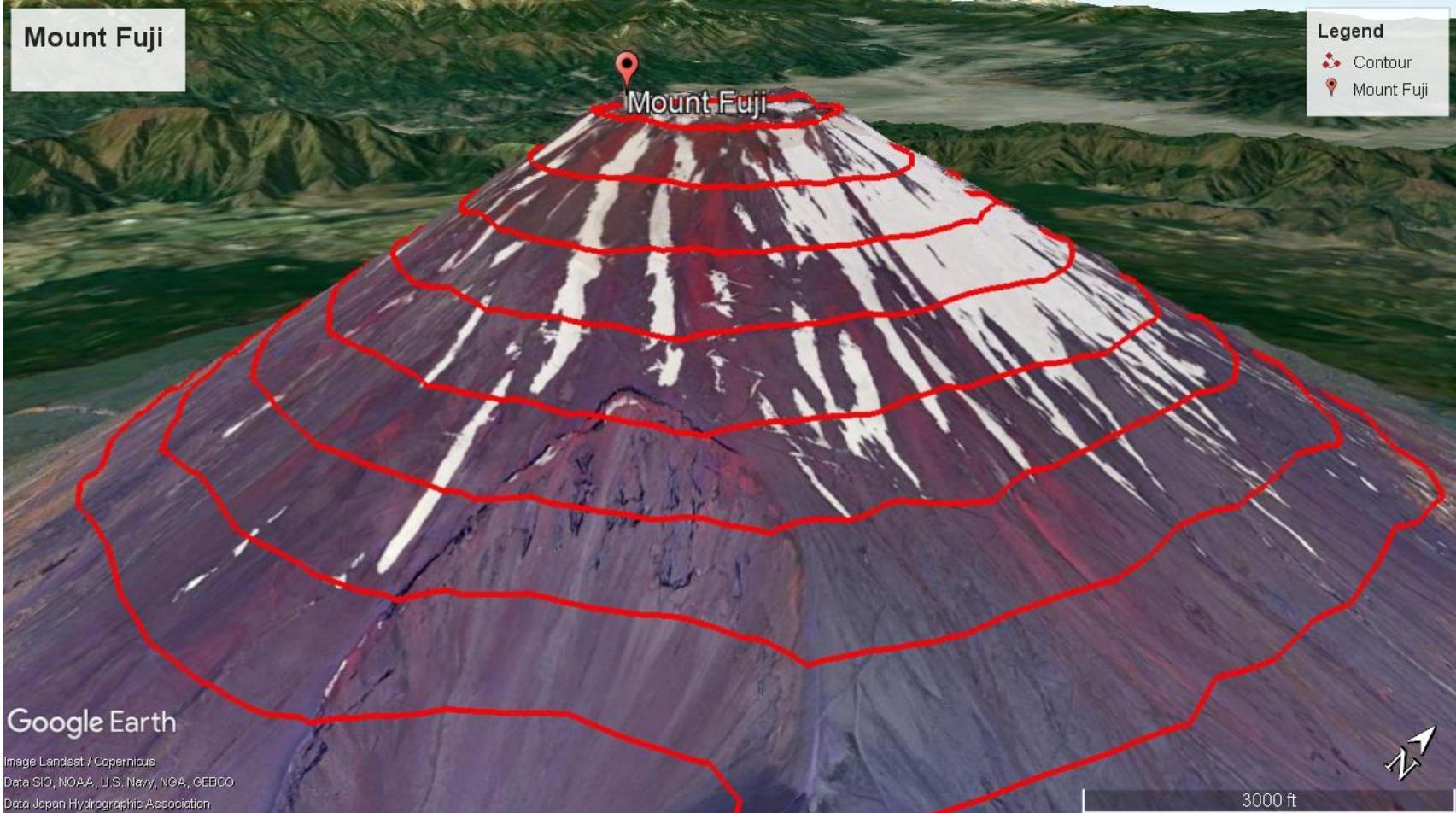
Alternative  
contour plot

# Mount Fuji

Mount Fuji

## Legend

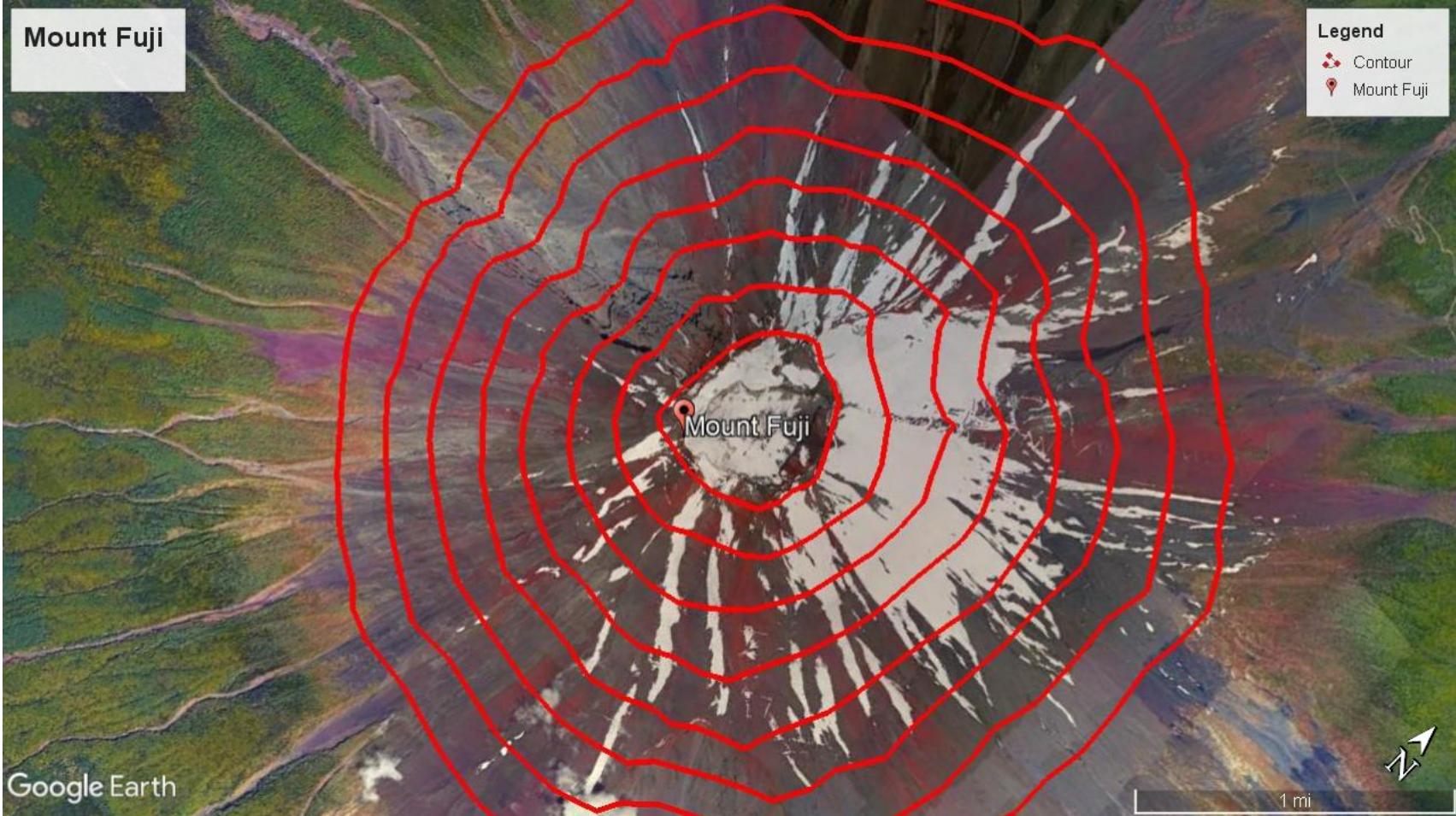
- Contour
- Mount Fuji

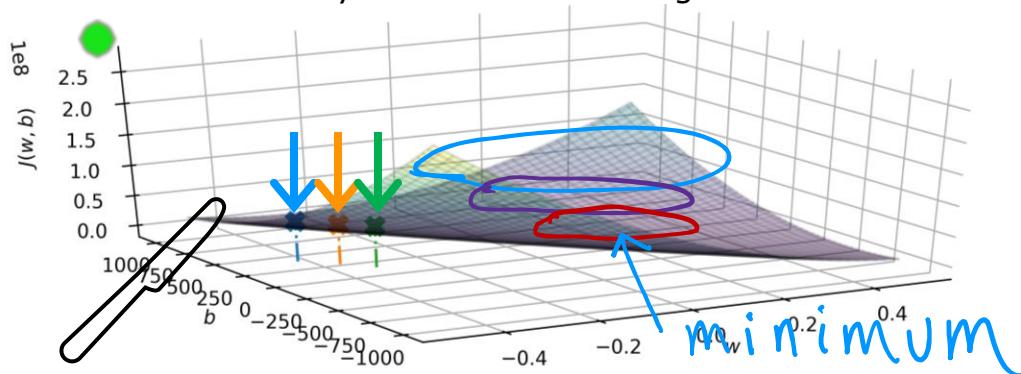
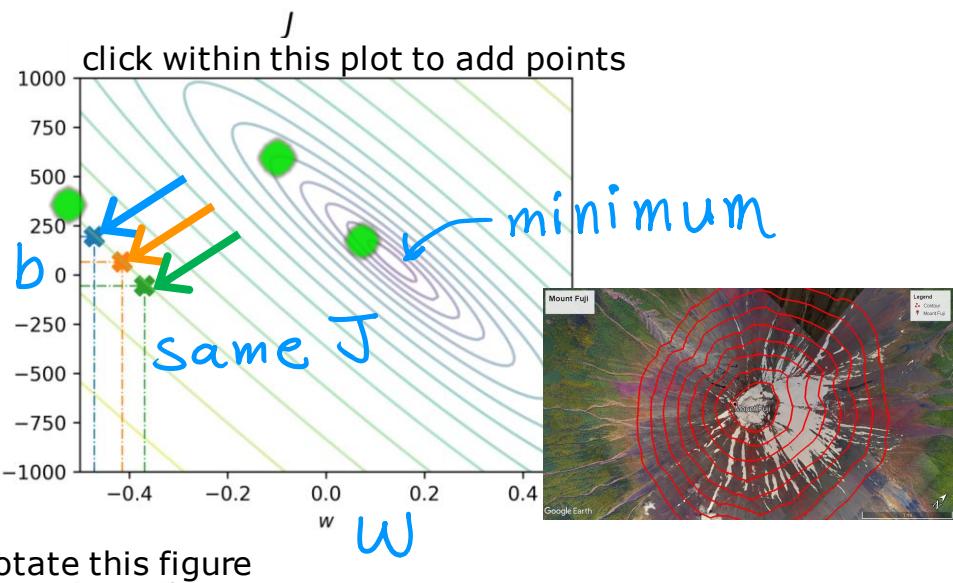
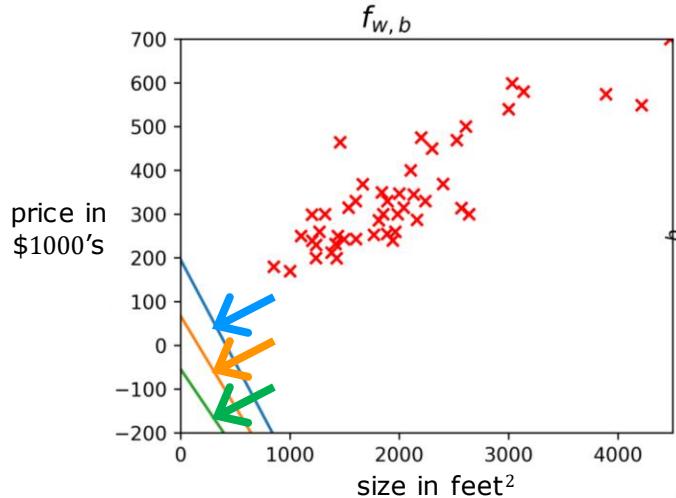


# Mount Fuji

Legend

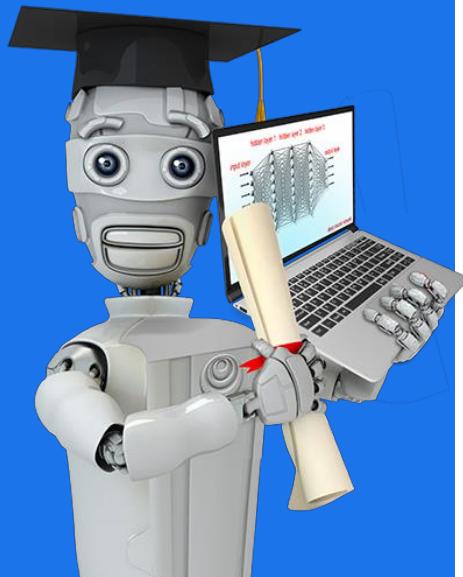
- Contour
- Mount Fuji





Stanford  
ONLINE

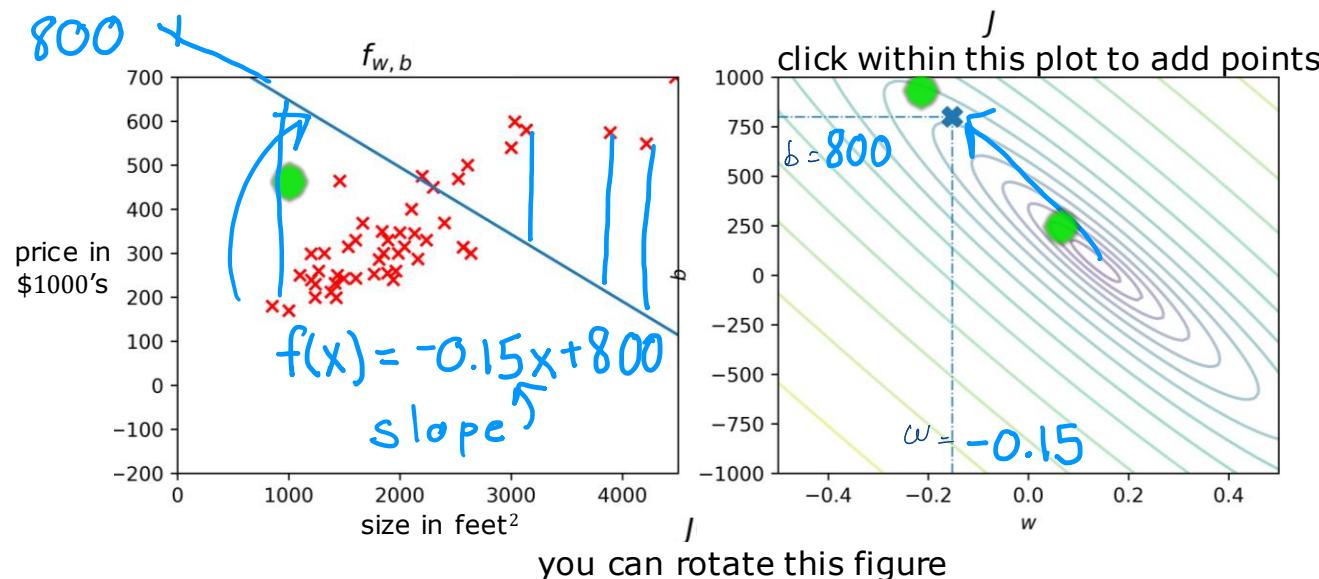
DeepLearning.AI



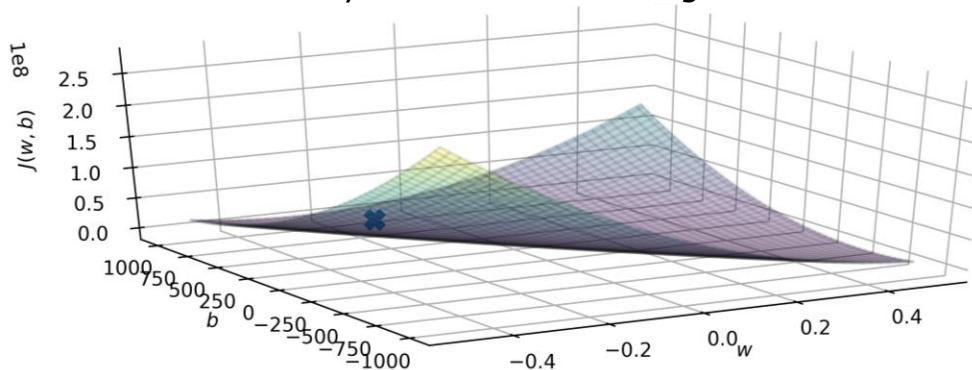
# Linear Regression with One Variable

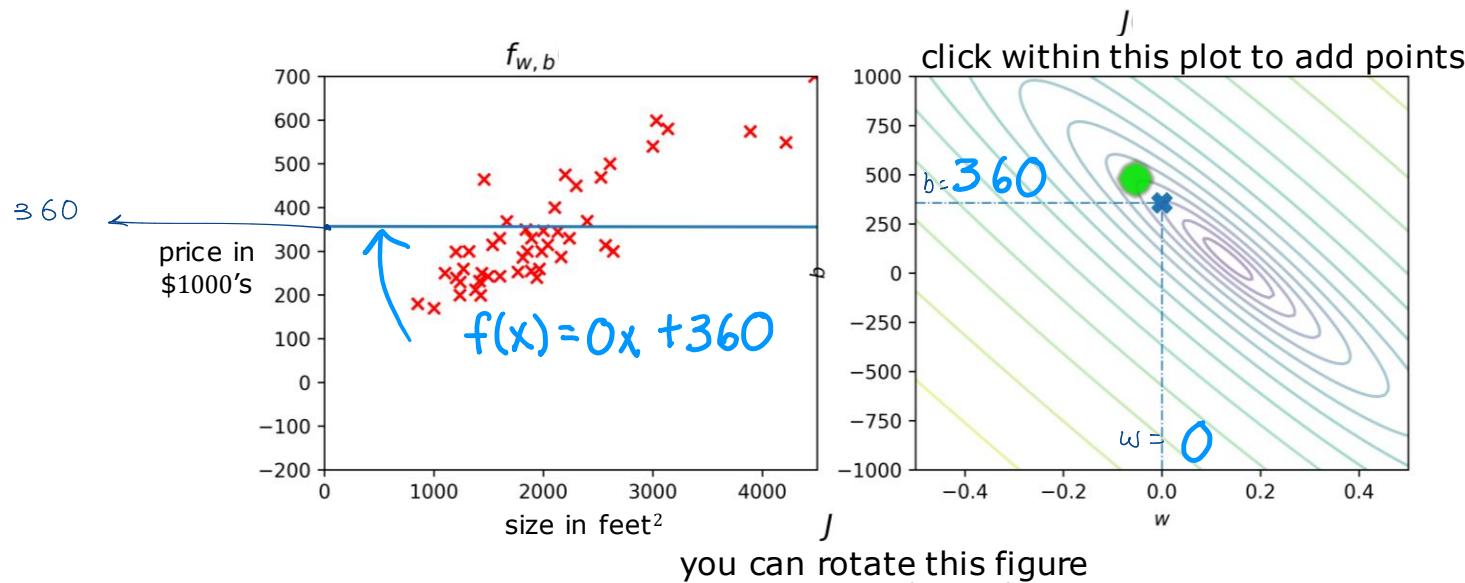
---

## Visualization examples

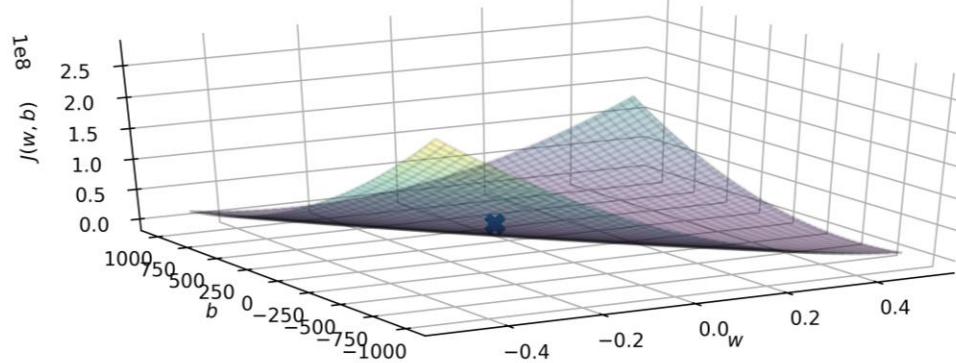


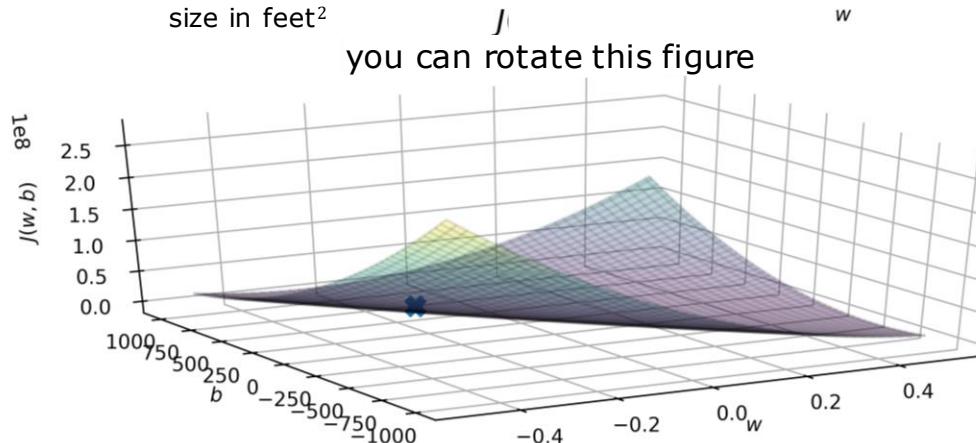
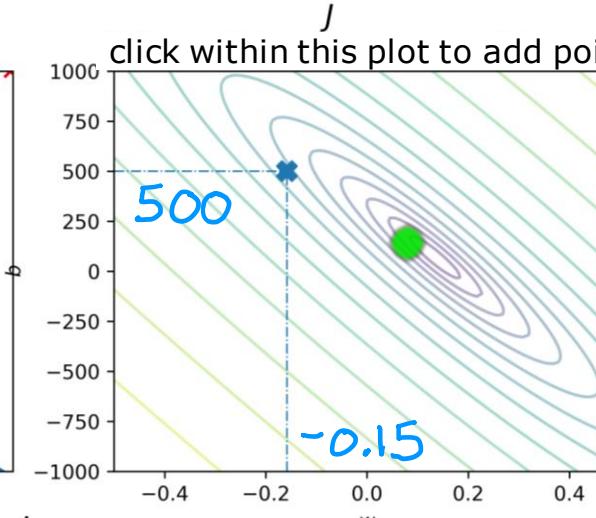
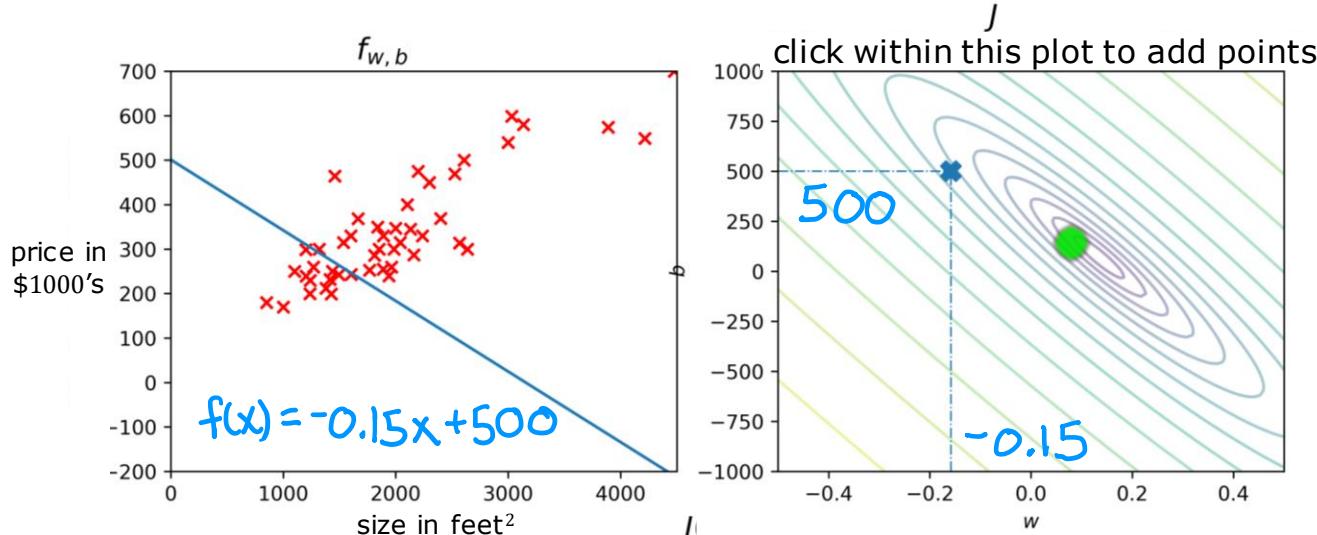
you can rotate this figure

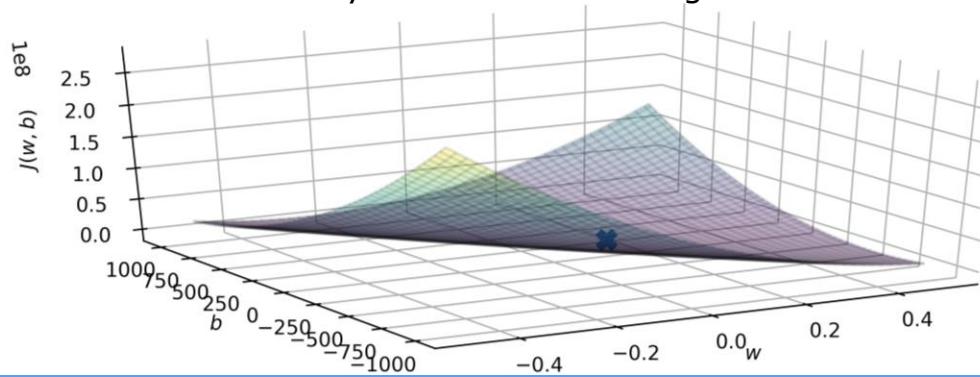
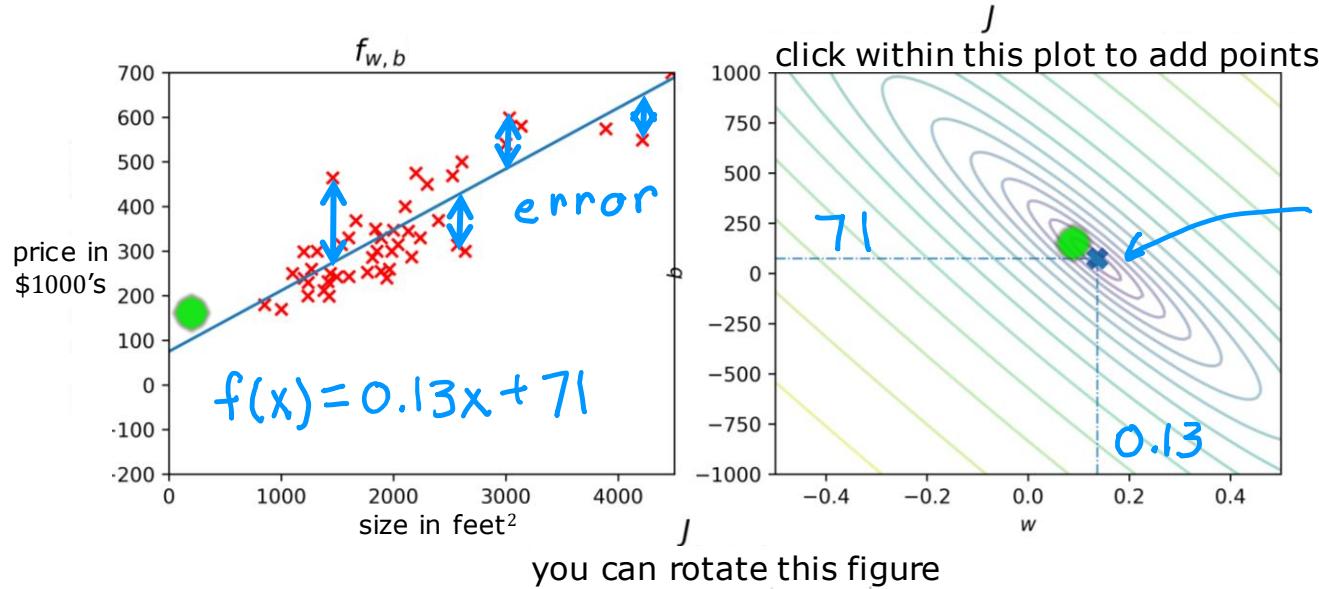




you can rotate this figure

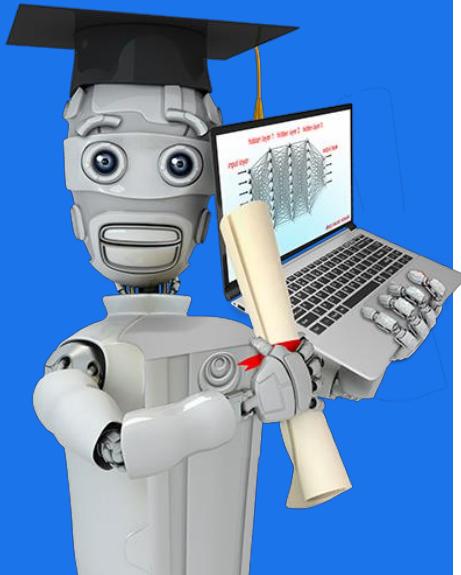






Stanford  
ONLINE

DeepLearning.AI



# Training Linear Regression

---

## Gradient Descent

Have some function  $\underline{J(w, b)}$  for linear regression  
or any function

Want  $\min_{w, b} \underline{J(w, b)}$   $\min_{w_1, \dots, w_n, b} \underline{J(w_1, w_2, \dots, w_n, b)}$

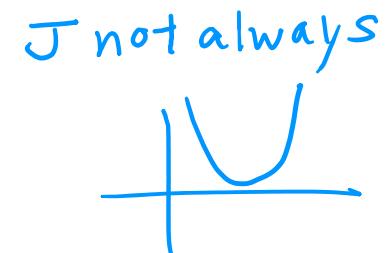
Outline:

Start with some  $\underline{w, b}$  (set  $w=0, b=0$ )

Keep changing  $w, b$  to reduce  $J(w, b)$

Until we settle at or near a minimum

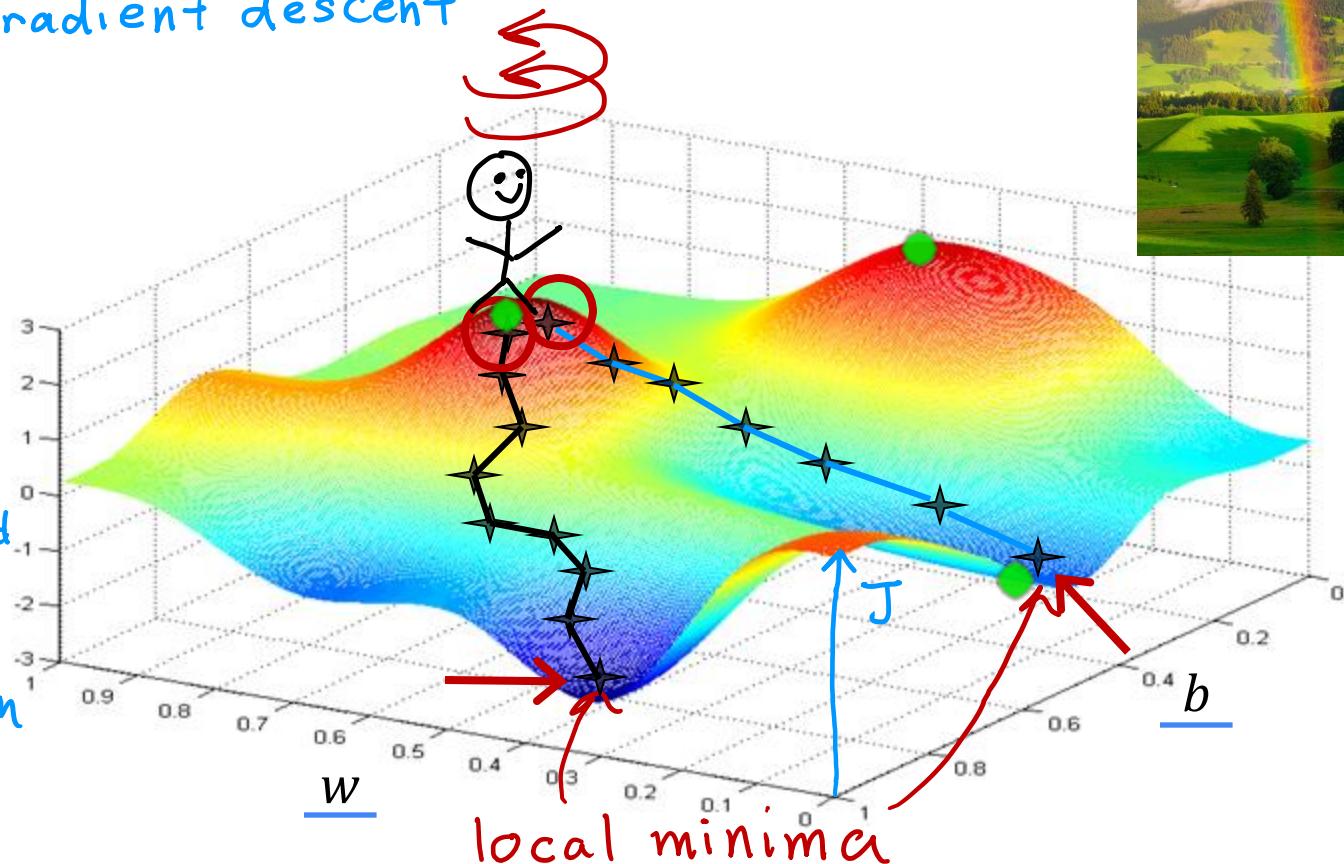
may have >1 minimum



gradient descent

$$J(w, b)$$

not squared  
error cost  
not linear  
regression



Stanford  
ONLINE

DeepLearning.AI



# Training Linear Regression

---

## Implementing Gradient Descent

# Gradient descent algorithm

Repeat until convergence

$$\left\{ \begin{array}{l} w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ b = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{array} \right.$$

$0 < \alpha < 1$   
Learning rate  
Derivative

Simultaneously  
update w and b

Assignment

$$a = c$$

$$a = a + 1$$

Code

Truth assertion

$$a = c$$

equal

$$a = a + 1$$

Math

$$a == c$$

Correct: Simultaneous update

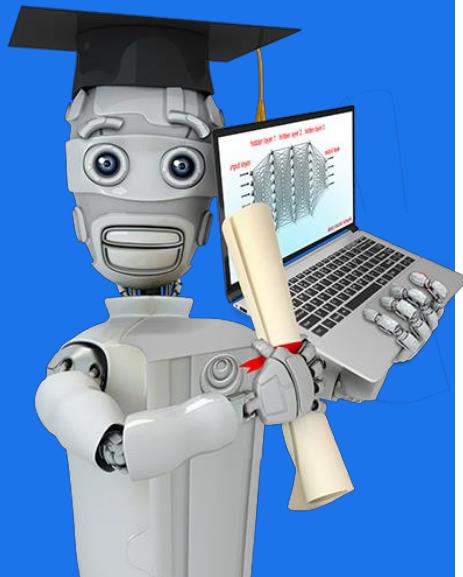
$$\left. \begin{array}{l} tmp\_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp\_b = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w = tmp\_w \\ b = tmp\_b \end{array} \right\}$$

Incorrect

$$\left. \begin{array}{l} tmp\_w = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w = tmp\_w \\ tmp\_b = b - \alpha \frac{\partial}{\partial b} J(w, b) \\ b = tmp\_b \end{array} \right\}$$

Stanford  
ONLINE

DeepLearning.AI



# Training Linear Regression

---

## Gradient Descent Intuition

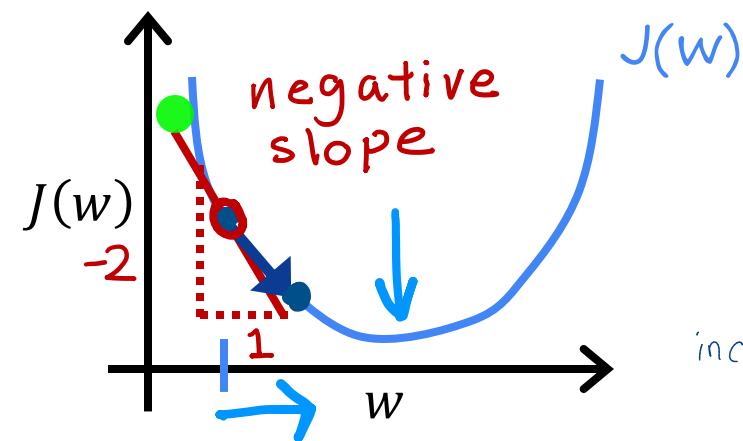
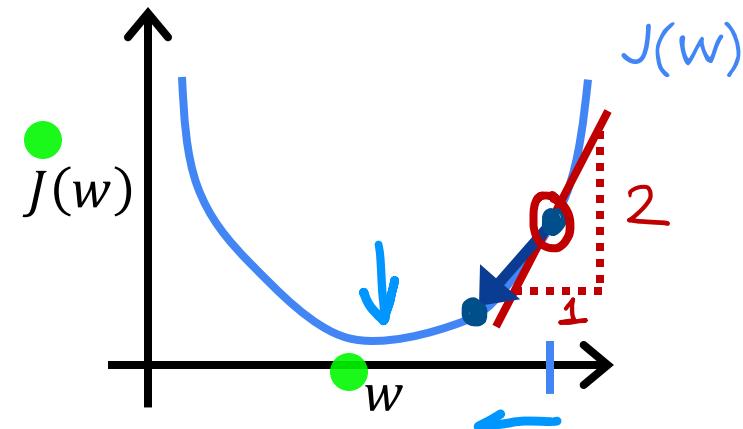
# Gradient descent algorithm

- repeat until convergence {  
  learning rate  $\alpha$   
     $w = w - \alpha \frac{\partial}{\partial w} J(w, b)$  *derivative*  
     $b = b - \alpha \frac{\partial}{\partial b} J(w, b)$

$$J(w)$$

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

$$\min_w J(w)$$



$$w = w - \alpha \frac{\frac{d}{dw} J(w)}{> 0}$$

always +

$w = w - \underline{\alpha} \cdot$  (positive number)

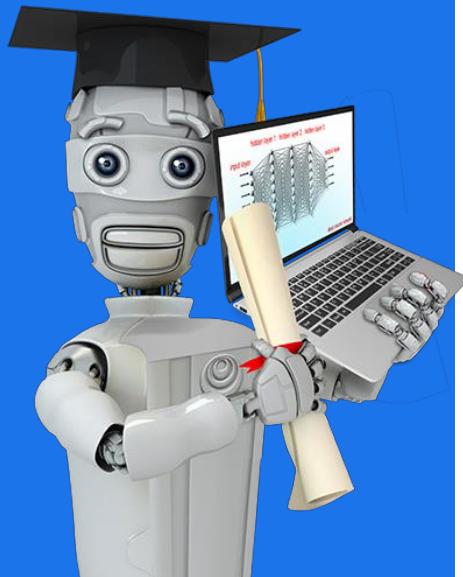
$$\frac{d}{dw} J(w) < 0$$

decrease  $w$

increase  $w$   $\leftarrow w = w - \alpha \cdot$  (negative number)

Stanford  
ONLINE

DeepLearning.AI



# Training Linear Regression

---

## Learning Rate

$$w = w - \alpha \frac{d}{dw} J(w)$$

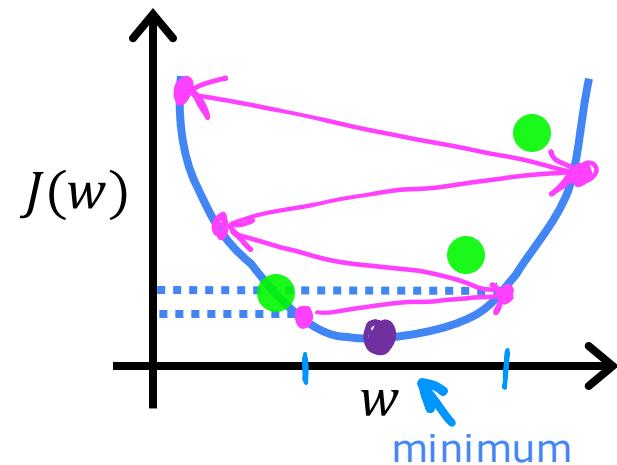
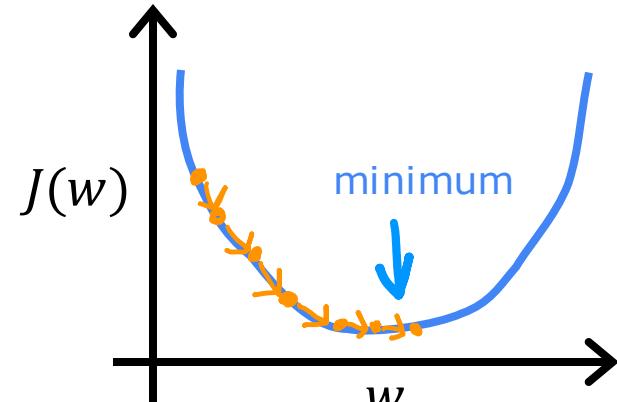
If  $\alpha$  is too small...

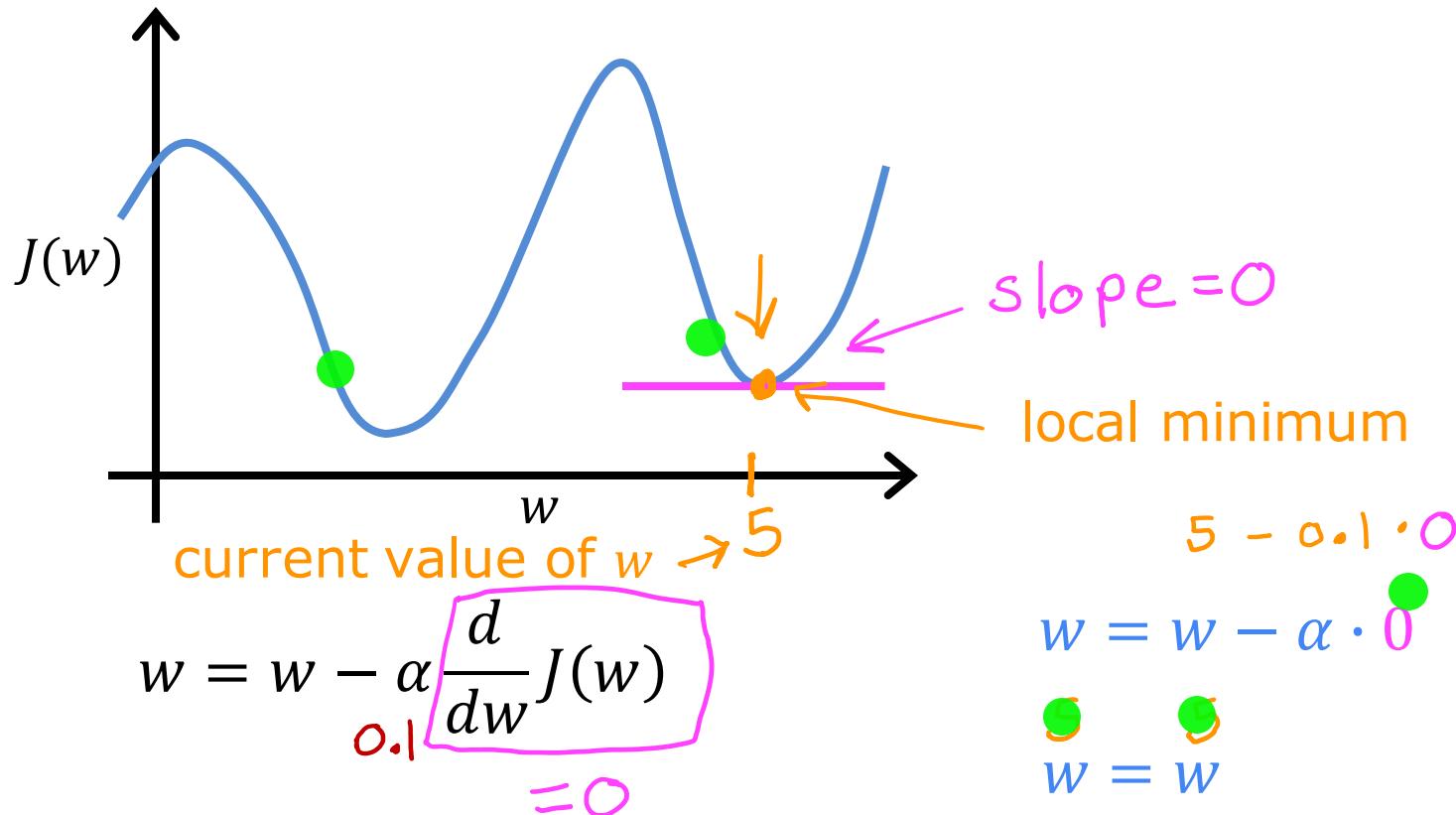
Gradient descent may be slow.

If  $\alpha$  is too large...

Gradient descent may:

- Overshoot, never reach minimum
- Fail to converge, diverge





# Can reach local minimum with fixed learning rate $\alpha$

$$w = w - \alpha \frac{d}{dw} J(w)$$

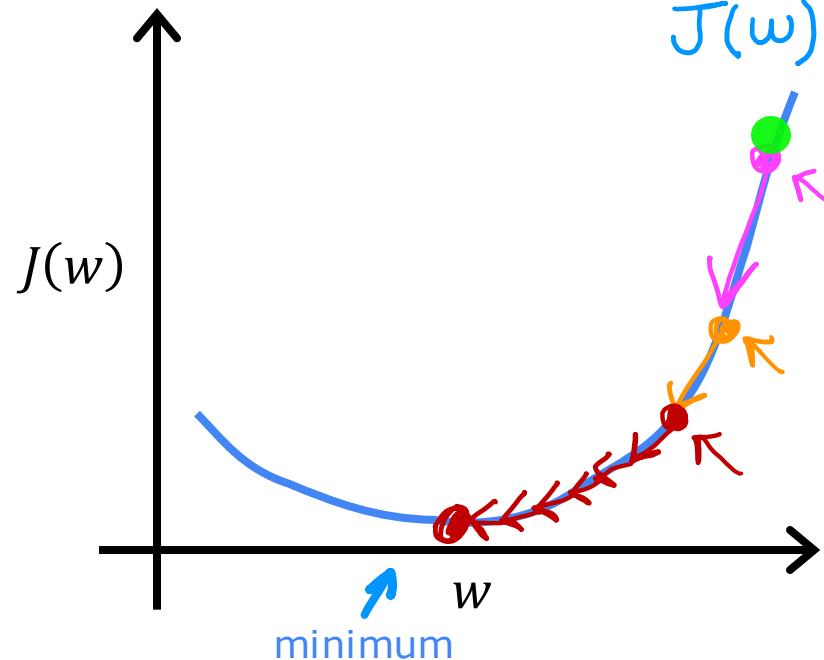
Diagram illustrating the effect of different learning rates  $\alpha$  on the update step:

- smaller**: A small blue step.
- not as large**: An orange step.
- large**: A large red step.

Near a local minimum,

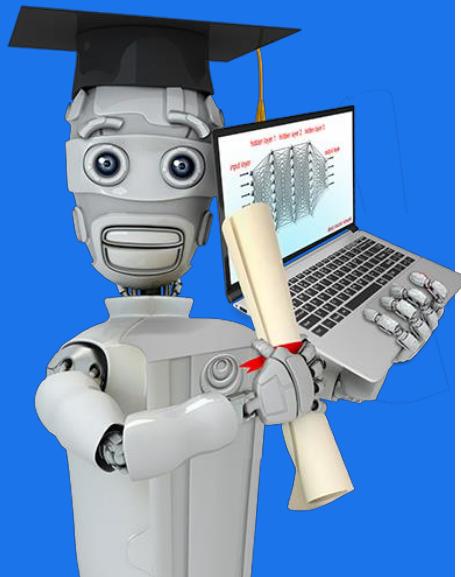
- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate  $\alpha$



Stanford  
ONLINE

DeepLearning.AI



# Training Linear Regression

---

Gradient Descent  
for Linear Regression

## Linear regression model

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

## Cost function

$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$

## Gradient descent algorithm

repeat until convergence {

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

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

}

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

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

next slide  
is optional!

(Optional)

$$\frac{\partial}{\partial w} J(w, b) = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial w} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)})^2$$

$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2x^{(i)}} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})x^{(i)}}$$

$$\frac{\partial}{\partial b} J(w, b) = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2 = \frac{\partial}{\partial b} \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)})^2$$

$$= \cancel{\frac{1}{2m}} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2} = \boxed{\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})}$$

no  $x^{(i)}$

# Gradient descent algorithm

repeat until convergence {

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

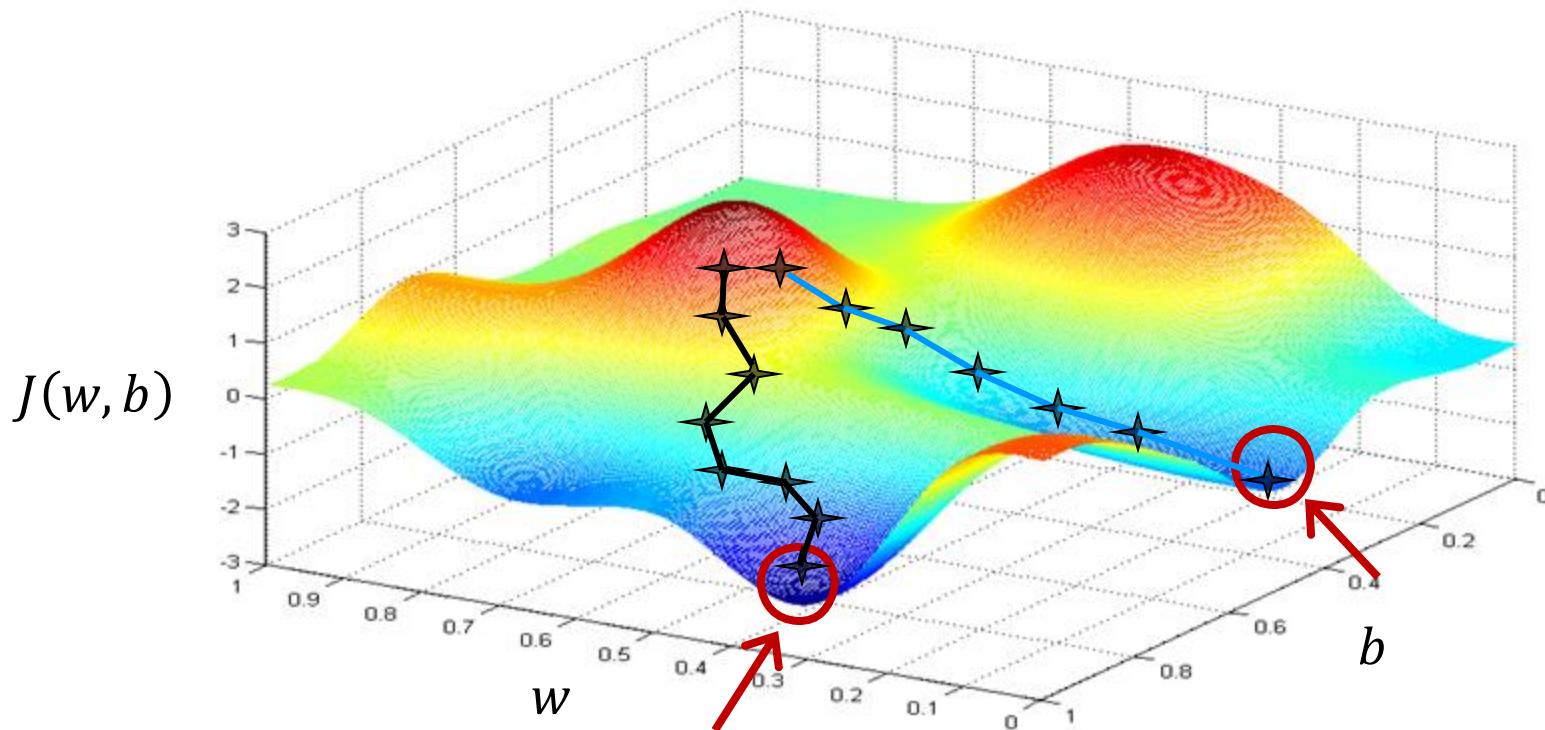
}

Update w and b simultaneously

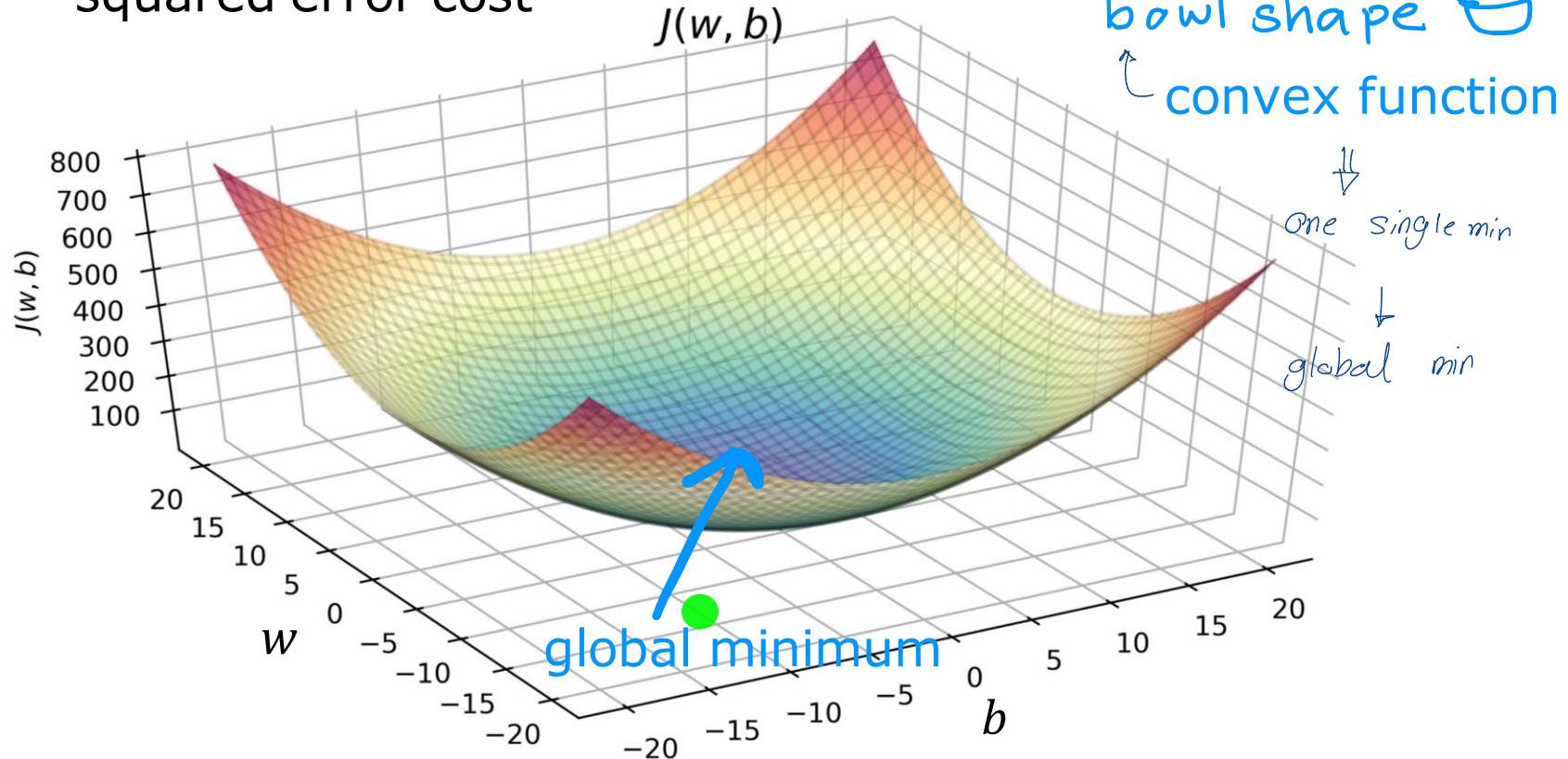
$\frac{\partial}{\partial w} J(w, b)$        $\frac{\partial}{\partial b} J(w, b)$

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

## More than one local minimum

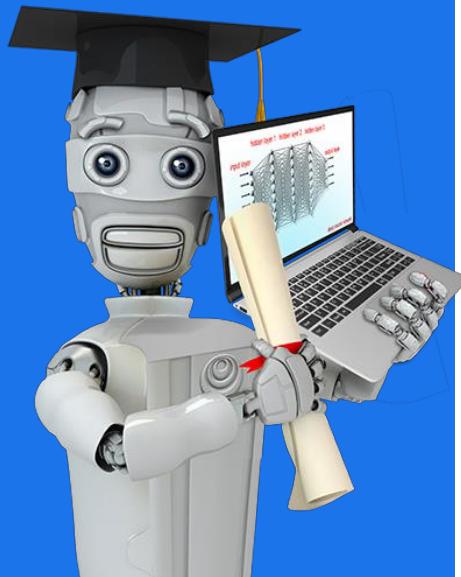


squared error cost



Stanford  
ONLINE

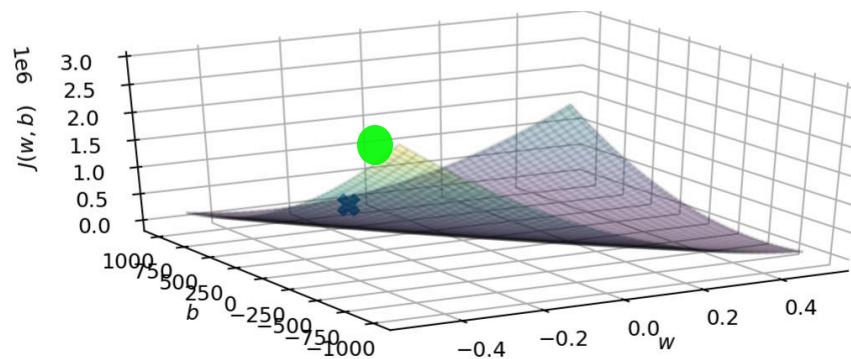
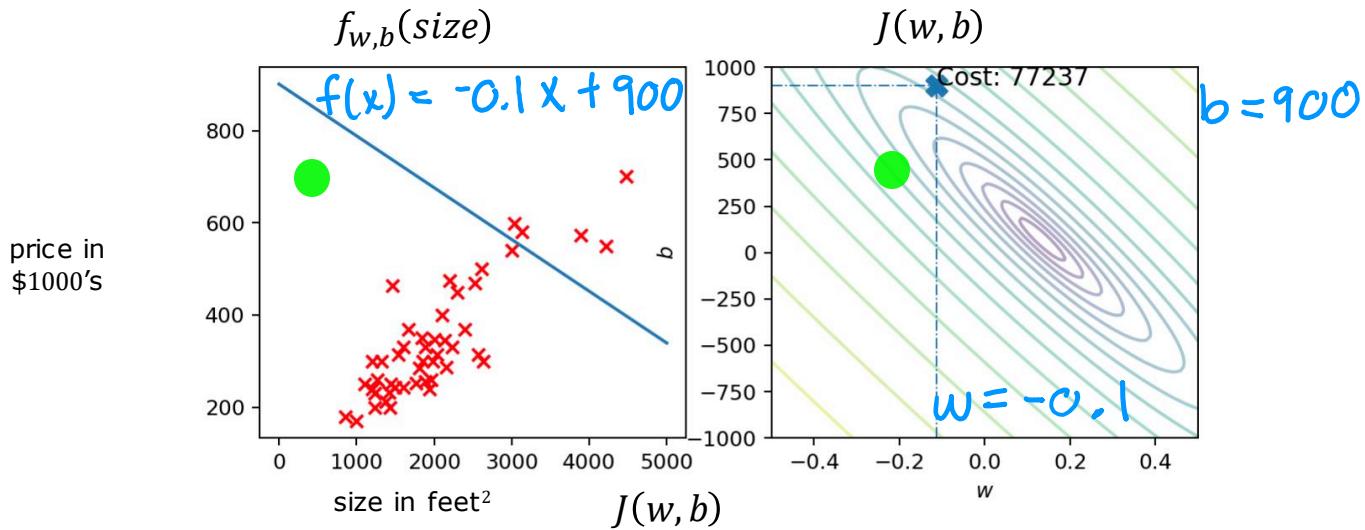
DeepLearning.AI

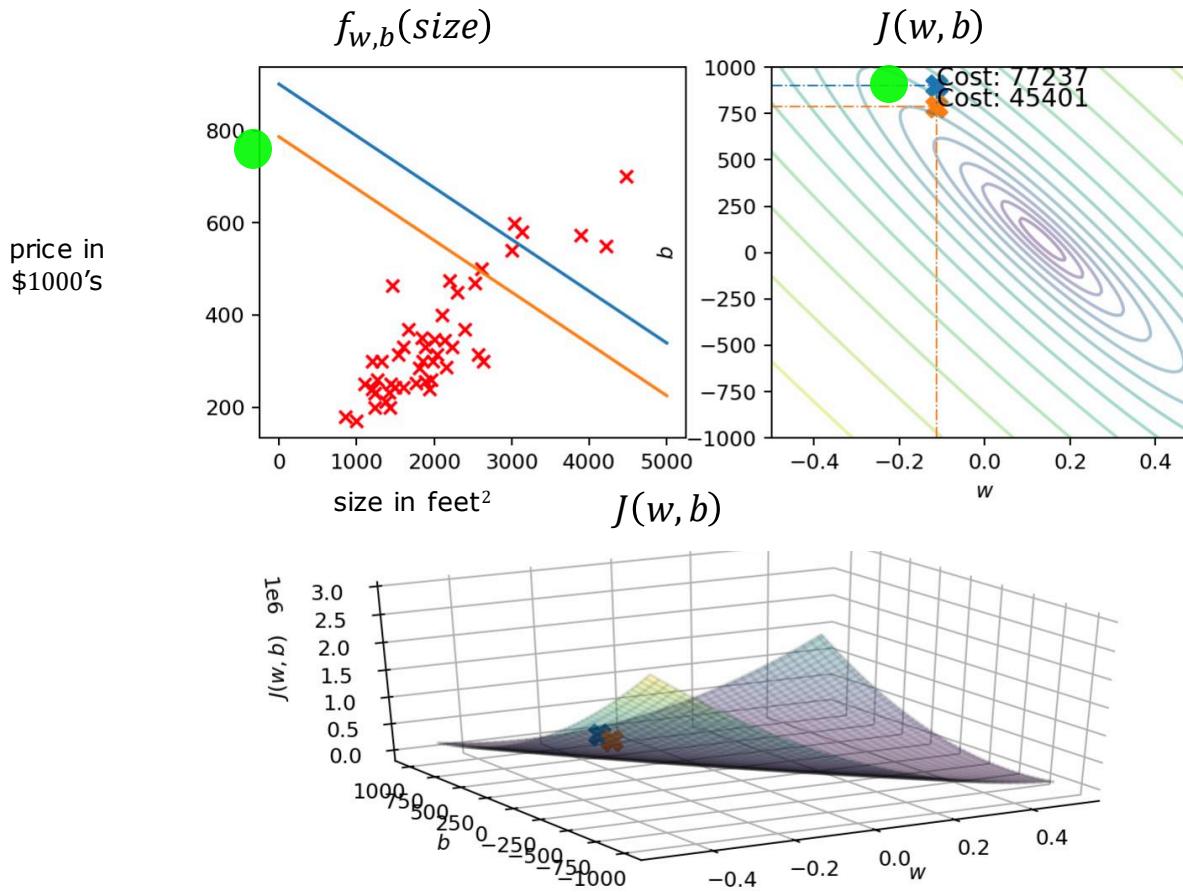


# Training Linear Regression

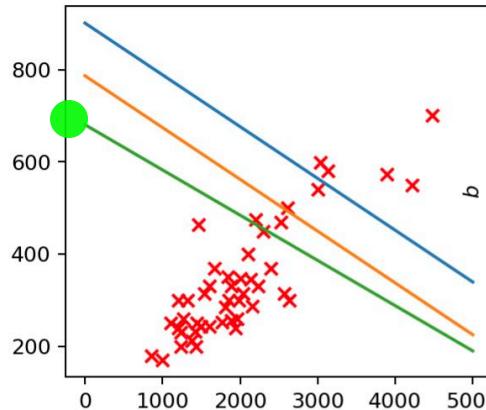
---

## Running Gradient Descent

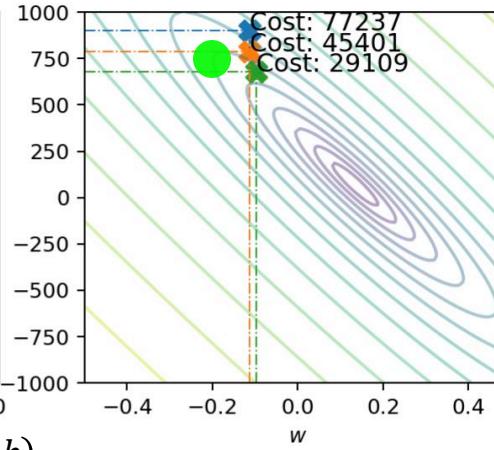




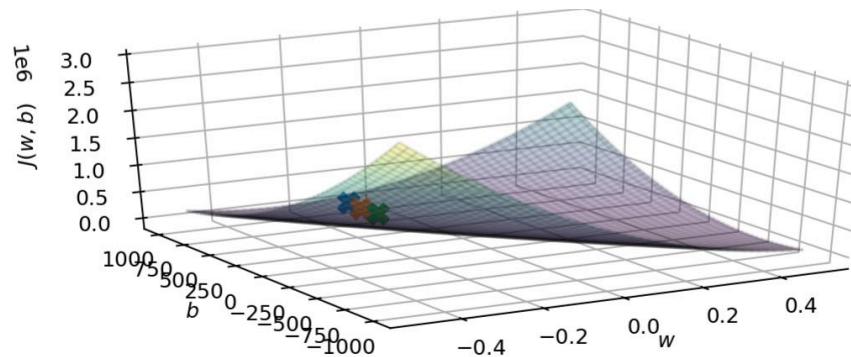
$f_{w,b}(\text{size})$

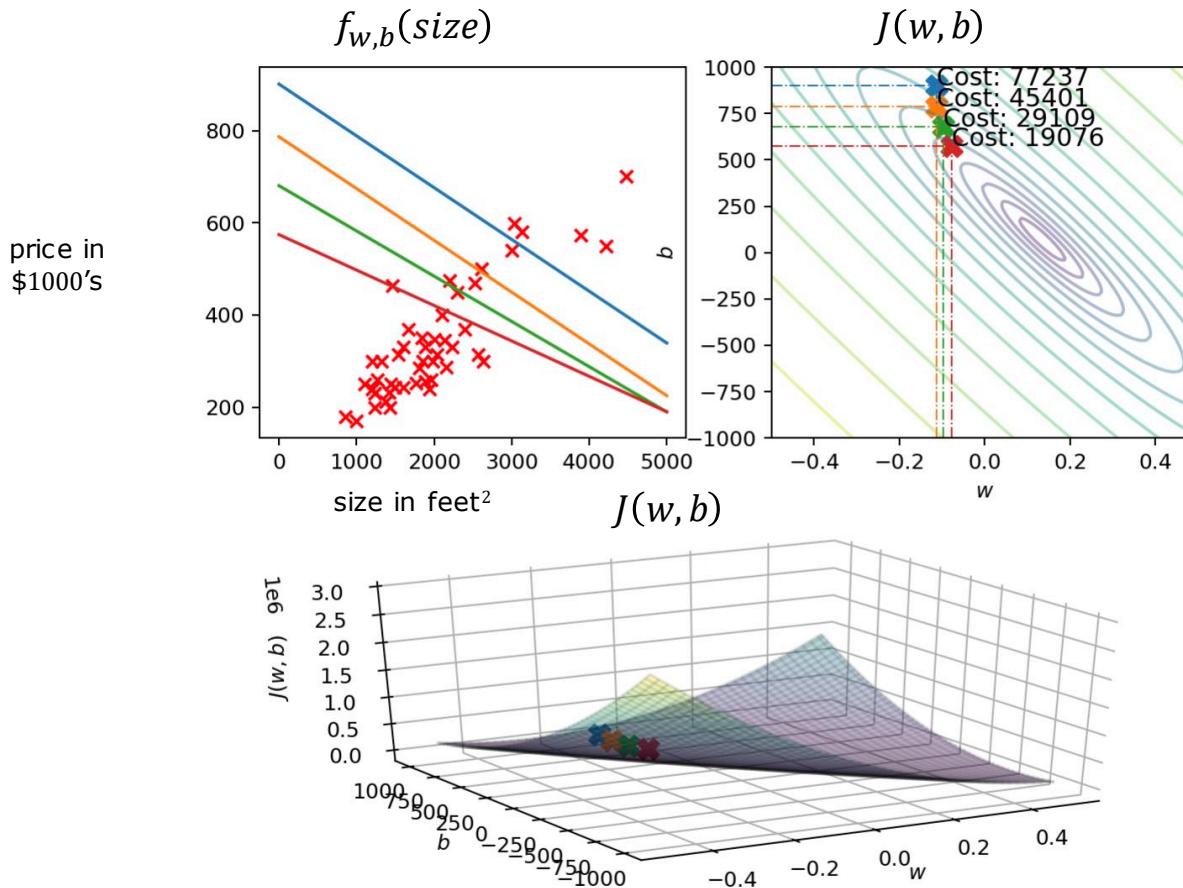


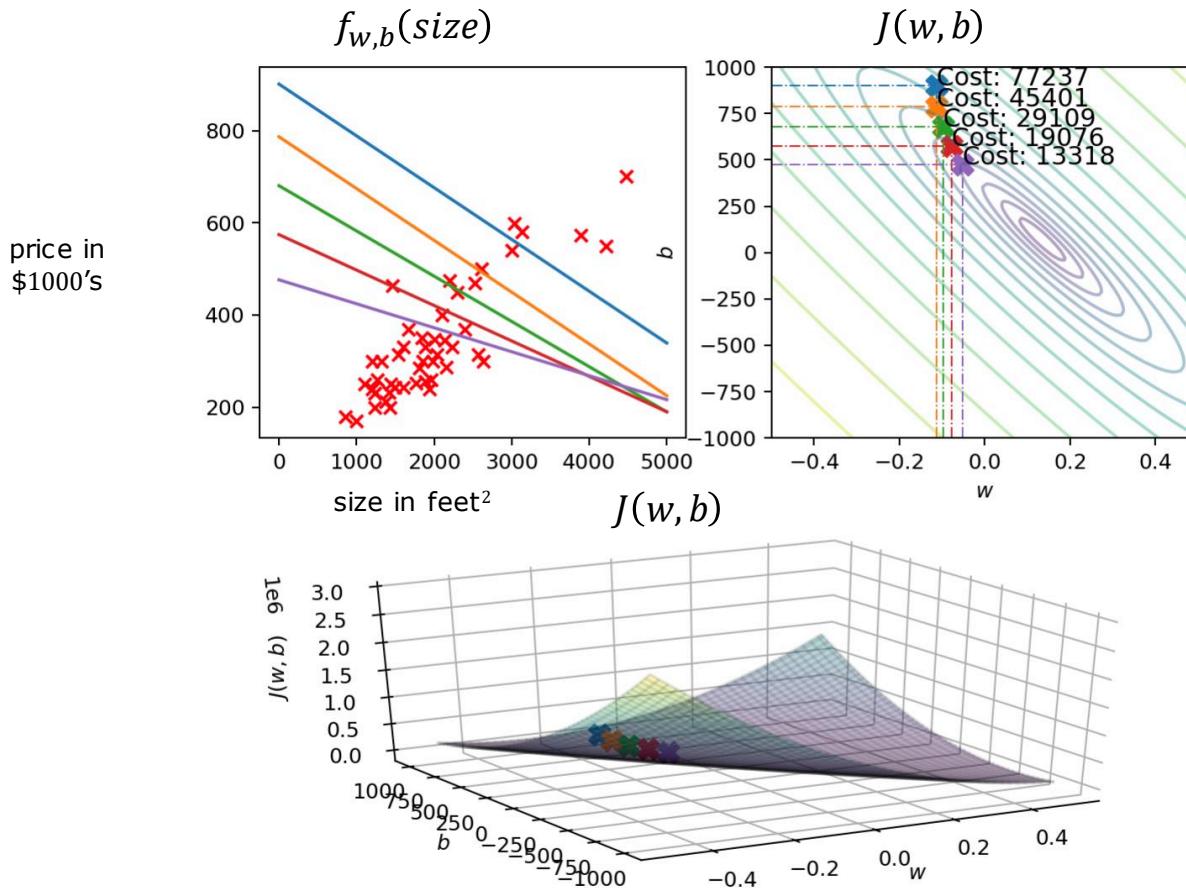
$J(w, b)$

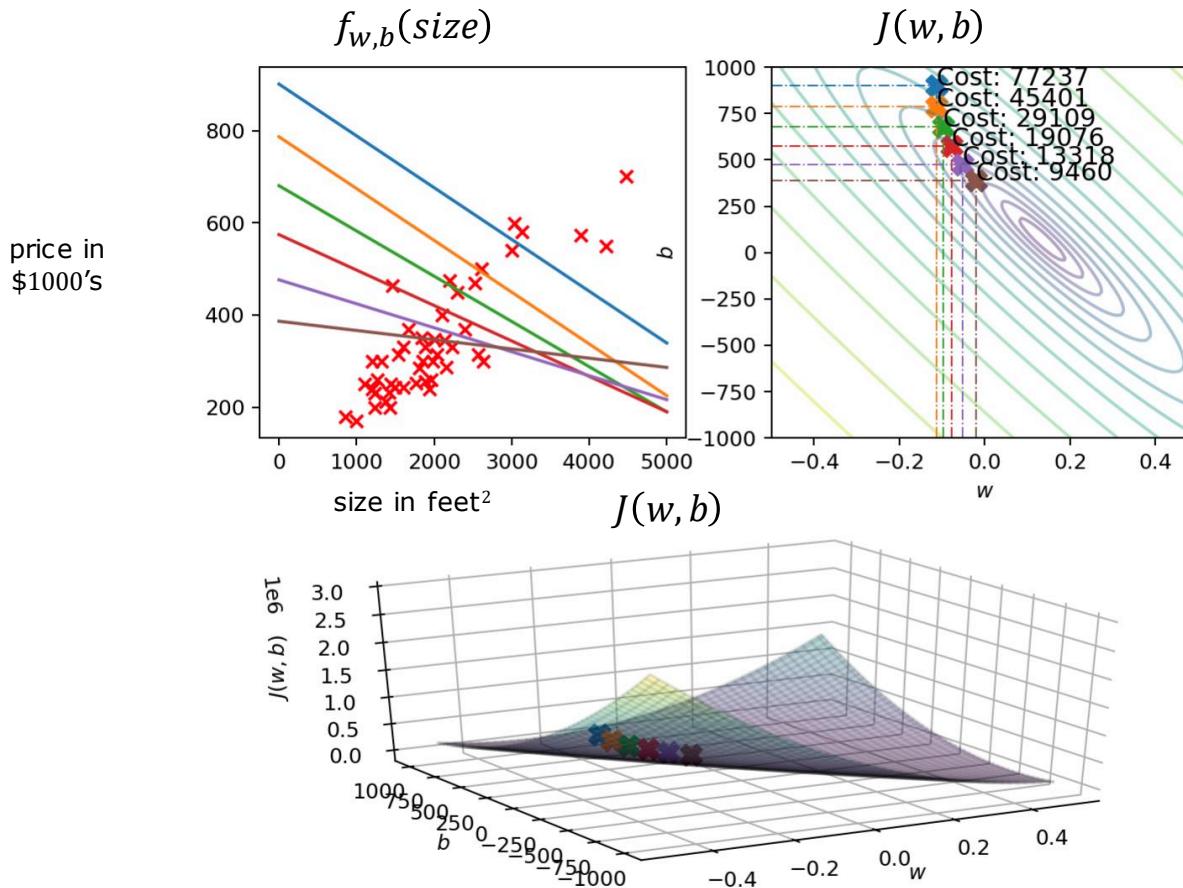


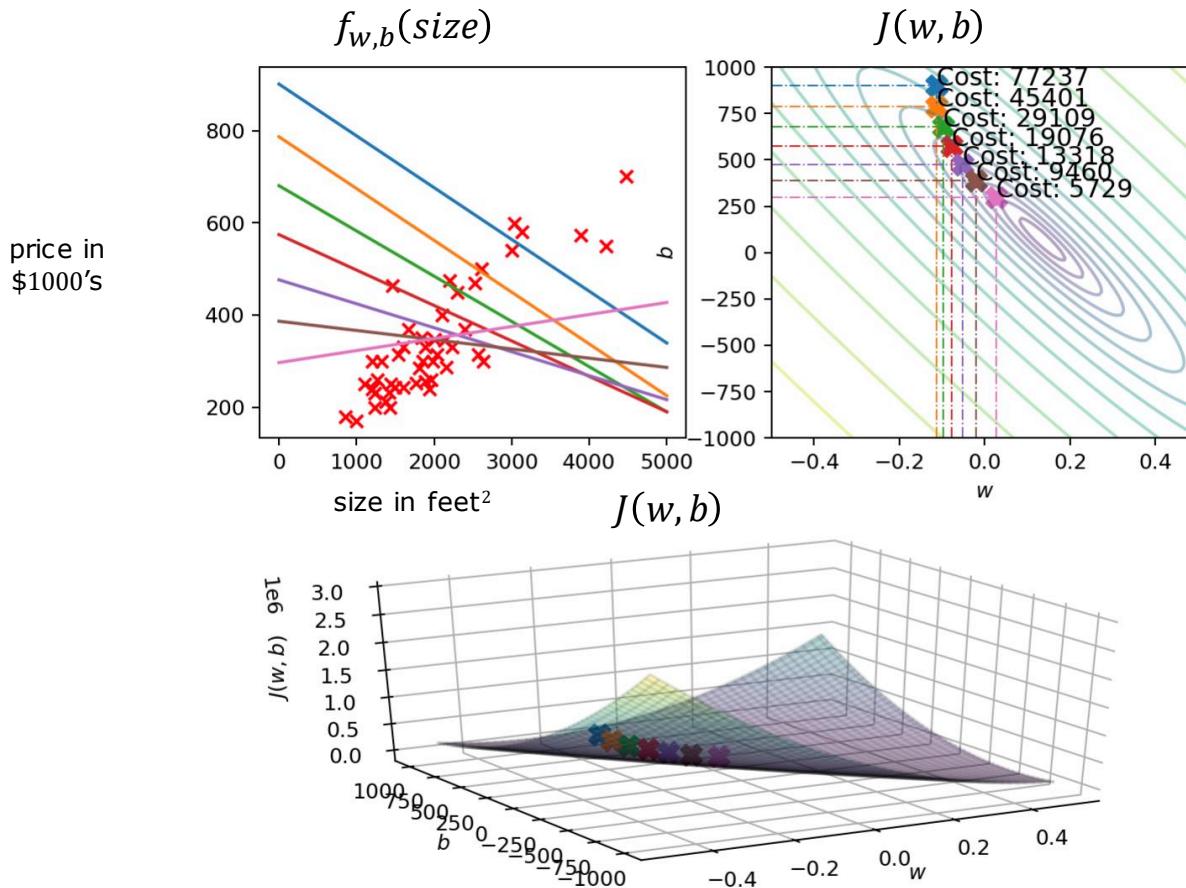
$J(w, b)$



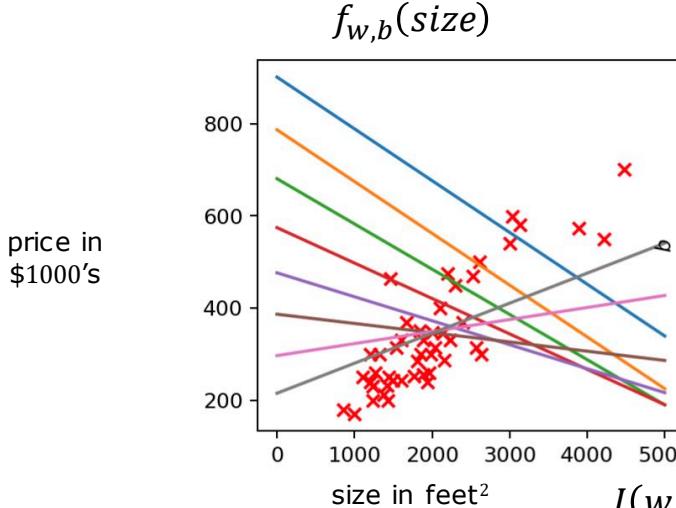




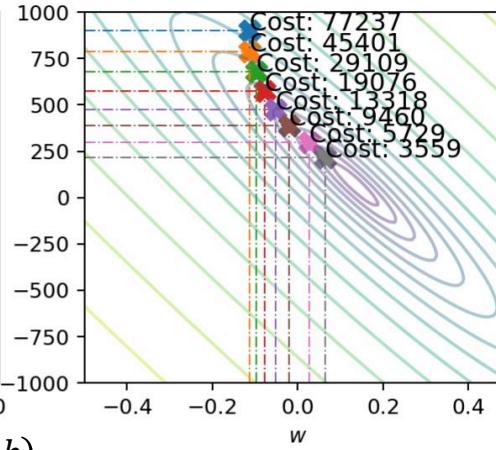




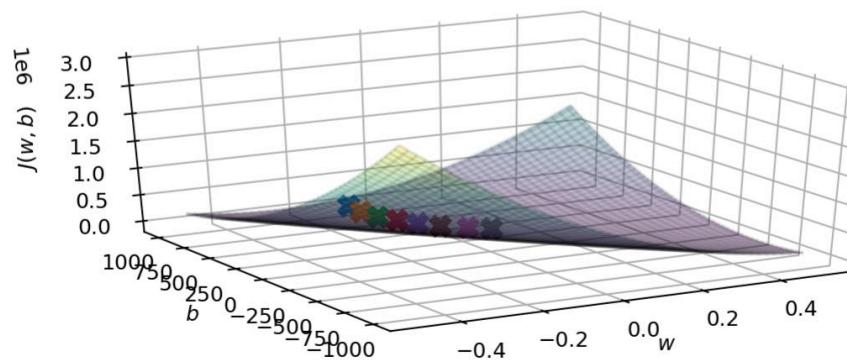
$$f_{w,b}(\text{size})$$

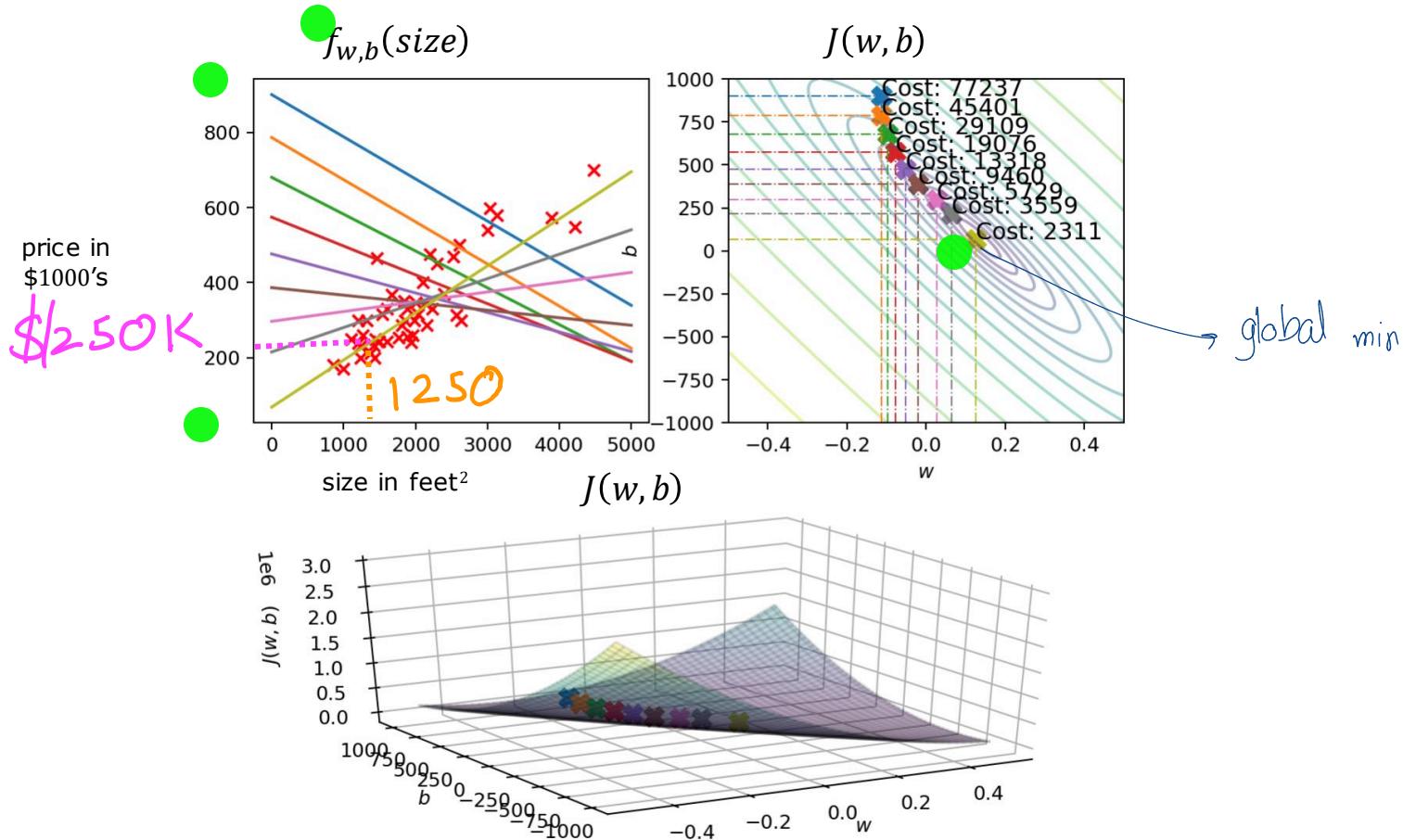


$$J(w, b)$$



$$J(w, b)$$

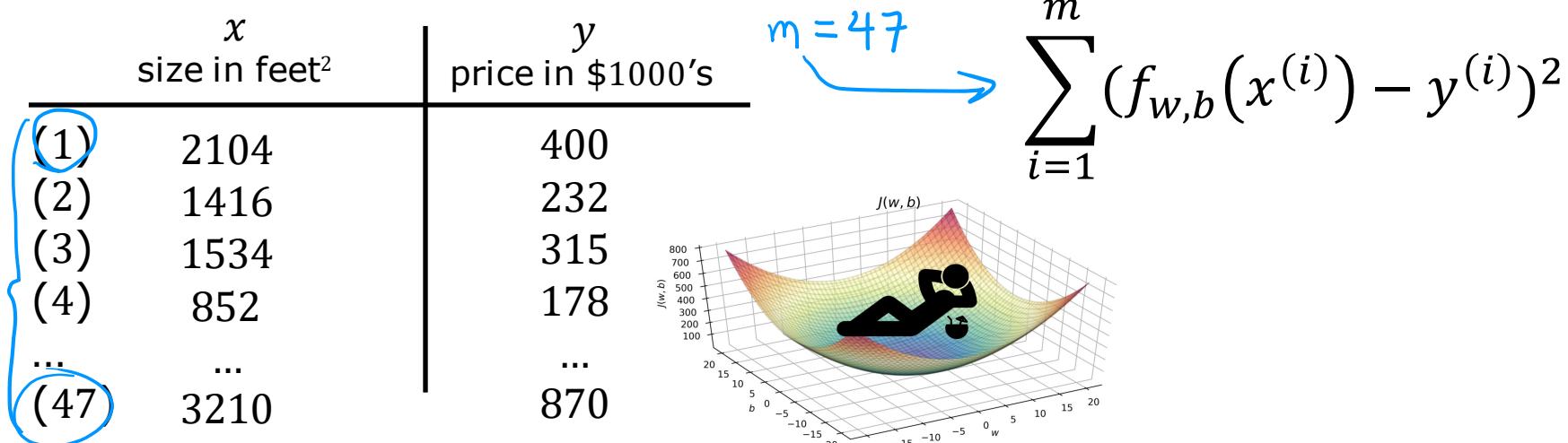




# “Batch” gradient descent

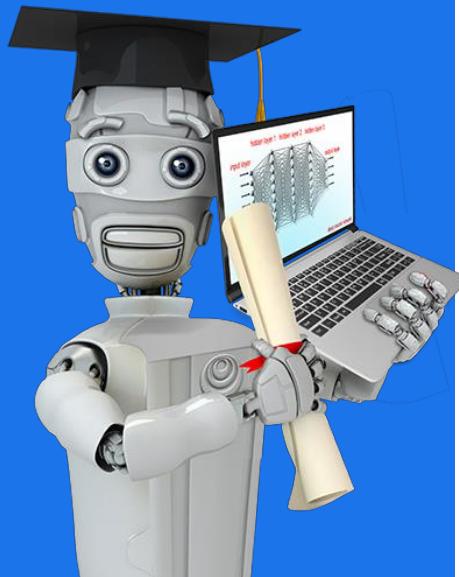
“Batch”: Each step of gradient descent uses all the training examples.

other gradient descent: subsets



Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with Multiple Variables

Week 2

## Multiple Features

# Multiple features (variables)

One  
feature



Size in feet <sup>2</sup> ( $x$ )	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 <sup>2</sup> $x_1$	Number of bedrooms $x_2$	Number of floors $x_3$	Age of home in years $x_4$	Price (\$) in \$1000's
$i=2$	2104	5	1	45	460
	1416	3	2	40	232
	1534	3	2	30	315
	852	2	1	36	178
	...	...	...	...	...

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

$x_j = j^{th}$  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

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

$x_3^{(2)} = 2$

# 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

$$f_{w,b}(x) = 0.1 \underset{\text{size}}{\uparrow} x_1 + 4 \underset{\text{\#bedrooms}}{\uparrow} x_2 + 10 \underset{\text{\#floors}}{\uparrow} x_3 + -2 \underset{\text{years}}{\uparrow} x_4 + 80 \underset{\text{base price}}{\uparrow}$$

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

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

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

b is a number

parameters  
of the model

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

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

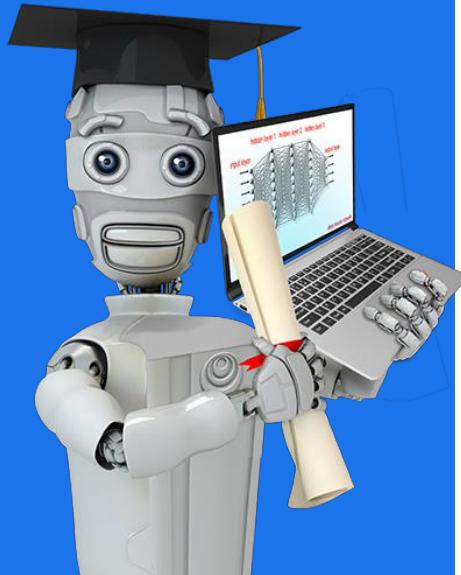
dot product

multiple linear regression

(not multivariate regression)

Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with Multiple Variables

---

## Vectorization Part 1

## 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

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

```
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 $n=100,000$

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

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



## 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  
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
```



Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with Multiple Variables

---

Vectorization  
Part 2

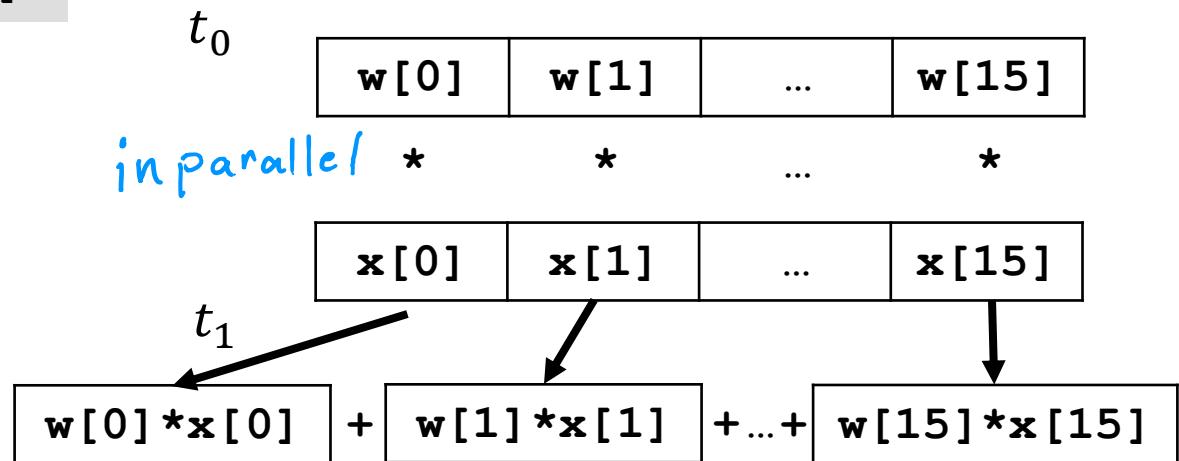
## Without vectorization

```
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)
```



efficient → scale to large datasets

Gradient descent

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

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

Without vectorization

$$w_1 = w_1 - 0.1d_1$$

$$w_2 = w_2 - 0.1d_2$$

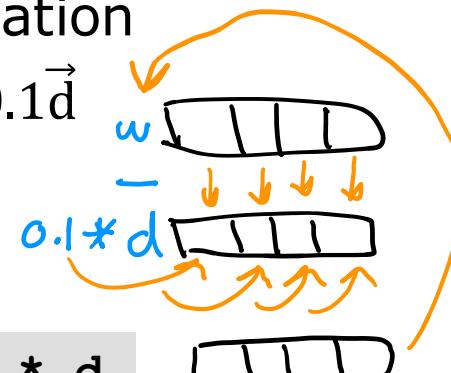
:

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

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

With vectorization

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

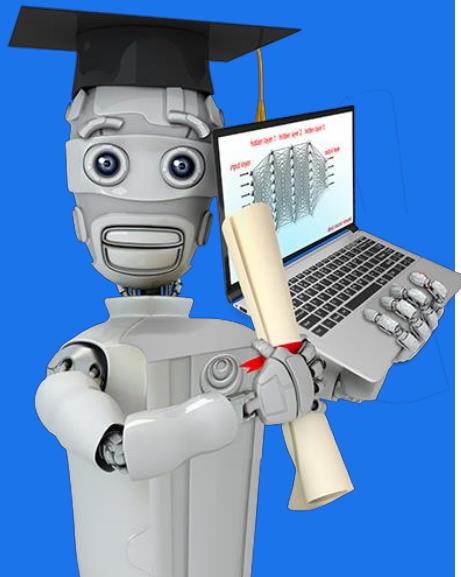


in code :

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

Stanford  
ONLINE

DeepLearning.AI



# Linear Regression with Multiple Variables

---

## Gradient Descent for Multiple Regression

## Previous notation

Parameters

$$w_1, \dots, w_n$$

$$\bullet \quad b$$

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

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

Gradient descent

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

## Vector notation

$\vec{w}$  *vector of length n*  
 $\vec{w} = [w_1 \dots w_n]$   
 $b$  still a number  
 $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$   
 $J(\vec{w}, b)$  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 {

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

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

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

simultaneously update  $w, b$

}

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

repeat {

$$\begin{aligned} j &= 1 & \underline{w}_1 &= w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\overrightarrow{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \underline{x}_1^{(i)} \\ &\vdots & & \hookrightarrow \frac{\partial}{\partial \underline{w}_1} J(\overrightarrow{w}, b) \\ j &= n & w_n &= w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\overrightarrow{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)} \end{aligned}$$

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

simultaneously update  
 $w_j$  (for  $j = 1, \dots, n$ ) 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  $w, b$



Stanford  
ONLINE

DeepLearning.AI



# Practical Tips for Linear Regression

enable gradient descend to run much faster

## Feature Scaling

### Part 1

# Feature and parameter values

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

↓      ↓  
size   #bedrooms

$x_1$ : size (feet<sup>2</sup>)  
range: 300 – 2,000

$x_2$ : # bedrooms  
range: 0 – 5



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

one training example

size of the parameters  $w_1, w_2$ ?



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



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

small      large

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

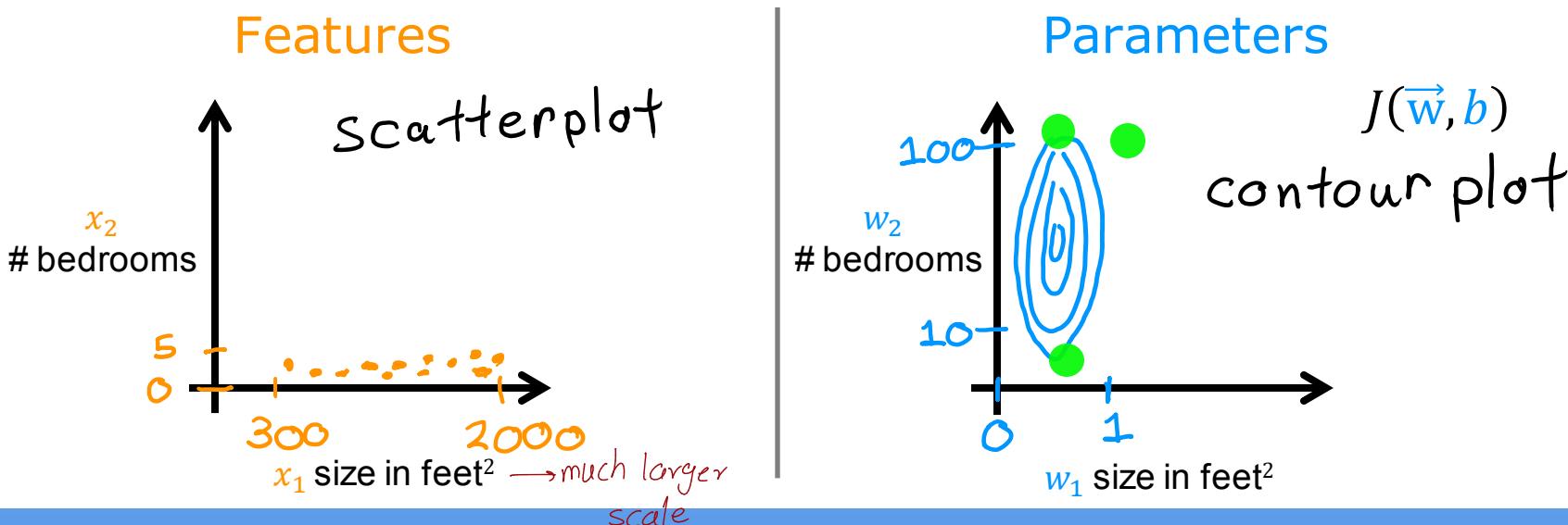
$$\widehat{\text{price}} = \$500k \quad \text{more reasonable}$$

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

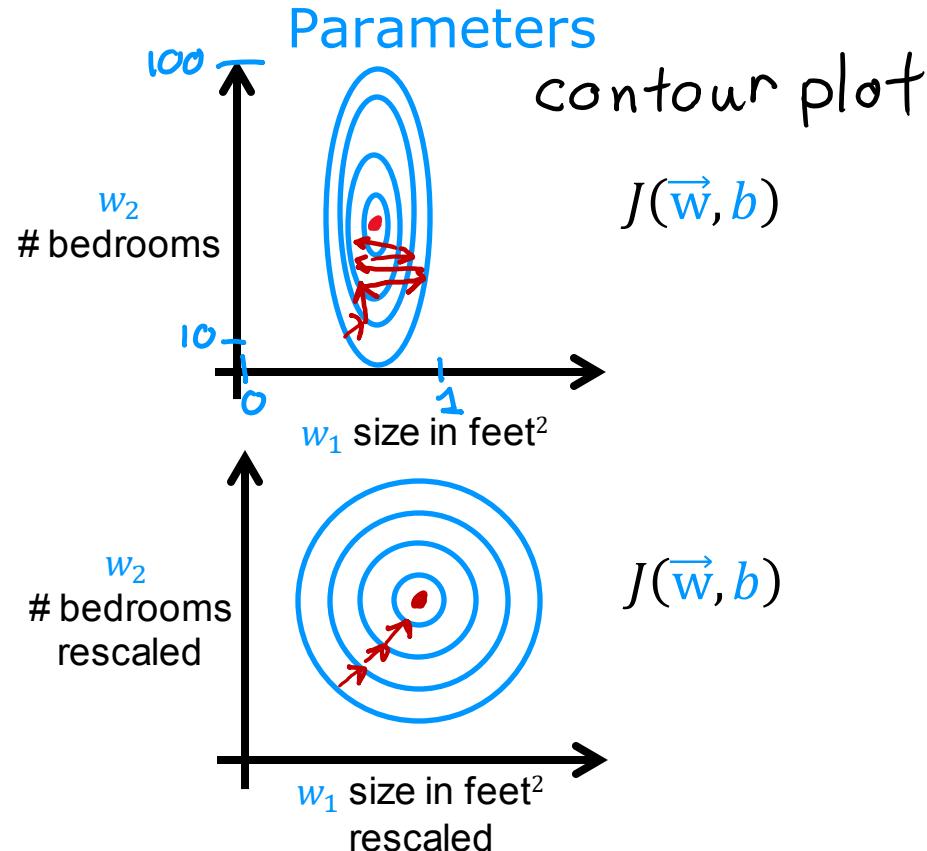
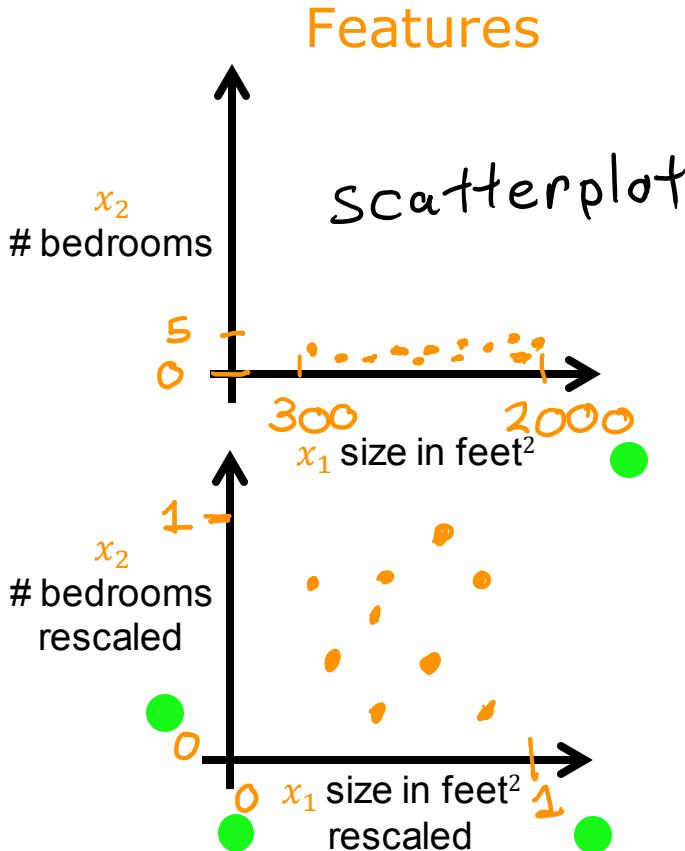
not good

# Feature size and parameter size

	size of feature $x_j$	size of parameter $w_j$
size in feet <sup>2</sup>	↔	↔
#bedrooms	↔	↔↔

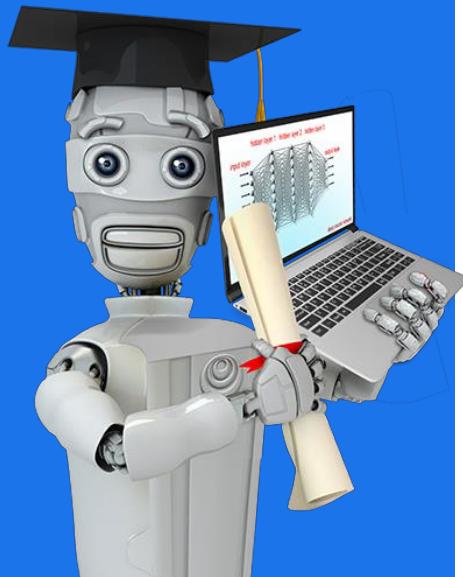


# Feature size and gradient descent



Stanford  
ONLINE

DeepLearning.AI

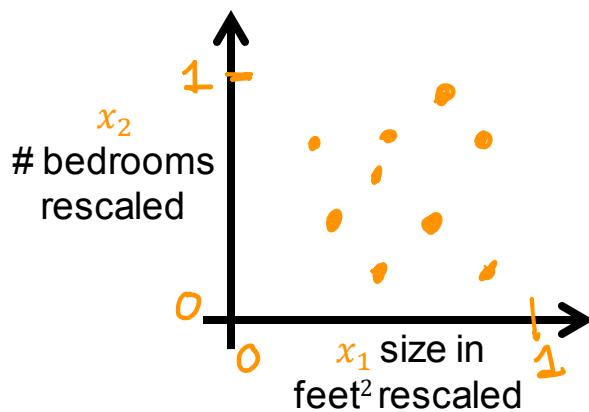
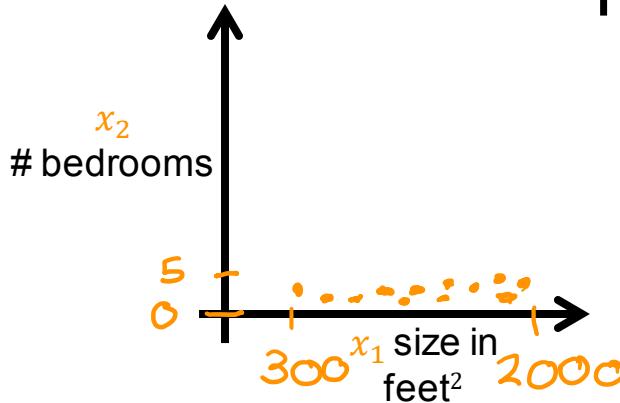


# Practical Tips for Linear Regression

---

## Feature Scaling Part 2

# Feature scaling



$$300 \leq x_1 \leq 2000$$

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

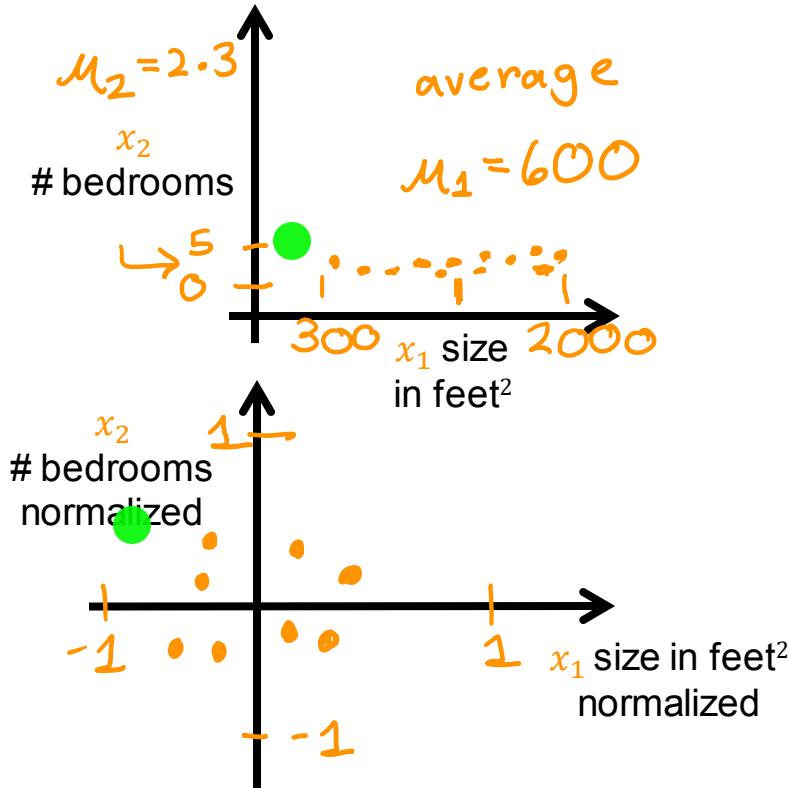
$$0.15 \leq x_{1,scaled} \leq 1$$

$$0 \leq x_2 \leq 5$$

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

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

# Mean normalization



$$300 \leq x_1 \leq 2000$$

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

max-min

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

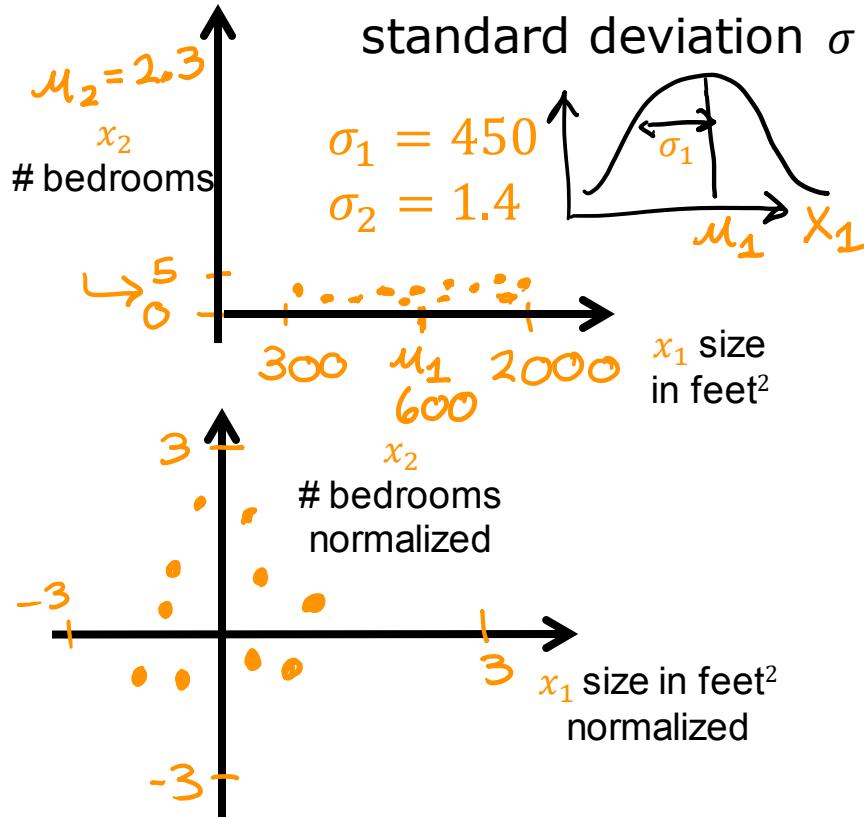
$$0 \leq x_2 \leq 5$$

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

max-min

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

# Z-score normalization



$$300 \leq x_1 \leq 2000$$

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

$$-0.67 \leq x_1 \leq 3.1$$

$$0 \leq x_2 \leq 5$$

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

$$-1.6 \leq x_2 \leq 1.9$$

# Feature scaling

aim for about  $-1 \leq x_j \leq 1$  for each feature  $x_j$   
 $-3 \leq x_j \leq 3$   
 $-0.3 \leq x_j \leq 0.3$

$\left. \begin{matrix} -3 \leq x_j \leq 3 \\ -0.3 \leq x_j \leq 0.3 \end{matrix} \right\}$  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

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

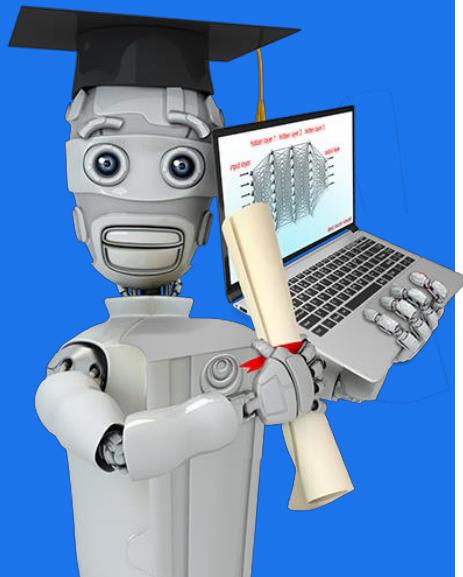
too small → rescale

$$98.6 \leq x_5 \leq 105$$

too large → rescale

Stanford  
ONLINE

DeepLearning.AI



# Practical Tips for Linear Regression

---

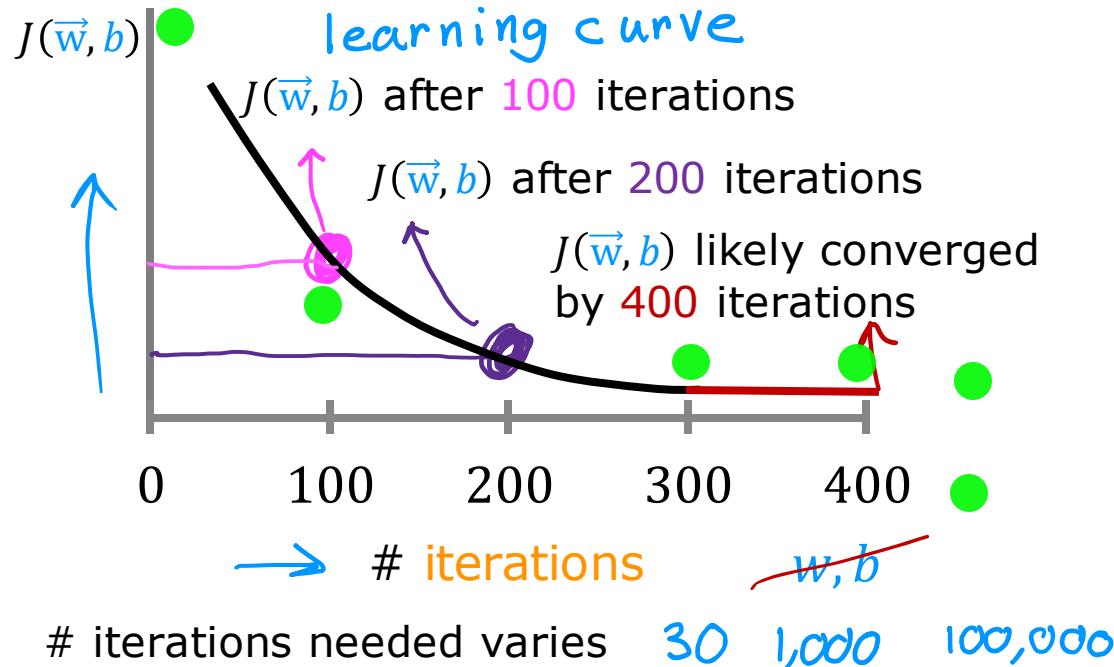
## Checking Gradient Descent for Convergence

# Gradient descent

$$\left\{ \begin{array}{l} w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b) \\ b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b) \end{array} \right.$$

# 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



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)

Stanford  
ONLINE

DeepLearning.AI

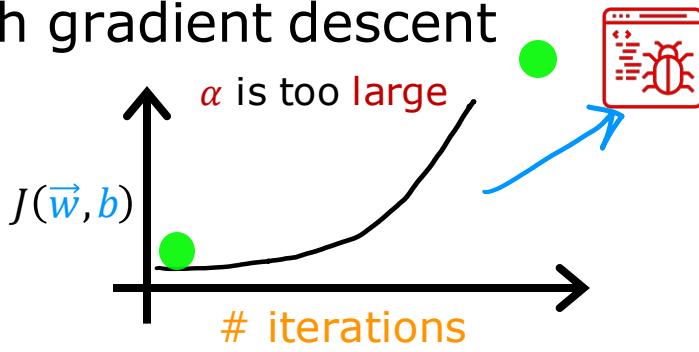
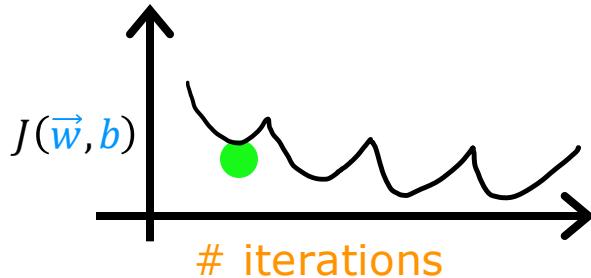


# Practical Tips for Linear Regression

---

## Choosing the Learning Rate

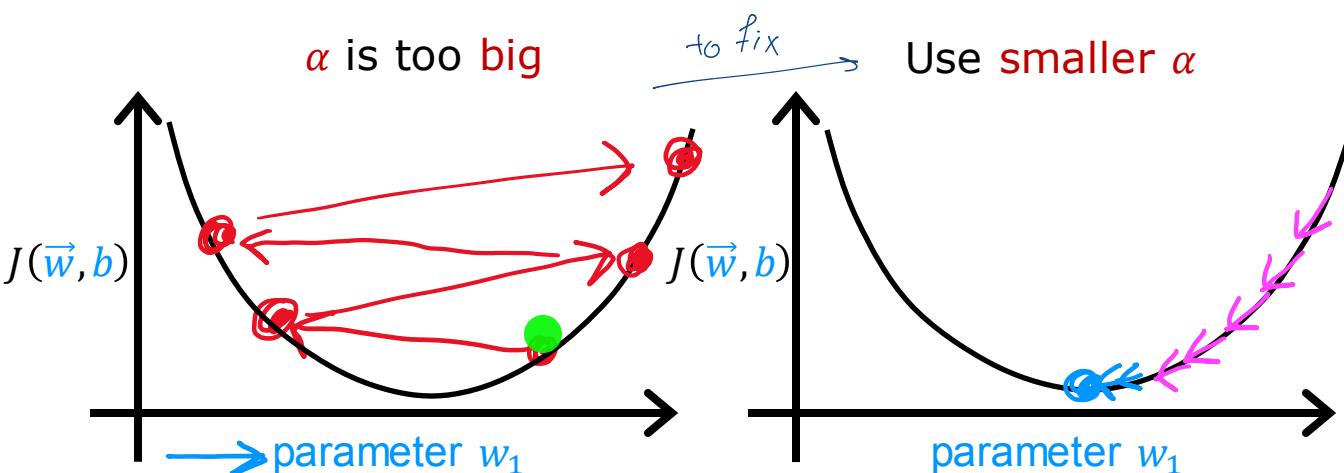
# Identify problem with gradient descent



or learning rate is too large

$w_1 = w_1 + \alpha d_1$  increase  
use a minus sign  
 $w_1 = w_1 - \alpha d_1$

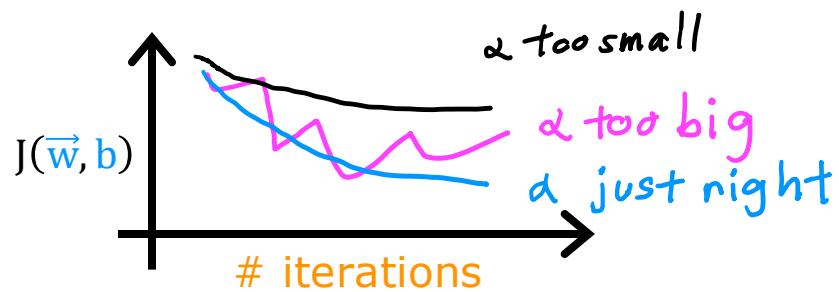
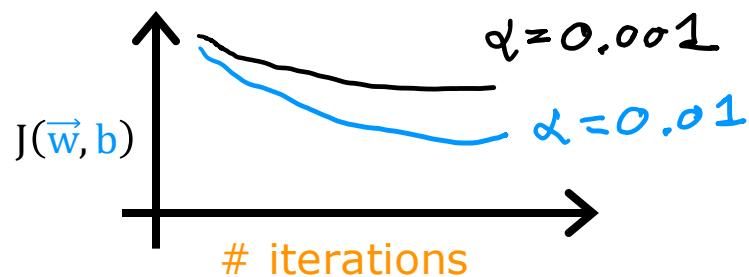
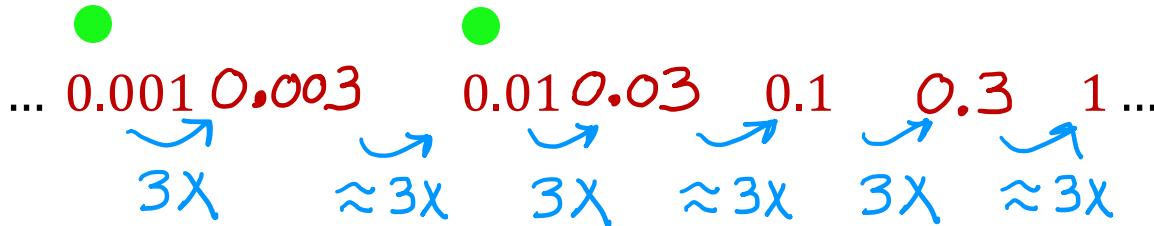
## Adjust learning rate



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:



Stanford  
ONLINE

DeepLearning.AI



# Practical Tips for Linear Regression

---

## Feature Engineering

# Feature engineering

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

frontage      depth

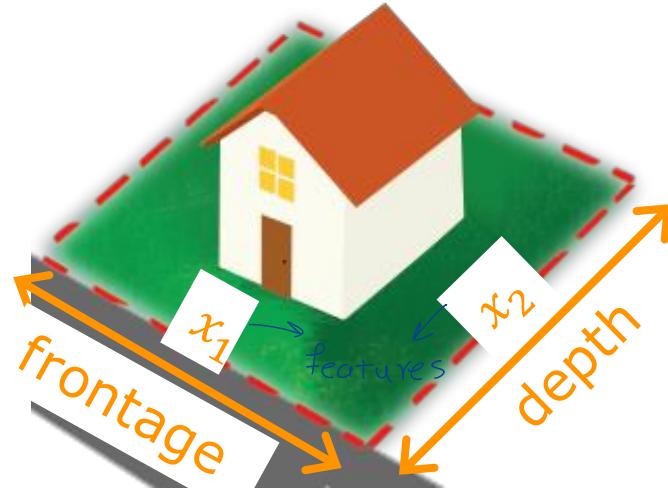
$$\text{area} = \text{frontage} \times \text{depth}$$

$$x_3 = x_1 x_2$$

new feature

Creating new feature  
↓  
Feature Engineering

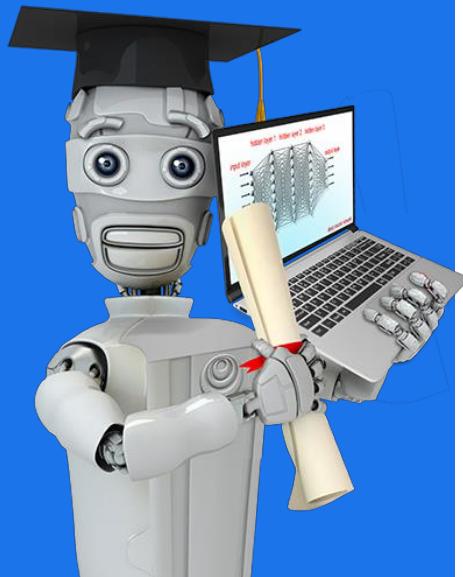
$$f_{\vec{w}, b}(\vec{x}) = \underline{w_1 x_1} + \underline{w_2 x_2} + \underline{w_3 x_3} + b$$



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

Stanford  
ONLINE

DeepLearning.AI



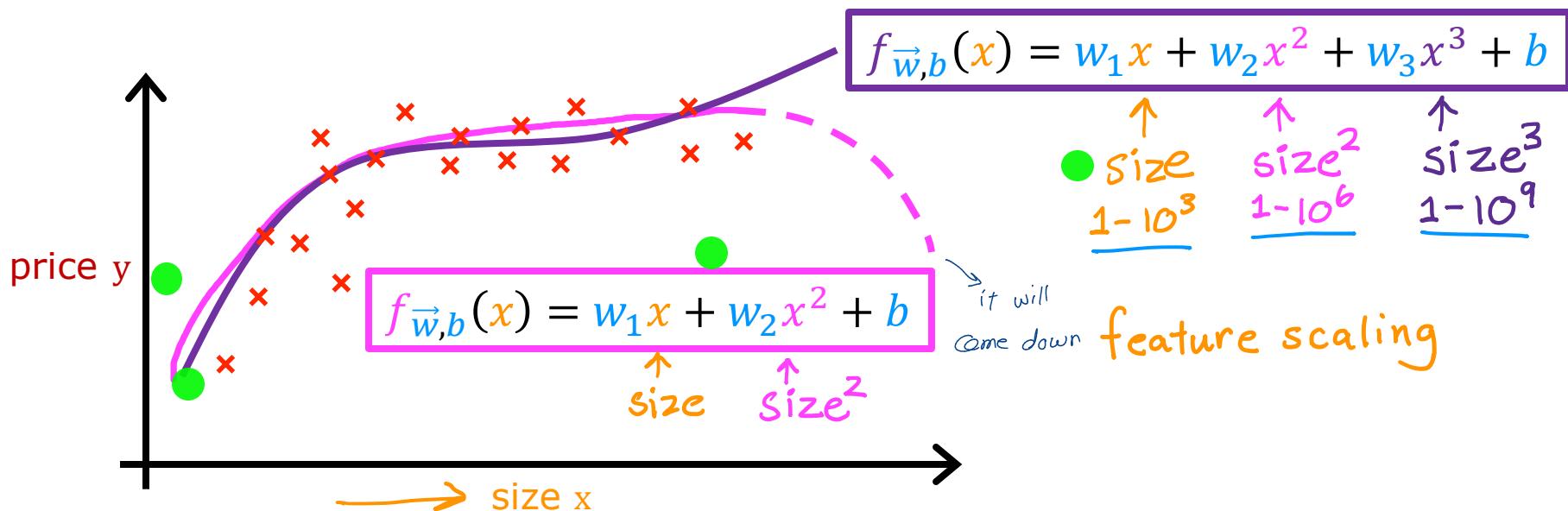
# Practical Tips for Linear Regression

multiple linear regression + feature engineering

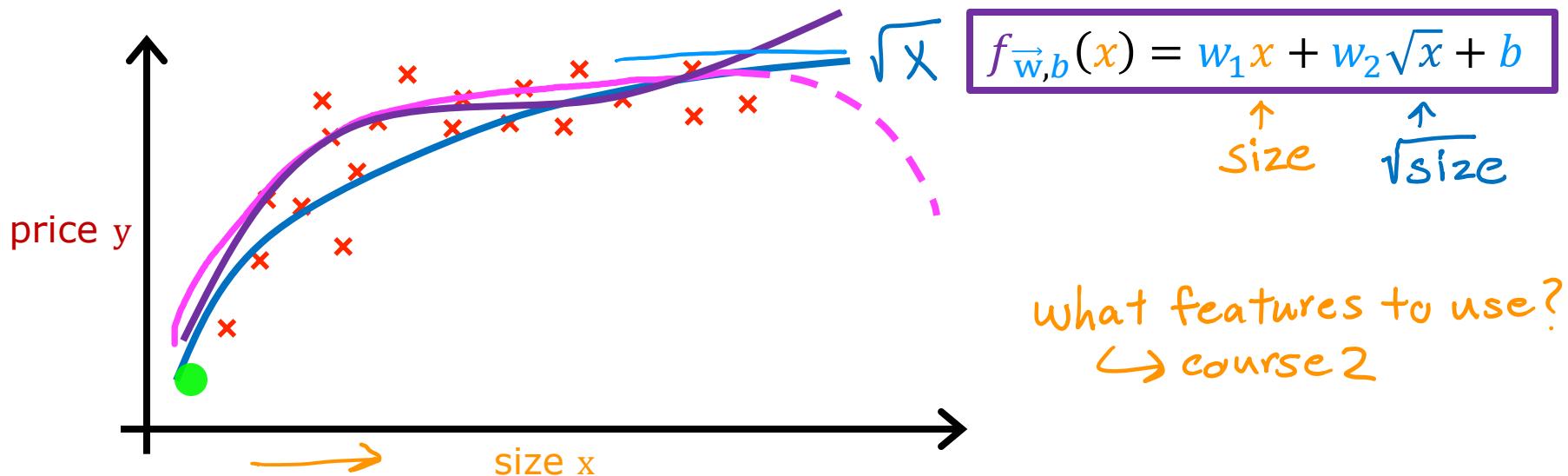
## Polynomial Regression

let you fit curves non-linear functions, to your data

# Polynomial regression



# Choice of features



Week 3

Stanford  
ONLINE

DeepLearning.AI



# Classification

---

## Motivations

# Classification

Question

Is this email spam?

Is the transaction fraudulent?

Is the tumor malignant?

Answer " $y$ "

no	yes
no	yes
no	yes

$y$  can only be one of **two** values

"binary classification"

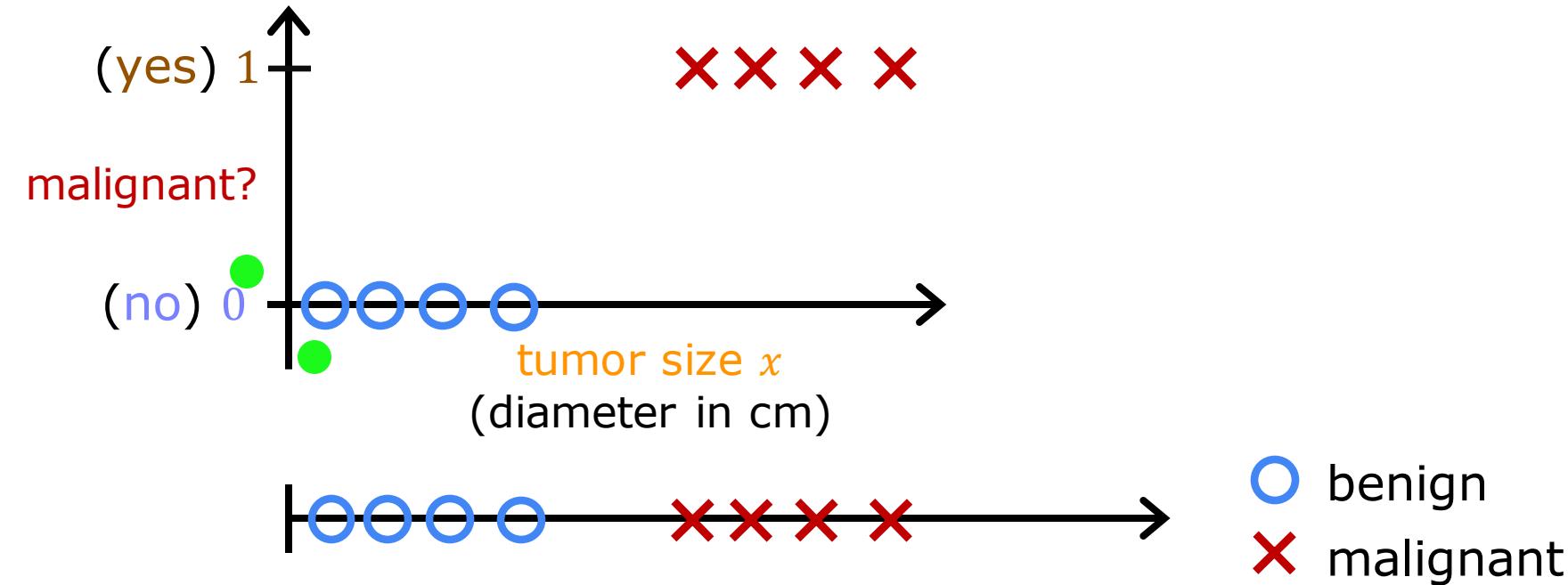
class = category

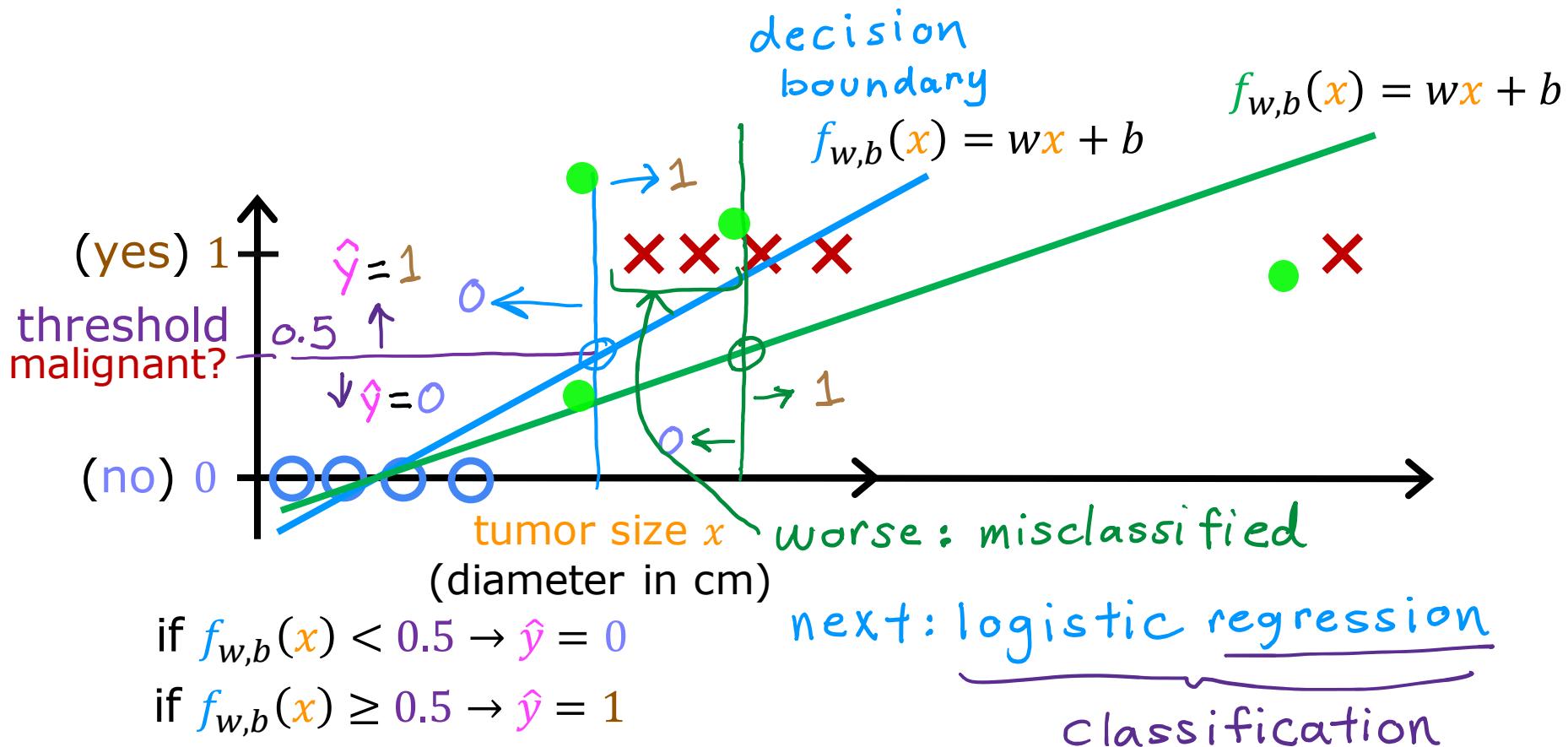
"negative class"  
 $\neq$  "bad"  
absence

false      true  
0            1

useful for  
classification

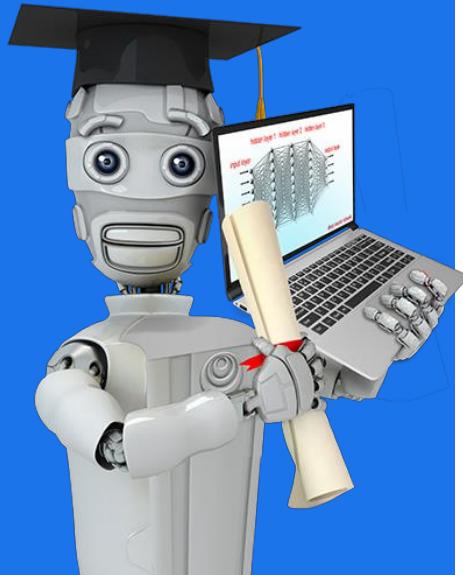
"positive class"  
 $\neq$  "good"  
presence





Stanford  
ONLINE

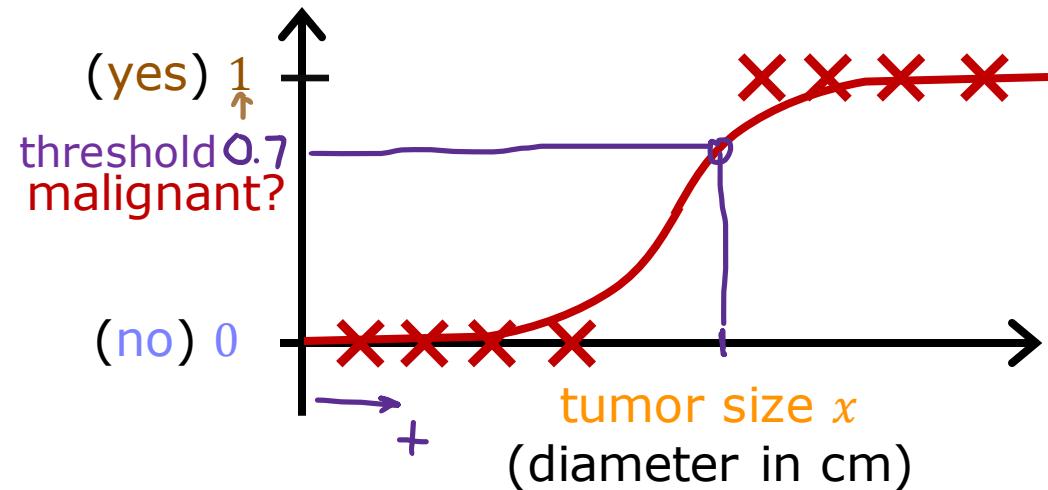
DeepLearning.AI



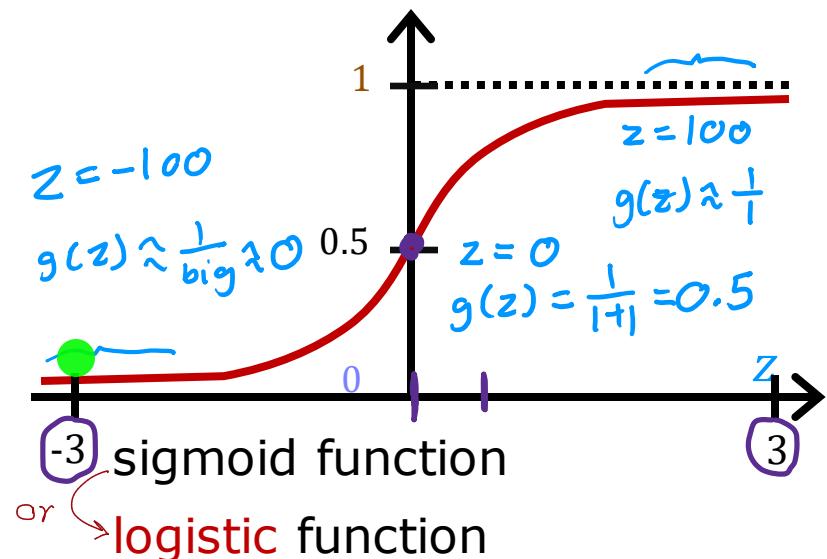
# Classification

Binary classification

Logistic Regression



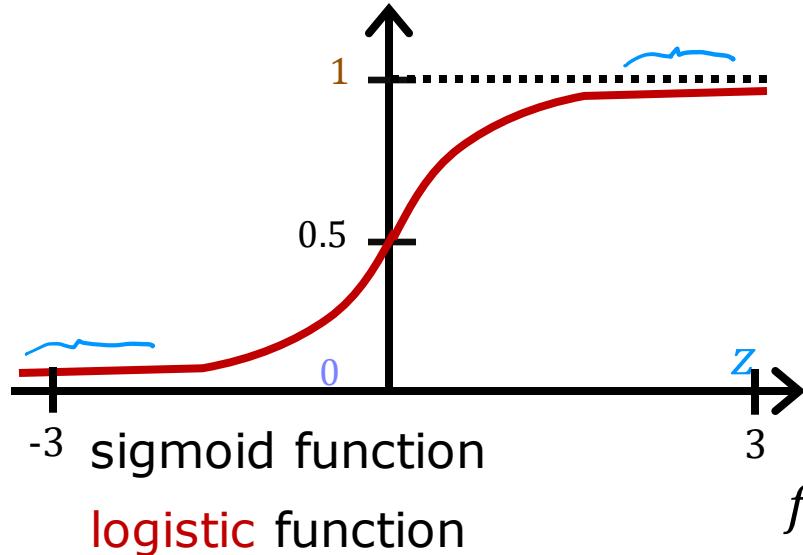
Want outputs between 0 and 1



outputs between 0 and 1

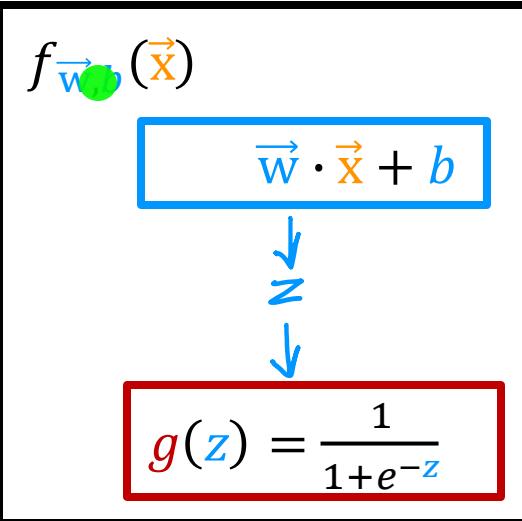
$$\text{Sigmoid Func} \rightarrow g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$

Want outputs between 0 and 1



outputs between 0 and 1

$$g(z) = \frac{1}{1+e^{-z}} \quad 0 < g(z) < 1$$



$$f_{\vec{w}, b}(\vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x} + b}_z) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

"logistic regression"

$$e \approx 2.7$$

# Interpretation of logistic regression output

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

“probability” that class is 1

$$f_{\vec{w}, b}(\vec{x}) = P(y = 1 | \vec{x}; \vec{w}, b)$$

Probability that  $y$  is 1,  
given input  $\vec{x}$ , parameters  $\vec{w}, b$

Example:

$x$  is “tumor size”

$y$  is 0 (not malignant)  
or 1 (malignant)

$$P(y = 0) + P(y = 1) = 1$$

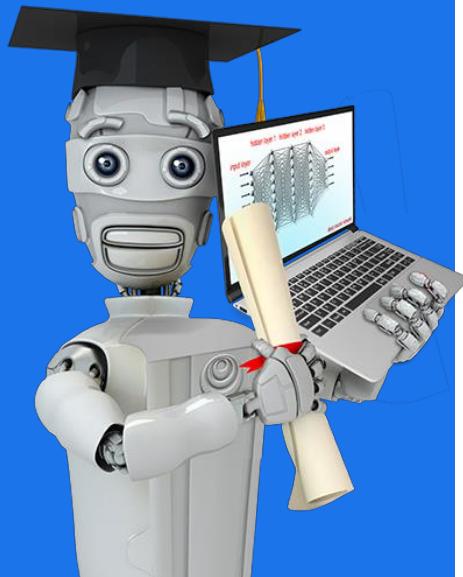
$$f_{\vec{w}, b}(\vec{x}) = 0.7$$

70% chance that  $y$  is 1



Stanford  
ONLINE

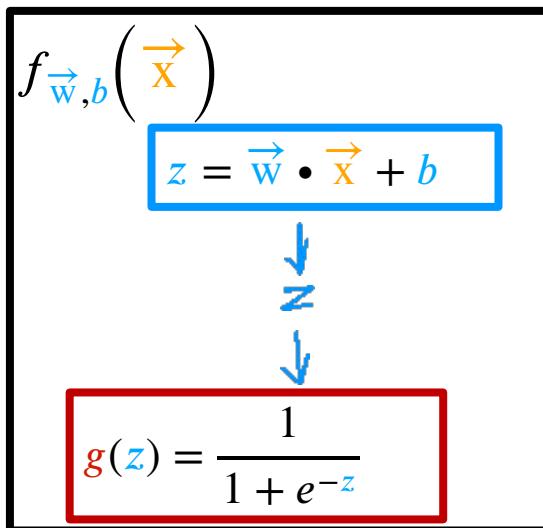
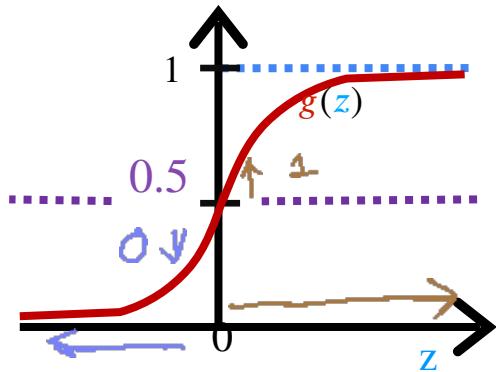
DeepLearning.AI



# Classification

---

## Decision Boundary



$$f_{\vec{w}, b}(\vec{x}) = g(\underbrace{\vec{w} \cdot \vec{x}}_z + \bar{b}) = P(y = 1 | \vec{x}; \vec{w}, b)$$

0 or 1? threshold

Is  $f_{\vec{w}, b}(\vec{x}) \geq \underline{0.5}$ ?

Yes:  $\hat{y} = 1$

No:  $\hat{y} = 0$

When is

$$f_{\vec{w}, b}(\vec{x}) \geq 0 \quad g(z) \geq 0.5$$

$$z \geq 0$$

$$z < 0$$

$$\vec{w} \cdot \vec{x} + b \geq 0$$

$$\vec{w} \cdot \vec{x} + b < 0$$

$$\hat{y} = 1$$

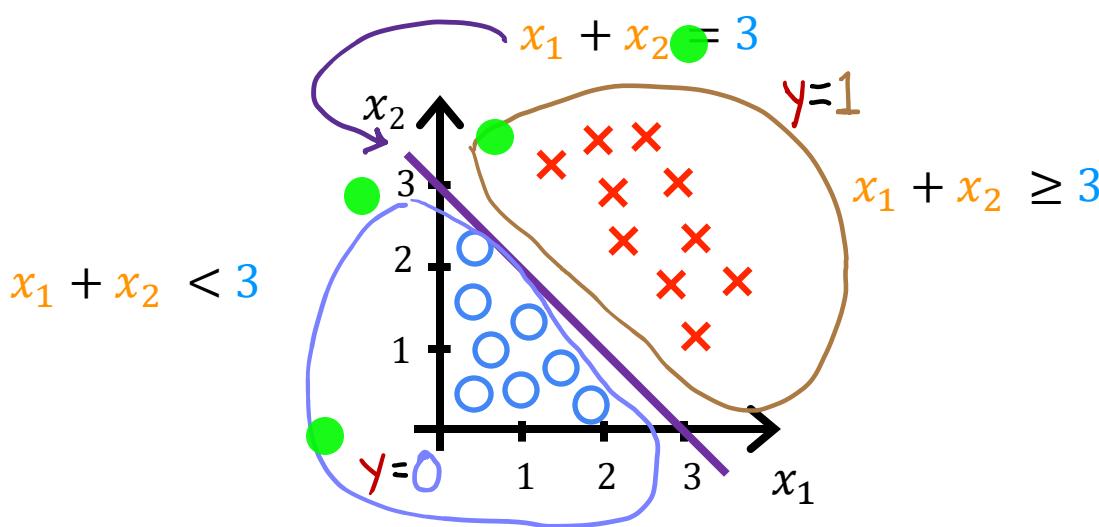
$$\hat{y} = 0$$

# Decision boundary

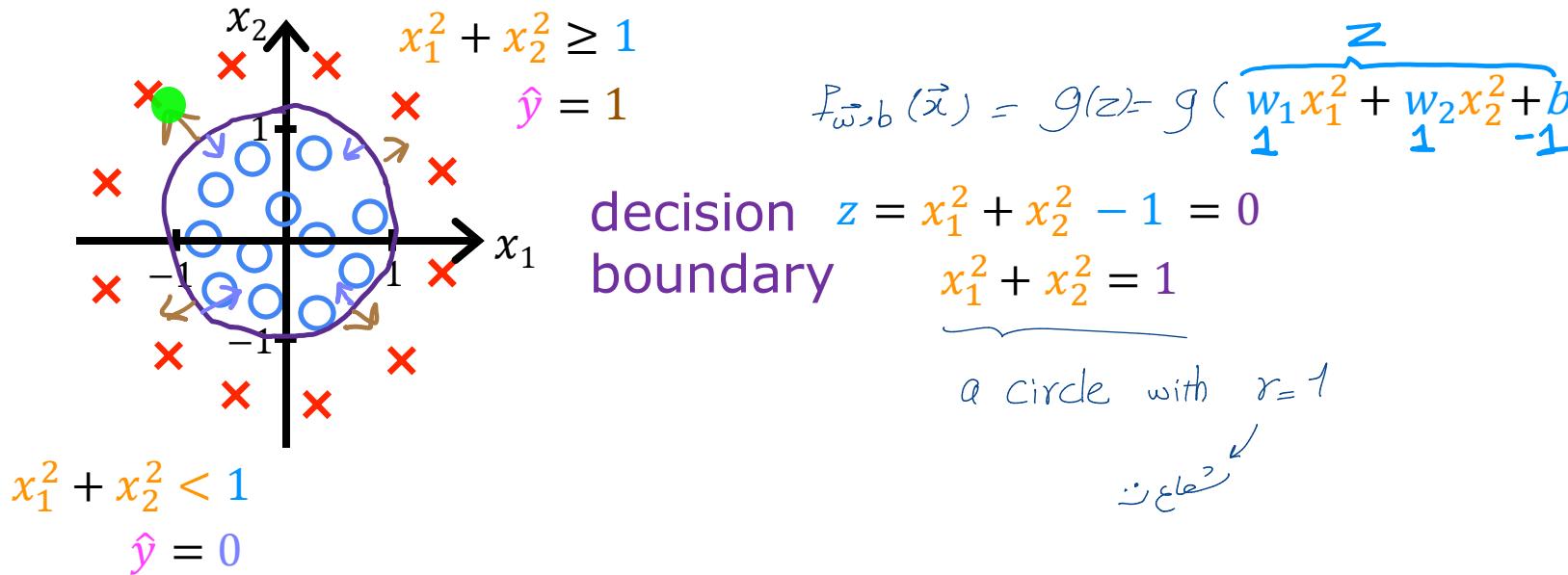
$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(\underbrace{w_1 x_1 + w_2 x_2}_{1} + \underbrace{b}_{-3})$$

Decision boundary  $z = \vec{w} \cdot \vec{x} + b = 0$

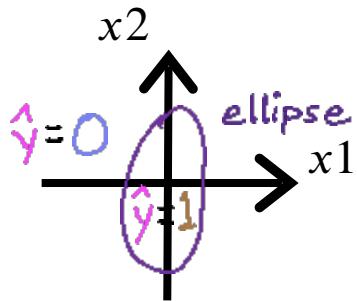
$$z = x_1 + x_2 - 3 = 0$$



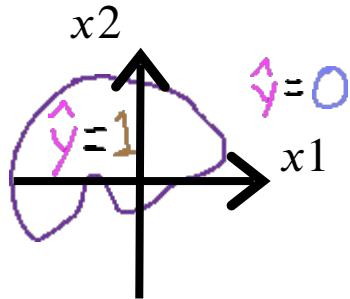
# Non-linear decision boundaries



# Non-linear decision boundaries



$$f_{\vec{w}, b}(\vec{x}) = g(z) = g(w_1 x_1 + w_2 x_2 + w_3 x_1^2 + w_4 x_1 x_2 + w_5 x_2^2 + w_6 x_1^3 + \dots + b)$$



Stanford  
ONLINE

DeepLearning.AI



# Cost Function

---

Cost Function for  
Logistic Regression

# Training set

each row = a patient

	tumor size (cm) $x_1$	...	patient's age $x_n$	malignant? $y$	binary classification $i = 1, \dots, m$ ↪ training examples
$i=1$	10		52	1	
:	2		73	0	
:	5		55	0	
$i=m$	12		49	1	
	...		...	...	

$j = 1, \dots, n$  ↪ features

target  $y$  is 0 or 1

$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$$

How to choose  $\vec{w} = [w_1 \ w_2 \ \cdots \ w_n]$  and  $b$ ?

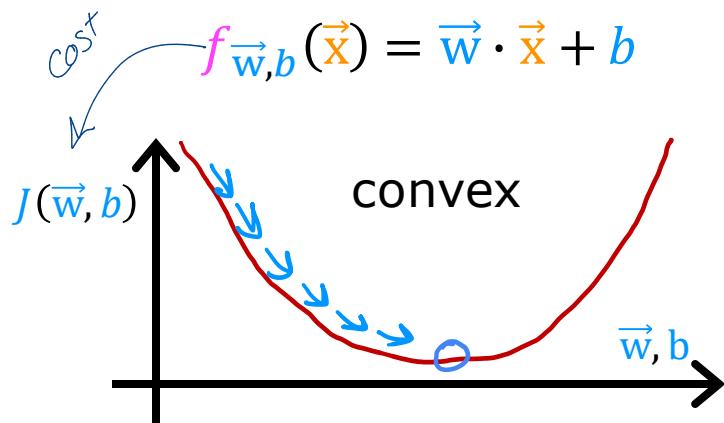
# Squared error cost

$$cost \quad J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2$$

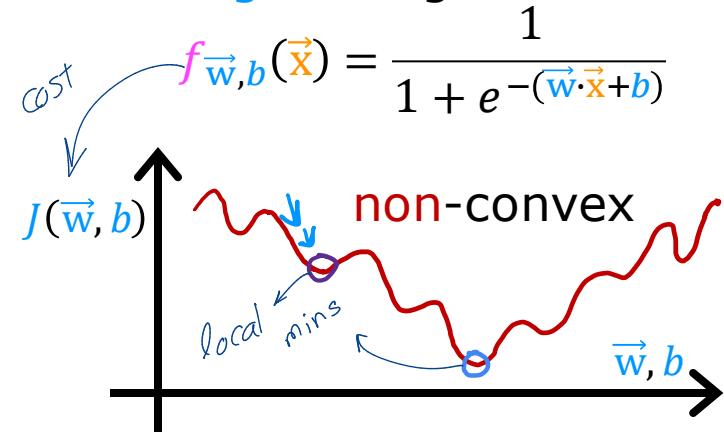
loss  $L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$

average of training set

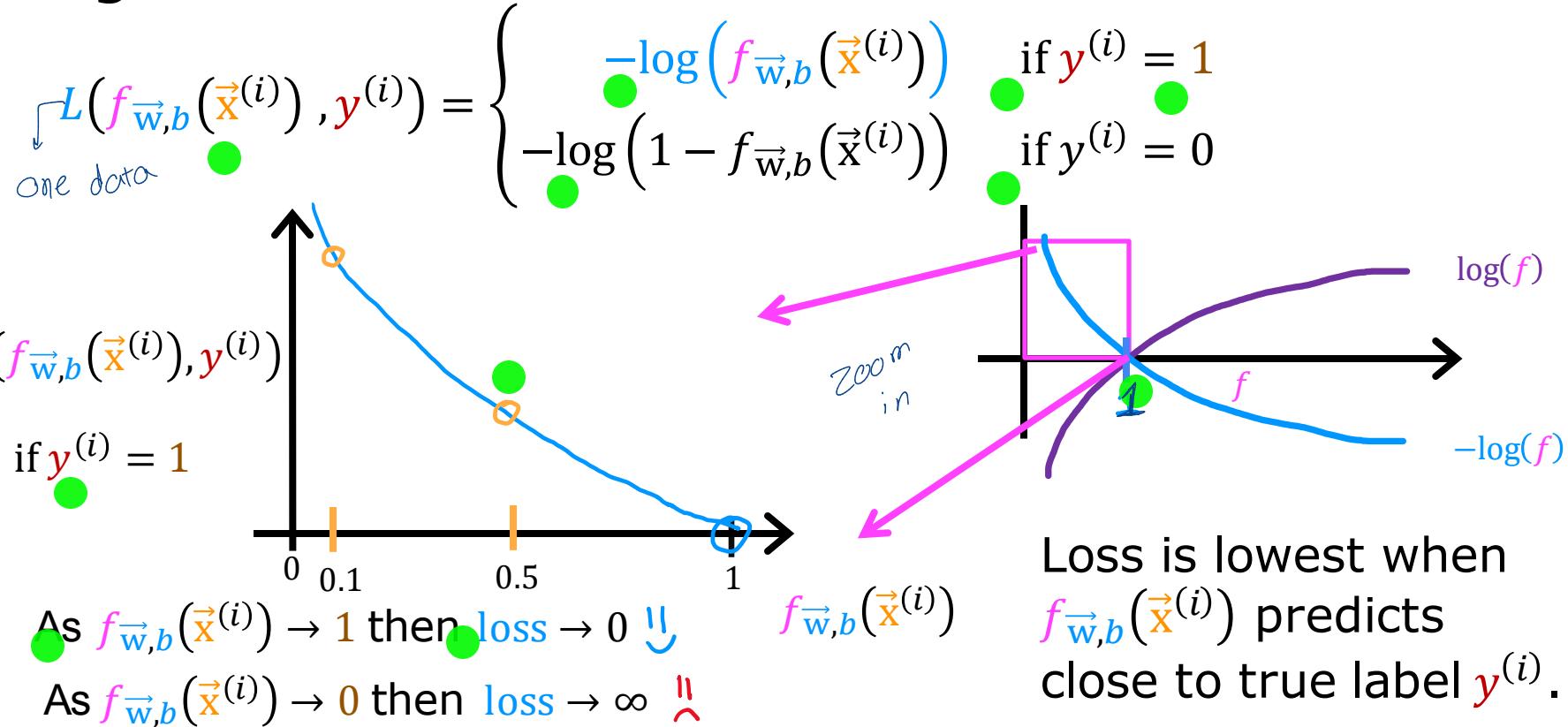
linear regression



logistic regression



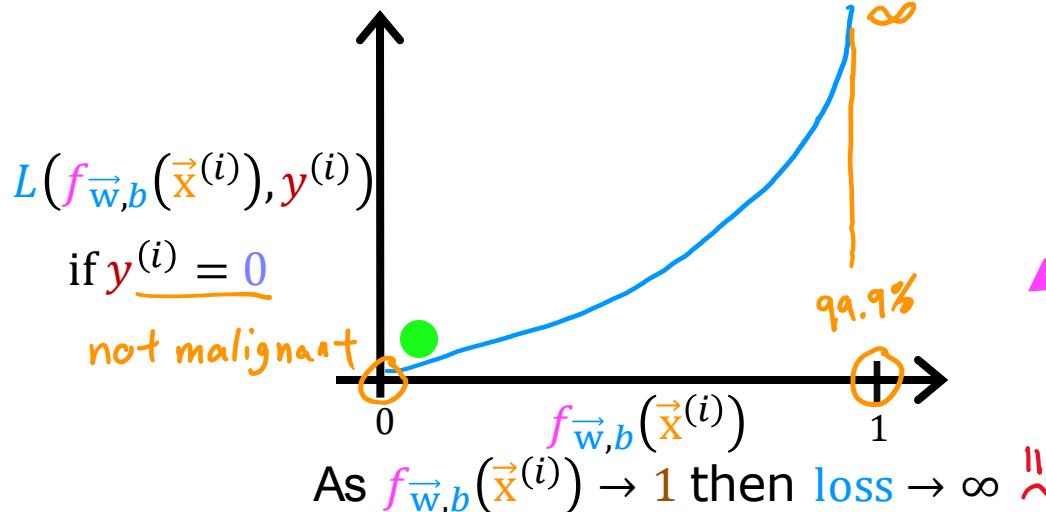
# Logistic loss function



# Logistic loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

As  $f_{\vec{w}, b}(\vec{x}^{(i)}) \rightarrow 0$  then loss  $\rightarrow 0$  



The further prediction  $f_{\vec{w}, b}(\vec{x}^{(i)})$  is from target  $y^{(i)}$ , the higher the loss.

# Cost

$$cost \quad J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})$$

*all training examples*

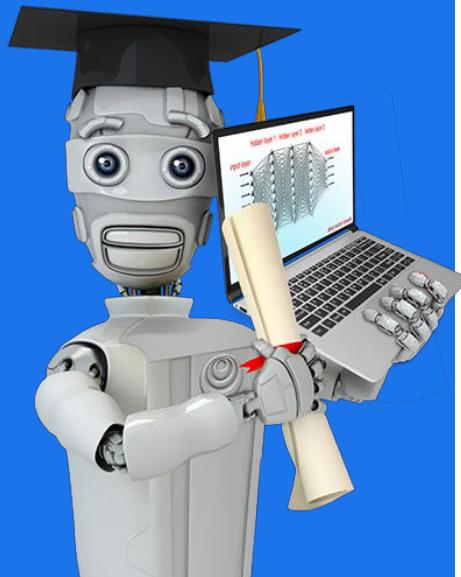
$$= \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

*Convex*  
*can reach a global minimum*

find  $w, b$  that minimize cost  $J$

Stanford  
ONLINE

DeepLearning.AI



# Cost Function

---

Simplified Cost  
Function for Logistic  
Regression

# Simplified loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

if  $y^{(i)} = 1$ :

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \underbrace{-1 \log(f(x))}_{\Theta}$$

# Simplified loss function

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = \begin{cases} -\log(f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 1 \\ -\log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) & \text{if } y^{(i)} = 0 \end{cases}$$

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

if  $y^{(i)} = 1$ :

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -1 \log(f(\vec{x}))$$

if  $y^{(i)} = 0$ :

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = - (1 - 0) \log(1 - f(\vec{x}))$$

# Simplified cost function

Simplified loss Func

$$L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)}) = -y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) - (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))$$

$$\text{cost} J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m [L(f_{\vec{w}, b}(\vec{x}^{(i)}), y^{(i)})]$$

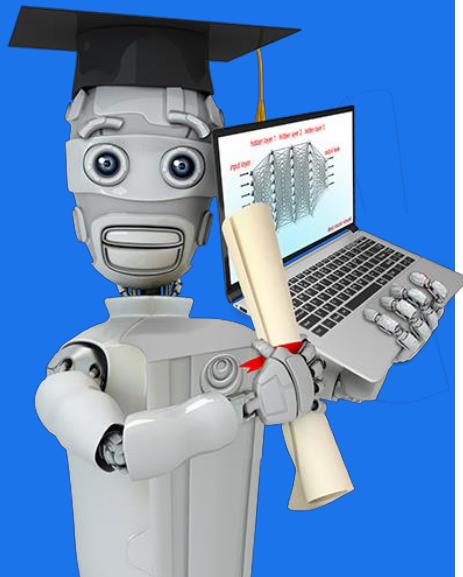
convex  
(single global minimum)

$$= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))]$$

maximum likelihood  
(don't worry about it!)

Stanford  
ONLINE

DeepLearning.AI



# Gradient Descent

---

## Gradient Descent Implementation

# Training logistic regression

Find  $\vec{w}, b$

Given new  $\vec{x}$ , output  $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1+e^{-(\vec{w} \cdot \vec{x} + b)}}$

$$P(y=1|\vec{x}; \vec{w}, b)$$

# Gradient descent

*cost*

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right]$$

repeat {  
   $j = 1 \dots n$   
   $w_j = w_j - \alpha \frac{\partial}{\partial w_j} J(\vec{w}, b)$   
   $b = b - \alpha \frac{\partial}{\partial b} J(\vec{w}, b)$   
}  
} simultaneous updates

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

# Gradient descent for logistic regression

repeat {

looks like linear regression!

(different in  $\vec{w}, b$  definition)

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

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

} simultaneous updates

Same concepts:

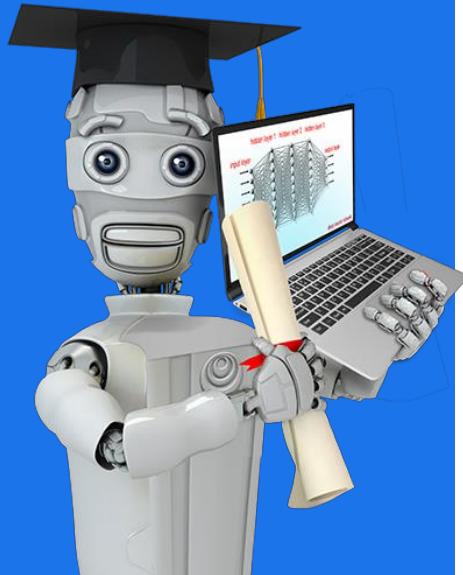
- Monitor gradient descent (learning curve)
- Vectorized implementation
- Feature scaling

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

Logistic regression  $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{(-\vec{w} \cdot \vec{x} + b)}}$

Stanford  
ONLINE

DeepLearning.AI

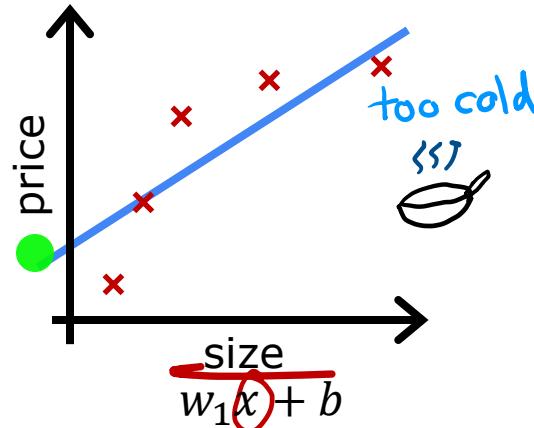


# Regularization to Reduce Overfitting

---

## The Problem of Overfitting

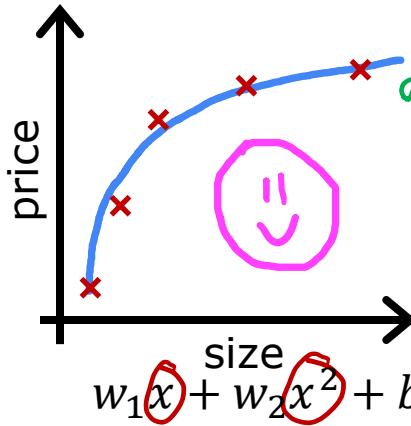
# Regression example



underfit

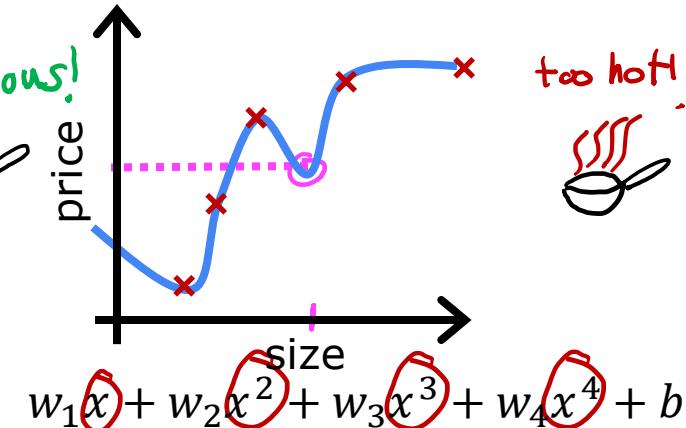
- Does not fit the training set well

high bias



just right

- Fits training set pretty well



overfit fits the data almost too well.

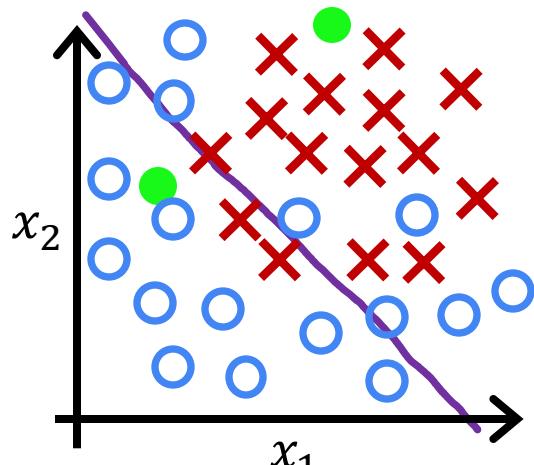
- Fits the training set extremely well

high variance

generalization  
the idea that you want your learning algo to do well, even on examples that are not

- \* **underfitting**: is a scenario in data science where a data model is unable to capture the relationship between the input and output variables accurately.
- \* **overfitting**: is the production of an analysis that corresponds too closely or exactly to a particular set of data , and may therefore fail to fit to additional data or predict future observations reliably .

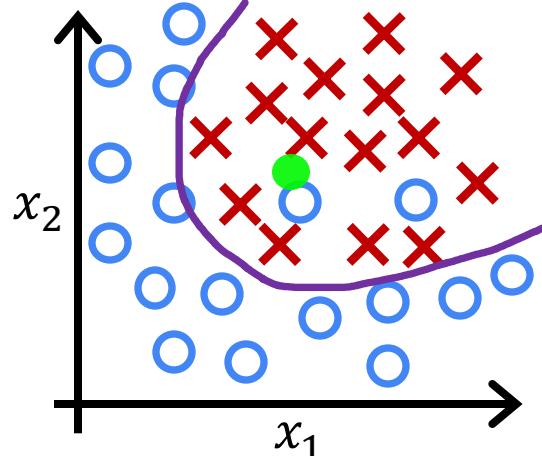
# Classification



$$z = w_1x_1 + w_2x_2 + b$$

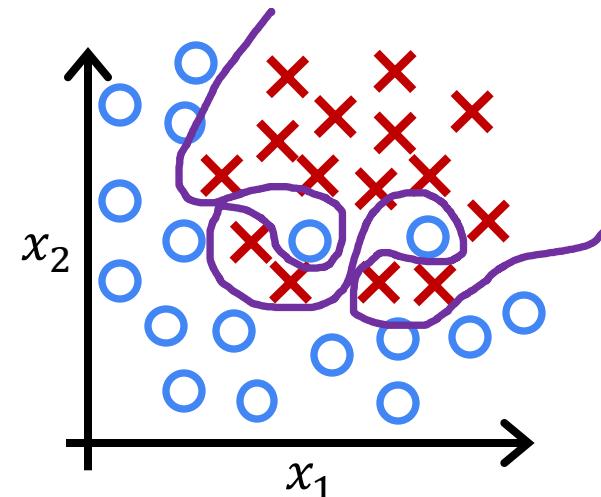
$$f_{\vec{w}, b}(\vec{x}) = g(z)$$

$g$  is the sigmoid function  
underfit      high bias



$$\begin{aligned} z = & w_1x_1 + w_2x_2 \\ & + w_3x_1^2 + w_4x_2^2 \\ & + w_5x_1x_2 + b \end{aligned}$$

just right

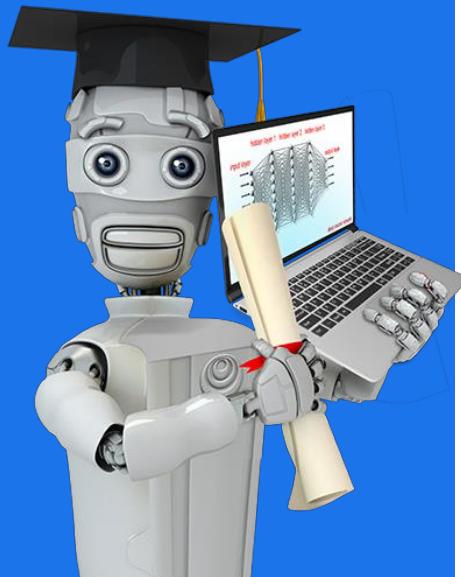


$$\begin{aligned} z = & w_1x_1 + w_2x_2 \\ & + w_3x_1^2 + w_4x_2^2 \\ & + w_5x_1^2x_2^3 + w_6x_1^3x_2 \\ & + \dots + b \end{aligned}$$

Overfit

Stanford  
ONLINE

DeepLearning.AI

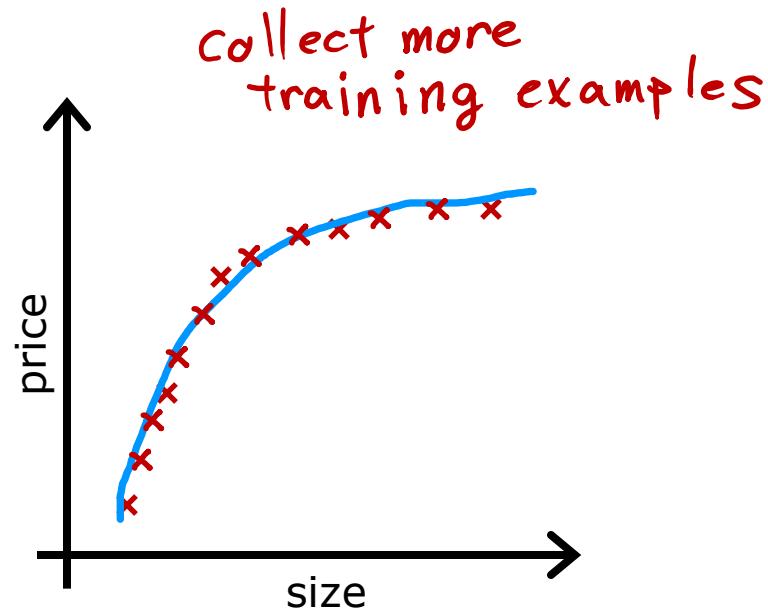
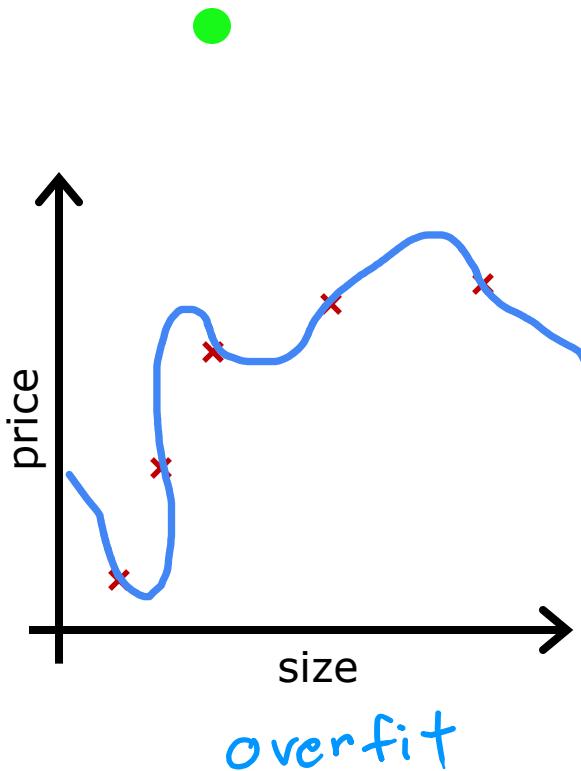


# Regularization to Reduce Overfitting

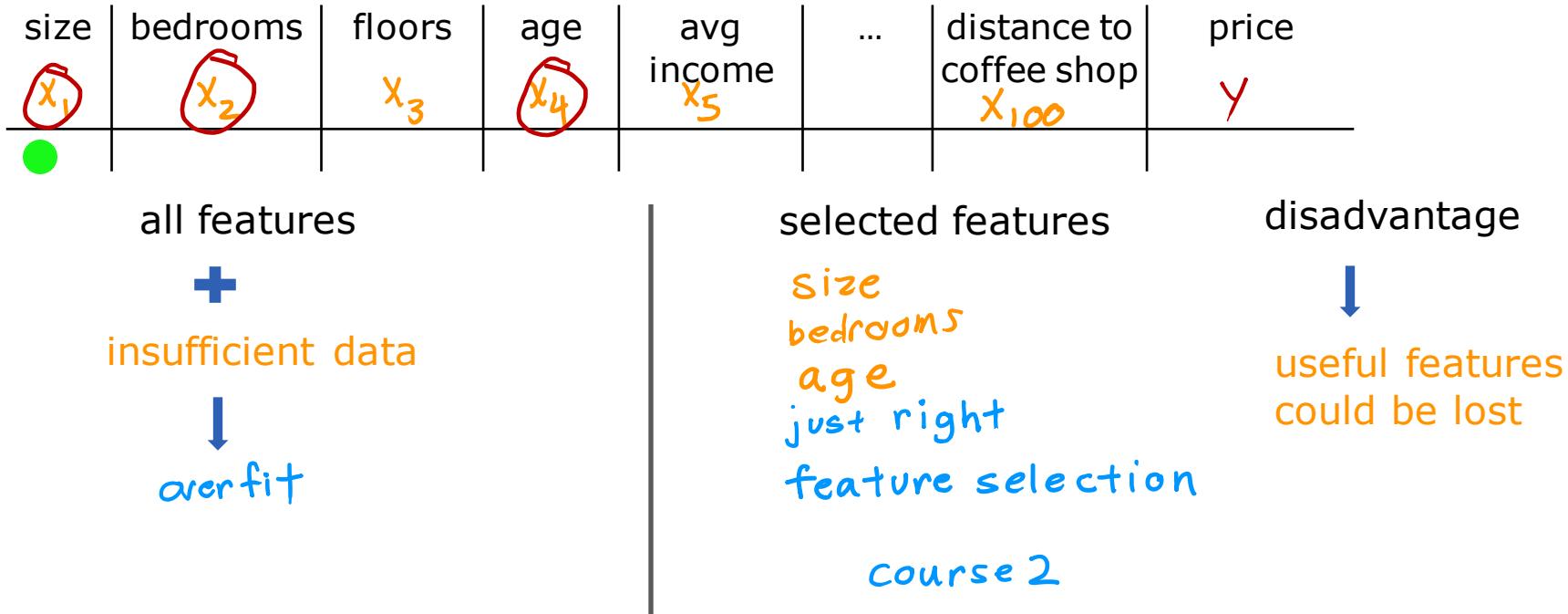
---

## Addressing Overfitting

# Collect more training examples

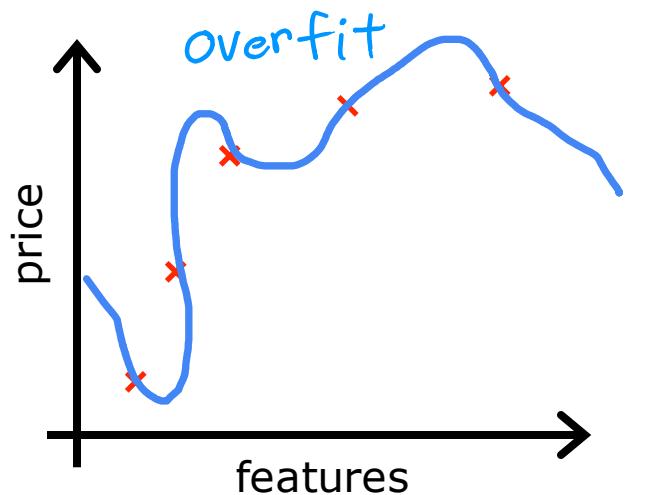


# Select features to include/exclude



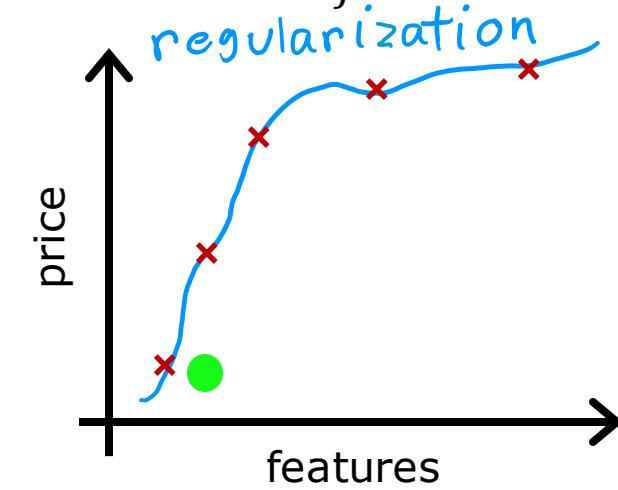
# Regularization

Reduce the size of parameters  $w_j$



$$f(x) = 28x - 385x^2 + 39x^3 - \cancel{174x^4} + 100$$

large values for  $w_j$       eliminate feature



$$f(x) = 13x - 0.23x^2 + 0.000014x^3 - \cancel{0.0001x^4} + 10$$

small values for  $w_j$  (small effect)

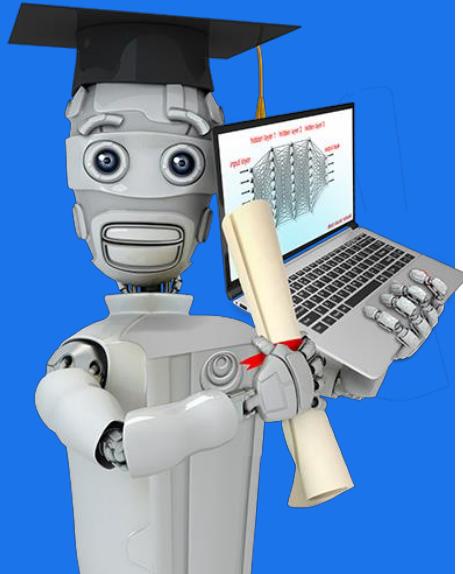
# Addressing overfitting

## Options

1. Collect more data
  2. Select features
    - Feature selection *in course 2*
  3. Reduce size of parameters
    - “Regularization” *next videos!*
- 

Stanford  
ONLINE

DeepLearning.AI

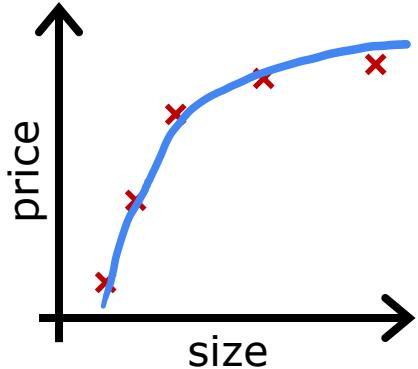


# Regularization to Reduce Overfitting

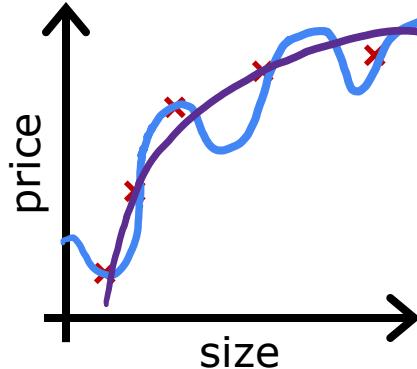
---

## Cost Function with Regularization

# Intuition



$$w_1x + w_2x^2 + b$$



$$w_1x + w_2x^2 + \cancel{w_3x^3} + \cancel{w_4x^4} + b$$

$\approx 0$        $\approx 0$

make  $w_3, w_4$  really small ( $\approx 0$ )

$$\min_{\vec{w}, b} \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + 1000 \underbrace{w_3^2}_{0.001} + 1000 \underbrace{w_4^2}_{0.002}$$

# Regularization

small values  $w_1, w_2, \dots, w_n, b$

simpler model

$$w_3 \approx 0$$

less likely to overfit

$$w_4 \approx 0$$

size	bedrooms	floors	age	avg income	...	distance to coffee shop	price
$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	...	$x_{100}$	$y$

$n$  features       $n = 100$

$w_1, w_2, \dots, w_{100}, b$

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m \left( f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)} \right)^2 + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{"lambda"} \quad \text{regularization parameter}} + \underbrace{\frac{\lambda}{2m} b^2}_{\lambda > 0}$$

can include or exclude  $b$

# Regularization

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[ \underbrace{\frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2}_{\text{mean squared error}} + \underbrace{\frac{\lambda}{2m} \sum_{j=1}^n w_j^2}_{\text{regularization term}} \right]$$

fit data

Keep  $w_j$  small

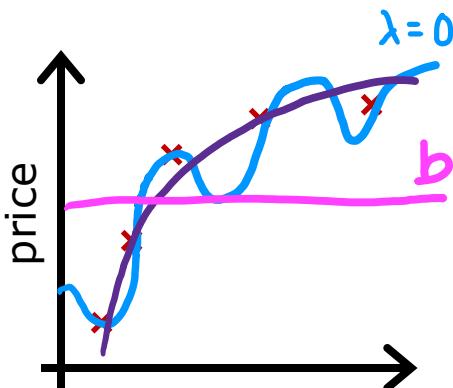
$\lambda$  balances both goals

choose  $\lambda = 10^{10}$

$$f_{\vec{w}, b}(\vec{x}) = \underbrace{w_1 x}_\approx + \underbrace{w_2 x^2}_\approx + \underbrace{w_3 x^3}_\approx + \underbrace{w_4 x^4}_\approx + b$$

$$f(x) = b$$

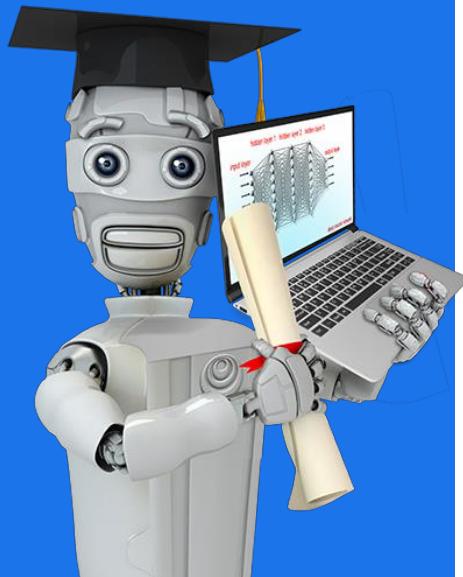
Choose  $\lambda$



- recap
- if  $\lambda = 0$   
↓  
over fit
  - ~~~~~
  - if  $\lambda = \infty$  large num  
↓  
under fit

Stanford  
ONLINE

DeepLearning.AI



# Regularization to Reduce Overfitting

---

## Regularized Linear Regression

# Regularized linear regression

$$\min_{\vec{w}, b} J(\vec{w}, b) = \min_{\vec{w}, b} \left[ \frac{1}{2m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right]$$

Gradient descent

repeat {

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

$j = 1, \dots, n$

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

} simultaneous update

$$\begin{aligned} &= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j \\ &= \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \end{aligned}$$

don't have to regularize  $b$

# Implementing gradient descent

repeat {

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

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

} simultaneous update  $j = 1, \dots, n$



# Implementing gradient descent

repeat {

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

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

} simultaneous update  $j = 1, \dots, n$

$$w_j = \underbrace{w_j - \alpha \frac{\lambda}{m} w_j}_{w_j \left( 1 - \alpha \frac{\lambda}{m} \right)} - \underbrace{\alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)}}_{\text{usual update}}$$

shrink  $w_j$

$$\alpha \frac{\lambda}{m}$$
$$0.01 \frac{1}{50} = 0.0002$$
$$w_j \left( 1 - 0.0002 \right)$$
$$0.9998$$

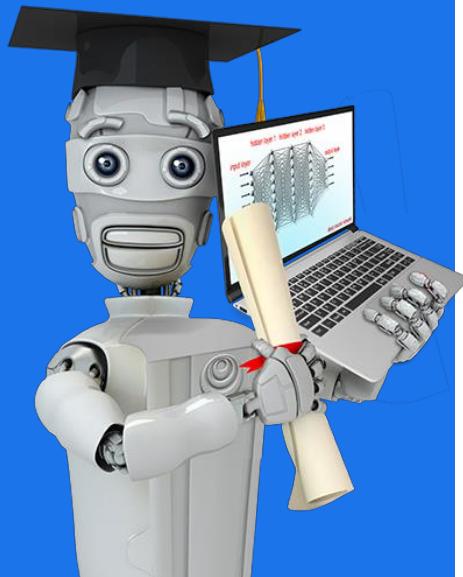


# How we get the derivative term (optional)

$$\begin{aligned} \frac{\partial}{\partial w_j} J(\vec{w}, b) &= \frac{\partial}{\partial w_j} \left[ \frac{1}{2m} \sum_{i=1}^m \underbrace{(\vec{w} \cdot \vec{x}^{(i)}) + b - y^{(i)}}_{f(\vec{x})}^2 + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2 \right] \\ &= \frac{1}{2m} \sum_{i=1}^m \left[ (\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)}) \cancel{\times} x_j^{(i)} \right] + \frac{\lambda}{2m} \cancel{\times} w_j \quad \text{No } \sum_{j=1}^n \\ &= \frac{1}{m} \sum_{i=1}^m \left[ \underbrace{(\vec{w} \cdot \vec{x}^{(i)} + b - y^{(i)})}_{f(\vec{x})} x_j^{(i)} \right] + \frac{\lambda}{m} w_j \\ &= \frac{1}{m} \sum_{i=1}^m \left[ (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j \end{aligned}$$

Stanford  
ONLINE

DeepLearning.AI

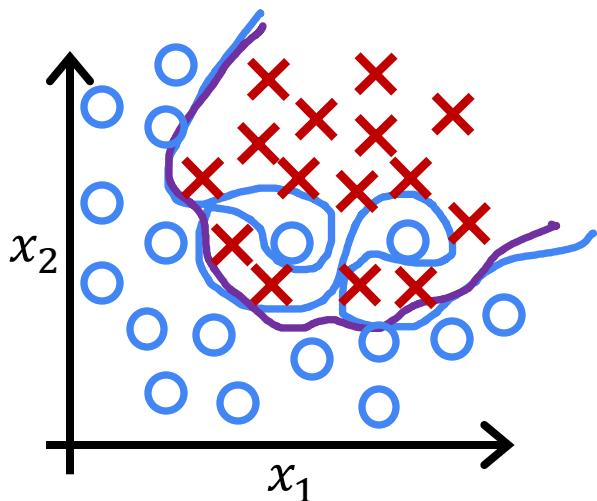


# Regularization to Reduce Overfitting

---

## Regularized Logistic Regression

# Regularized logistic regression



$$z = w_1x_1 + w_2x_2 + w_3x_1^2x_2 + w_4x_1^2x_2^2 + w_5x_1^2x_2^3 + \dots + b$$
$$f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-z}}$$

## Cost function

$$J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)}))] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

$\min_{\vec{w}, b} J(\vec{w}, b) \rightarrow w_j \downarrow$

# Regularized logistic regression

$$\min_{\vec{w}, b} J(\vec{w}, b) = -\frac{1}{m} \sum_{i=1}^m \left[ y^{(i)} \log(f_{\vec{w}, b}(\vec{x}^{(i)})) + (1 - y^{(i)}) \log(1 - f_{\vec{w}, b}(\vec{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

Gradient descent

repeat {

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

j = 1...n

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

}

Looks same as  
for linear regression!

$$= \frac{1}{m} \sum_{i=1}^m \left[ (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right] + \frac{\lambda}{m} w_j$$

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

logistic regression

don't have to  
regularize b