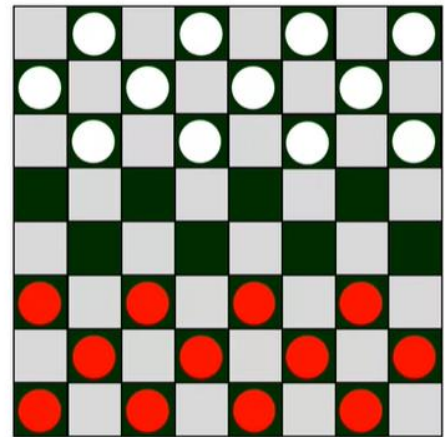


1. What is machine learning

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

2. Supervised learning part 1

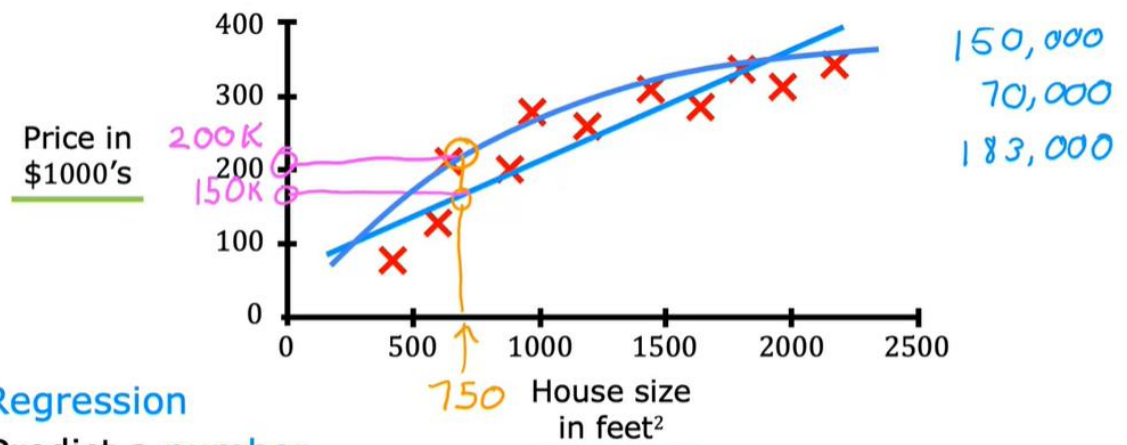
Supervised learning



Learns from being given “right answers”

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



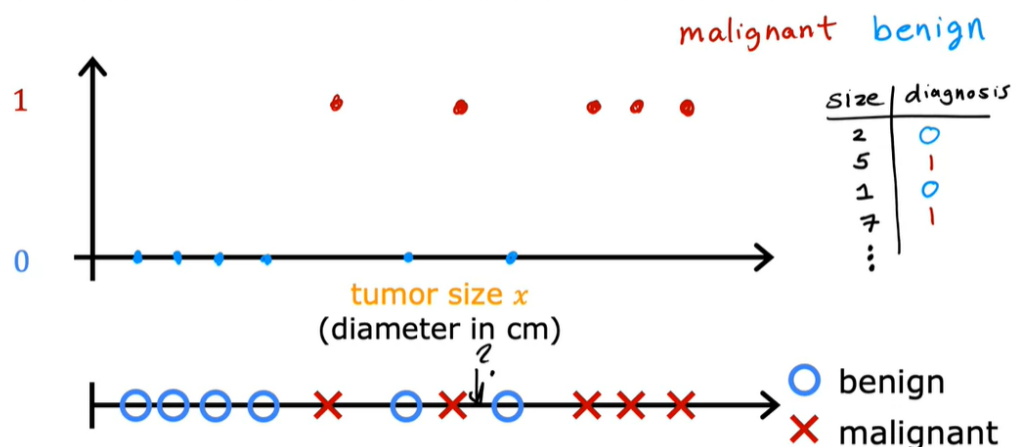
Regression

Predict a **number**

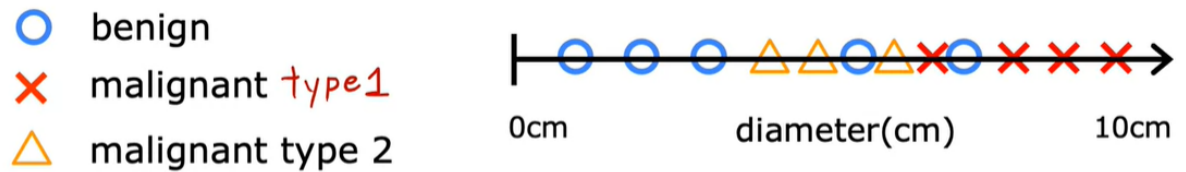
infinitely many possible outputs

3. Supervised learning part 2

Classification: Breast cancer detection



Classification: Breast cancer detection



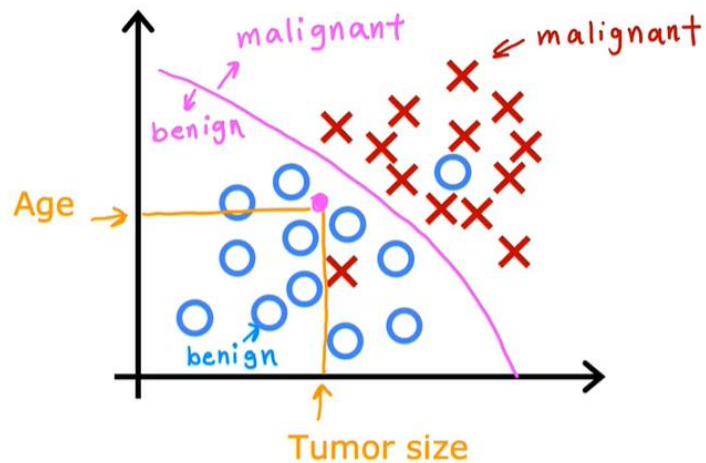
class category

Classification

predict categories cat dog benign malignant 0, 1, 2

small number of possible outputs

Two or more inputs



Supervised learning is when we give our learning algorithm the right answer y for each example to learn from. Which is an example of supervised learning?

☐ Calculating the average age of a group of customers.

☒ Spam filtering.

✓ Correct

For instance, emails labeled as "spam" or "not spam" are examples used for training a supervised learning algorithm. The trained algorithm will then be able to predict with some degree of accuracy whether an unseen email is spam or not.

Supervised learning

Learns from being given "right answers"

Regression

Predict a number

infinitely many possible outputs

Classification

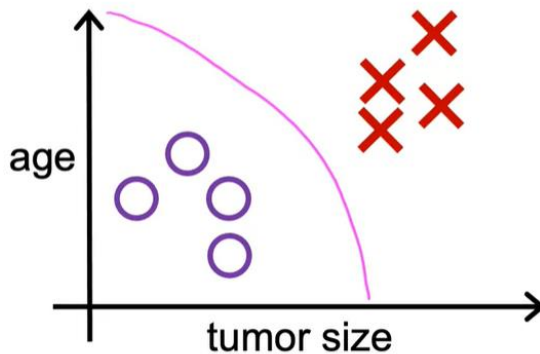
predict categories

small number of possible outputs

4. Unsupervised learning part 1

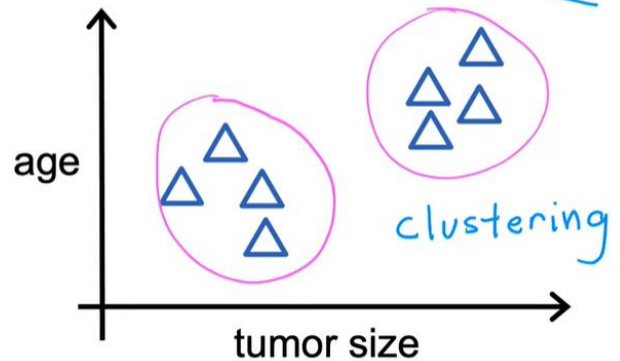
Supervised learning

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

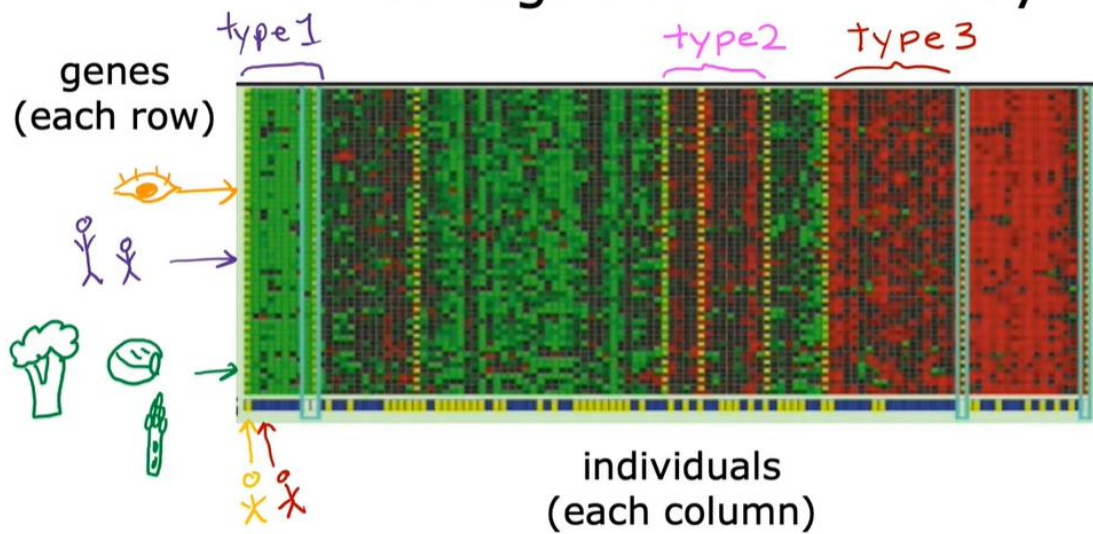
• Twin Panda Cubs Born at Tokyo's Ueno Zoo

PEOPLE · 6 hours ago

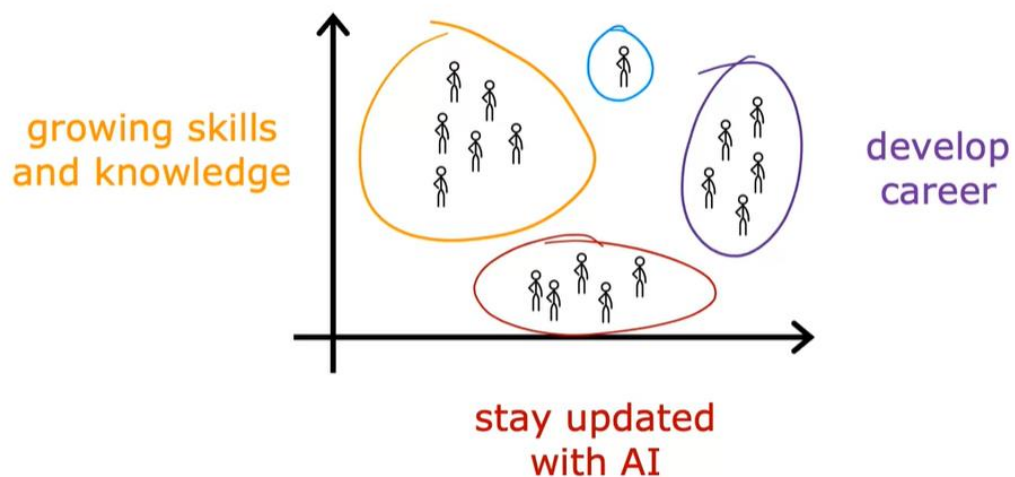
[View Full Coverage](#)



Clustering: DNA microarray



Clustering: Grouping customers



5. 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.

Clustering

Group similar data points together.

Dimensionality reduction

Compress data using fewer numbers.

Anomaly detection

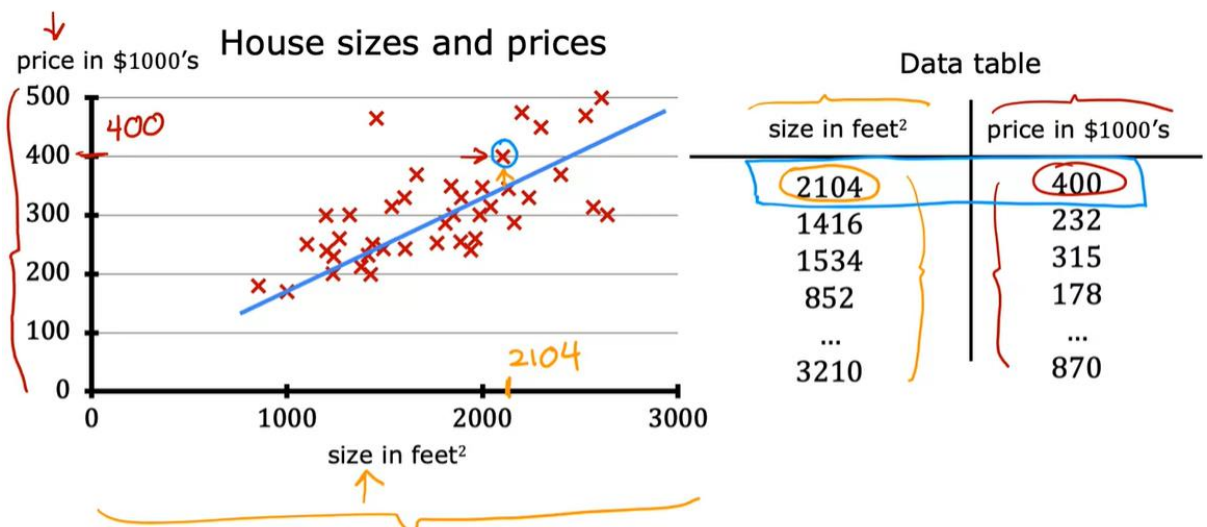
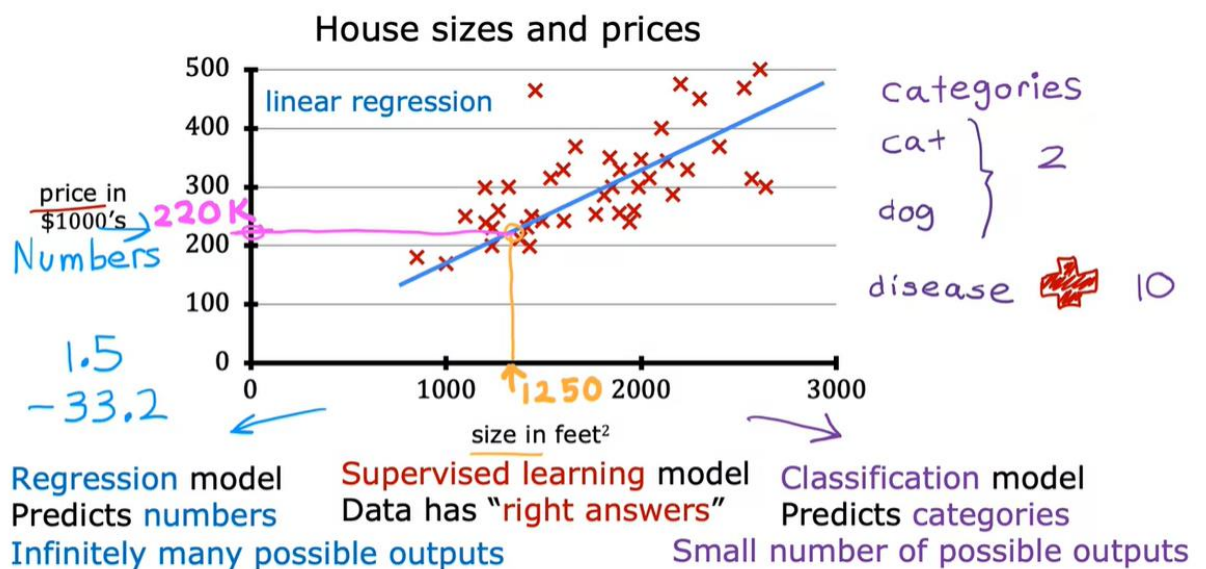
Find unusual data points.

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.
- ☒ ☒ Given a set of news articles found on the web, group them into sets of articles about the same story.
- ☒ ☒ 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

6. Linear regression model part 1



Terminology

Training set: Data used to train the model

	x size in feet ²	y price in \$1000's
(1)	2104	400
(2)	1416	232
(3)	1534	315
(4)	852	178
...
(47)	3210	870

$m = 47$

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

$x^{(2)} = 1416$ $x^{(2)} \neq x^2$ not exponent

Notation:

x = "input" variable
feature

y = "output" variable
"target" variable

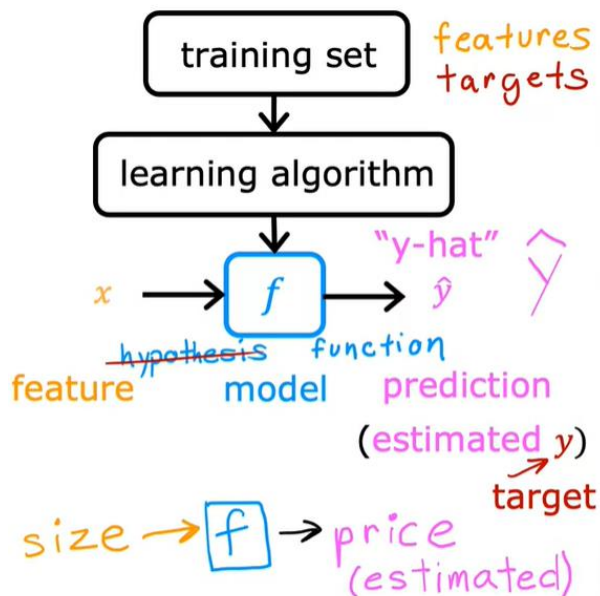
m = number of training examples

(x, y) = single training example

$(x^{(i)}, y^{(i)})$

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

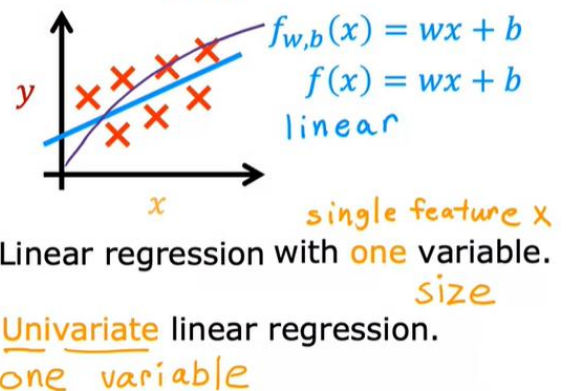
7. Linear regression model part 2



How to represent f ?

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

$f(x)$



8. Cost function formula

Training set

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

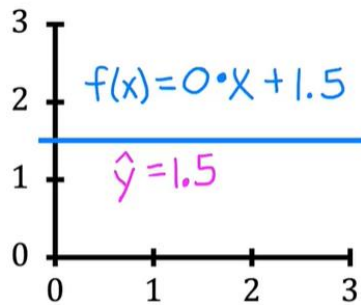
Model: $f_{w,b}(x) = wx + b$

w, b : parameters
coefficients
weights

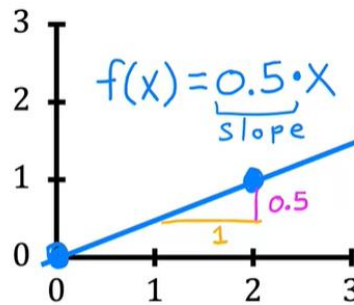
What do w, b do?

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

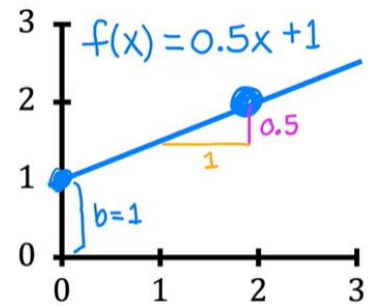
$f(x)$



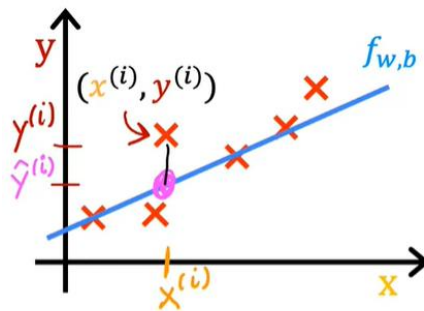
→ $w = 0$
→ $b = 1.5$
 y-intercept



→ $w = 0.5$
→ $b = 0$



→ $w = 0.5$
→ $b = 1$



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

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

Cost function: Squared error cost function

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

error

m = number of training examples

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

Find w, b :

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

9. Cost function: intuition

model:

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

parameters:

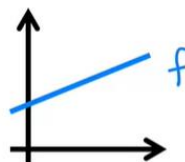
$$w, b$$

cost function:

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

goal:

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

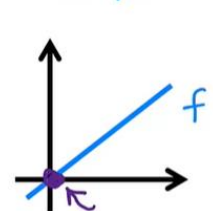


simplified

$$f_w(x) = wx$$

w

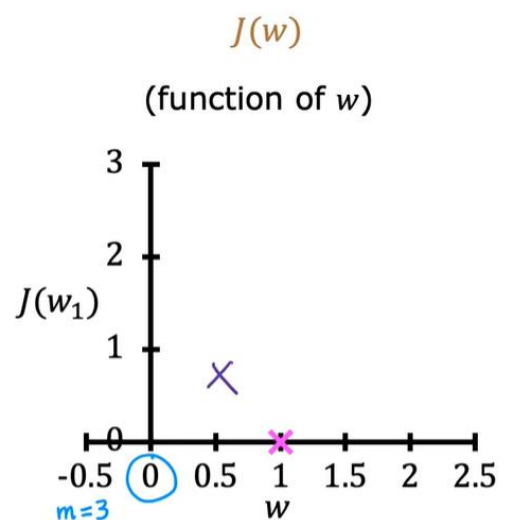
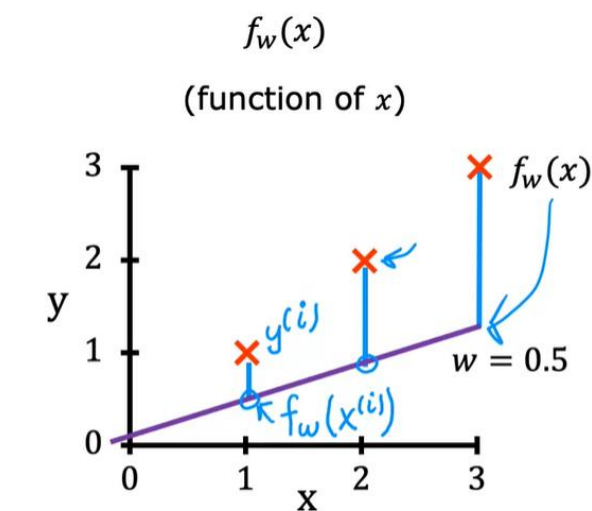
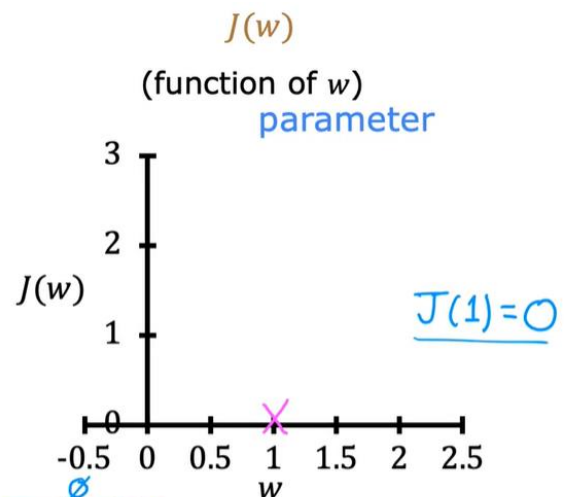
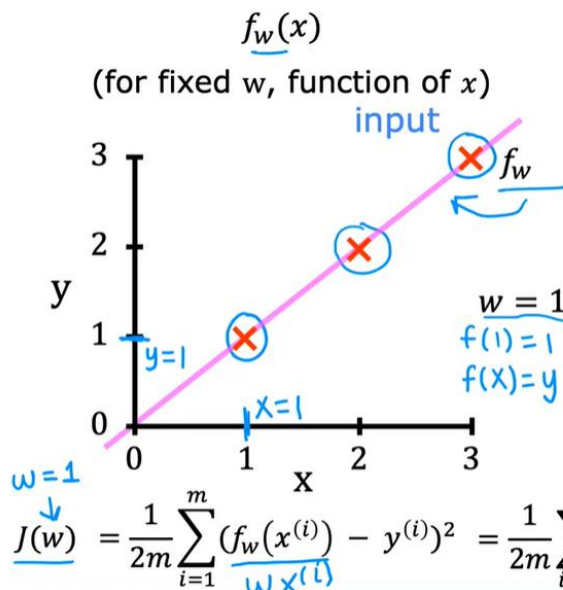
$$b = \emptyset$$



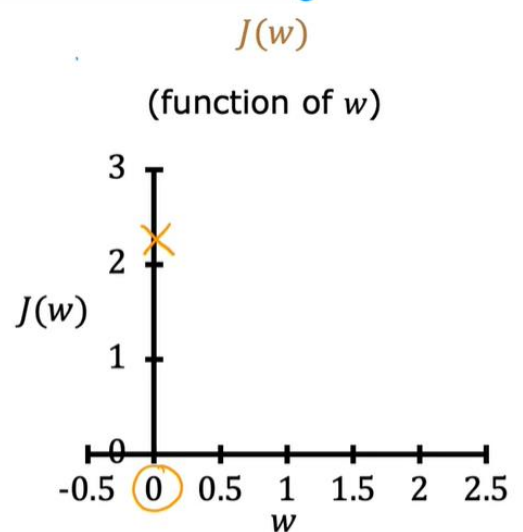
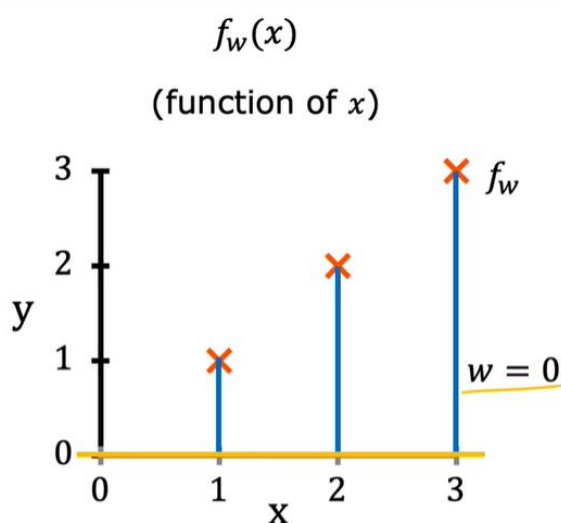
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

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

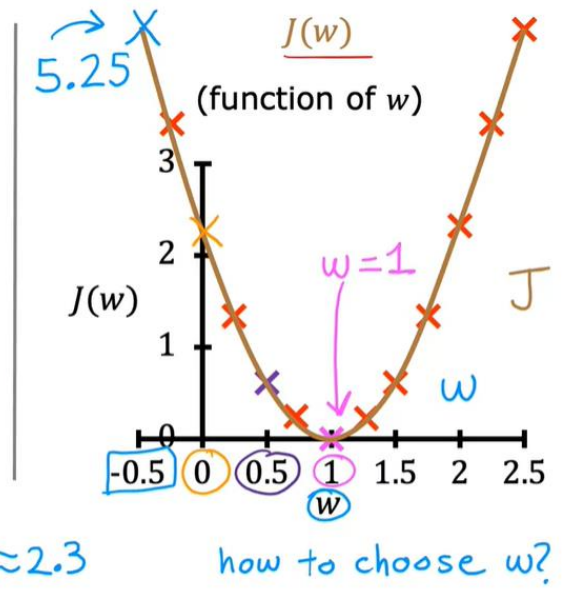
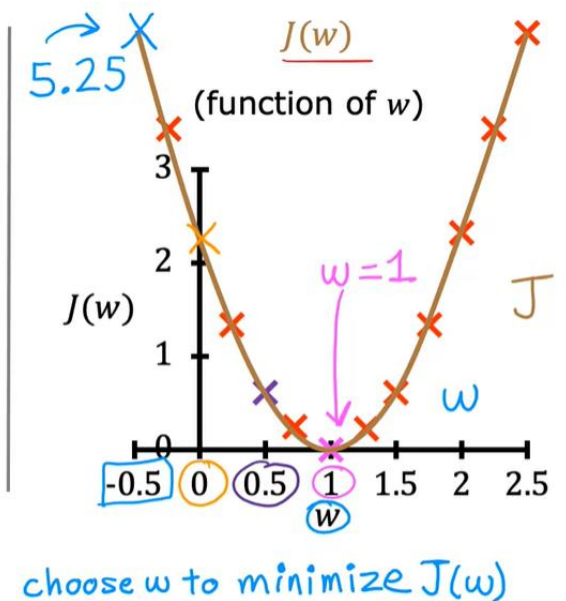
$w x^{(i)}$



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



$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2) = \frac{1}{6} [14] \approx 2.3$$

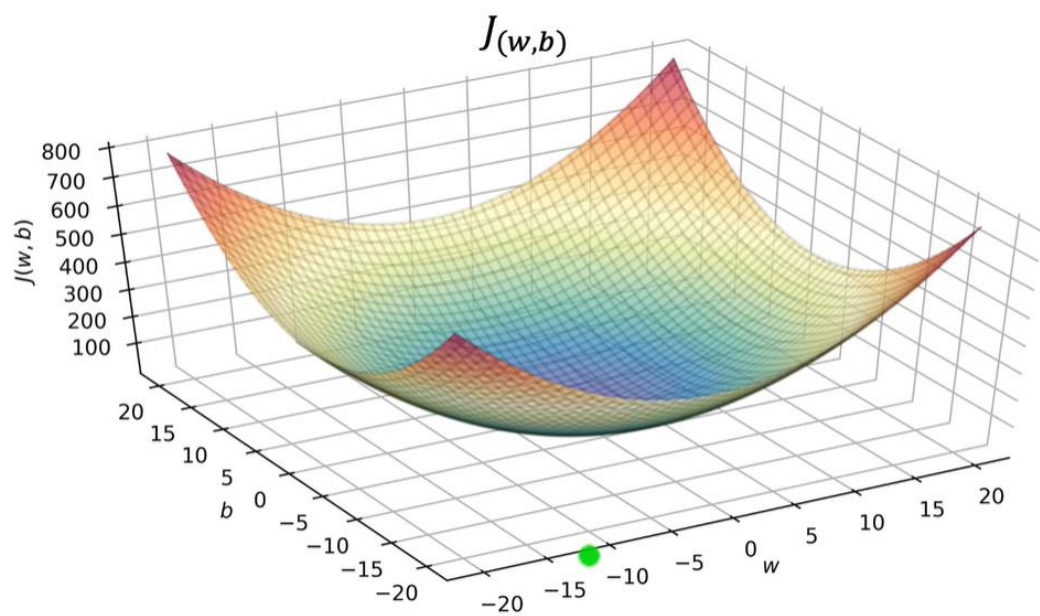
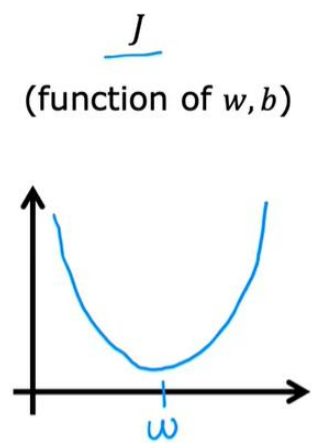
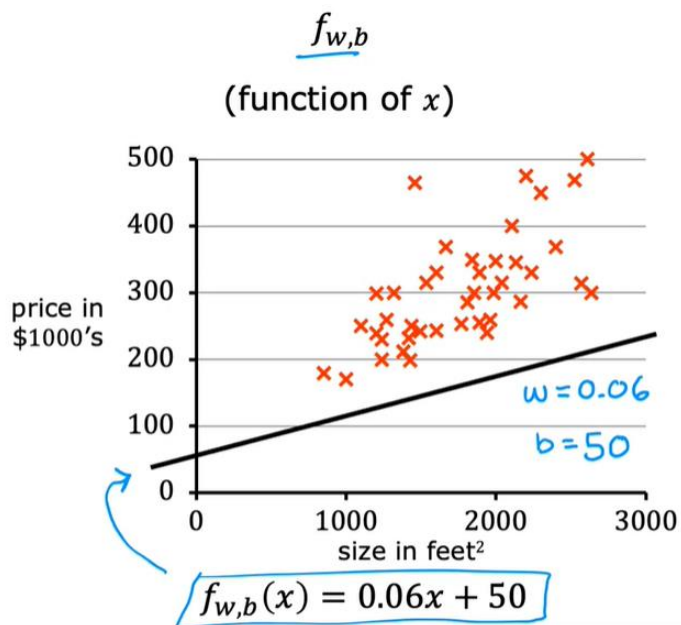

$$\underset{w}{\text{minimize}} J(w)$$
$$\underset{w, b}{\text{minimize}} J(w, b)$$


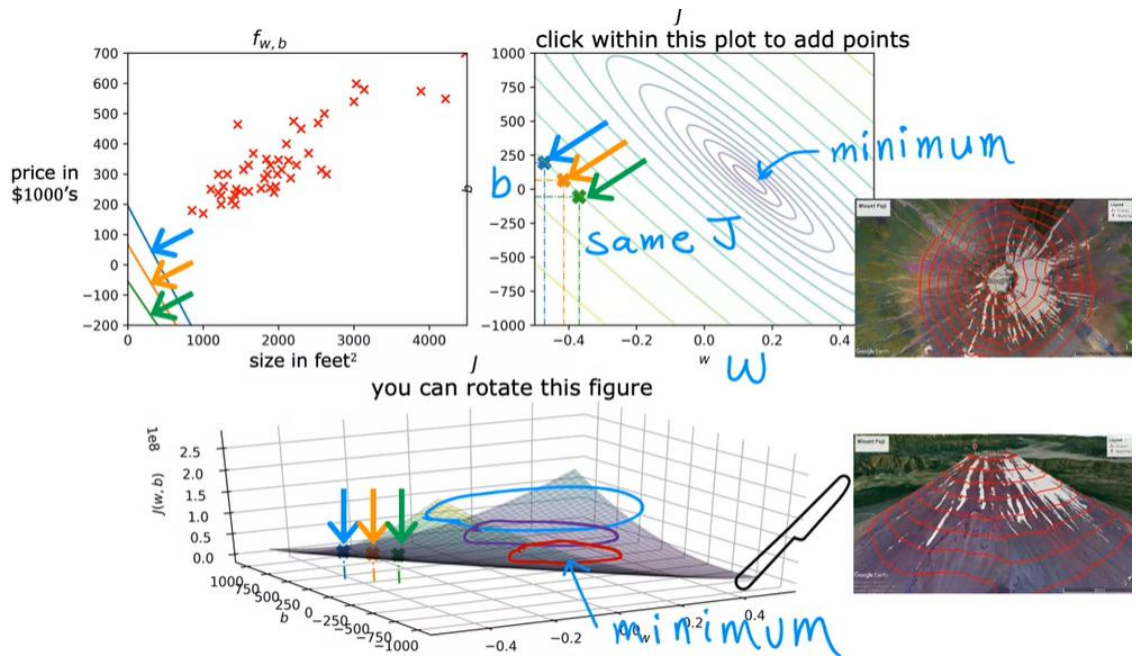
Model

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

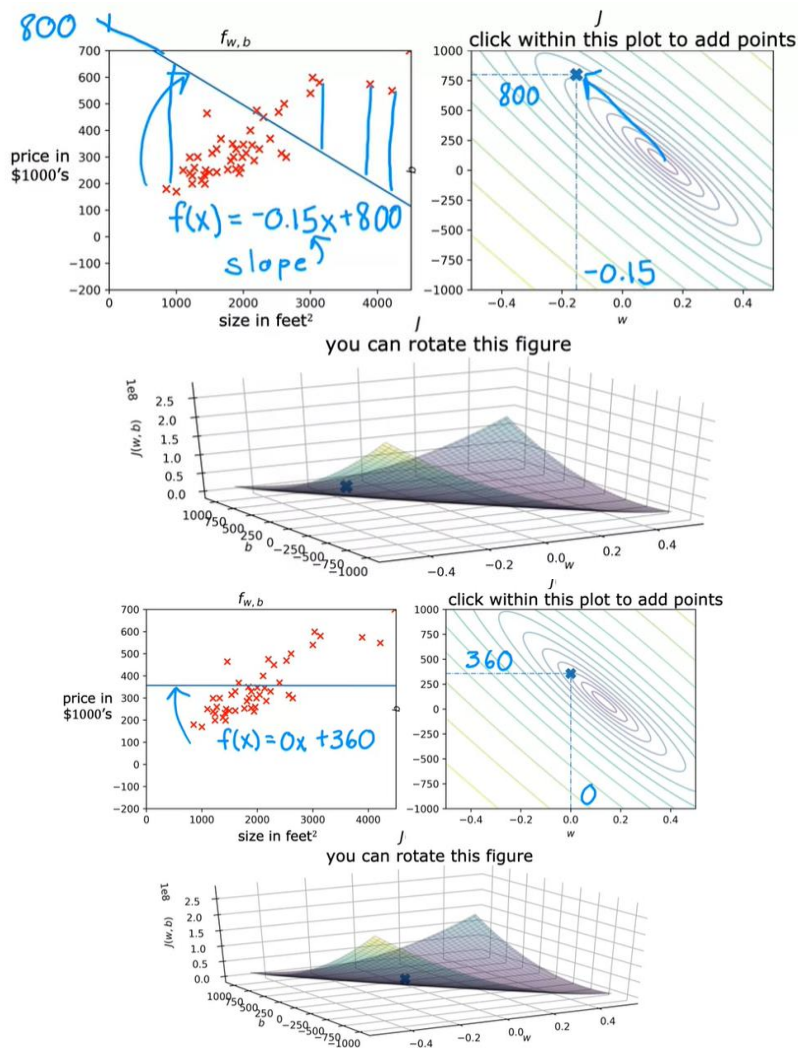
w, b ~~before: $b=0$~~

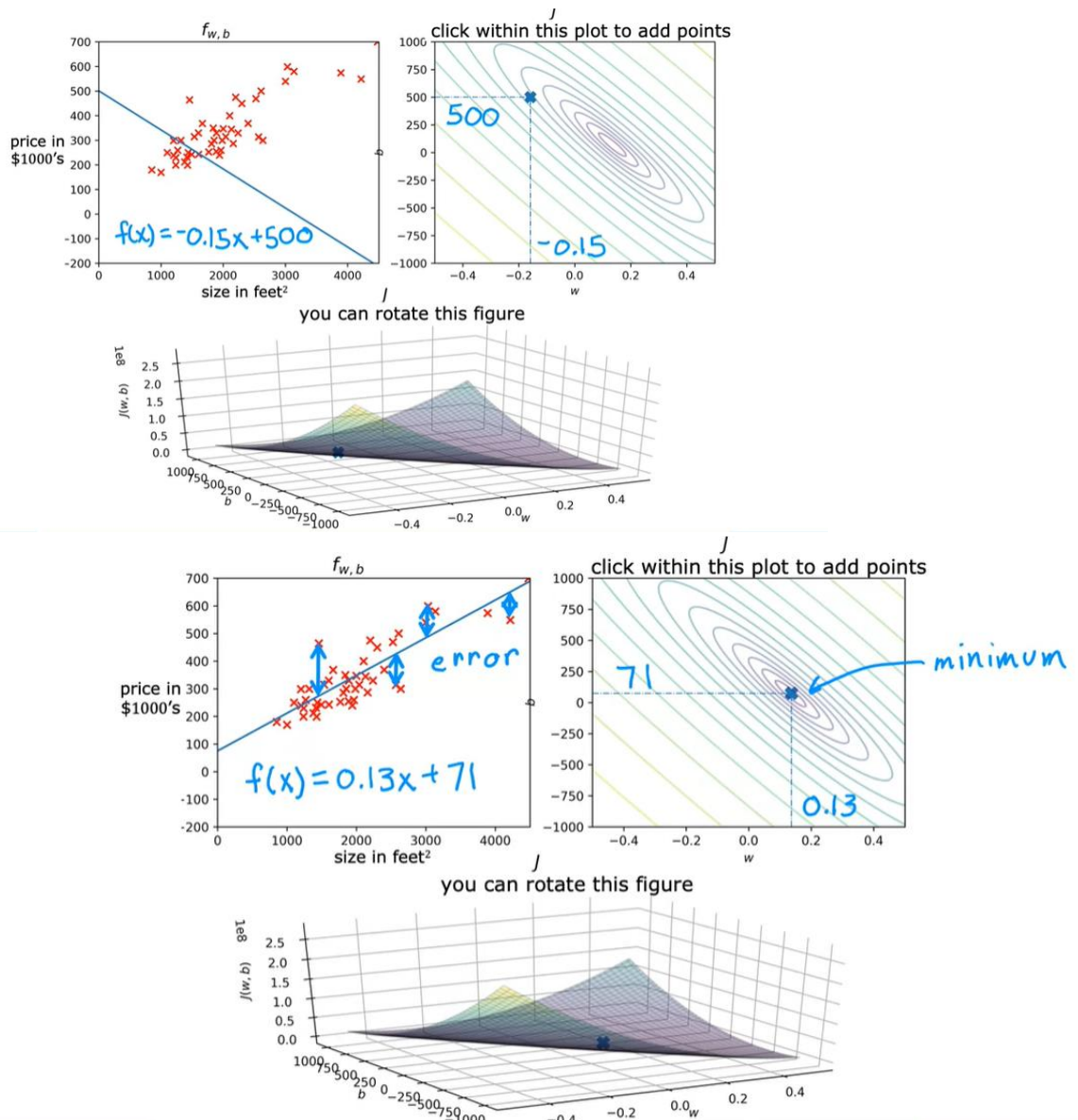
$$J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$$
$$\underset{w,b}{\text{minimize}} J(w,b)$$





11. Visualization examples





12. Practice quiz

1.

For linear regression, the model is $f_{w,b}(x) = wx + b$.

Which of the following are the inputs, or features, that are fed into the model and with which the model is expected to make a prediction?

- ☐ m
- ☐ w and b .
- ☒ x
- ☐ (x, y)

✓ Correct

The x , the input features, are fed into the model to generate a prediction $f_{w,b}(x)$

2. For linear regression, if you find parameters w and b so that $J(w, b)$ is very close to zero, what can you conclude?
- ☐ This is never possible -- there must be a bug in the code.
 - ☐ The selected values of the parameters w and b cause the algorithm to fit the training set really poorly.
 - ☒ The selected values of the parameters w and b cause the algorithm to fit the training set really well.

✓ Correct

When the cost is small, this means that the model fits the training set well.

13. Gradient descent

Have some function $J(w, b)$ for linear regression
or any function

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

Outline:

Start with some w, b (set $w=0, b=0$)

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

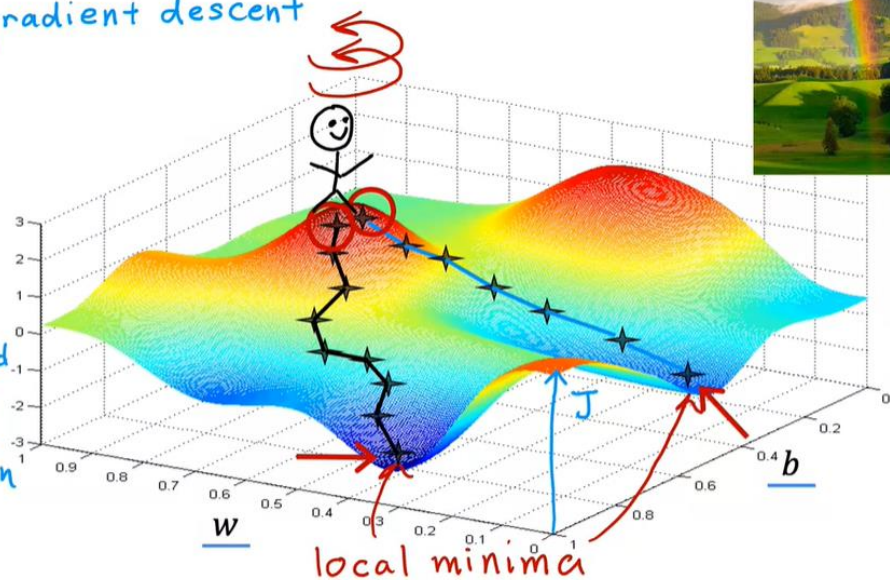
Until we settle at or near a minimum

may have >1 minimum

J not always

gradient descent

$J(w, b)$
not squared
error cost
not linear
regression



14. Implementing gradient descent

Gradient descent algorithm

Repeat until convergence

$$\begin{cases} \underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{cases}$$

Learning rate
Derivative

Simultaneously
update w and b

Assignment

$$a = c$$

$$a = a + 1$$

Code

Truth assertion

$$a = c$$

$$a = a + 1$$

Math

$$a == c$$

Correct: Simultaneous update

$$\begin{aligned} tmp_w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ tmp_b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ w &= tmp_w \\ b &= tmp_b \end{aligned}$$

Incorrect

$$\begin{aligned} tmp_w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w &= tmp_w \\ tmp_b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ b &= tmp_b \end{aligned}$$

Gradient descent is an algorithm for finding values of parameters w and b that minimize the cost function J . What does this update statement do? (Assume α is small.)

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

- ☐ Checks whether w is equal to $w - \alpha \frac{\partial J(w, b)}{\partial w}$
- ☒ Updates parameter w by a small amount

✓ Correct

This updates the parameter by a small amount, in order to reduce the cost J .

15. Gradient descent: intuition

Gradient descent algorithm

repeat until convergence {

$$\begin{cases} \underline{w} = w - \alpha \frac{\partial}{\partial w} J(w, b) \\ \underline{b} = b - \alpha \frac{\partial}{\partial b} J(w, b) \end{cases}$$

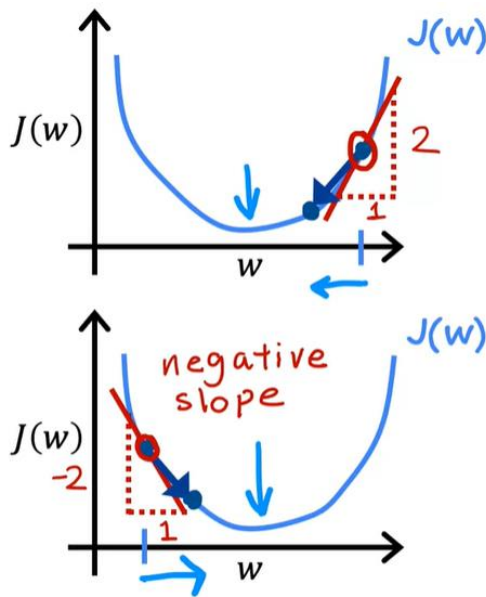
learning rate

derivative

$$J(w)$$

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

$$\underline{\min}_w J(w)$$



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

> 0

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

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

$$w = w - \alpha \cdot (\text{negative number})$$

Assume the learning rate α is a small positive number. When $\frac{\partial J(w,b)}{\partial w}$ is a positive number (greater than zero) -- as in the example in the upper part of the slide shown above -- what happens to w after one update step?

- ☐ w increases
- ☐ It is not possible to tell if w will increase or decrease.
- ☒ w decreases.
- ☐ w stays the same

✓ Correct

The learning rate α is always a positive number, so if you take w minus a positive number, you end up with a new value for w that is smaller

16. Learning rate

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

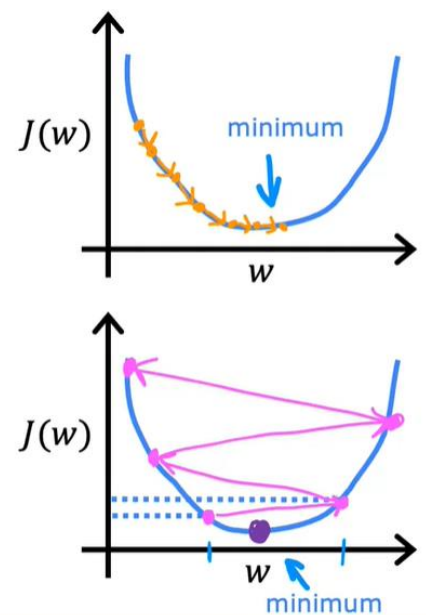
If α is too small...

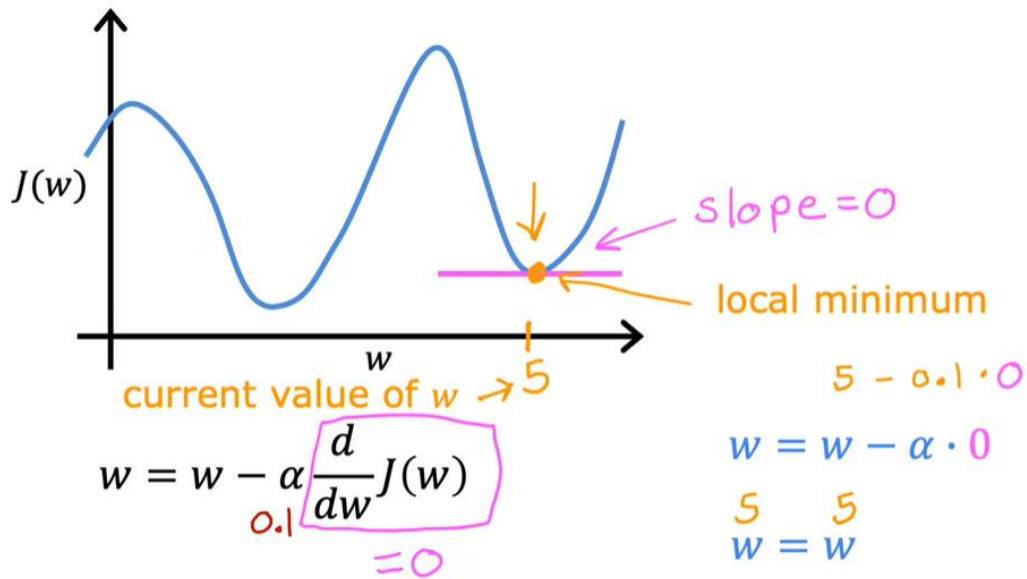
Gradient descent may be slow.

If α is too large...

Gradient descent may:

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





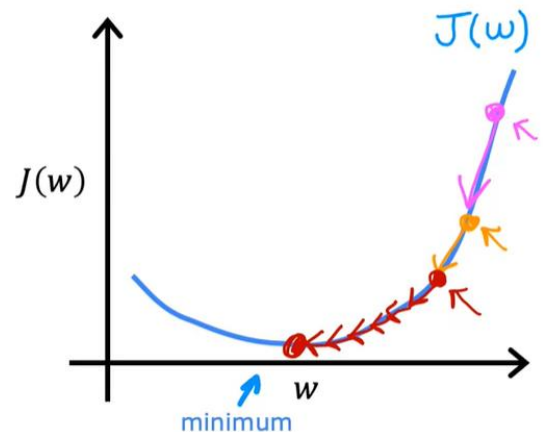
Can reach local minimum with fixed learning rate α

Annotations: α is smaller, not as large, large.

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

Near a local minimum,
 - Derivative becomes smaller
 - Update steps become smaller

Can reach minimum without decreasing learning rate α



17. Gradient descent for linear regression

Linear regression model

Cost function

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

Gradient descent algorithm

repeat until convergence {

$$\begin{aligned} w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)} \\ b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \rightarrow \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \end{aligned}$$

}

(Optional)

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

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

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

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

$\text{no } x^{(i)}$

Gradient descent algorithm

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

repeat until convergence {

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

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

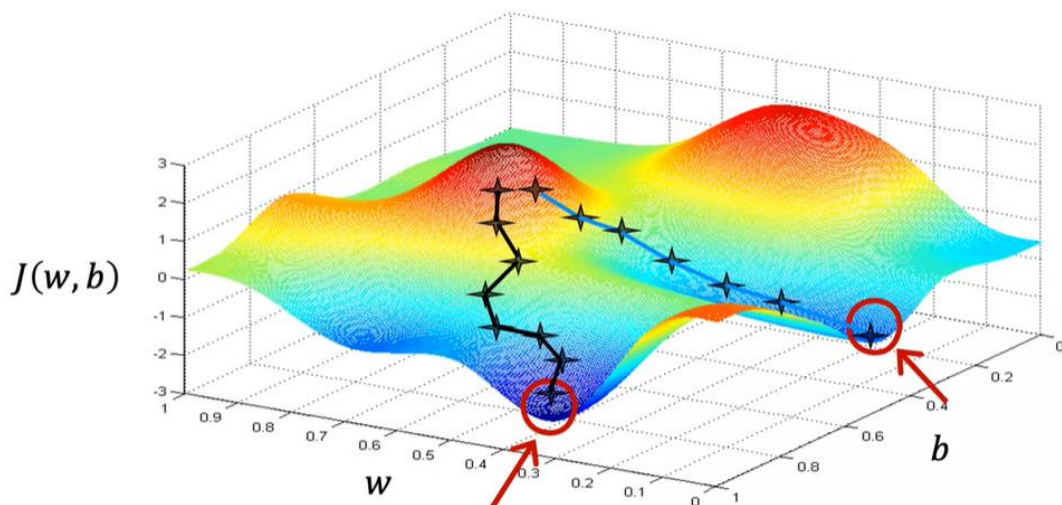
$\frac{\partial}{\partial b} J(w, b)$

Update w and b simultaneously

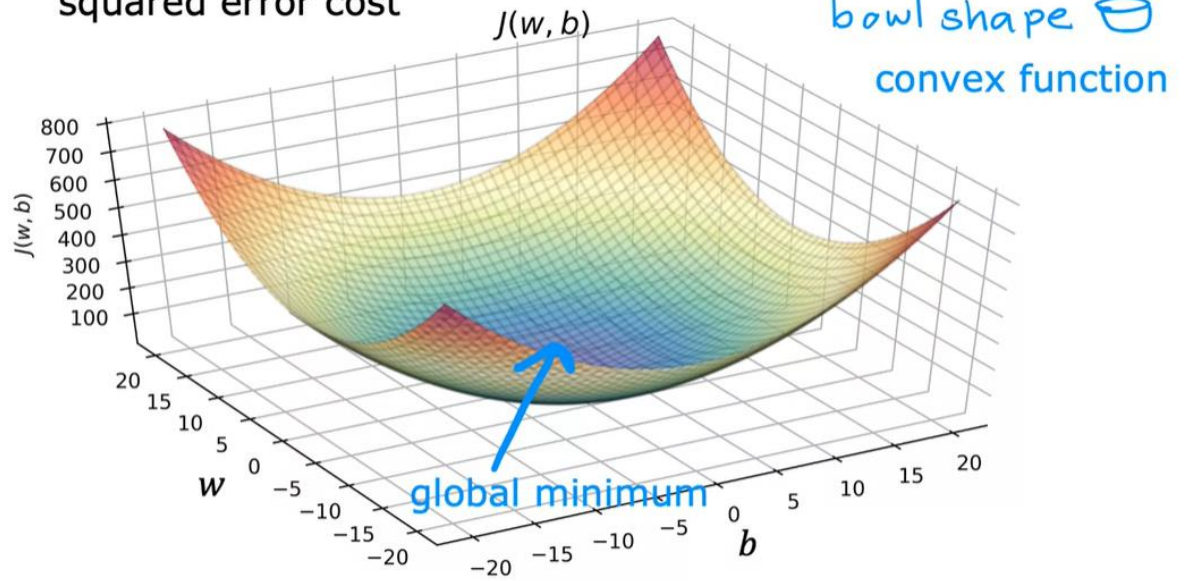
$f_{w,b}(x^{(i)}) = w x^{(i)} + b$

}

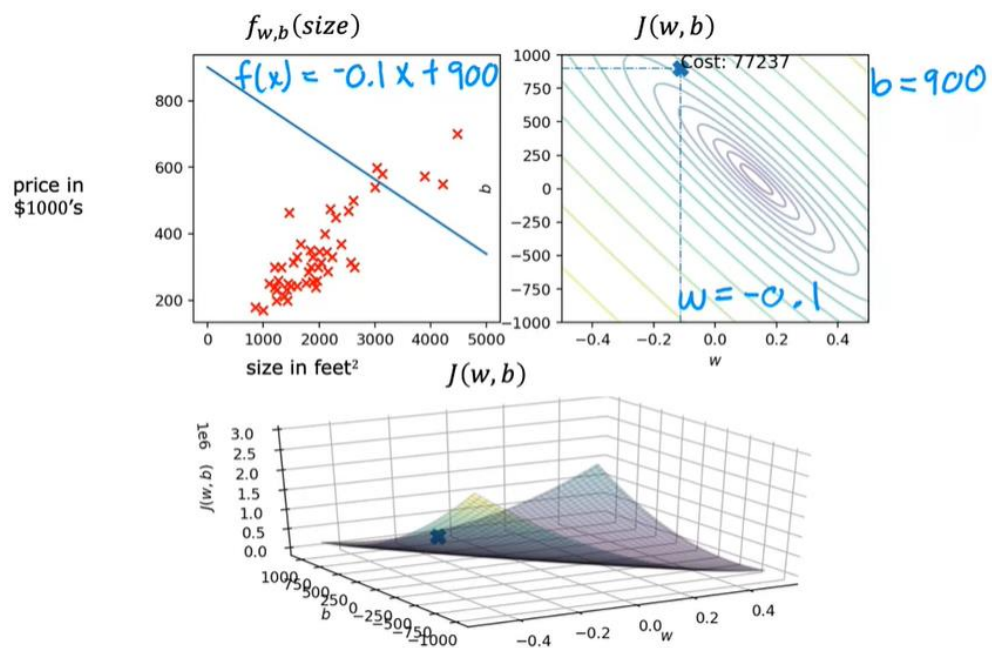
More than one local minimum

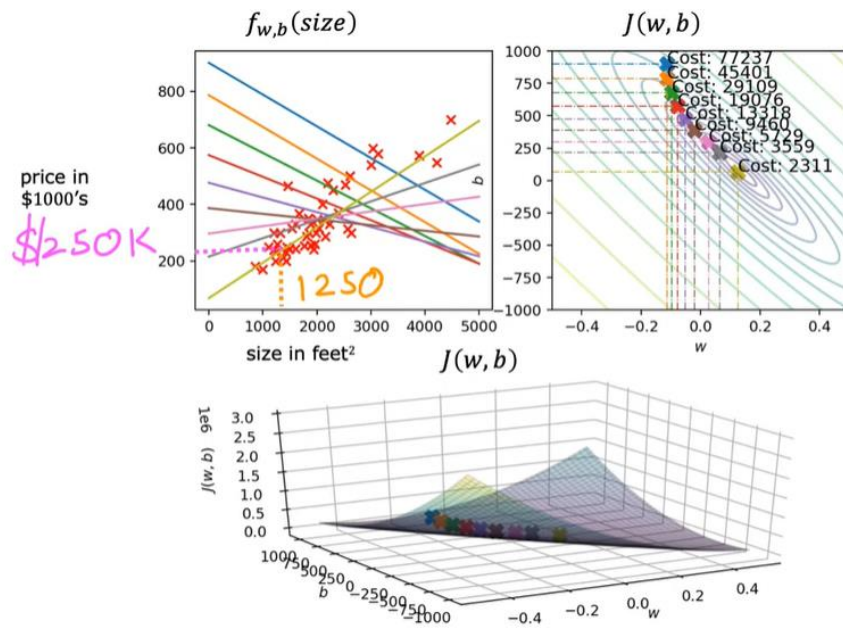


squared error cost



18. Running gradient descent





"Batch" gradient descent



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

other gradient descent: subsets

	x size in feet ²	y price in \$1000's
(1)	2104	400
(2)	1416	232
(3)	1534	315
(4)	852	178
...
(47)	3210	870

$m = 47$

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

1.

Gradient descent is an algorithm for finding values of parameters w and b that minimize the cost function J .

repeat until convergence {

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

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

When $\frac{\partial J(w, b)}{\partial w}$ is a negative number (less than zero), what happens to w after one update step?

- ☒ w increases.
- ☐ w stays the same
- ☐ w decreases
- ☐ It is not possible to tell if w will increase or decrease.

✓ **Correct**

The learning rate is always a positive number, so if you take w minus a negative number, you end up with a new value for w that is larger (more positive).

2.

For linear regression, what is the update step for parameter b ?

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

✓ **Correct**

The update step is $b = b - \alpha \frac{\partial J(w, b)}{\partial b}$ where $\frac{\partial J(w, b)}{\partial b}$ can be computed with this expression:
$$\sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})$$