

Week 1

Machine learning

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

Arthur Samuel (1959)

Machine learning algorithms

rapid advancements

used most in real-world applications

- Supervised learning ← course 1, 2

- Unsupervised learning ←

course 3

- Recommender systems

- Reinforcement learning

The other thing we're
going to spend a lot of

Supervised learning

Learns from being given “right answers”

Regression

Predict a number

infinitely many possible outputs

Classification

predict categories

small number of possible outputs

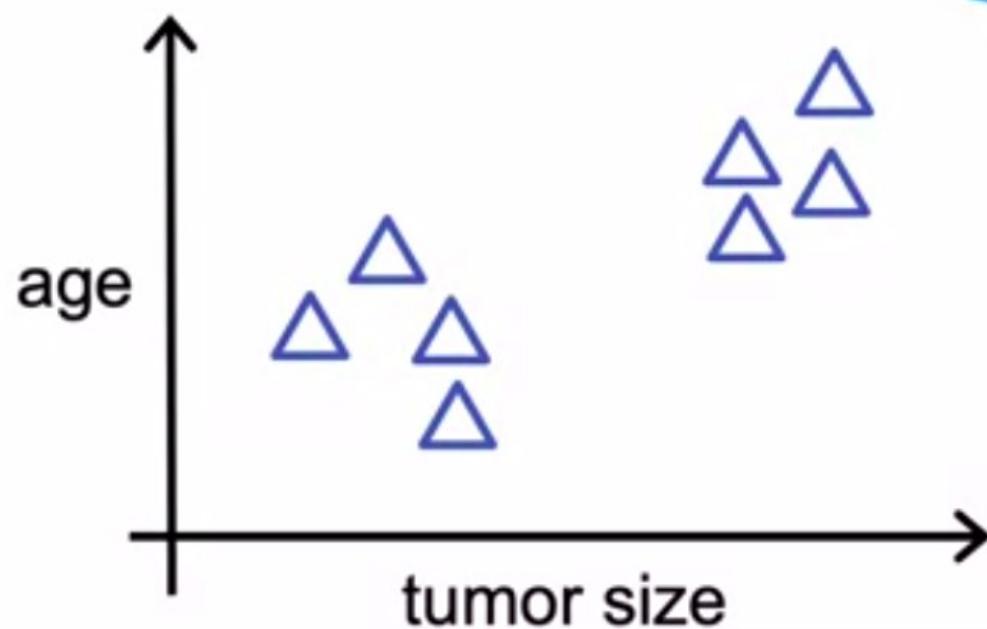
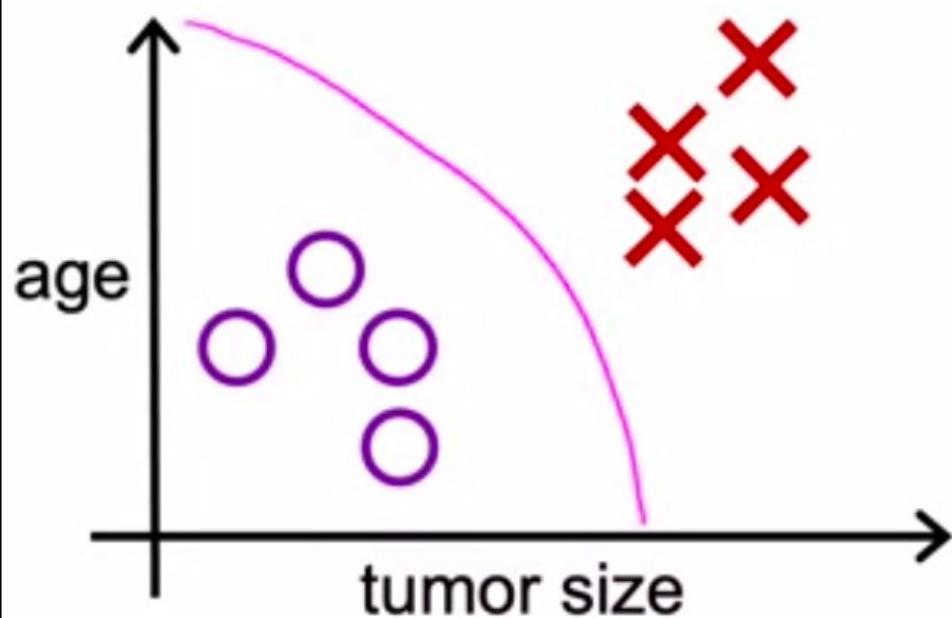
a category,

all of a small set of possible outputs.



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

Unsupervised learning
Find something interesting
in **unlabeled** data.



This is unsupervised learning,

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.

to make too much
sense to you yet.

[Back](#) Practice quiz: Supervised vs unsupervised learning
Graded Quiz • 15 min Due Jul 11, 12:29 PM IST

Practice quiz: Supervised vs unsupervised learning

Total points 2

1. Which are the two common types of supervised learning? (Choose two)

1 point

- Regression
- Classification
- Clustering

2.

1 point

Which of these is a type of unsupervised learning?

- Regression
- Clustering
- Classification

Terminology

Training Data used to train the model

set: x

size in feet²

y
price in \$1000's

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

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

$(1^{\text{st}}, 2^{\text{nd}}, 3^{\text{rd}} \dots)$

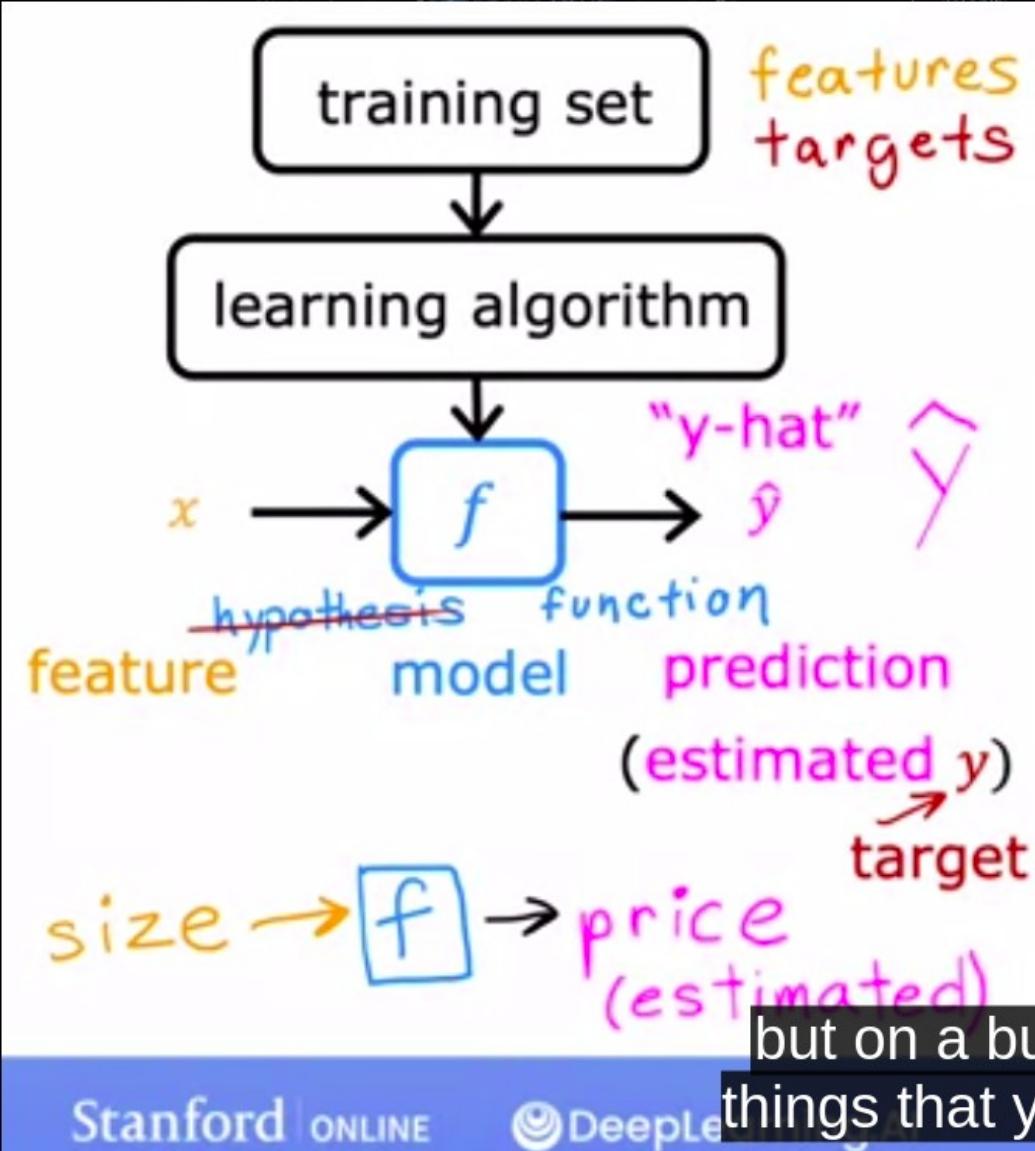
$$x^{(1)} = 2104$$

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

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

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

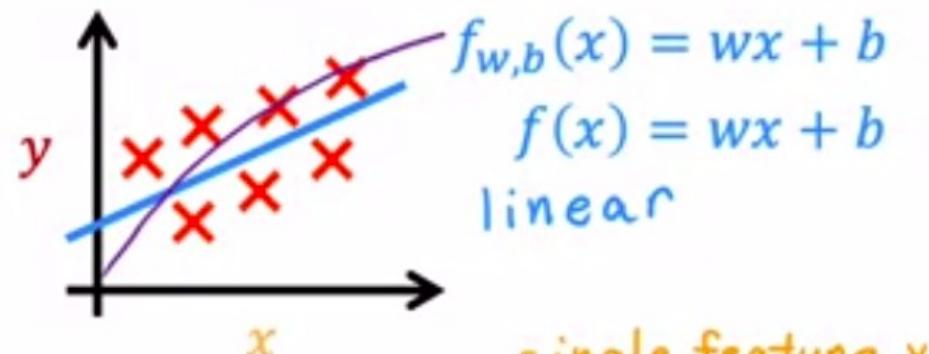
the training set and refers
to row i in the table.



How to represent f ?

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

$$f(x)$$



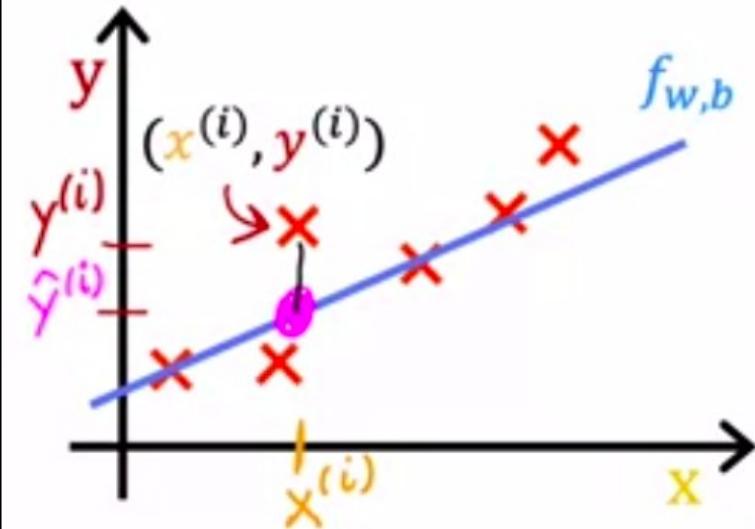
Linear regression with **one** variable.
size
Univariate linear regression.
one variable

Notation

Here is a summary of some of the notation you will encounter.

General	Notation	Description	Python (if applicable)
	a	scalar, non bold	
	\mathbf{a}	vector, bold	
Regression			
	\mathbf{x}	Training Example feature values (in this lab - Size (1000 sqft))	<code>x_train</code>
	y	Training Example targets (in this lab Price (1000s of dollars)).	<code>y_train</code>
	$x^{(i)}, y^{(i)}$	i_{th} Training Example	<code>x_i , y_i</code>
	m	Number of training examples	<code>m</code>
	w	parameter: weight,	<code>w</code>
	b	parameter: bias	<code>b</code>
	$f_{w,b}(x^{(i)})$	The result of the model evaluation at $x^{(i)}$ parameterized by w, b : $f_{w,b}(x^{(i)}) = wx^{(i)} + b$	<code>f_wb</code>

Cost function: Squared error cost function



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

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

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

m = number of training examples

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

intuition

Find w, b :

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

model:

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

parameters:

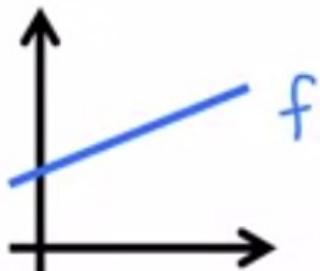
$$w, b$$

cost function:

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

goal:

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

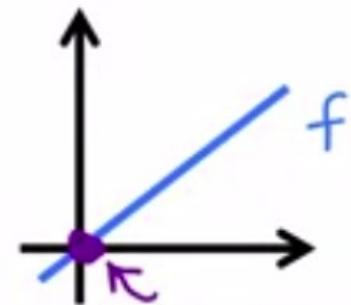


simplified

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

$$b = \emptyset$$

$$w$$



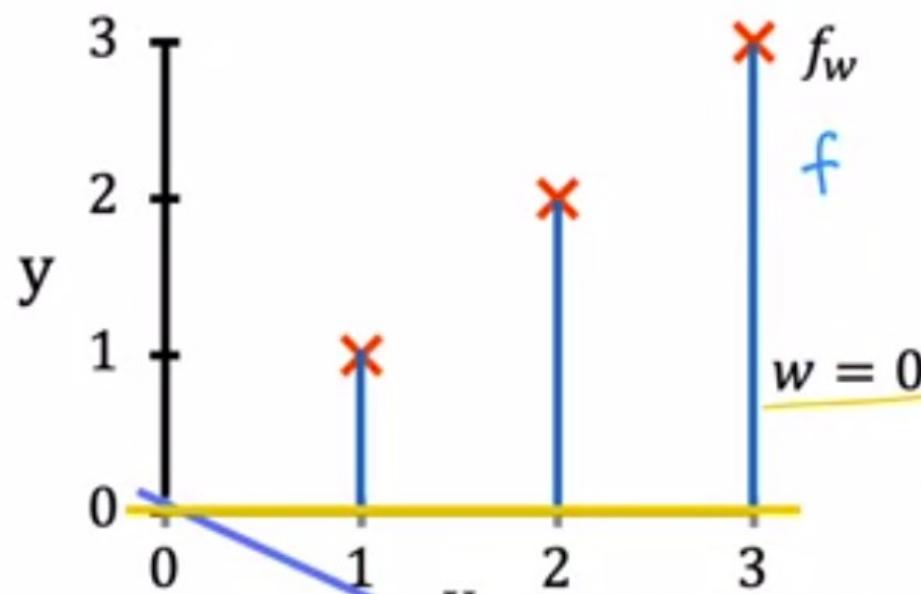
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2$$

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

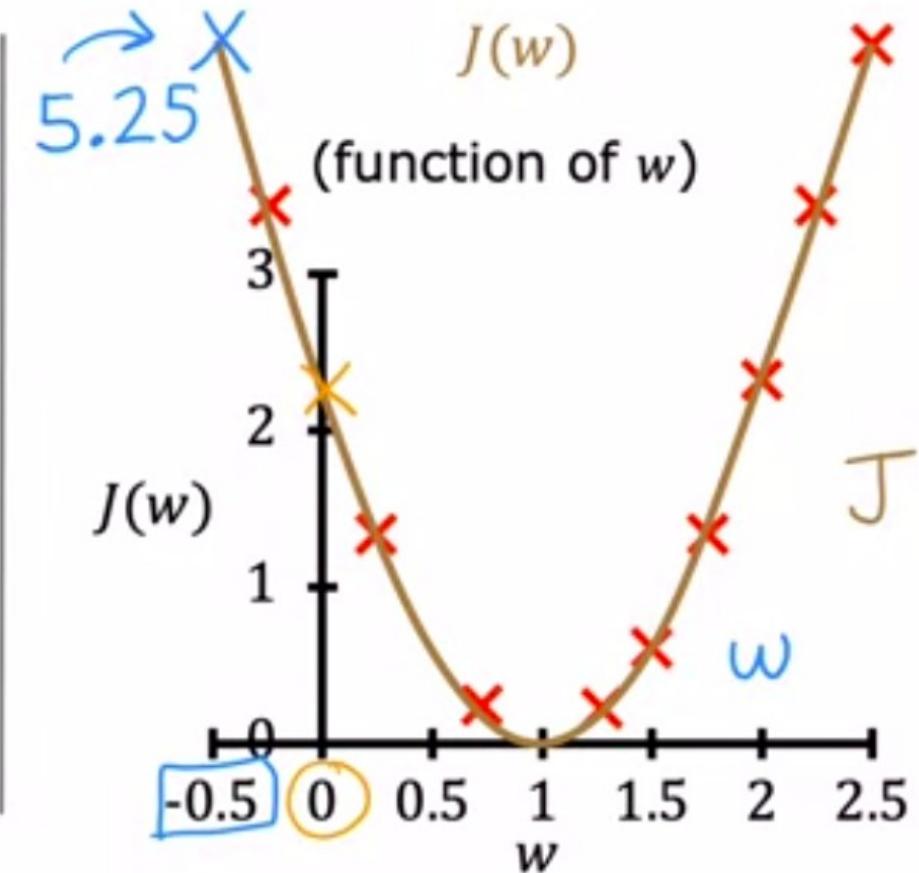
$$w x^{(i)}$$

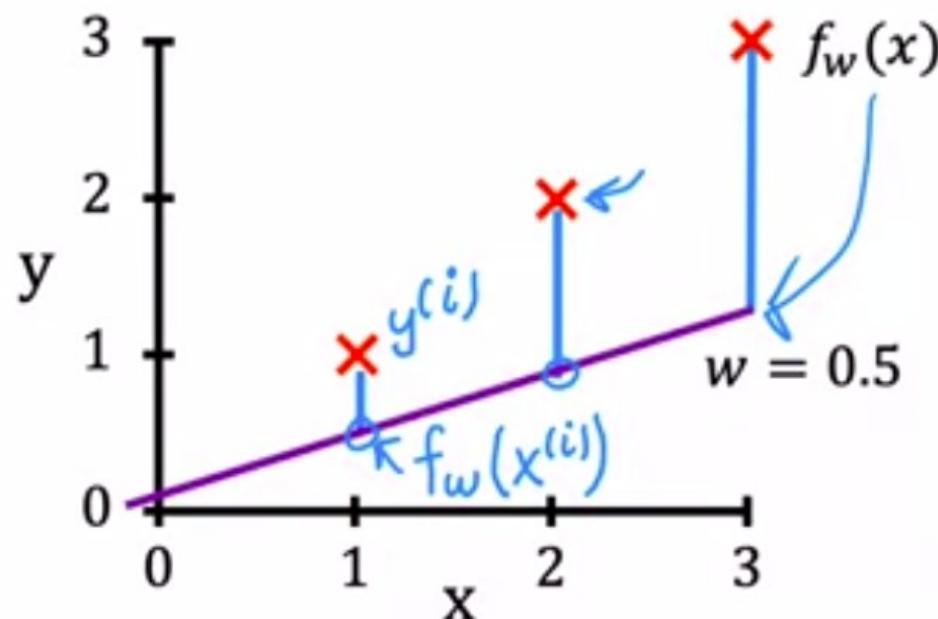
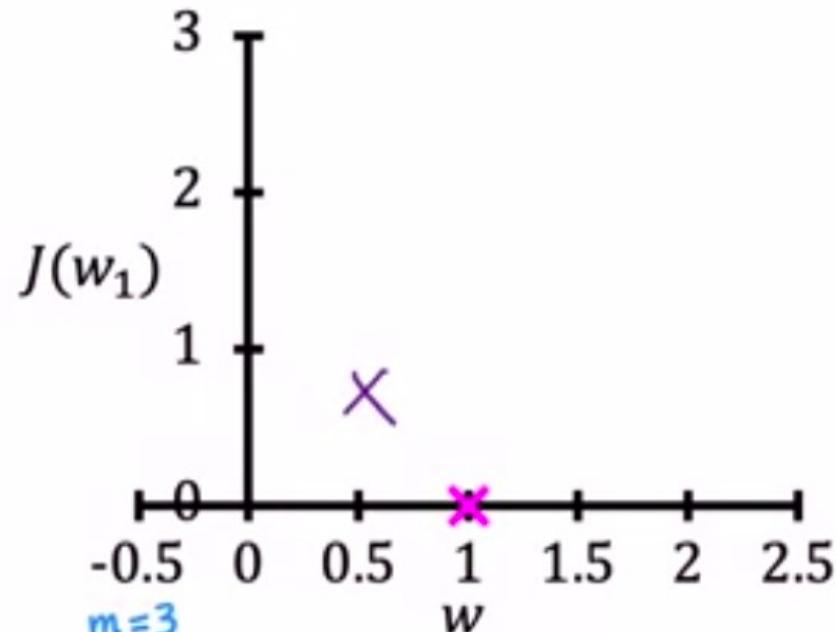
different values for the parameter w. In particular,

$f_w(x)$
(function of x)



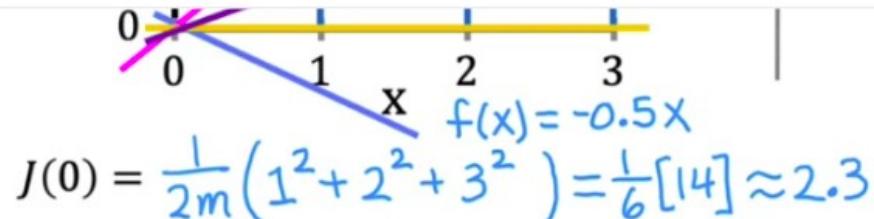
$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2) = \frac{1}{6}[14] \approx 2.3$$



$f_w(x)$ (function of x) $J(w)$ (function of w)

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

Question



When does the model fit the data relatively well, compared to other choices for parameter w?

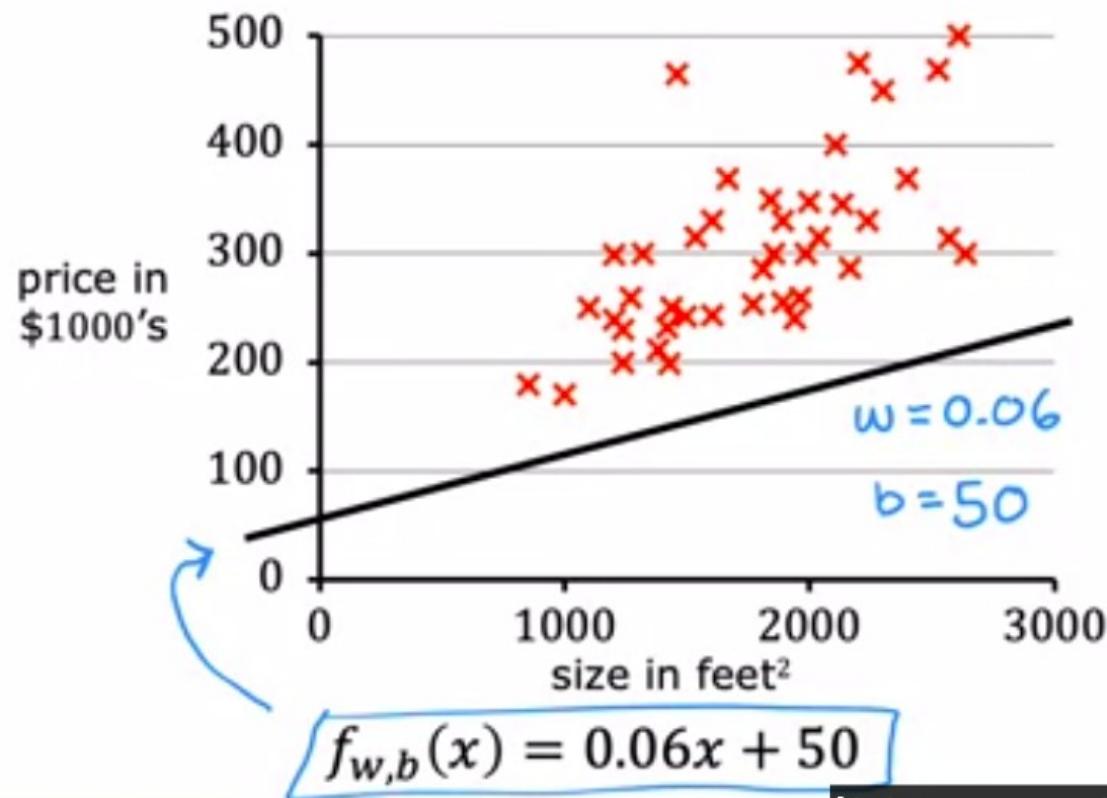
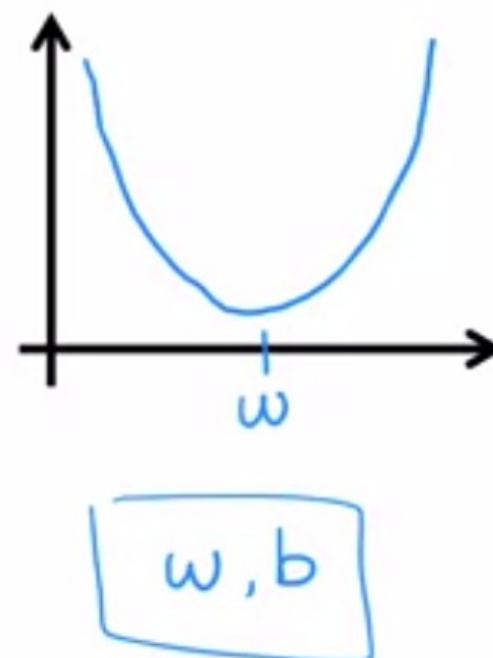
- When x is at or near a minimum.
- When w is close to zero.
- When the cost J is at or near a minimum.
- When $f_w(x)$ is at or near a minimum for all the values of x in the training set.

✓ **Correct**

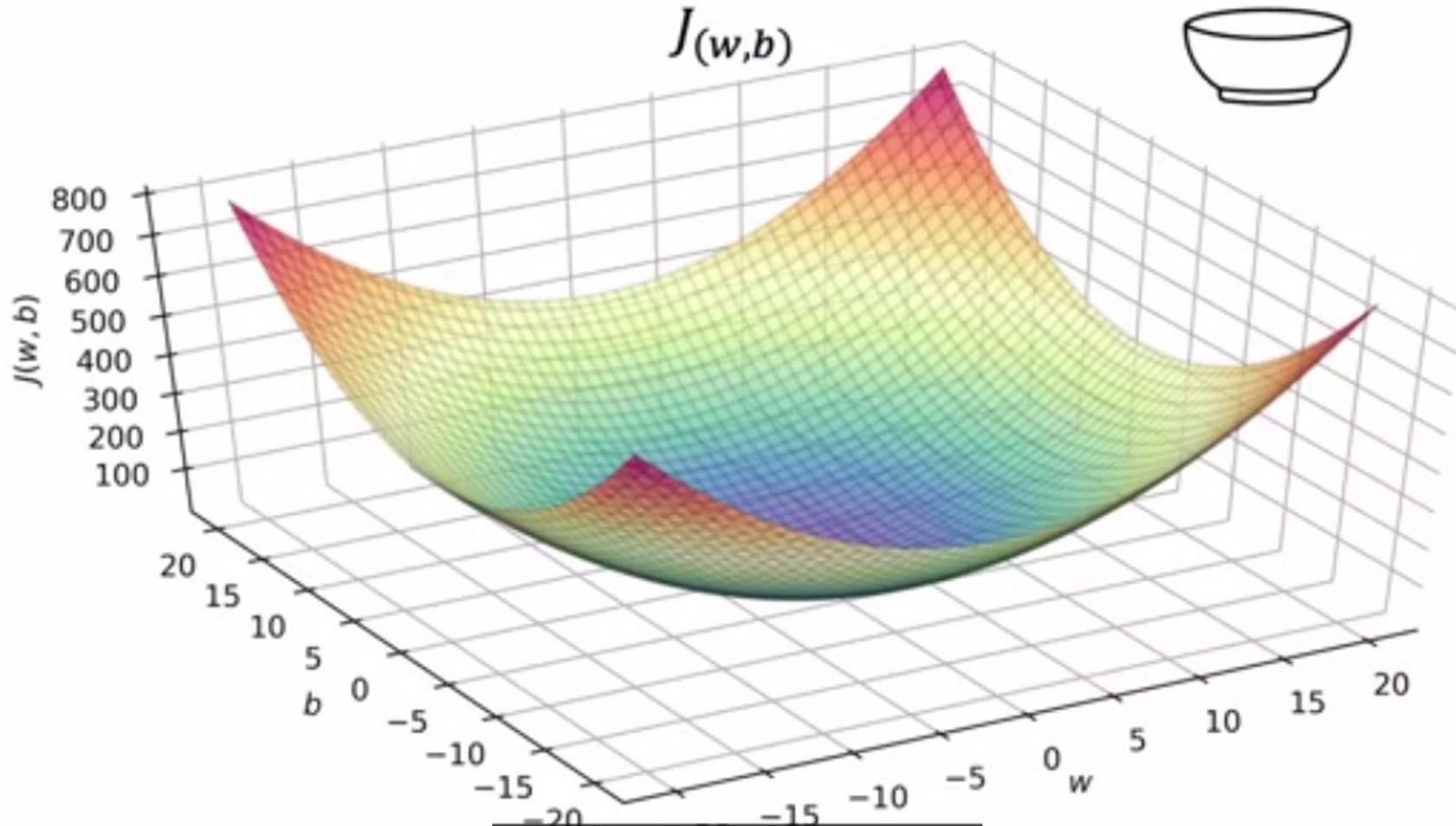
When the cost is relatively small, closer to zero, it means the model fits the data better compared to other choices for w and b .

Skip

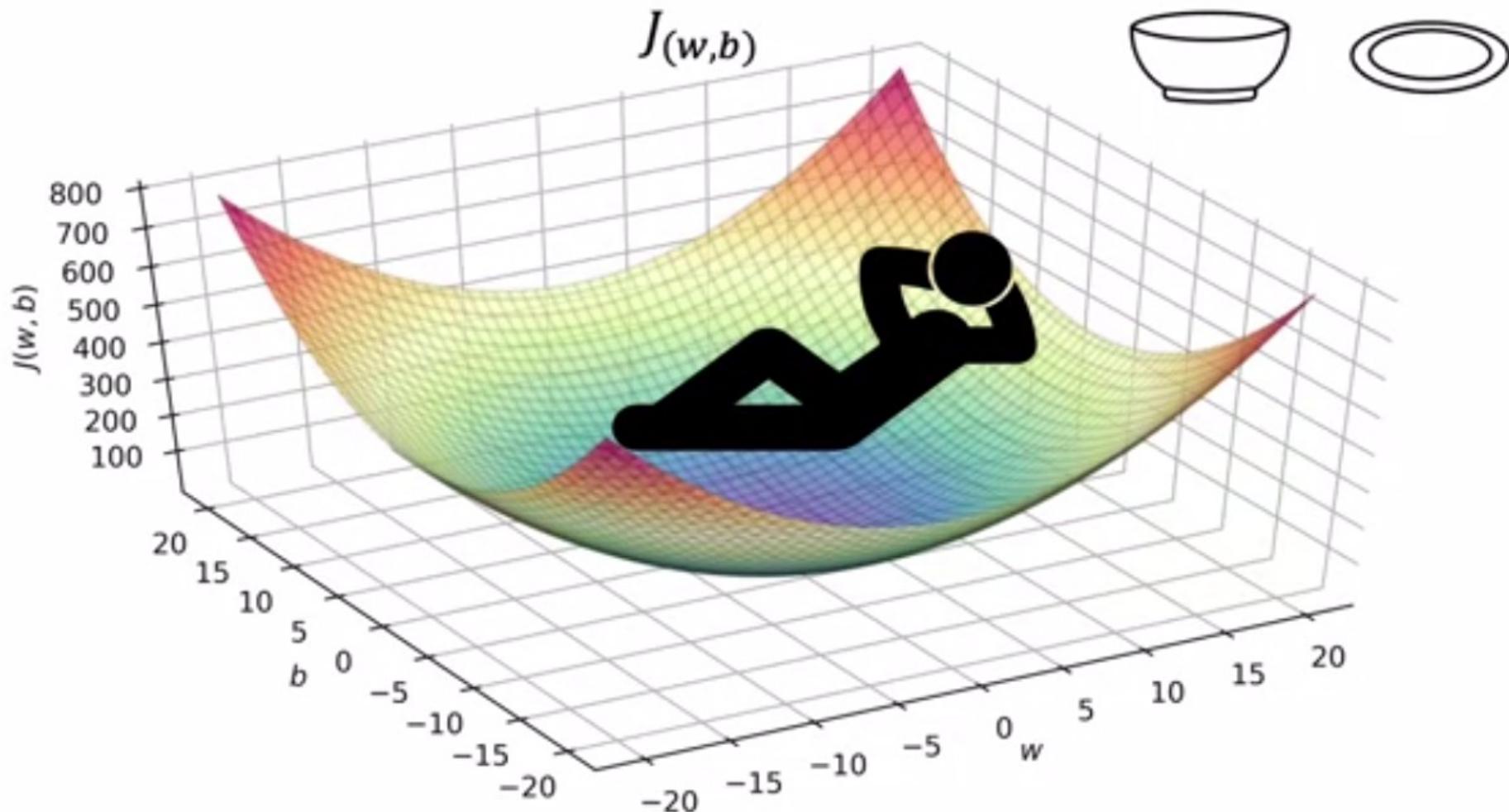
Continue

$f_{w,b}$ (function of x) J (function of w, b)

It turns out that
the cost function



maybe because I'm a
little bit hungry,

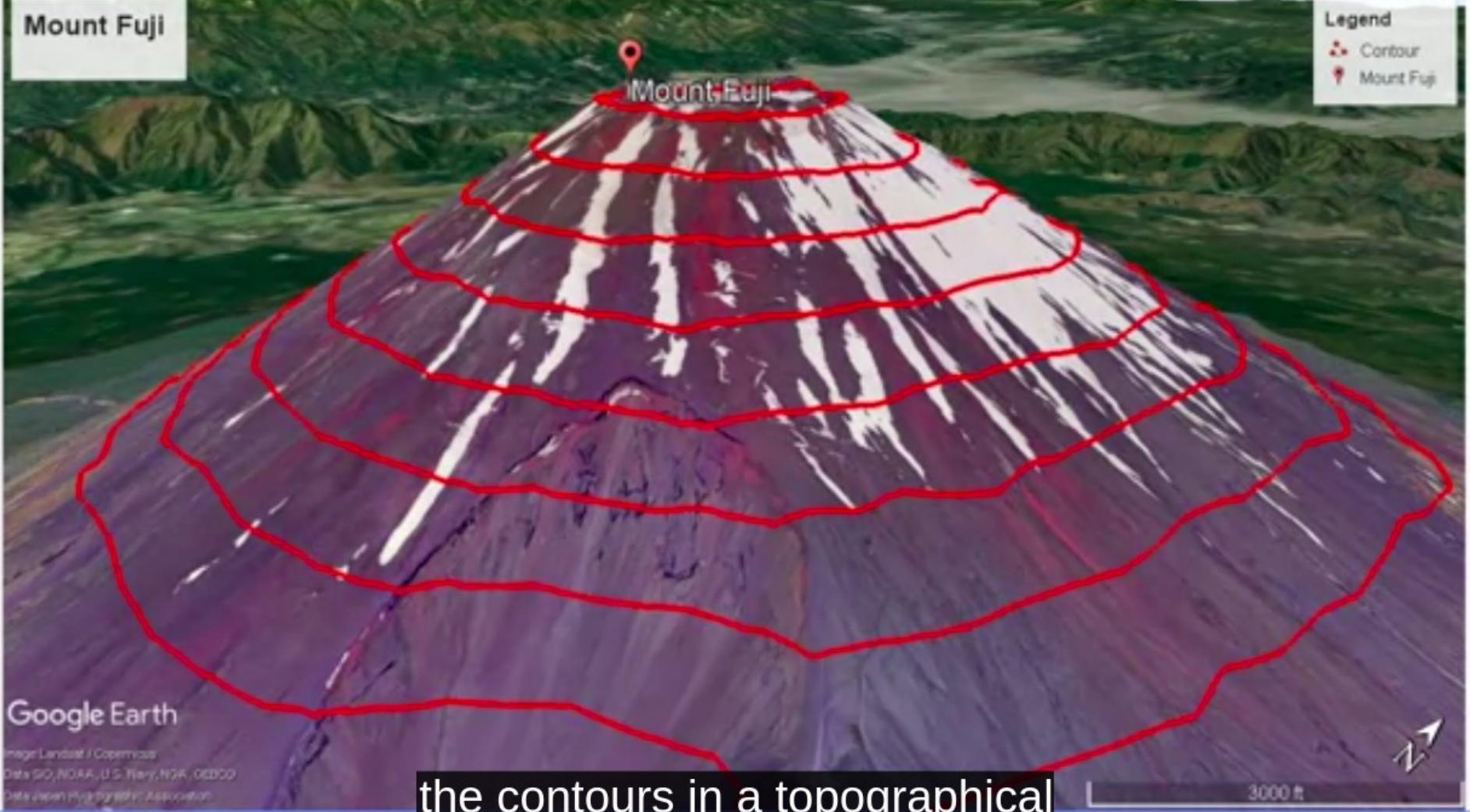


Mount Fuji

Legend

- Contour
- Mount Fuji

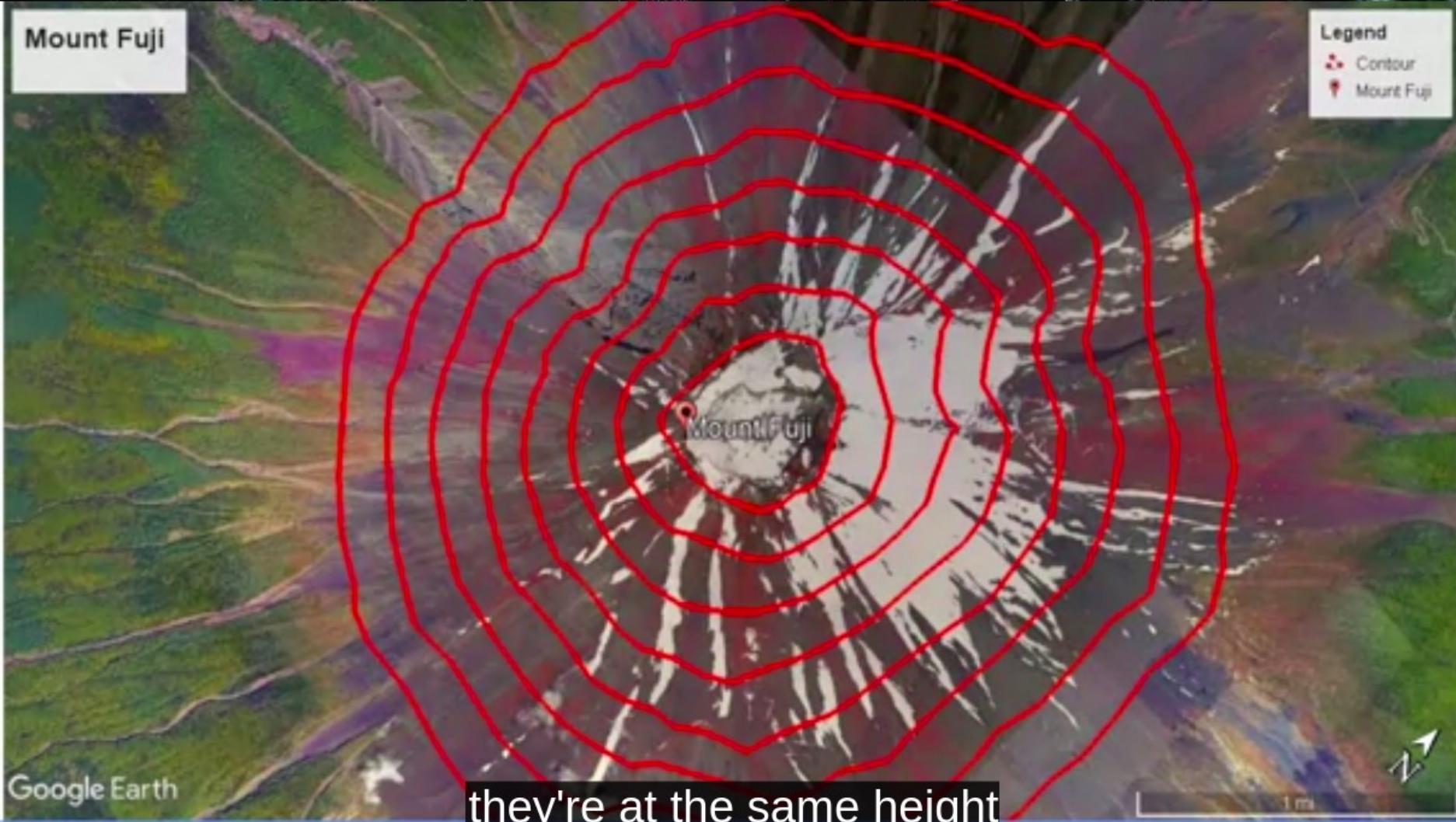
Mount Fuji



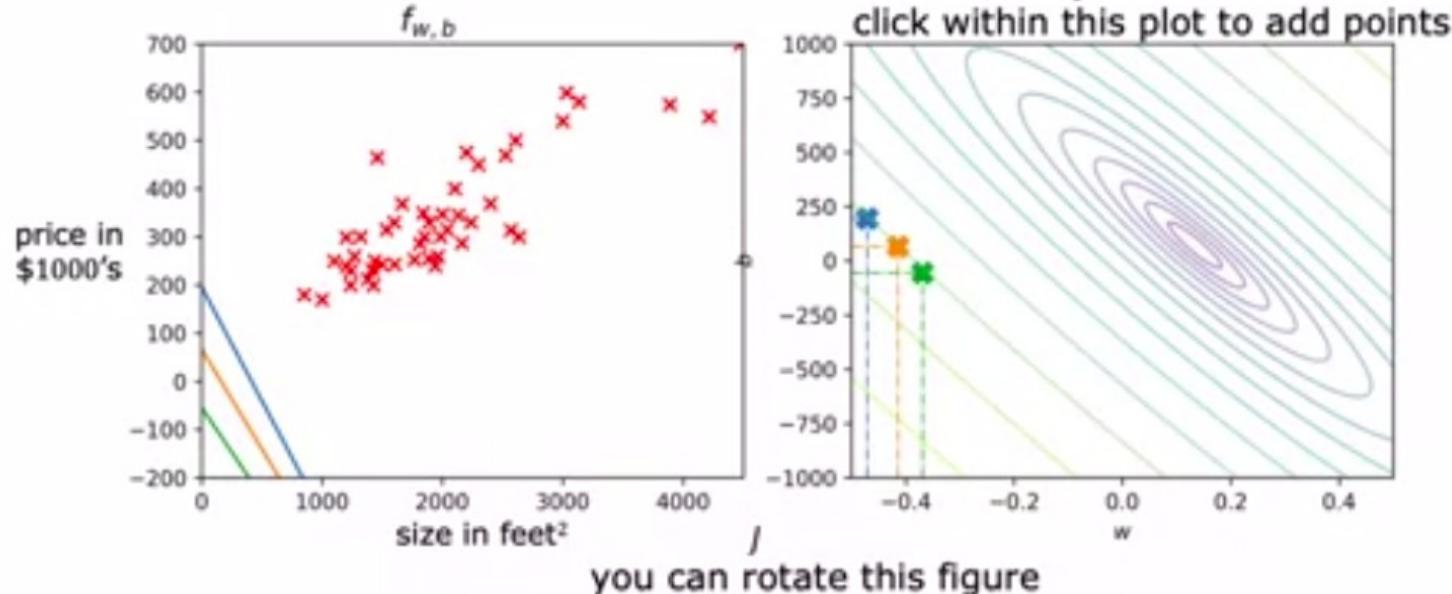
the contours in a topographical
map are basically

Mount Fuji

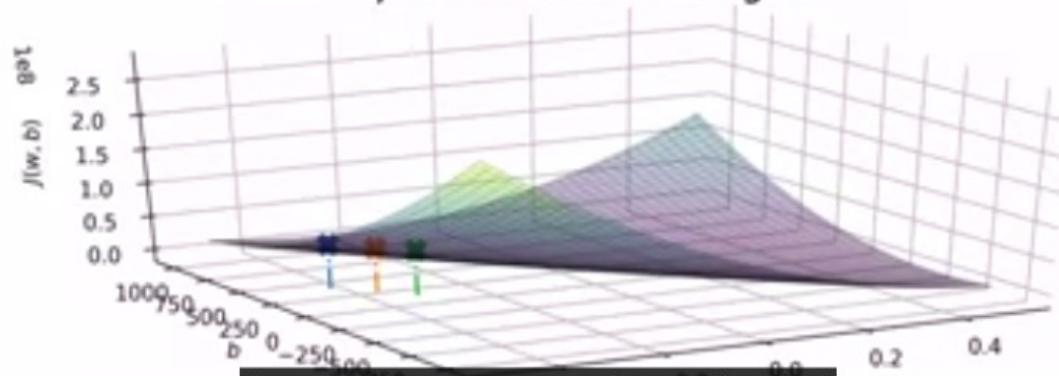
Legend
Contour
Mount Fuji



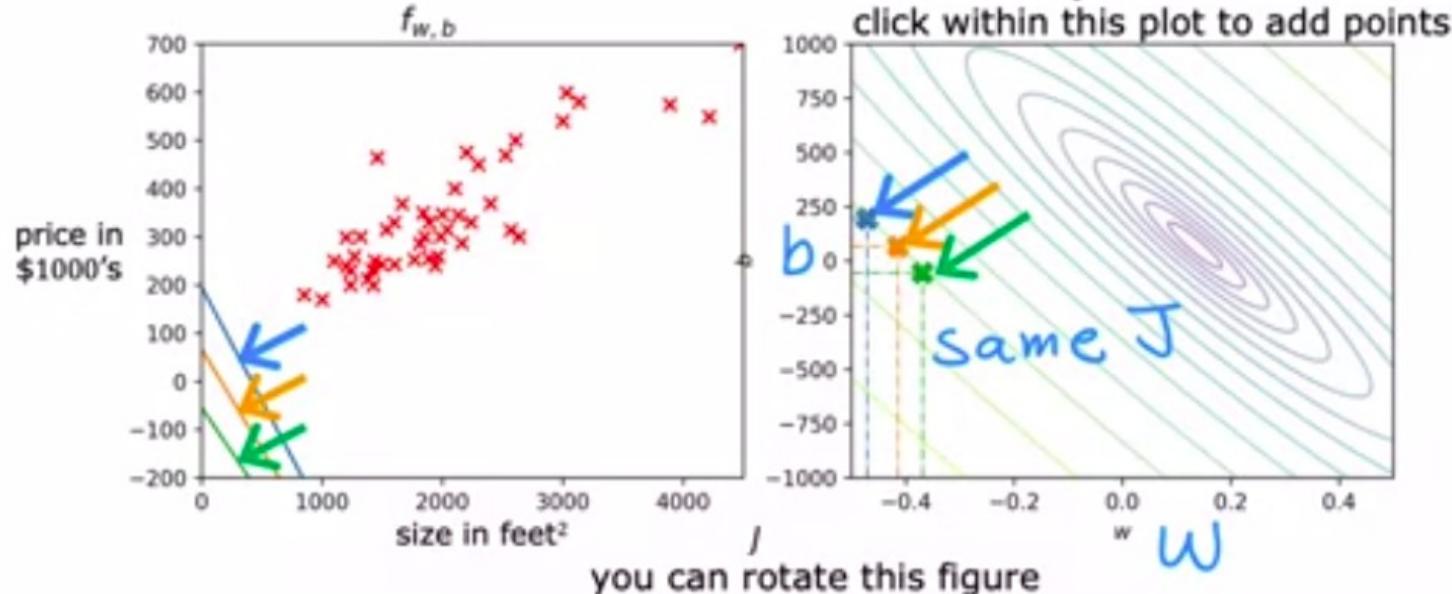
they're at the same height
for different heights.



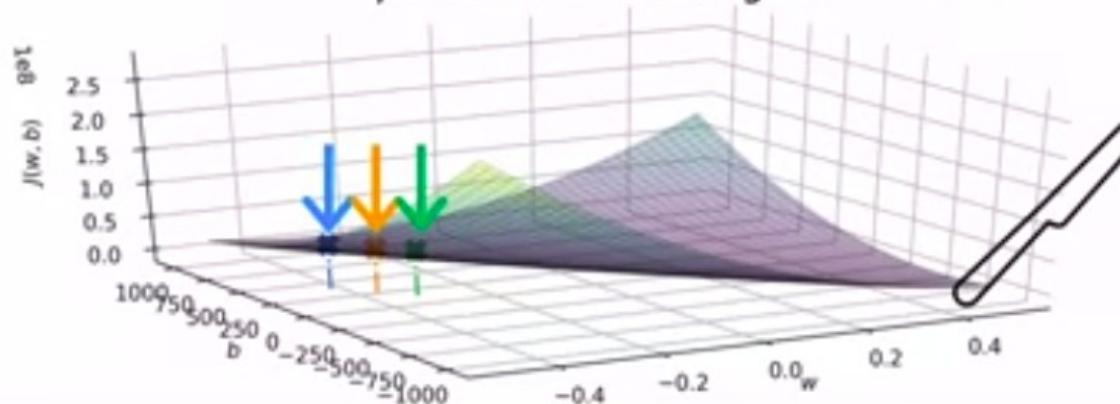
J
you can rotate this figure

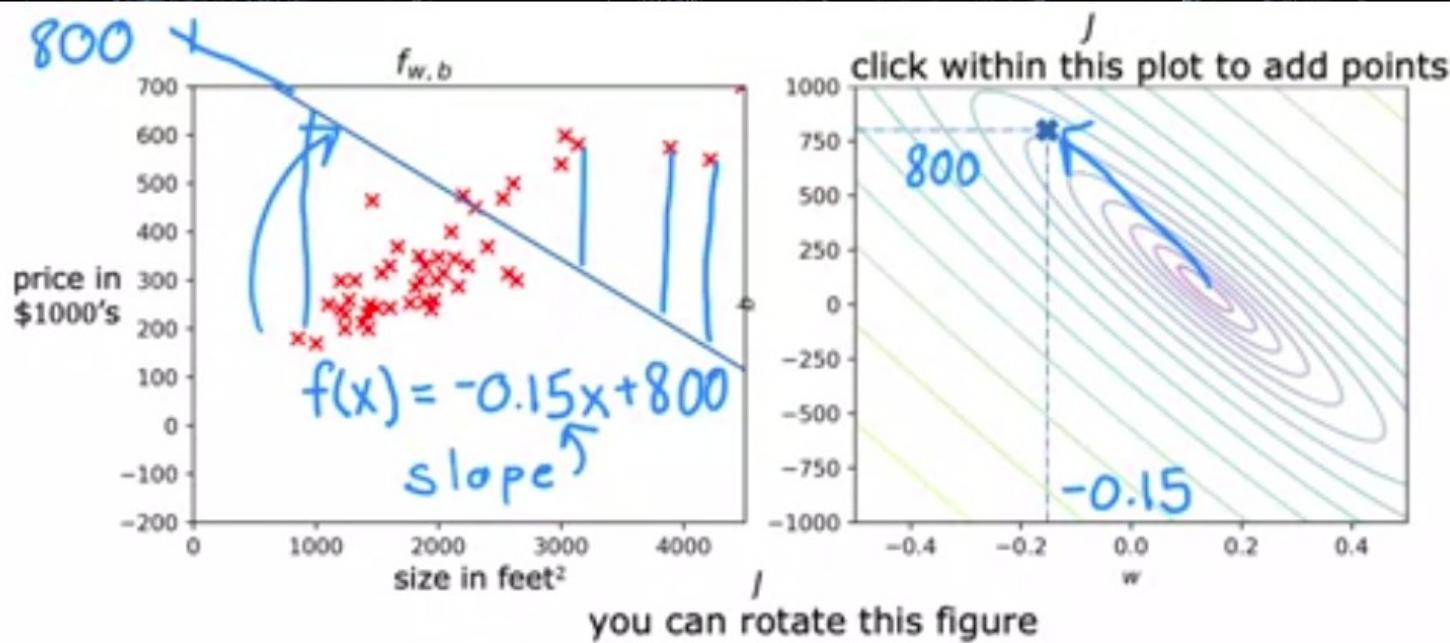


but it is actually a bowl
just very stretched out,

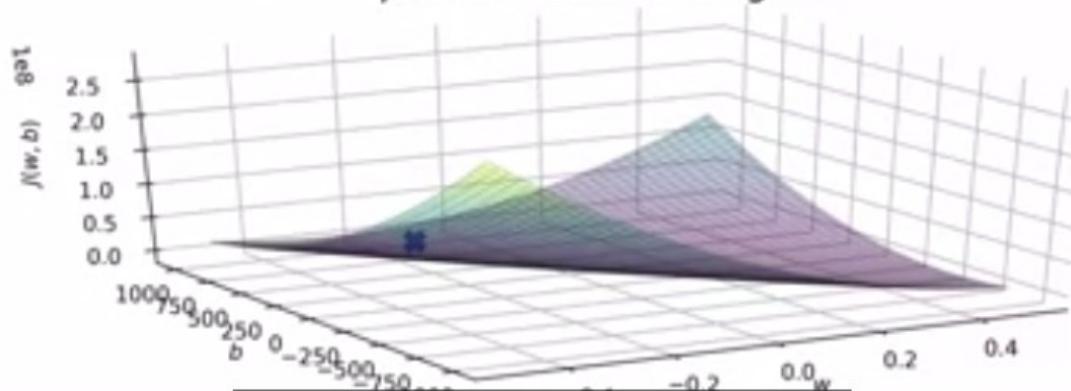


J
you can rotate this figure

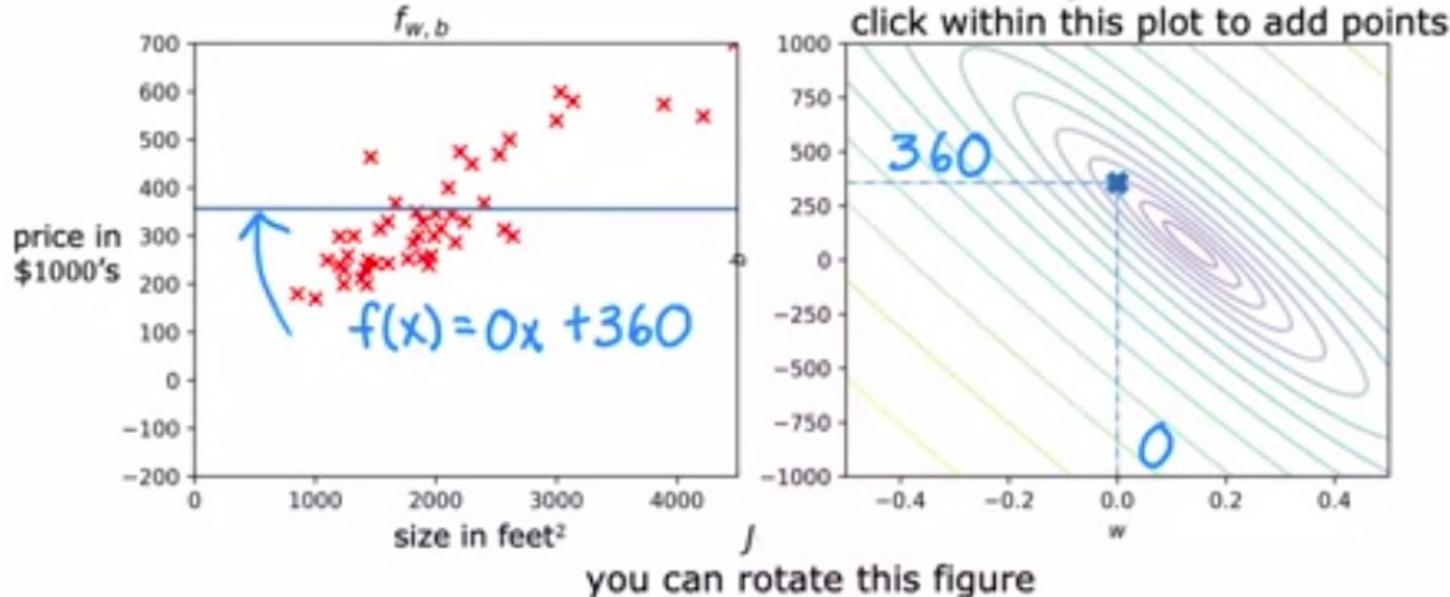




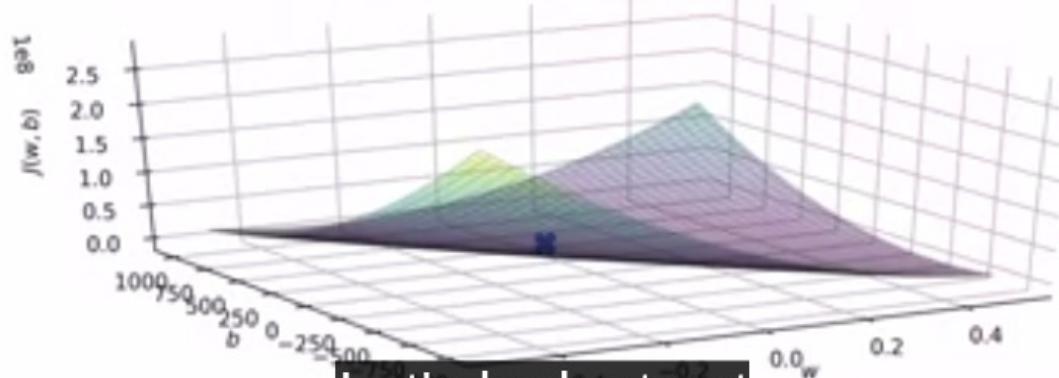
you can rotate this figure



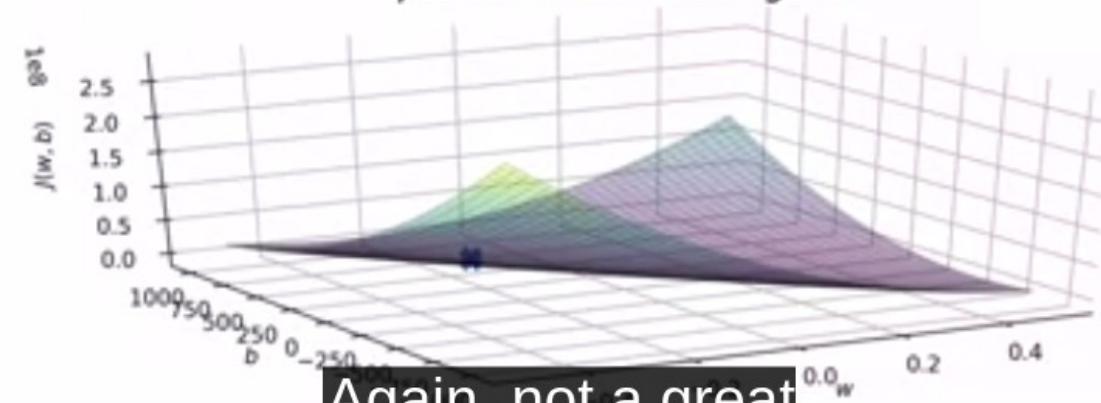
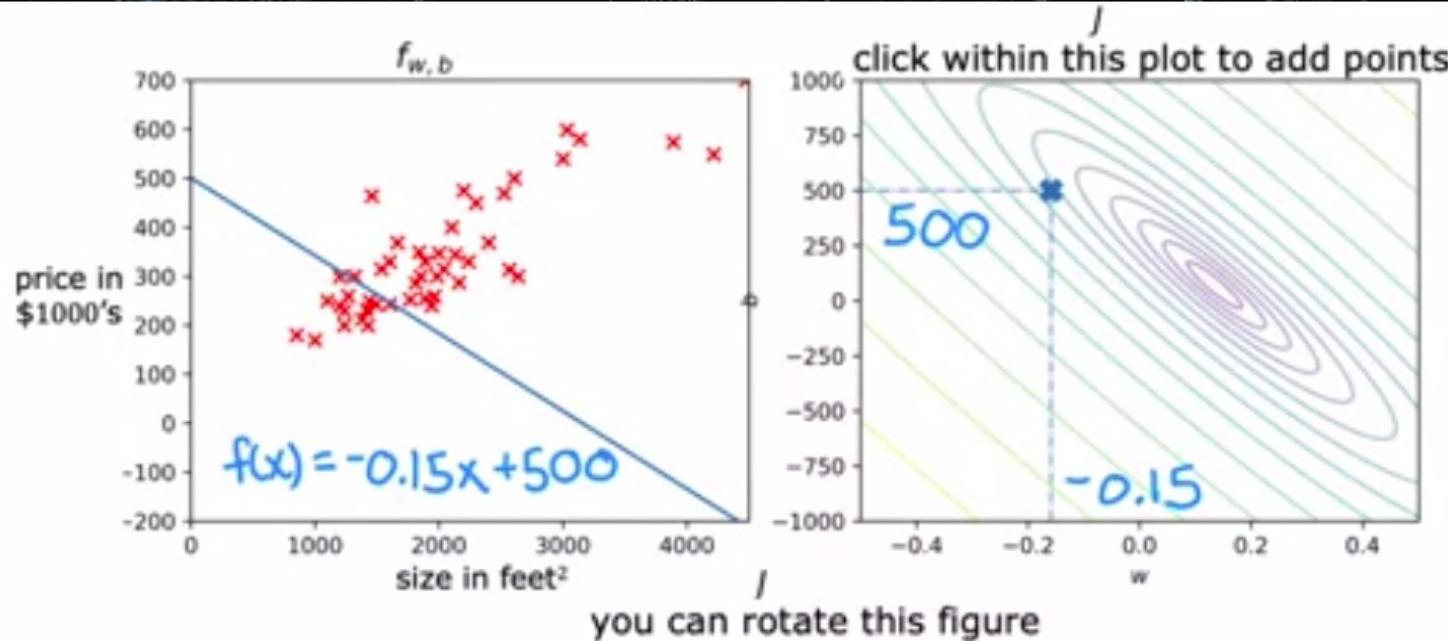
another example with a
different choice of w and b.



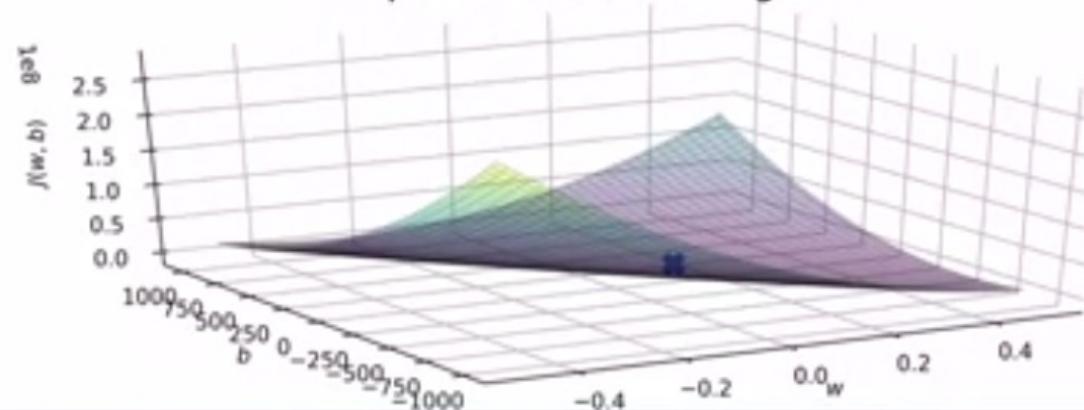
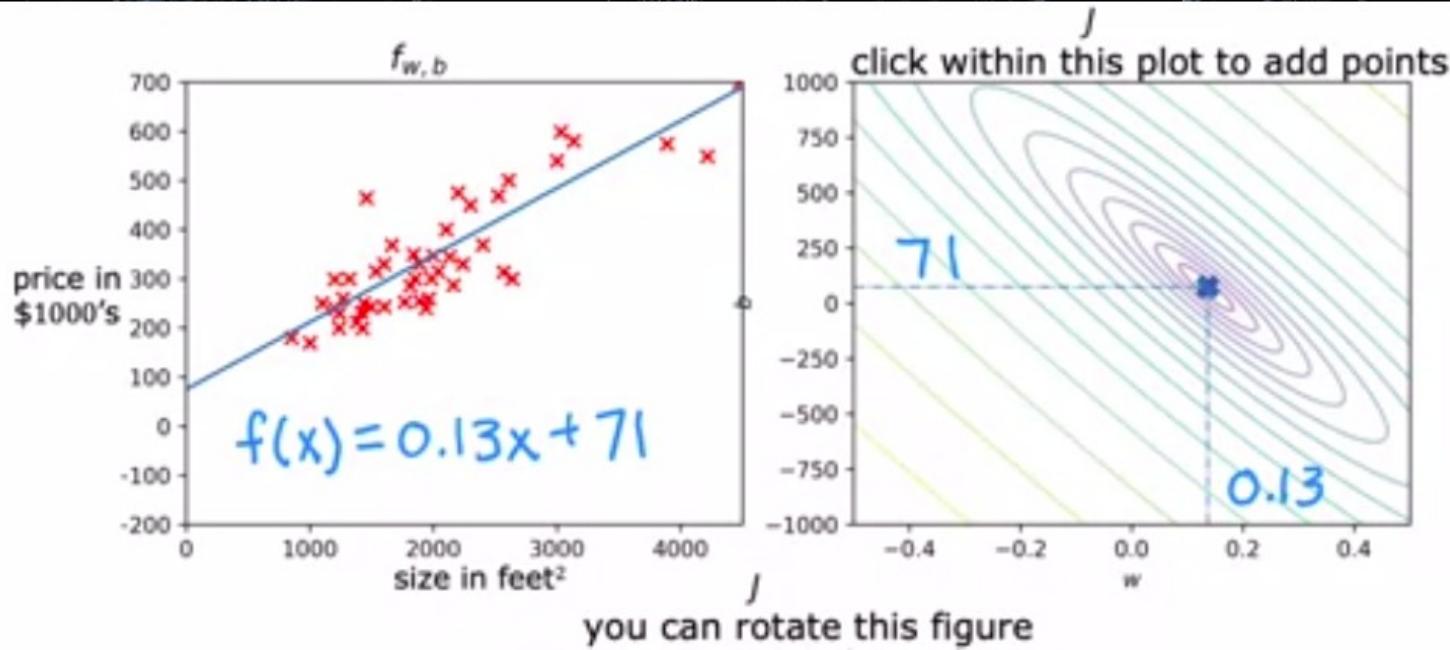
j
you can rotate this figure

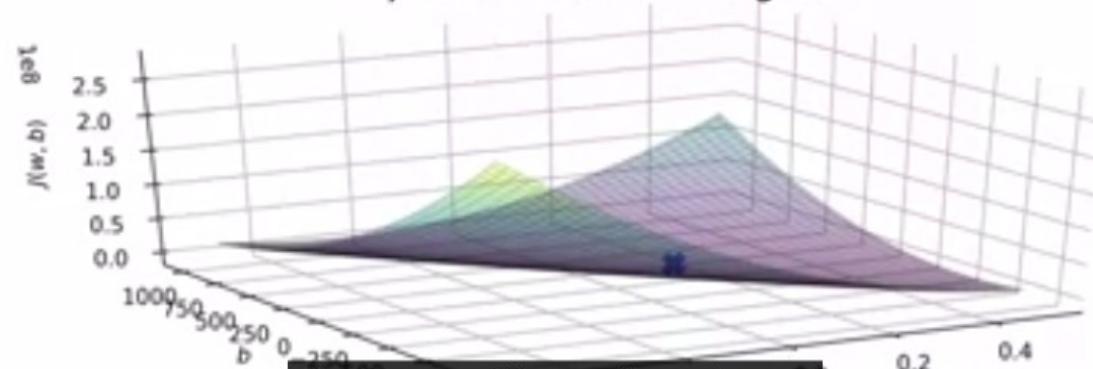
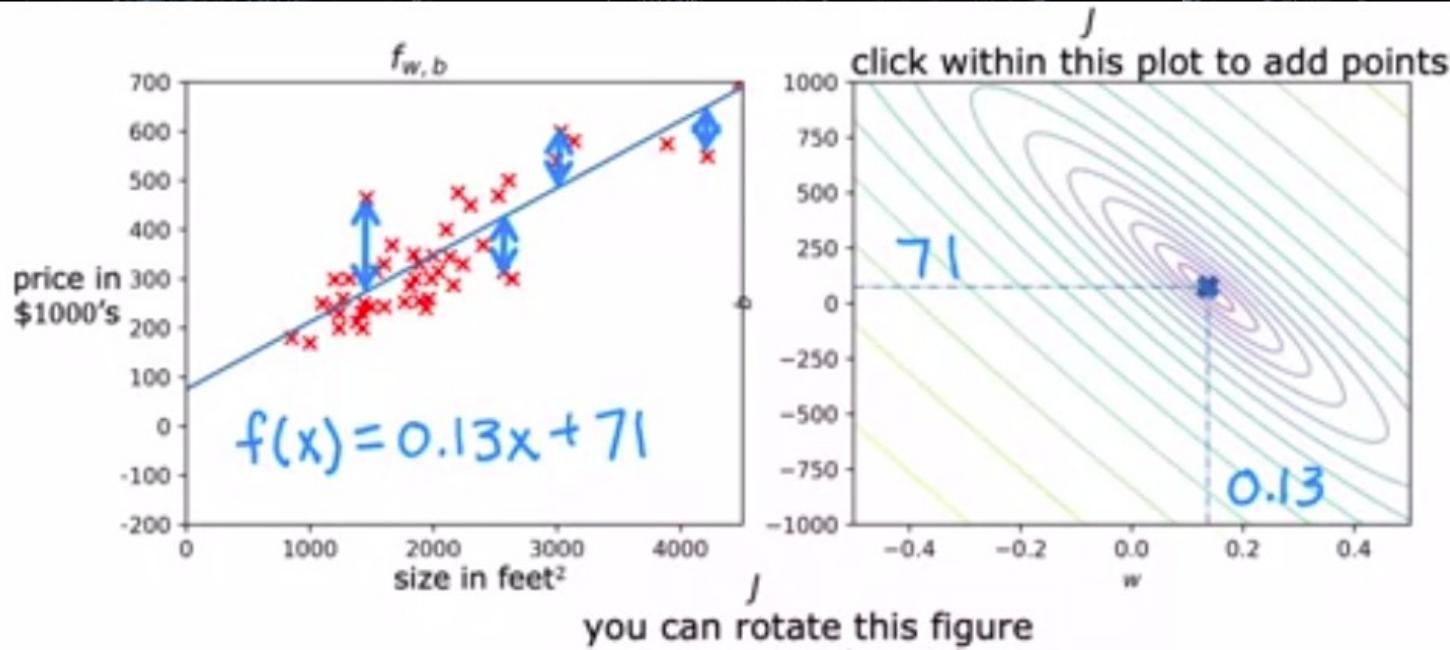


Let's look at yet another example.



Again, not a great
fit to the data,





the predicted values
on the straight line,

[Back](#) Practice quiz: Regression
Graded Quiz • 10 min

Due Jul 11, 12:29 PM IST

Practice quiz: Regression

Total points 2

1.

1 point

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?

- (x, y)
- x
- w and b .
- m

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?

1 point

- The selected values of the parameters w and b cause the algorithm to fit the training set really poorly.

[Back](#) Practice quiz: Regression

Due Jul 11, 12:29 PM IST

Graded Quiz • 10 min

- (x, y)
- x
- w and b .
- m



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?

1 / 1 point

- The selected values of the parameters w and b cause the algorithm to fit the training set really poorly.
- 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 well.



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

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

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

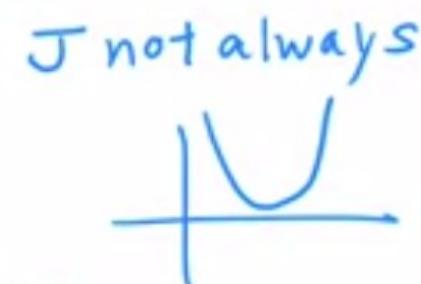
Outline:

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

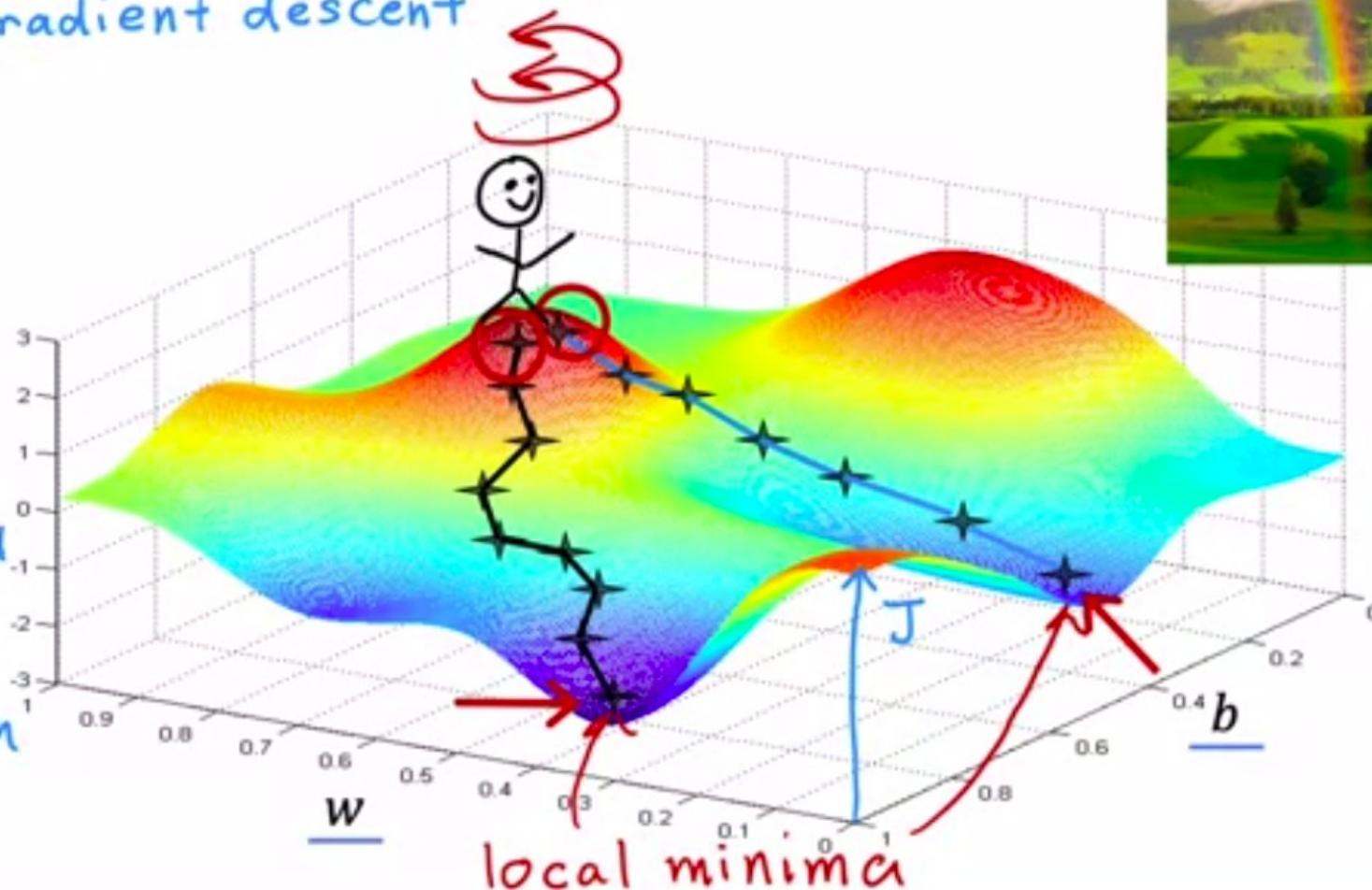


Let's take a look
at an example of

gradient descent

$J(w, b)$

not squared
error cost
not linear
regression



Gradient descent algorithm

Assignment | Truth assertion

Repeat

$$\left\{ \begin{array}{l} w = \\ b = \end{array} \right.$$

Question

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

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

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.)

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

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



Correct

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

Skip

Continue

Gradient descent algorithm

Repeat until convergence

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

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

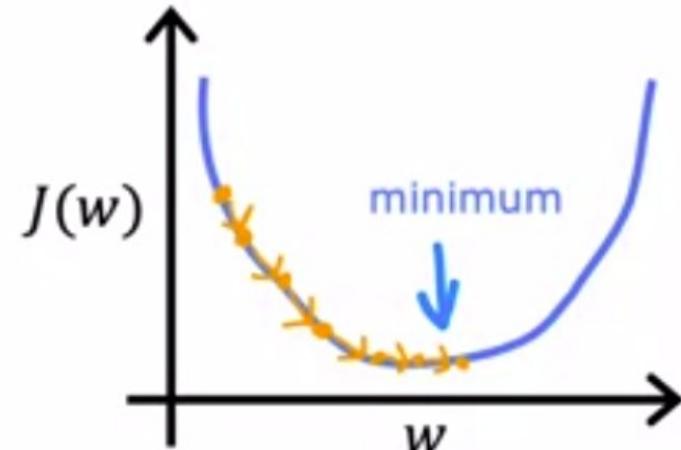
$$\begin{aligned} tmp_w &= w - \alpha \frac{\partial}{\partial w} J(w, b) \\ w &= tmp_w \\ tmp_b &= b - \alpha \frac{\partial}{\partial b} J(w, b) \\ b &= tmp_b \end{aligned}$$

and thus this updated value for

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

If α is too small...

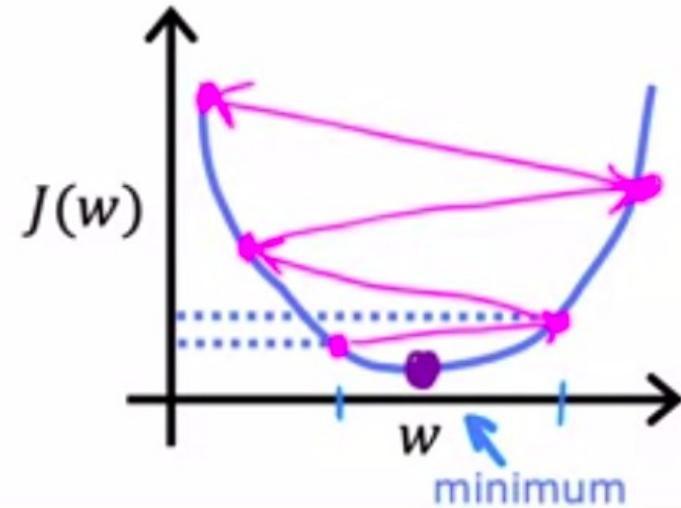
Gradient descent may be slow.

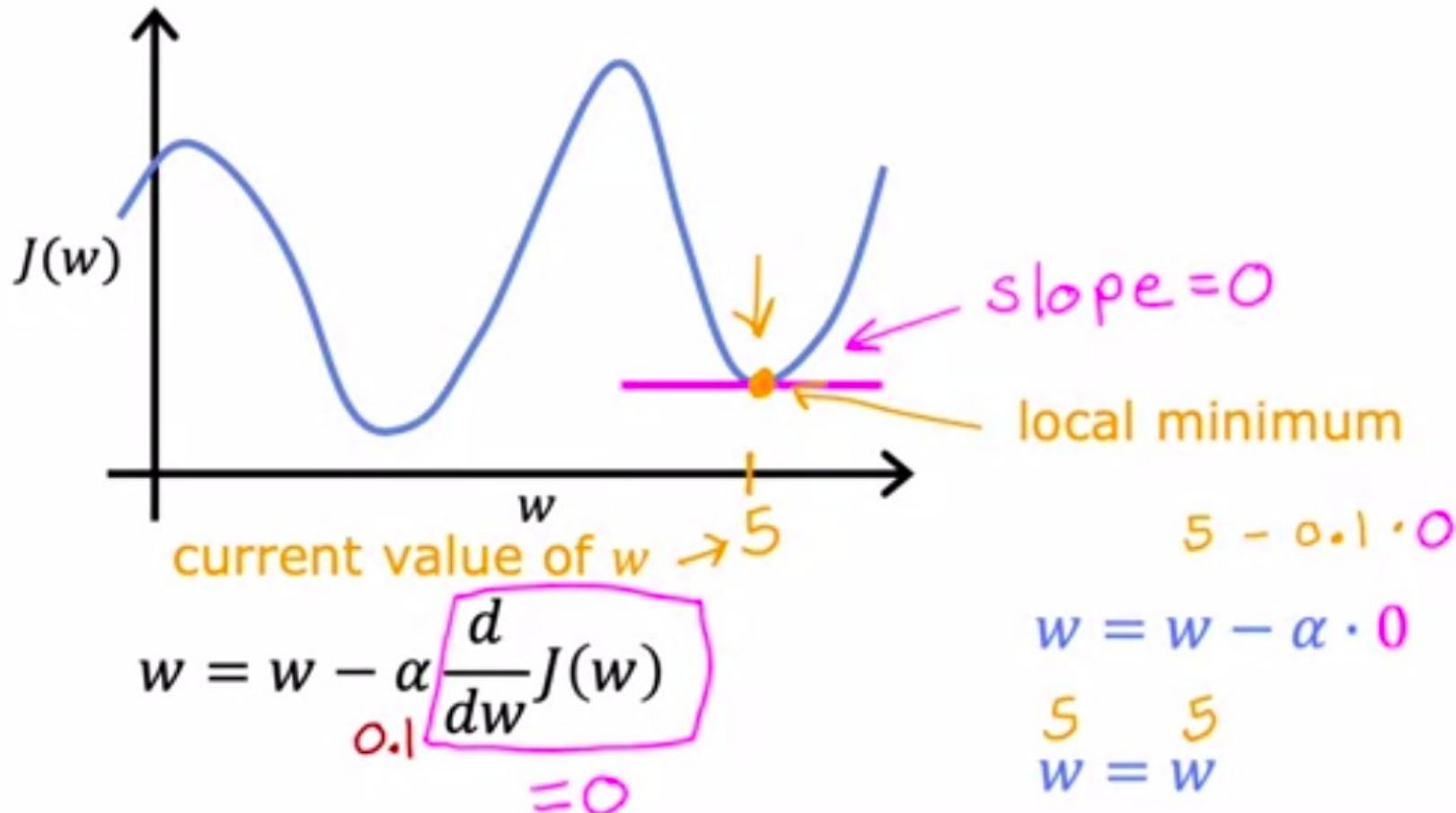


If α is too large...

Gradient descent may:

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





Can reach local minimum with fixed learning rate

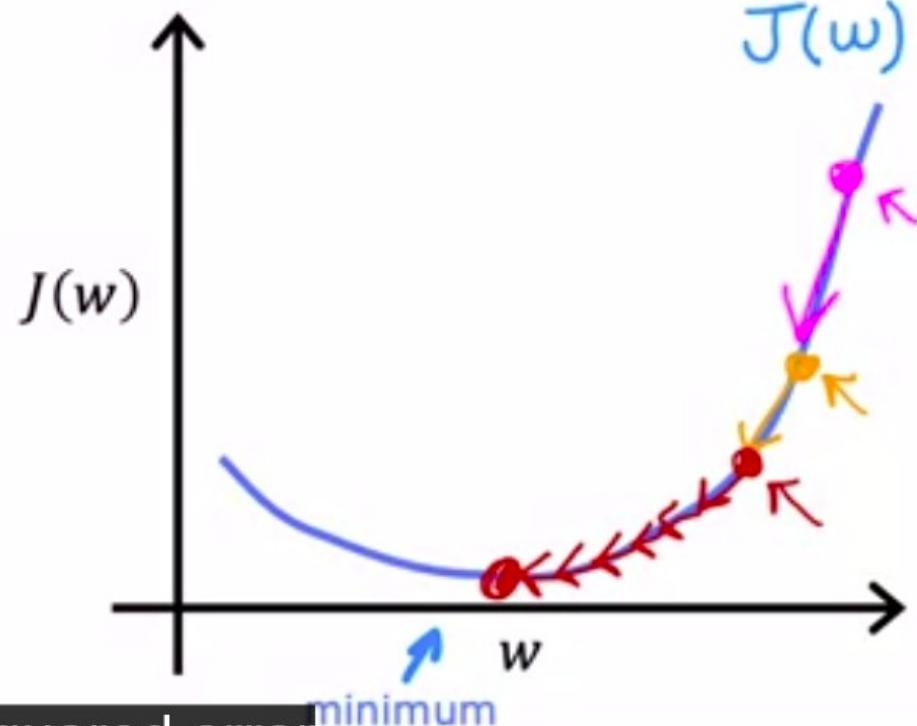
$$w = w - \alpha \frac{d}{dw} J(w)$$

smaller
not as large
large

Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate α



Not just the mean squared error cost function that we're using for

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 {

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

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

}

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 {

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

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

}

implements gradient descent

this way, it will work.

(Optional)

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

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

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

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

no $x^{(i)}$

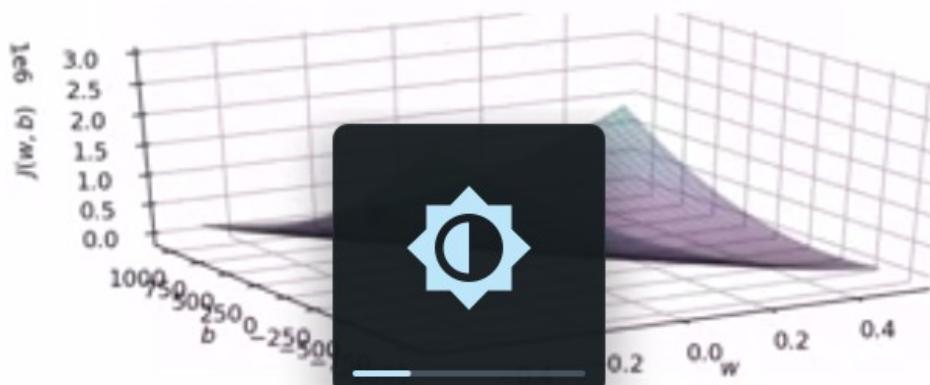
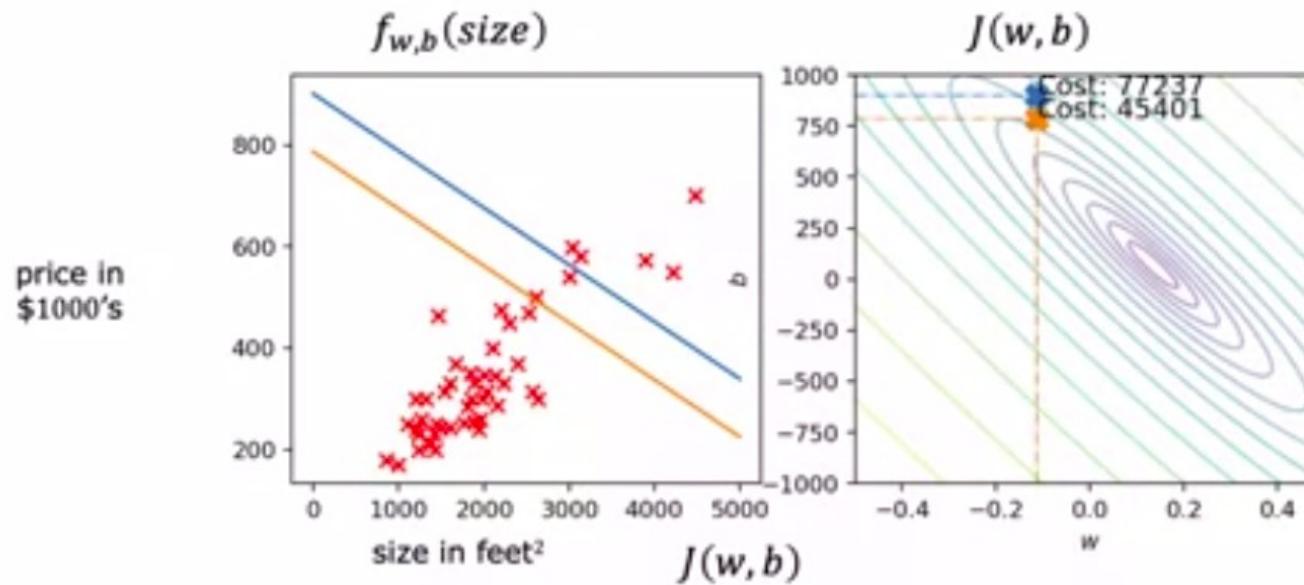
Gradient descent algorithm

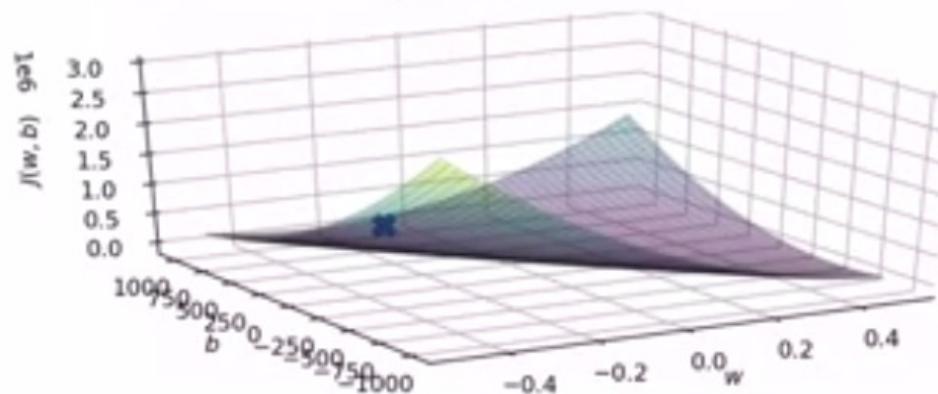
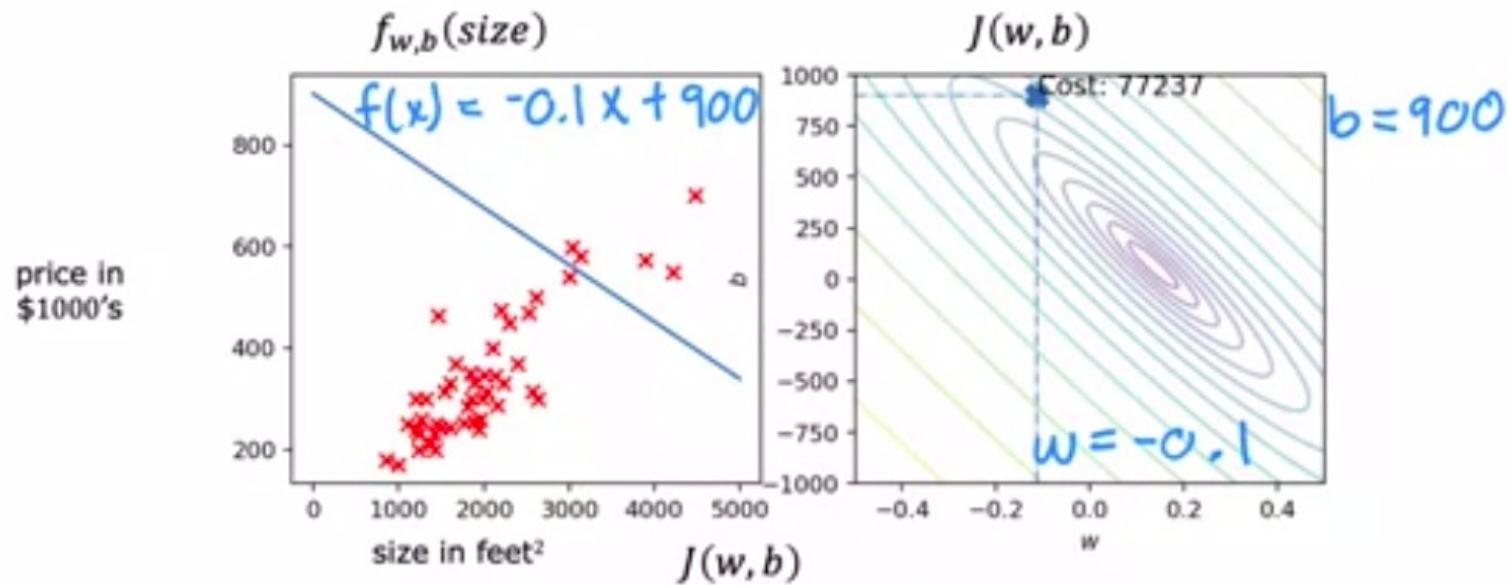
repeat until convergence {
 $w = w - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \right]$
 $b = b - \alpha \left[\frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) \right]$
}

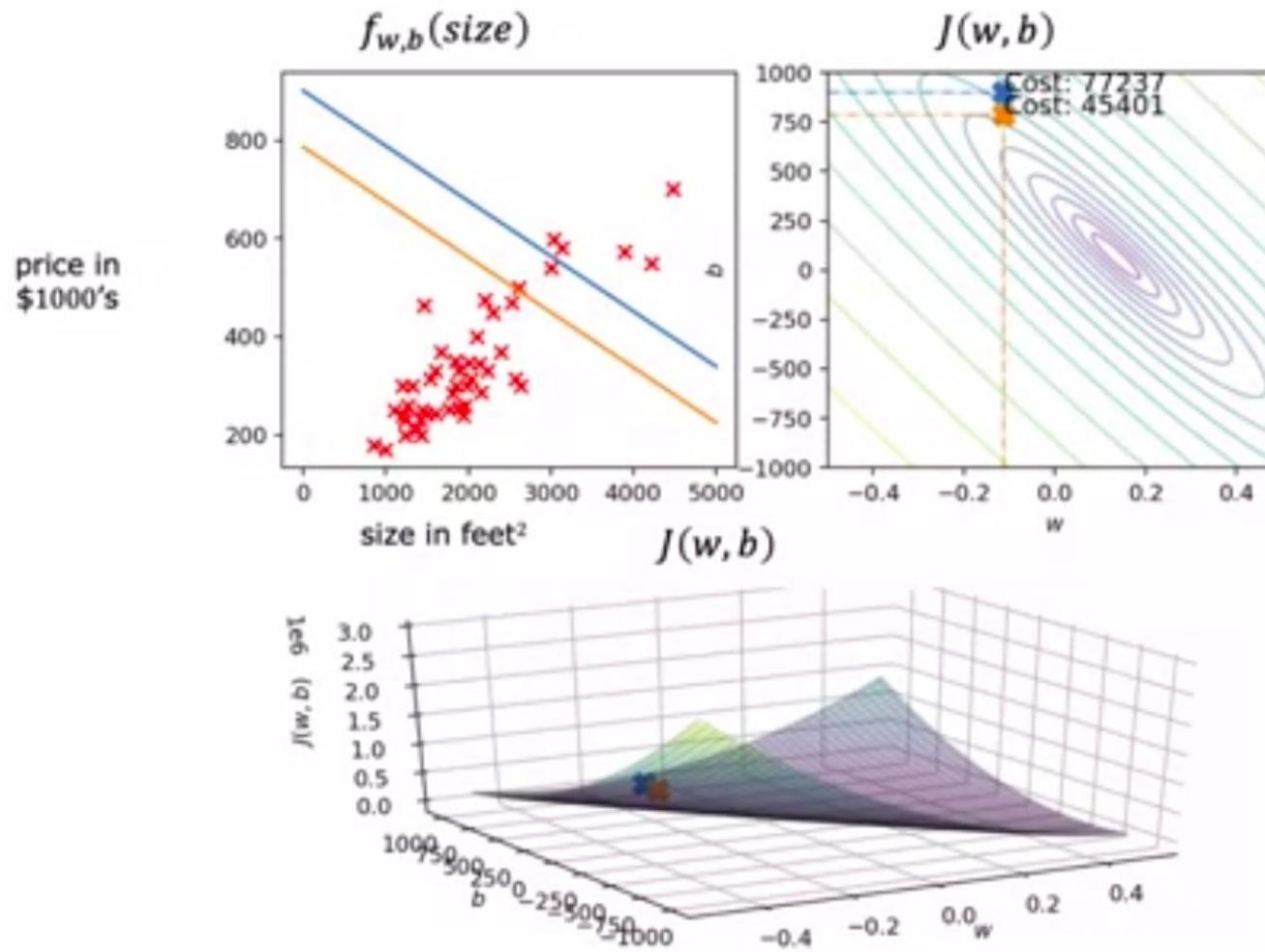
$\frac{\partial}{\partial w} J(w, b)$ $\frac{\partial}{\partial b} J(w, b)$

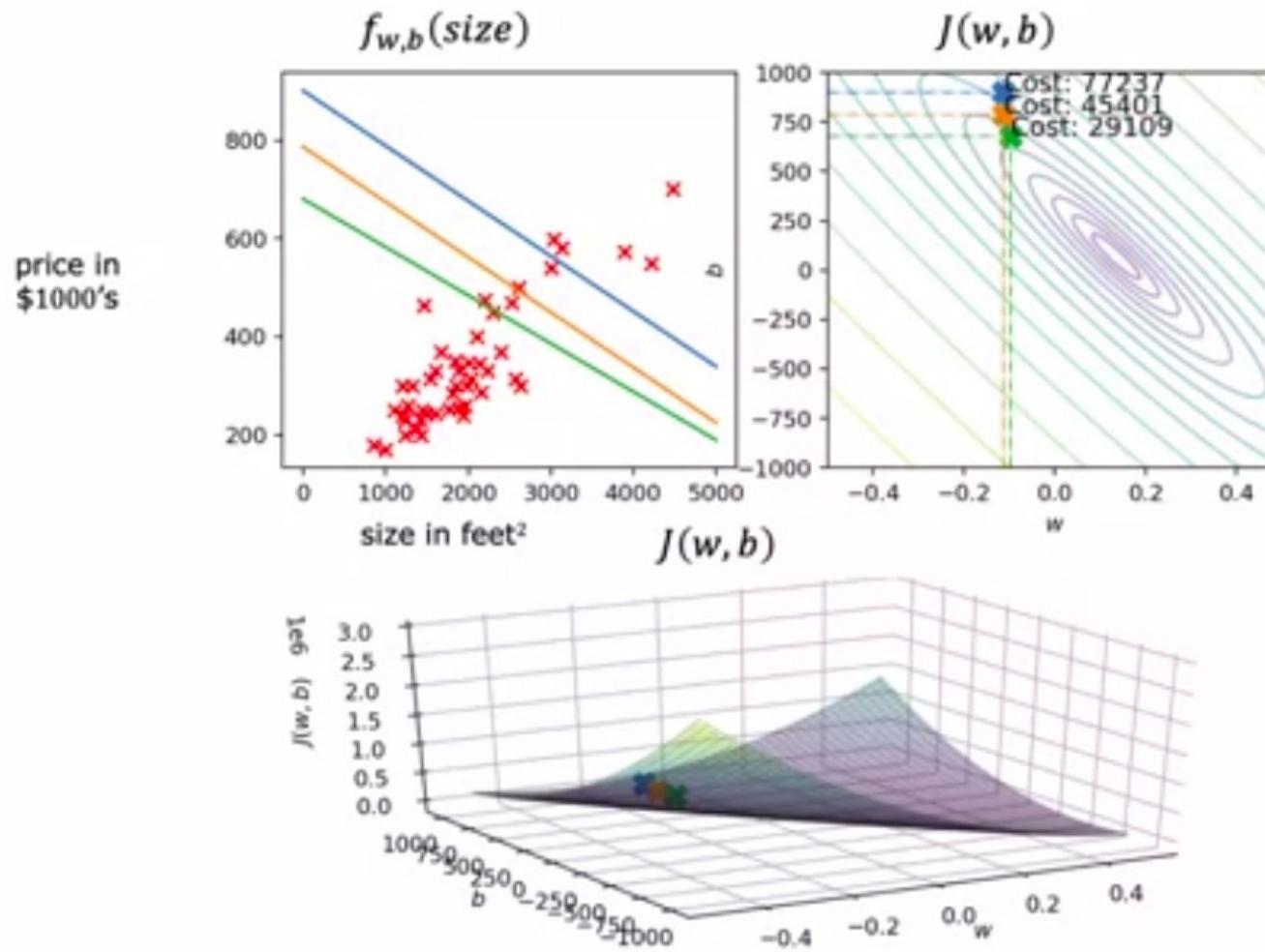
Update w and b simultaneously

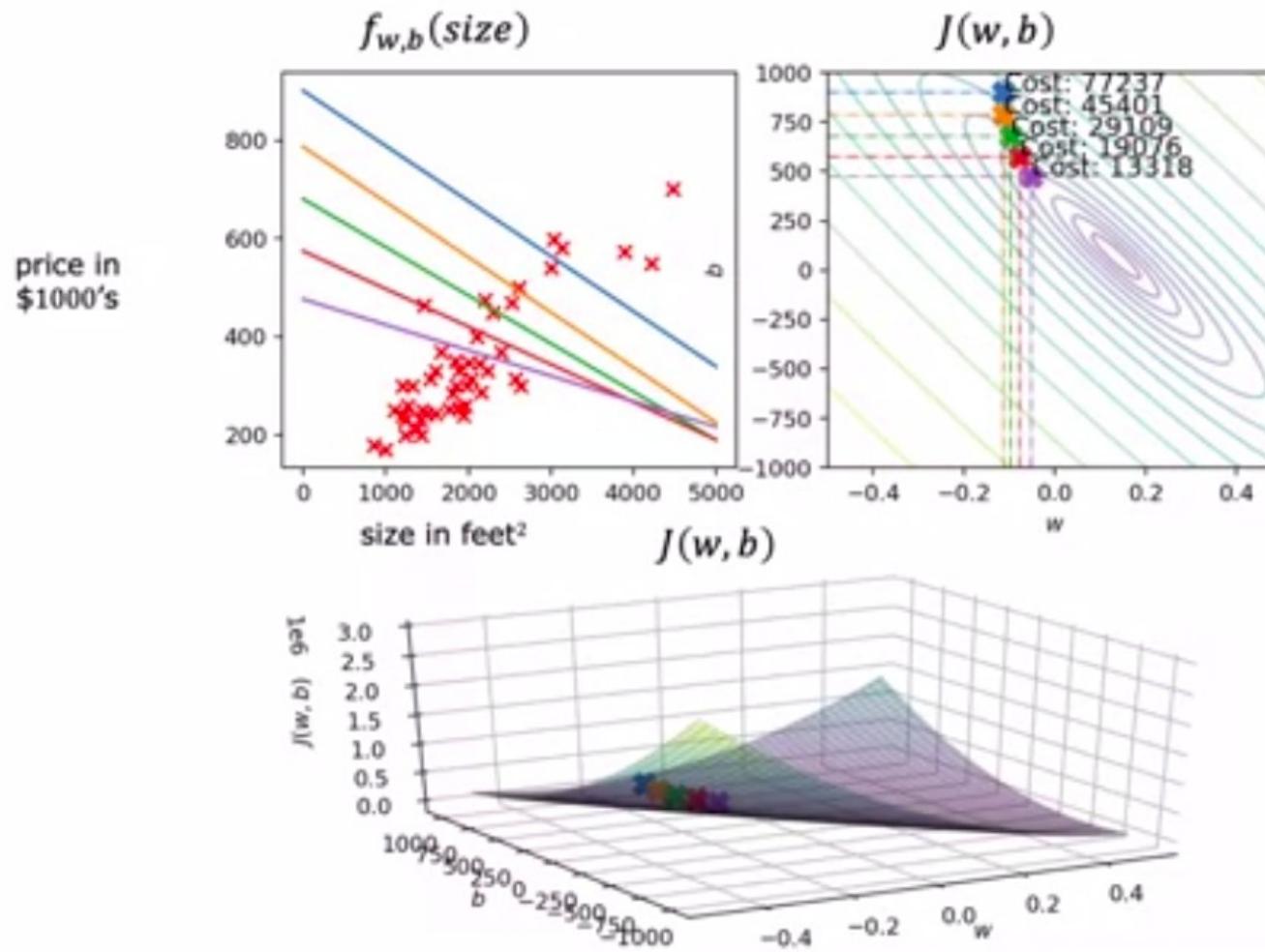
$f_{w,b}(x^{(i)}) = w x^{(i)} + b$

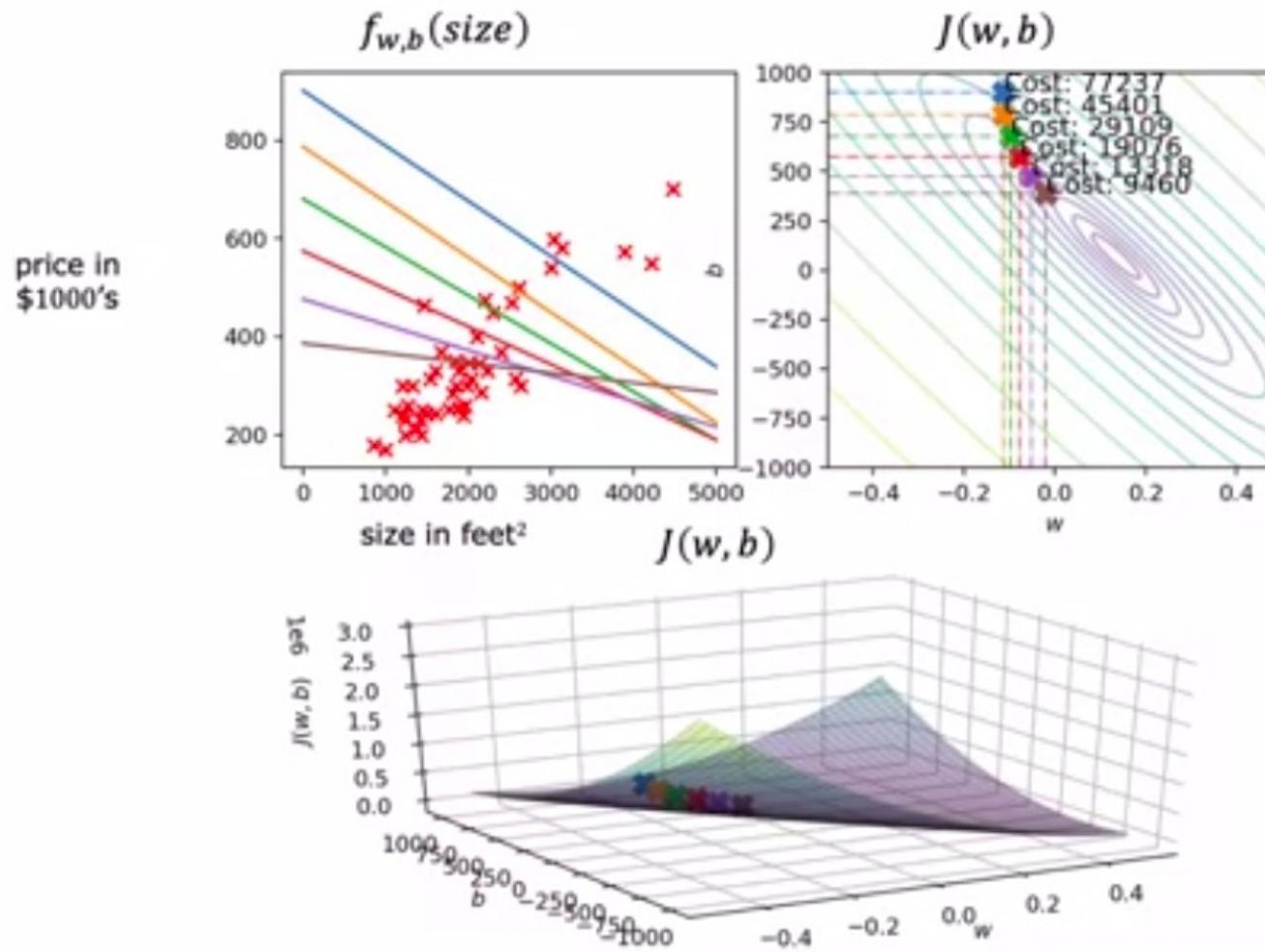


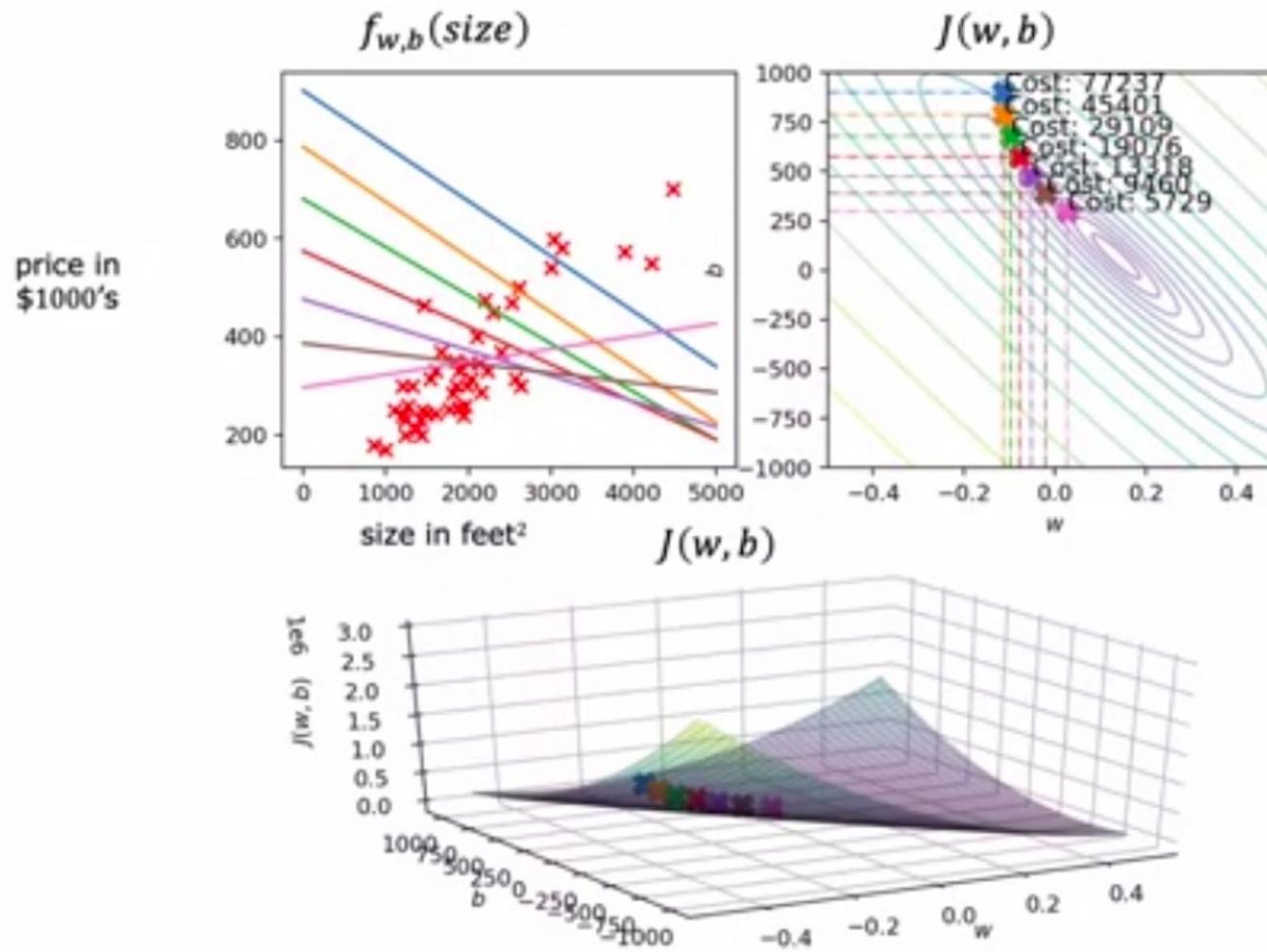


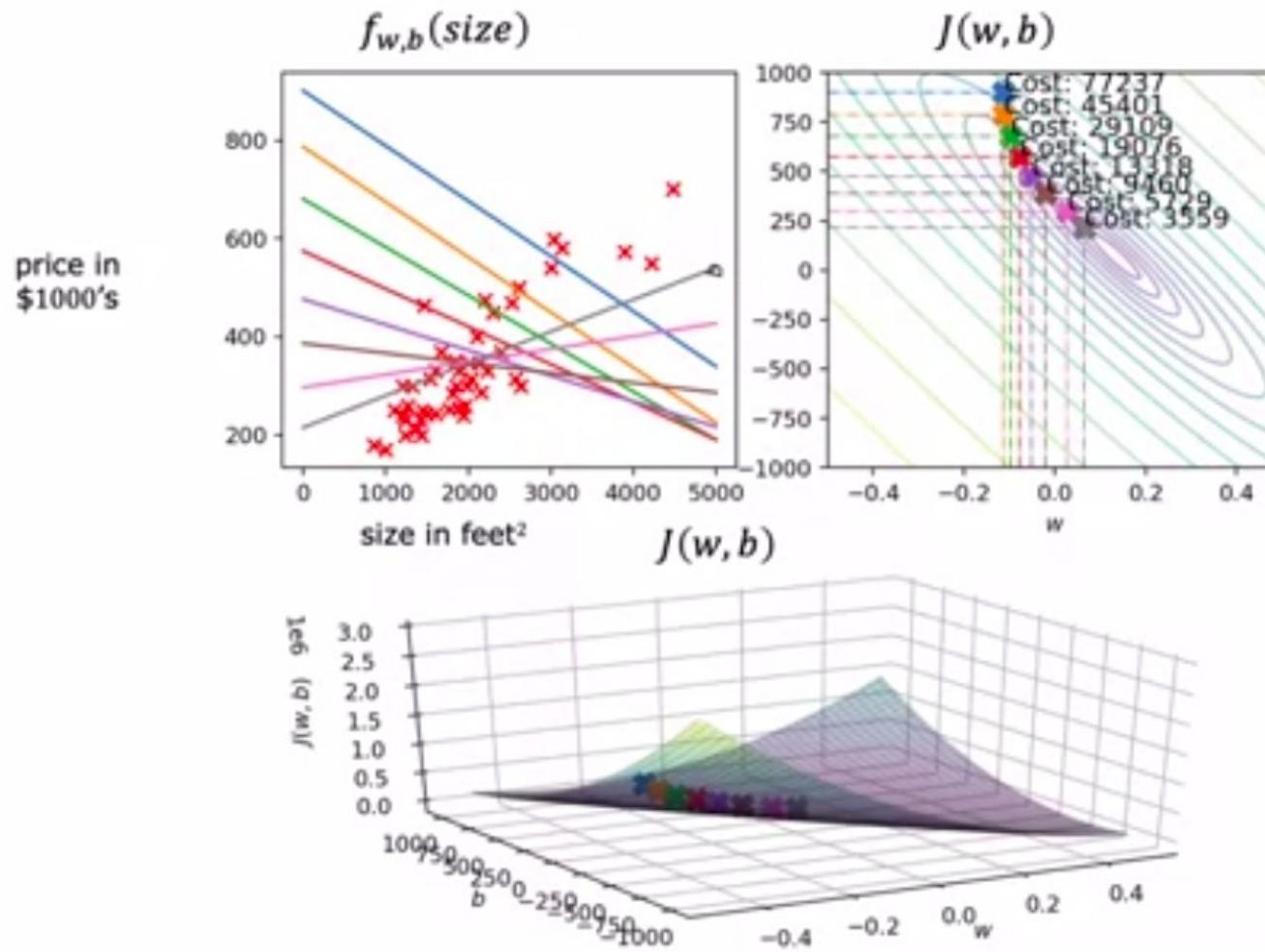


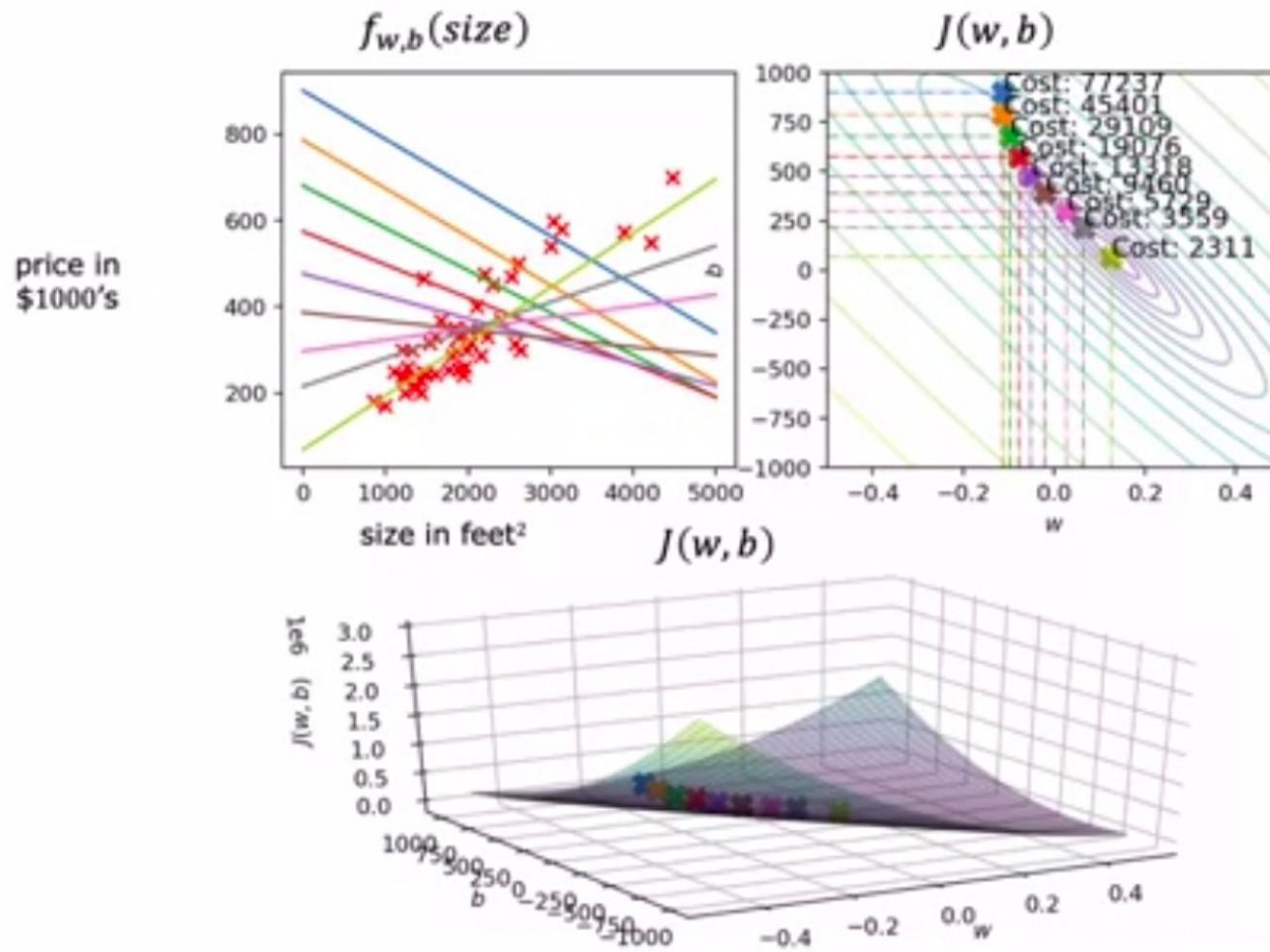








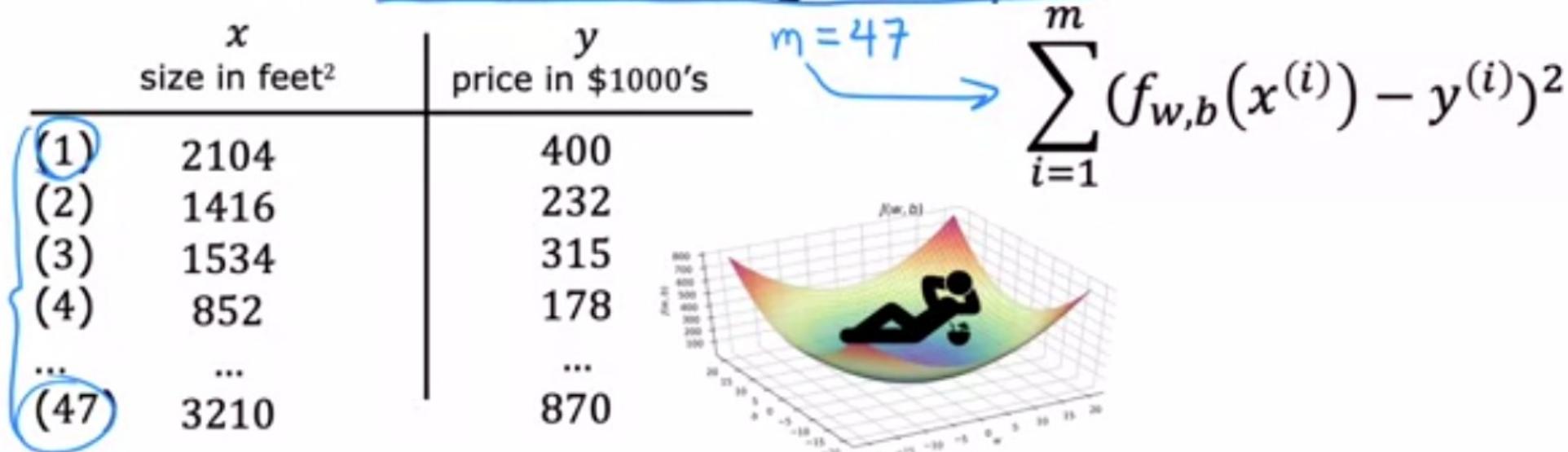




"Batch" gradient descent

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

other gradient descent: subsets



[Back](#)

Practice quiz: Train the model with gradient descent

Due Sep 4, 11:59 PM IST

Graded Quiz • 10 min

1.

1 / 1 point

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

[Back](#)

Practice quiz: Train the model with gradient descent

Due Sep 4, 11:59 PM IST

Graded Quiz • 10 min

- w stays the same
- 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.

1 / 1 point

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

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

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