

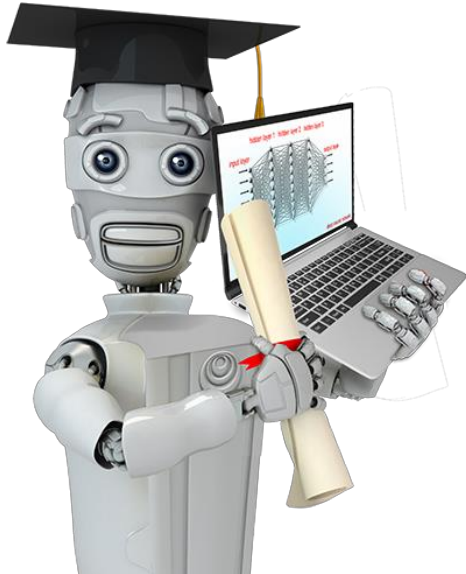
Copyright Notice

These slides are distributed under the Creative Commons License.

[DeepLearning.AI](#) makes these slides available for educational purposes. You may not use or distribute these slides for commercial purposes. You may make copies of these slides and use or distribute them for educational purposes as long as you cite [DeepLearning.AI](#) as the source of the slides.

For the rest of the details of the license, see <https://creativecommons.org/licenses/by-sa/2.0/legalcode>

Machine Learning Overview



Supervised Learning Part 1

Supervised learning



input

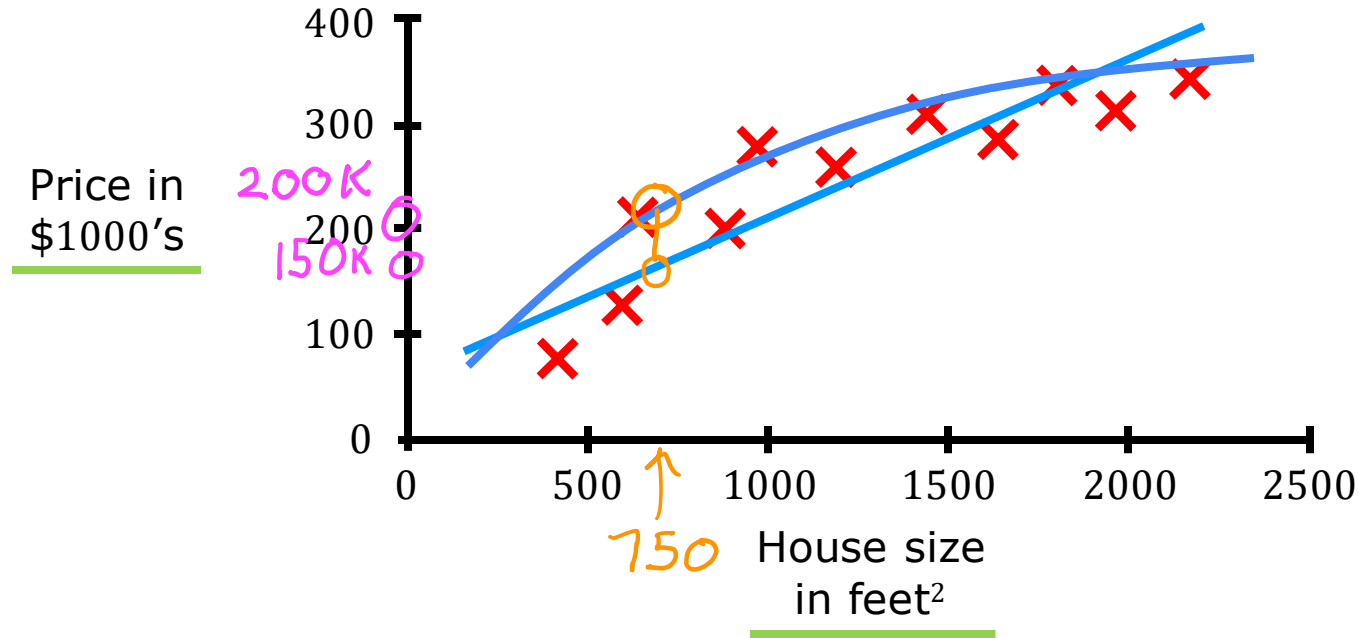


output label

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

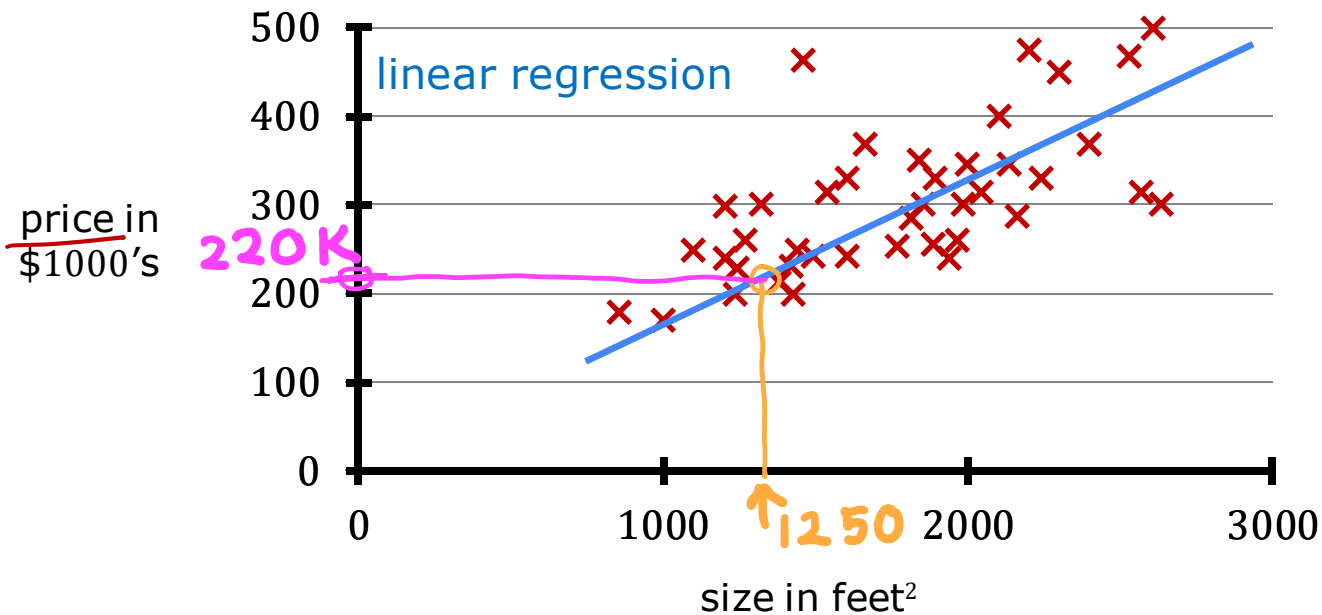
Regression: Housing price prediction



Linear Regression with One Variable

Linear Regression Model Part 1

House sizes and prices



Regression model

Predicts numbers

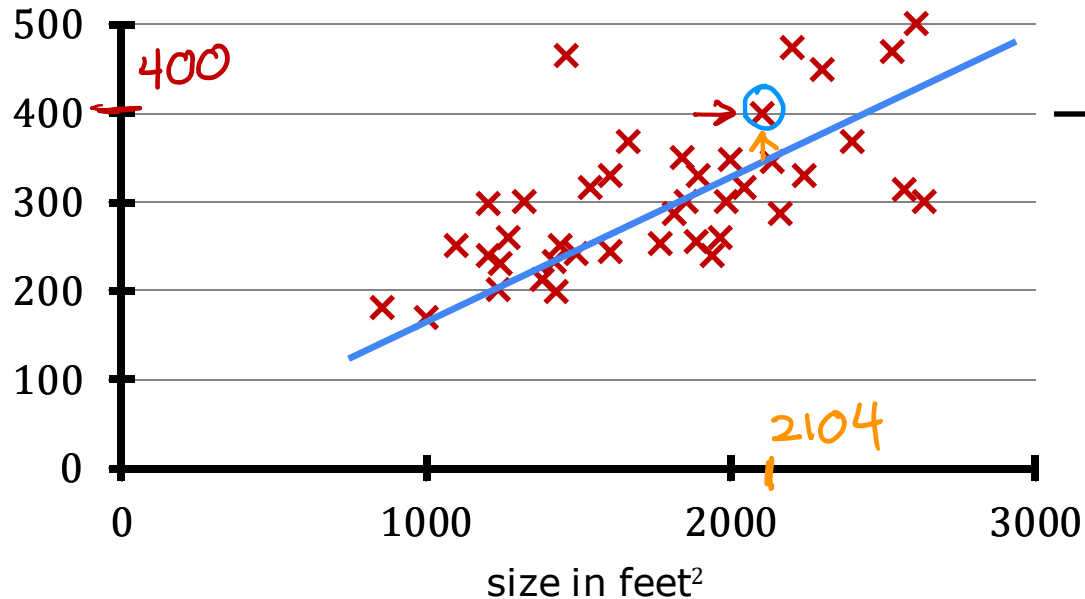
Infinitely many possible outputs

Supervised learning model

Data has "right answers"

House sizes and prices

price in \$1000's



Data table

size in feet ²	price in \$1000's
2104	400
1416	232
1534	315
852	178
...	...
3210	870

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

$m = 47$

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

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

Notation:

x = "input" variable
feature

y = "output" variable
"target" variable

m = number of training examples

(x, y) = single training example

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

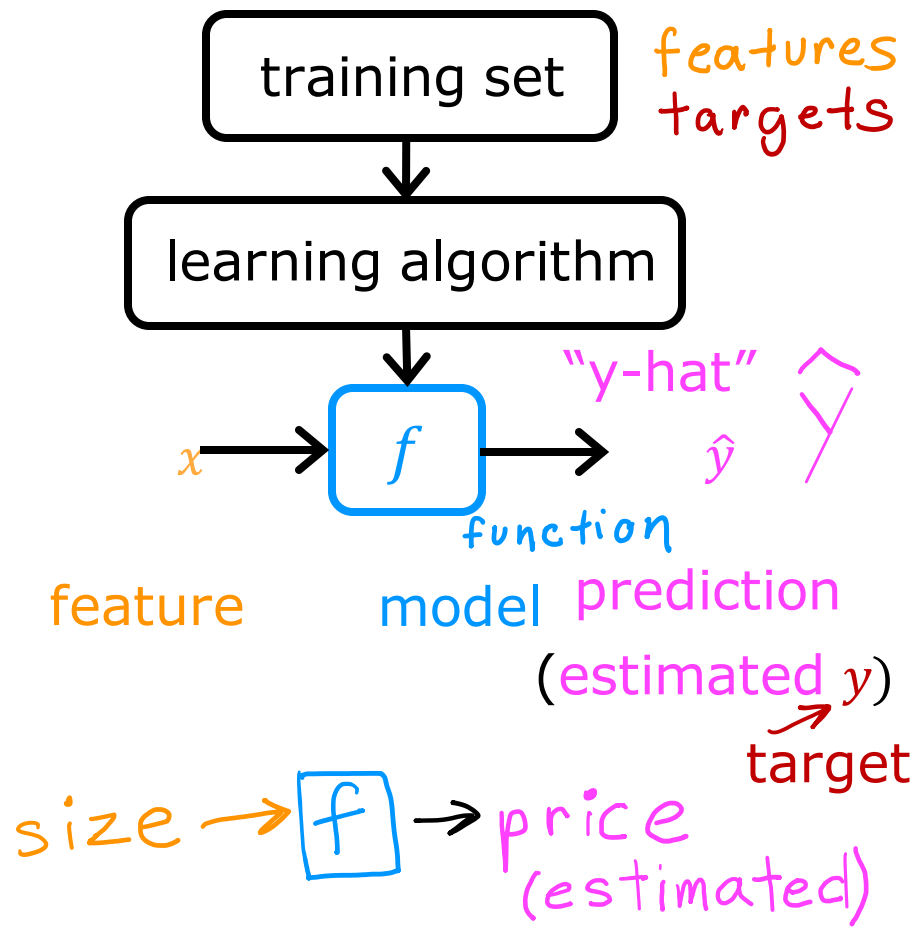
$(x^{(i)}, y^{(i)})$ = i^{th} training example

index

(1st, 2nd, 3rd ...)

Linear Regression with One Variable

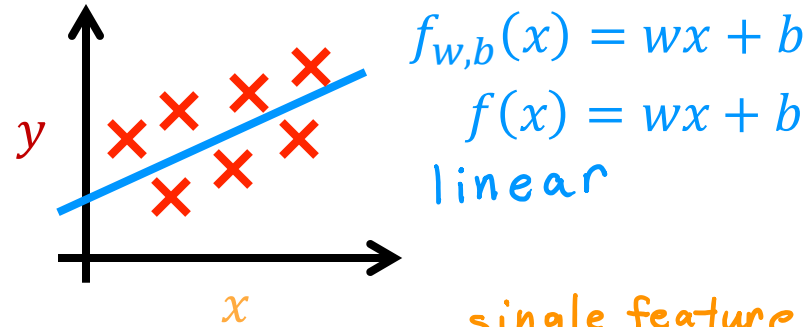
Linear Regression Model Part 2



How to represent f ?

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

$$f(x)$$



Linear regression with **one** variable.
size

Univariate linear regression.
one variable

Linear Regression with One Variable

Cost Function

Training set

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

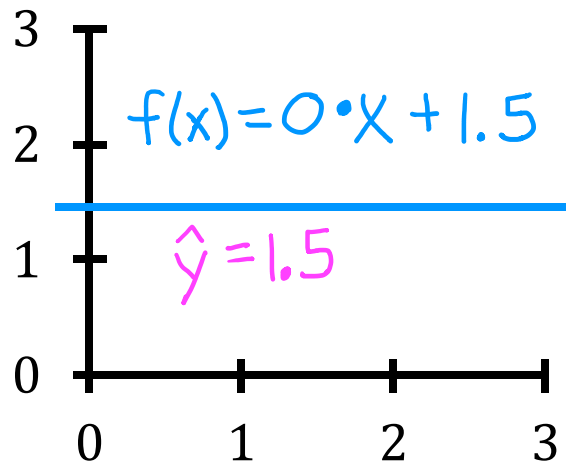
Model: $f_{w,b}(x) = wx + b$

w, b : parameters
coefficients
weights

What do w, b do?

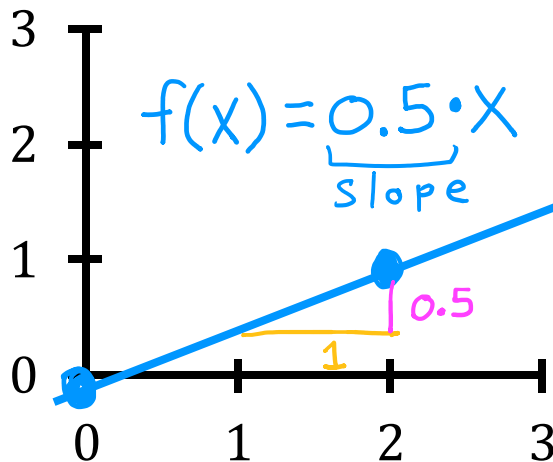
$$f_{w,b}(x) = wx + b$$

$$f(x)$$



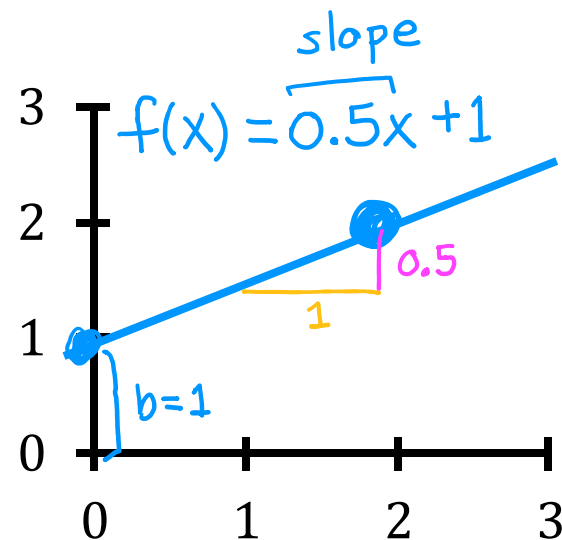
$$w = 0$$

$$b = 1.5$$



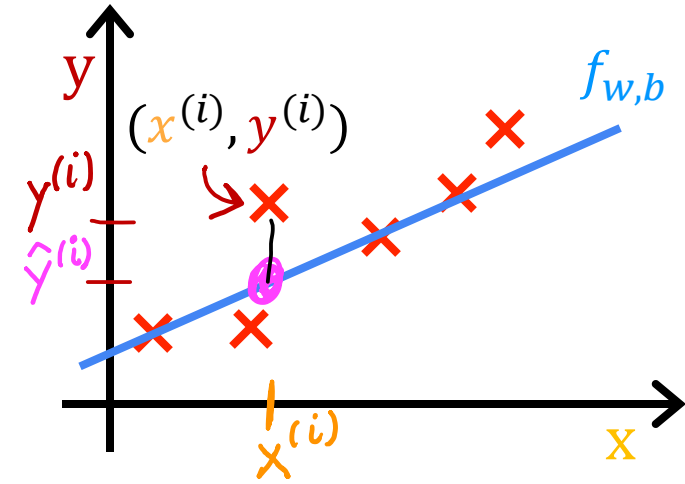
$$w = 0.5$$

$$b = 0$$



$$w = 0.5$$

$$b = 1$$



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

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

Maliyet
fonksiyonu

Kare hatası maliyet
fonksiyonu

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

error

m = eğitim örnekleri sayısı

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

Find w, b :

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

Linear Regression with One Variable

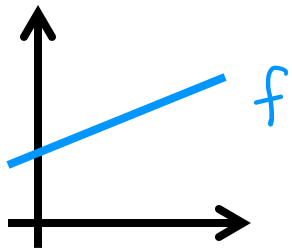
Cost Function
Intuition

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:

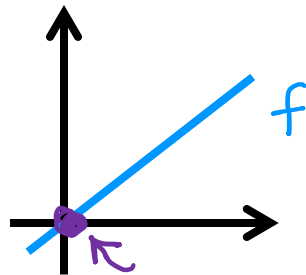
$$\underset{w, b}{\text{minimize}} J(w, b)$$

simplified

$$f_w(x) = wx$$

$$b = \emptyset$$

$$w$$



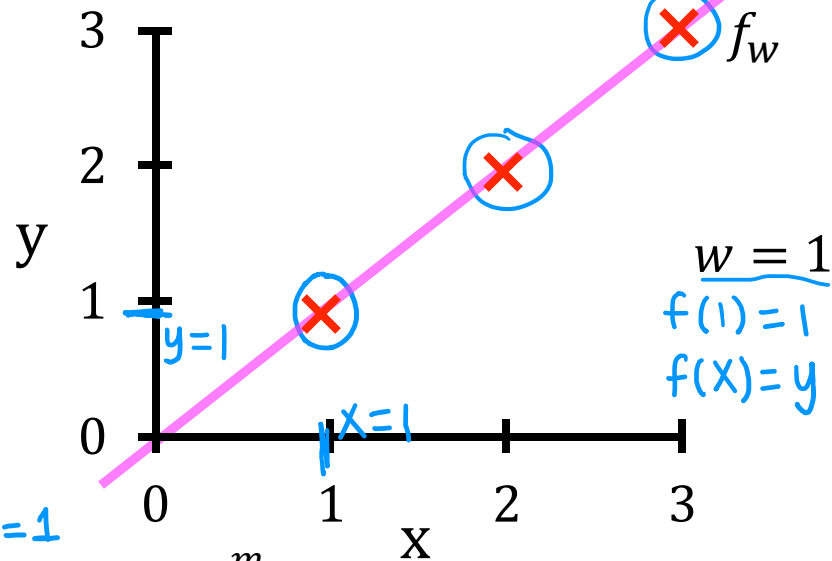
$$J(w) = \frac{1}{2m} \sum_{i=1}^m \underbrace{(f_w(x^{(i)}) - y^{(i)})^2}_{wx^{(i)}}$$

$$wx^{(i)}$$

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

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

input

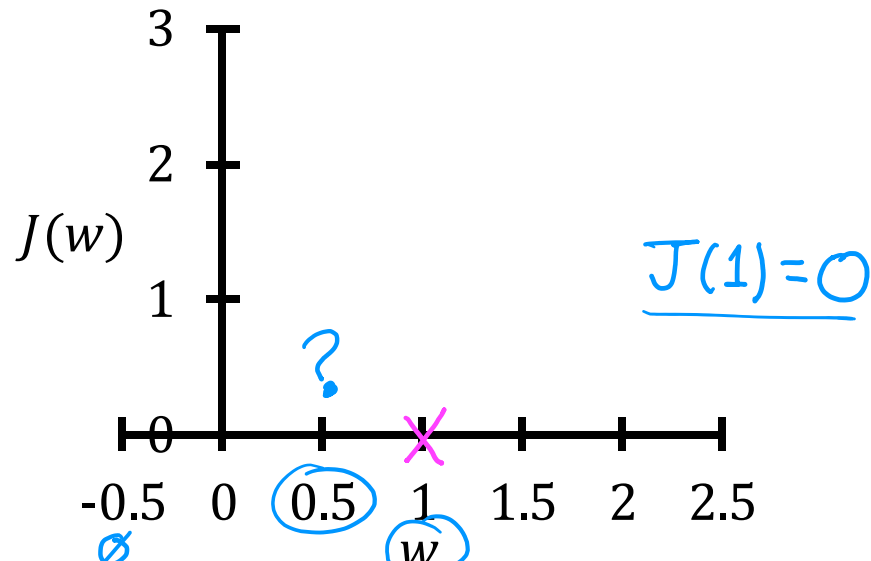


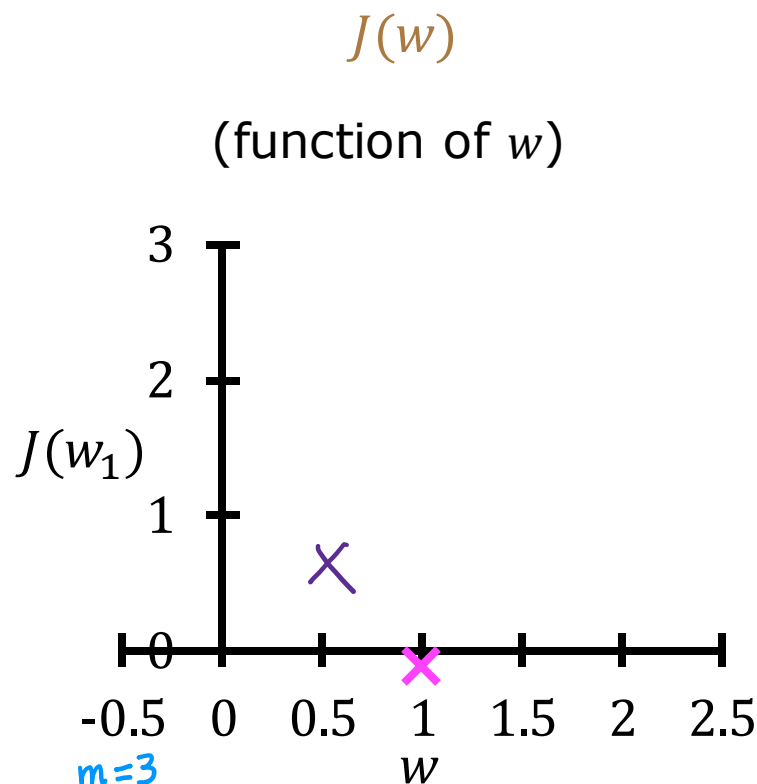
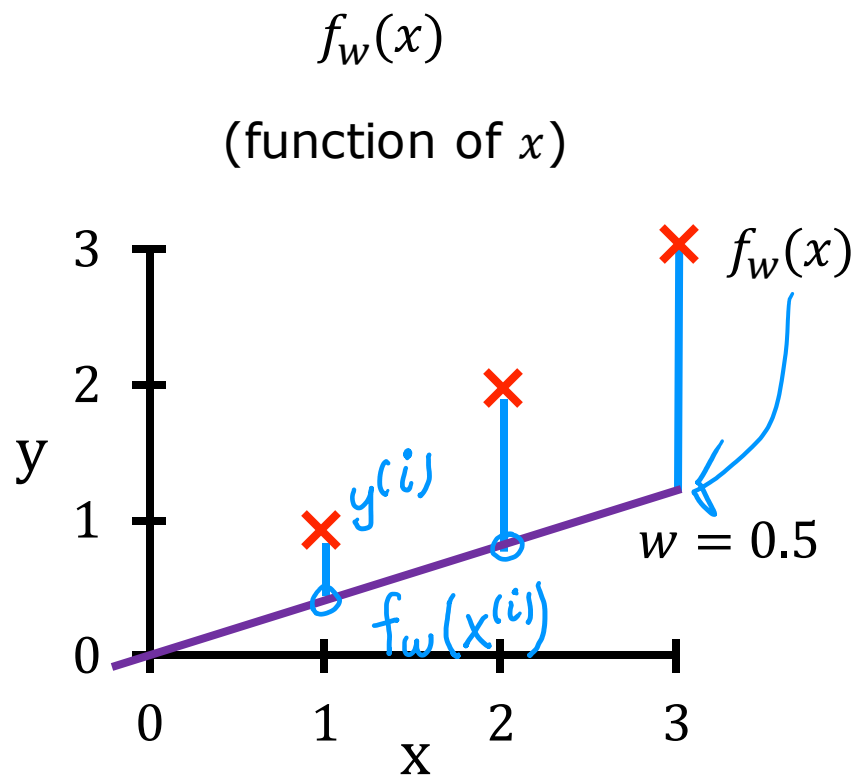
$$J(w) = \frac{1}{2m} \sum_{i=1}^m (f_w(x^{(i)}) - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m (wx^{(i)} - y^{(i)})^2 = \frac{1}{2m} (0^2 + 0^2 + 0^2) = 0$$

Handwritten blue annotations: $w=1 \downarrow$ (pointing to $J(w)$), $w x^{(i)}$ (under $f_w(x^{(i)})$), \emptyset (above the second sum), and w (circled around the w in the final term).

 $J(w)$ (function of w)

parameter

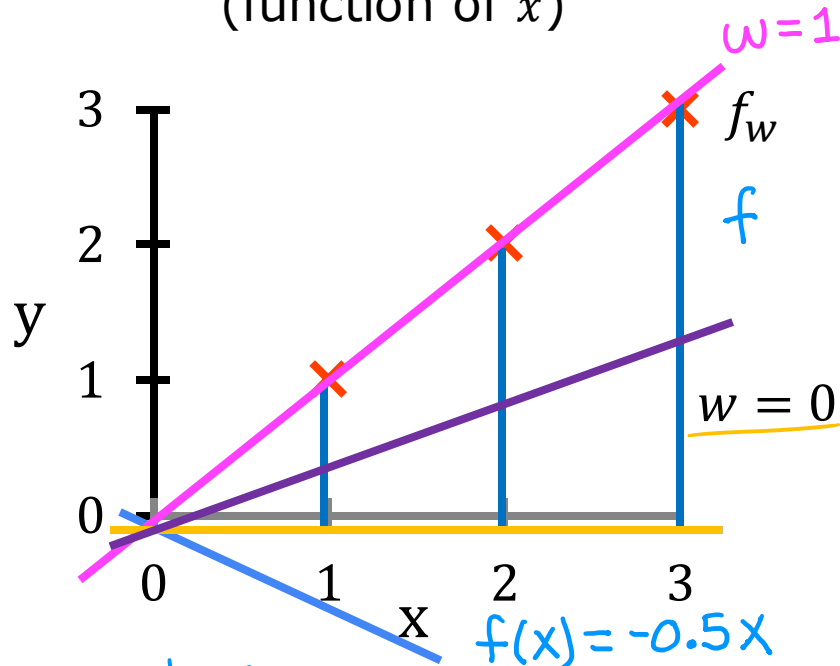




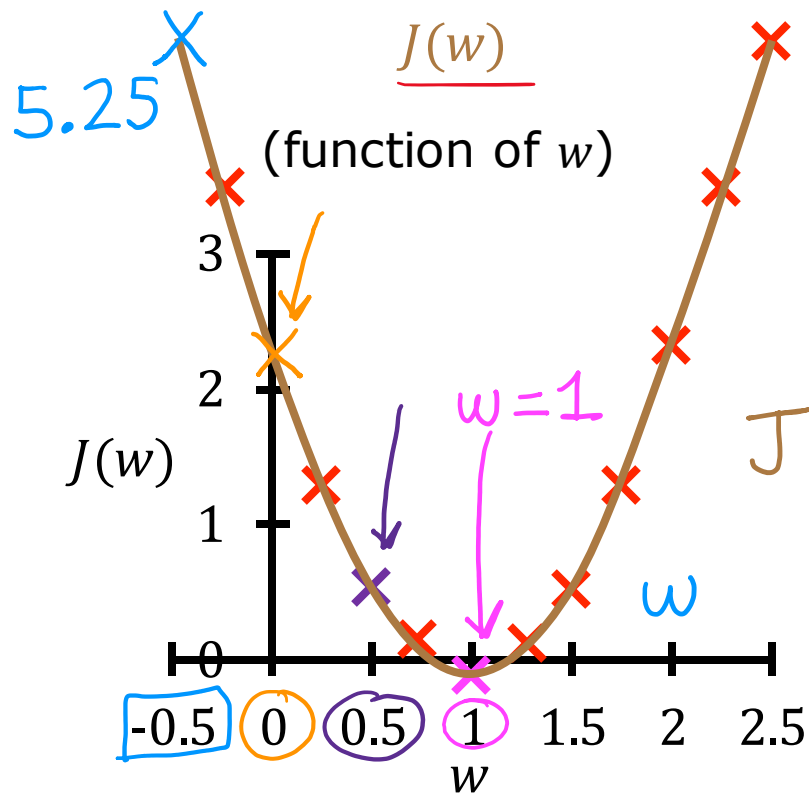
$$J(0.5) = \frac{1}{\underline{2m}} \left[(0.5-1)^2 + (1-2)^2 + (1.5-3)^2 \right] = \frac{1}{\underline{2 \times 3}} [3.5] = \frac{3.5}{6} \approx 0.58$$

$m=3$

$f_w(x)$
(function of x)



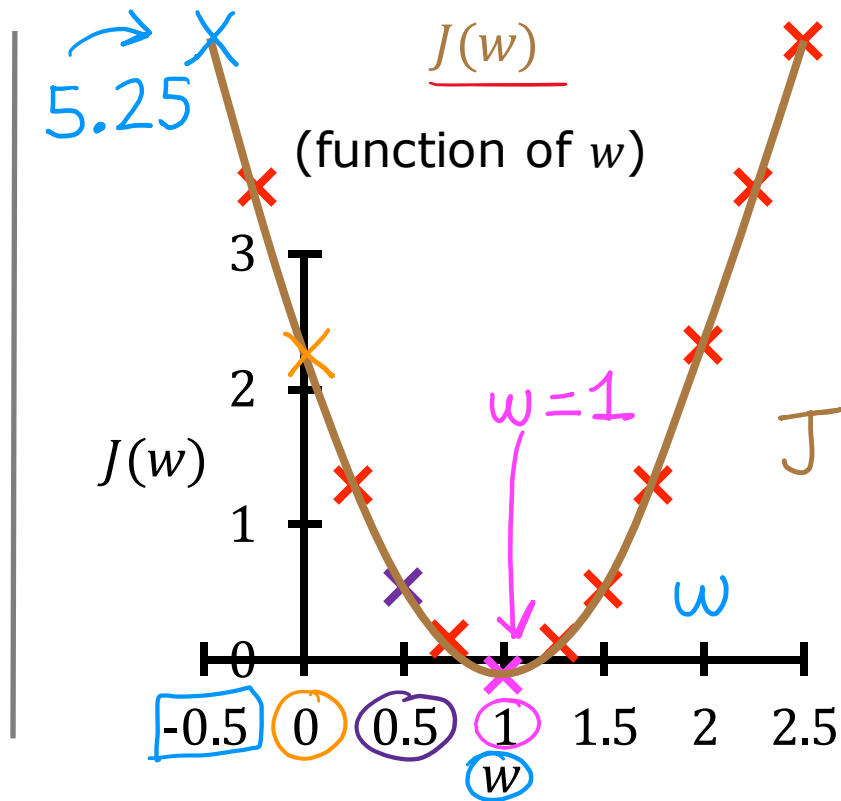
$$J(0) = \frac{1}{2m} (1^2 + 2^2 + 3^2) = \frac{1}{6} [14] \approx 2.3$$



how to choose w ?

goal of linear regression:

minimize $J(w)$



choose w to minimize $J(w)$

Linear Regression with One Variable

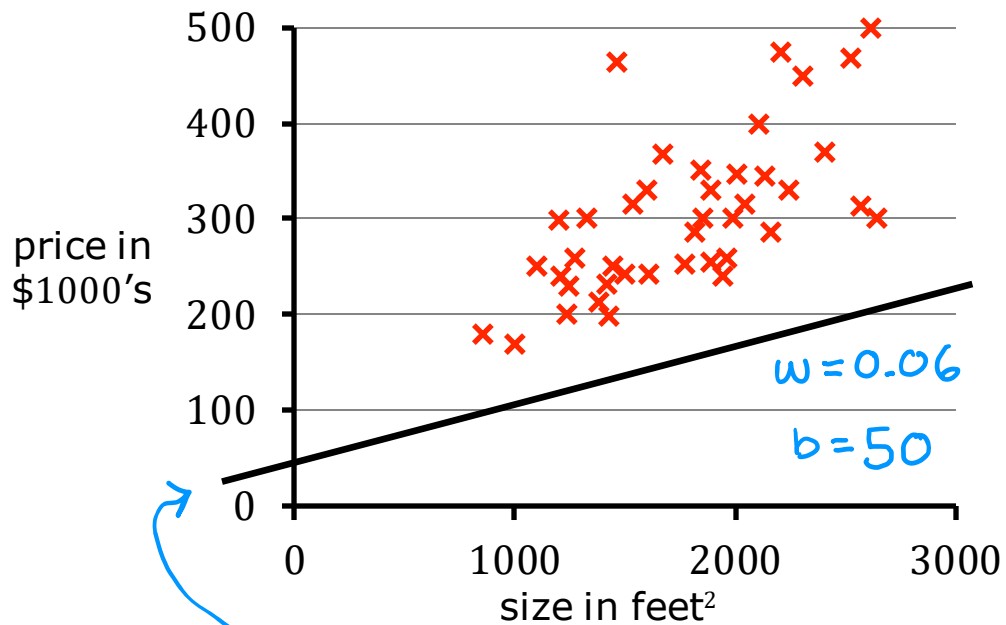
Visualizing
the Cost Function

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

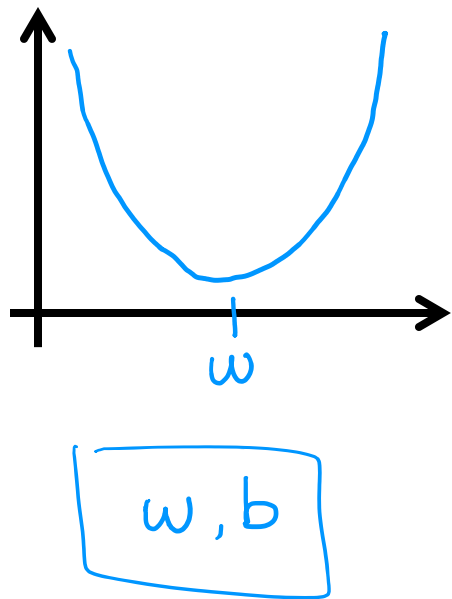
Parameters w, b

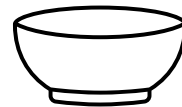
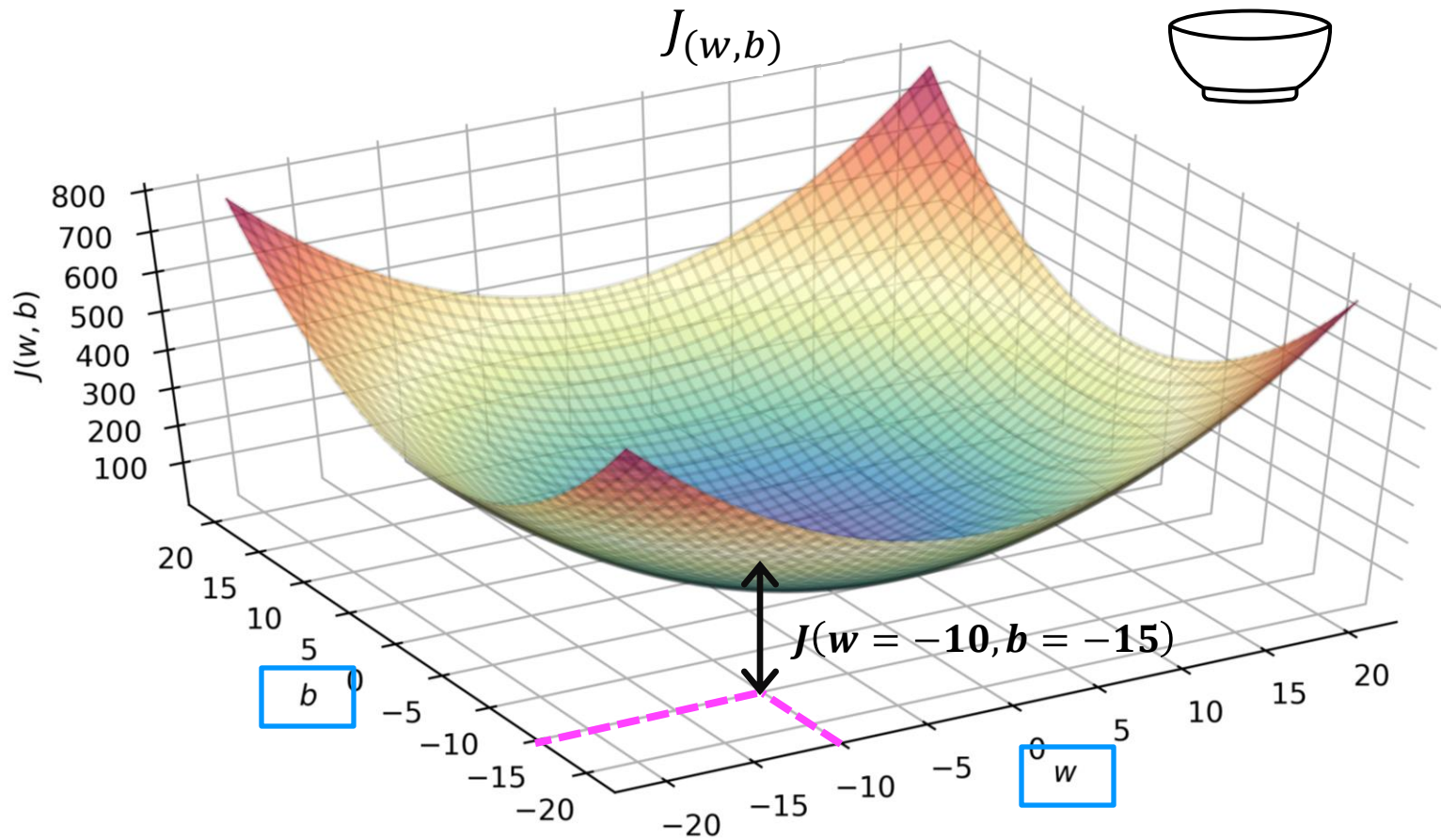
Cost Function $J(w, b) = \frac{1}{2m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)})^2$

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

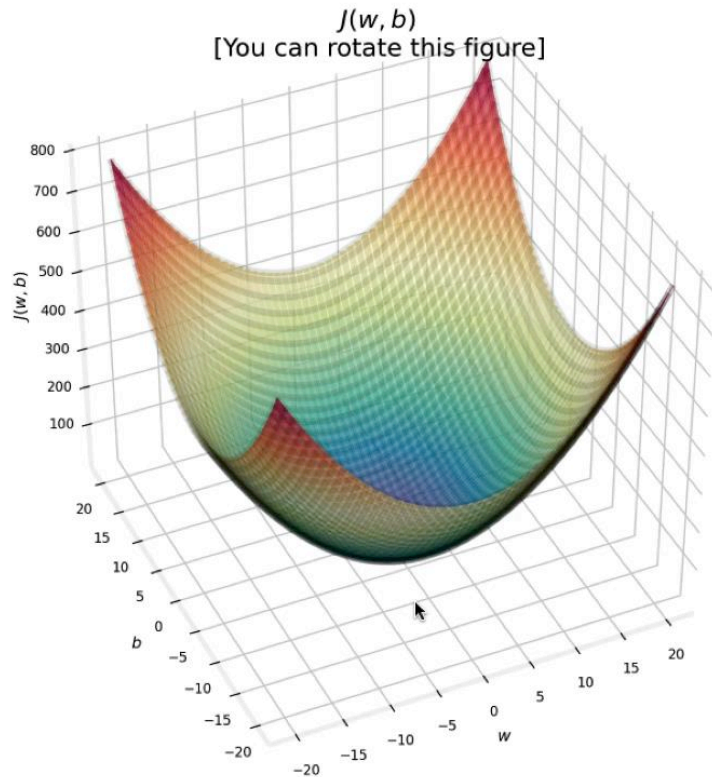
$f_{w,b}$ (function of x)

$$f_{w,b}(x) = 0.06x + 50$$

 J (function of w, b)



3D surface plot



Alternative
contour plot

Mount Fuji

Legend

- Contour
- Mount Fuji

Mount Fuji

Google Earth

Image Landsat / Copernicus

Data SIO, NOAA, U.S. Navy, NGA, GEBCO

Data Japan Hydrographic Association

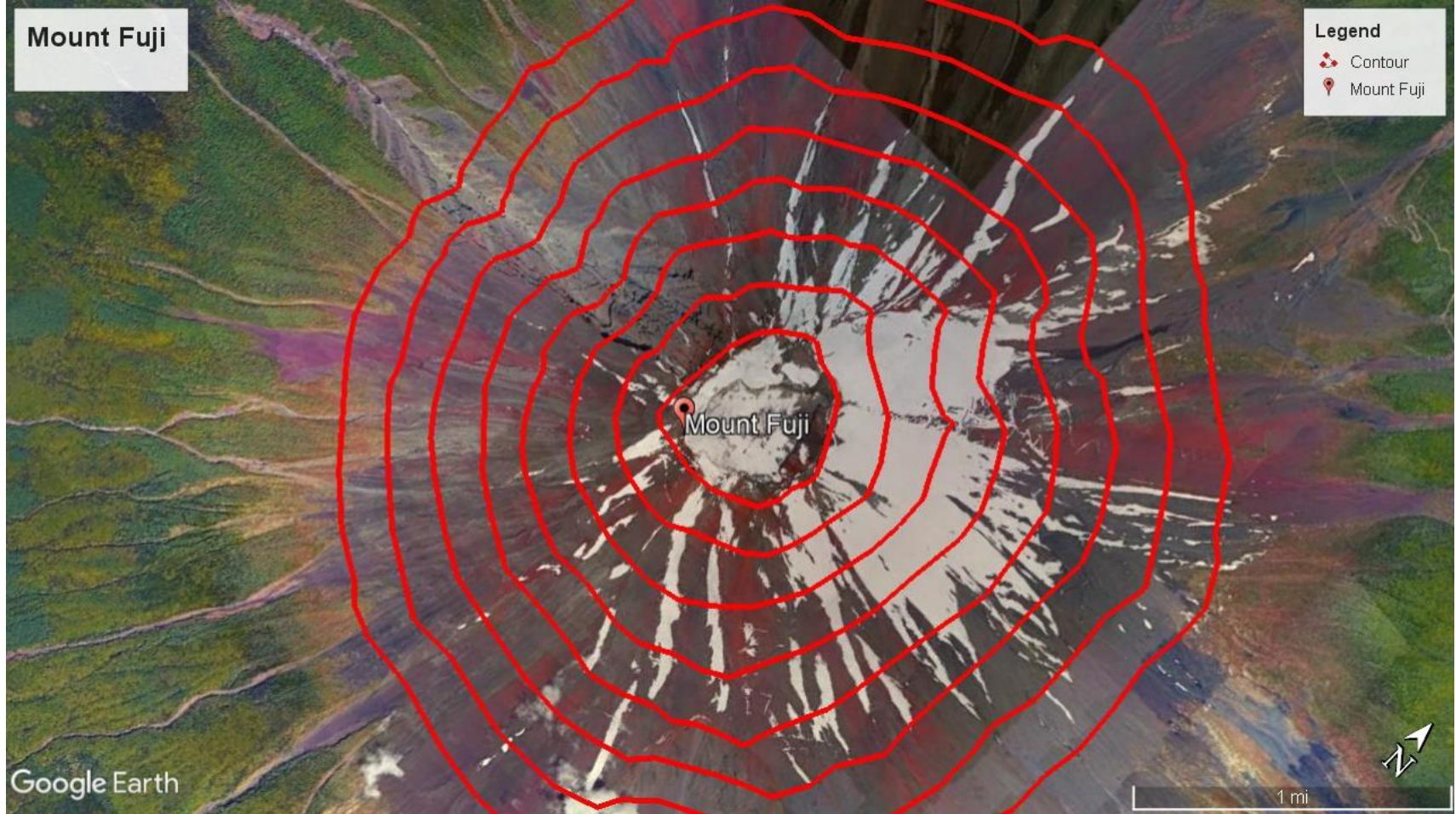
3000 ft

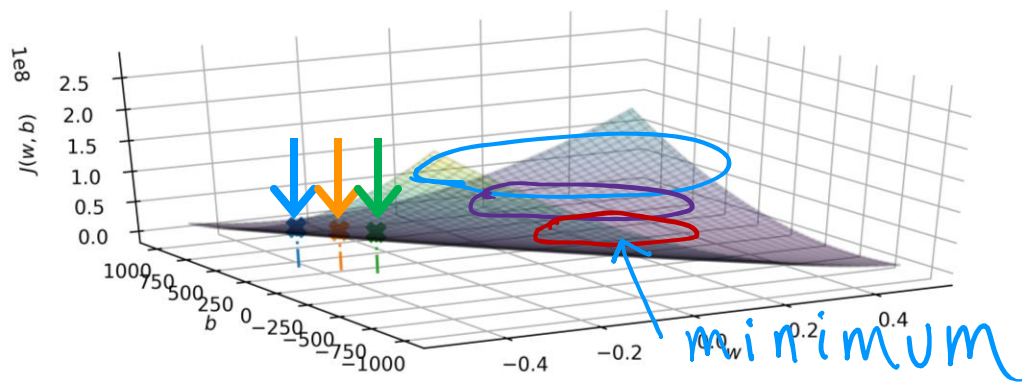
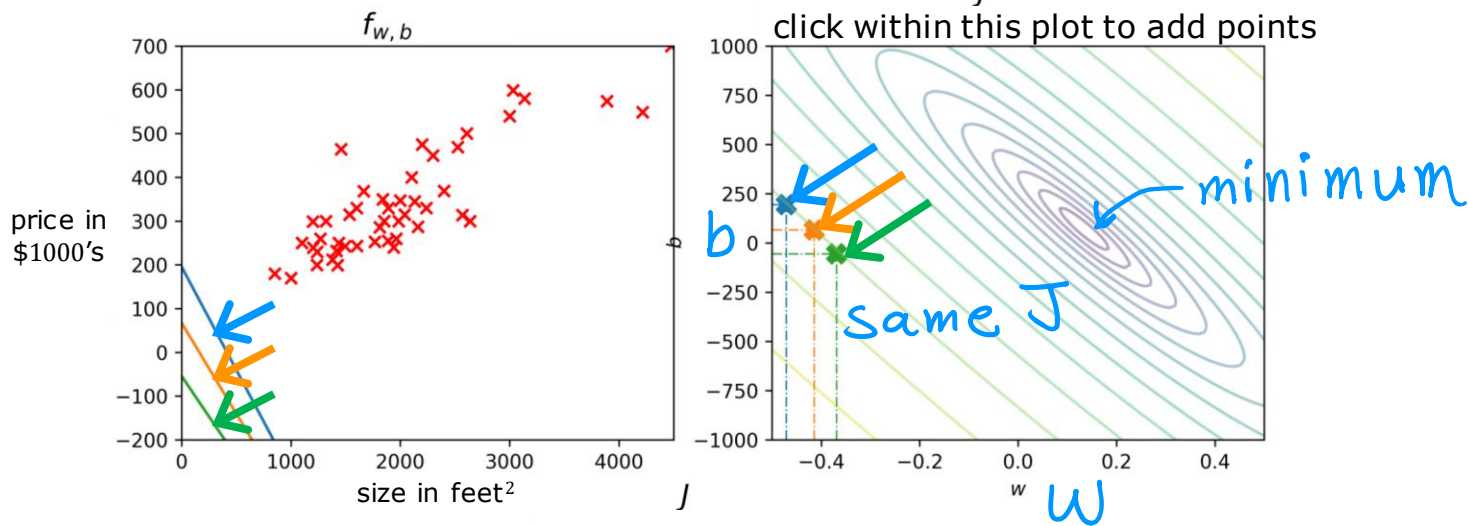


Mount Fuji

Legend

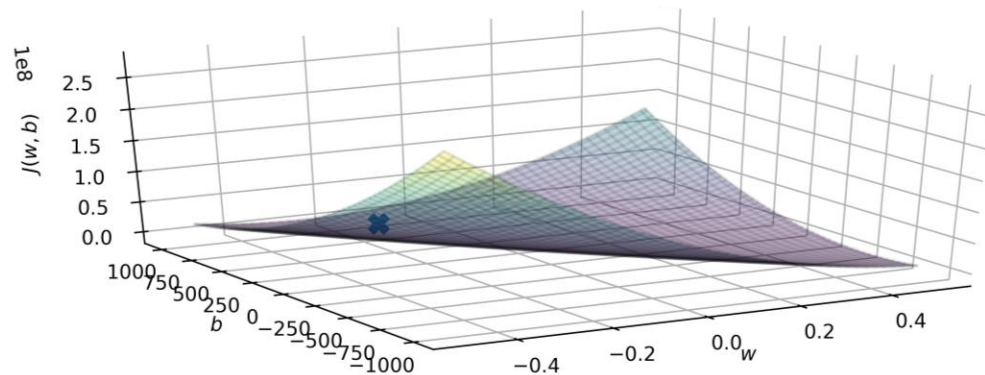
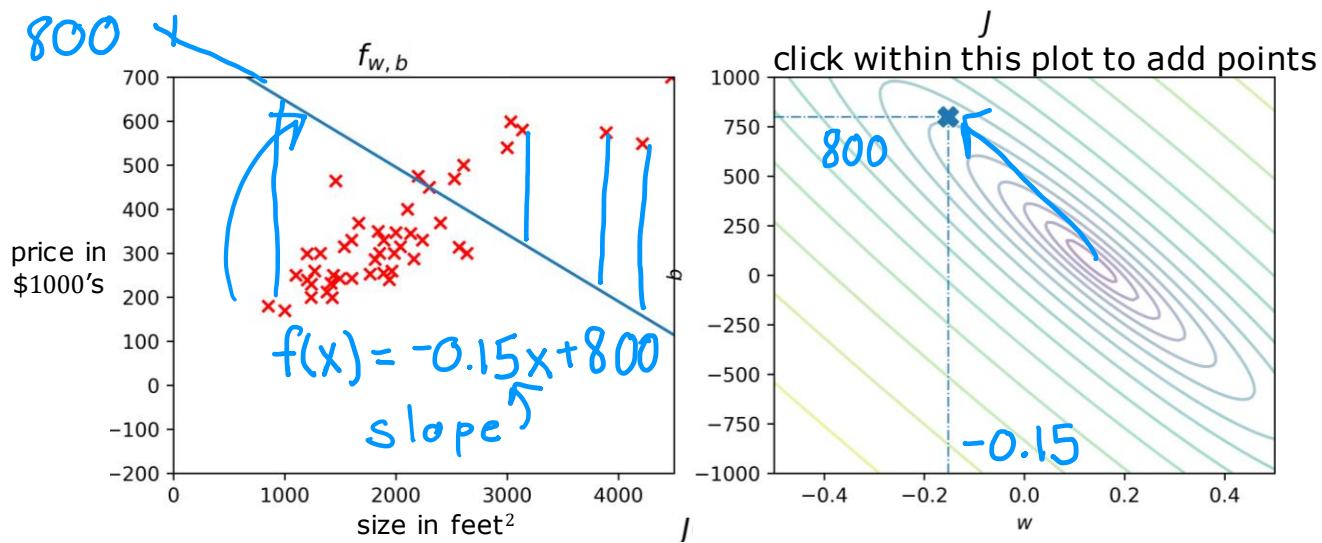
- Contour
- Mount Fuji

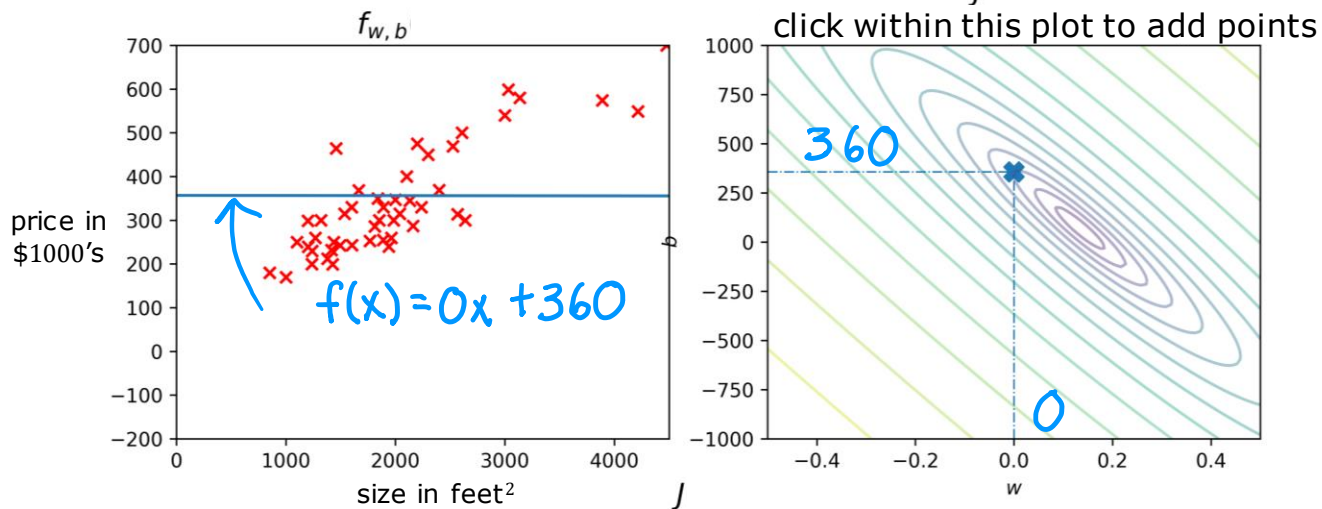




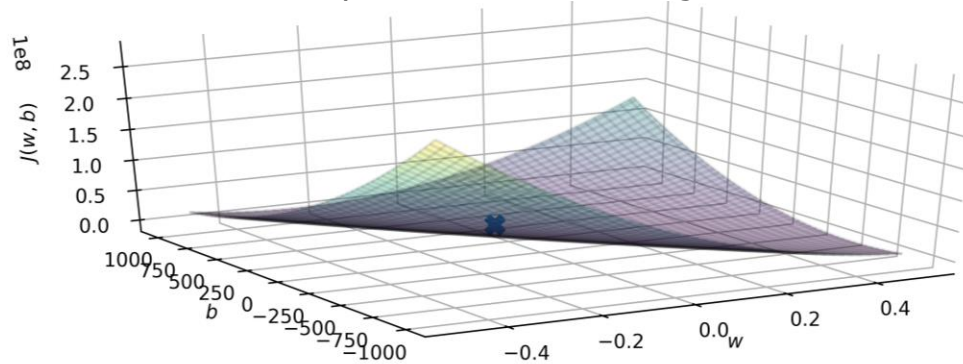
Linear Regression with One Variable

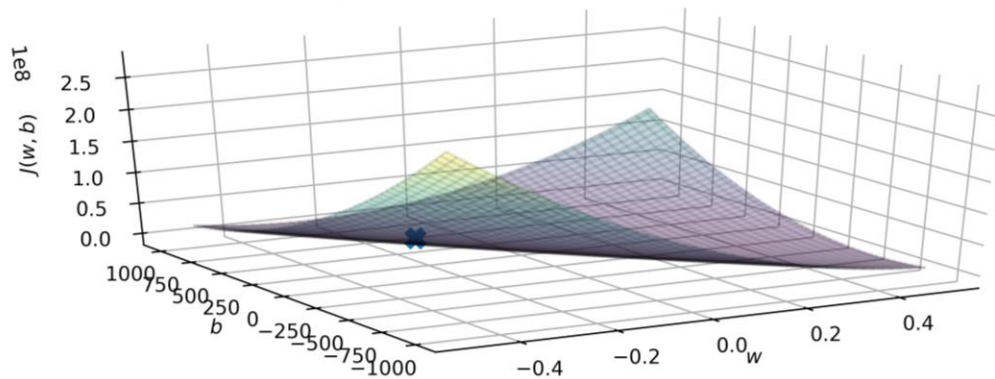
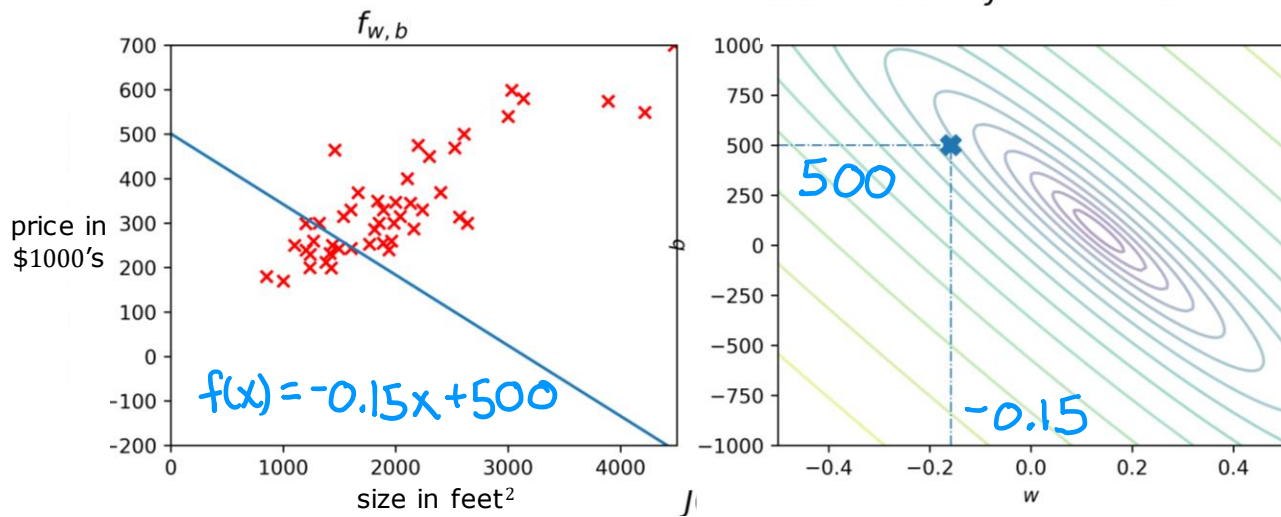
Visualization examples

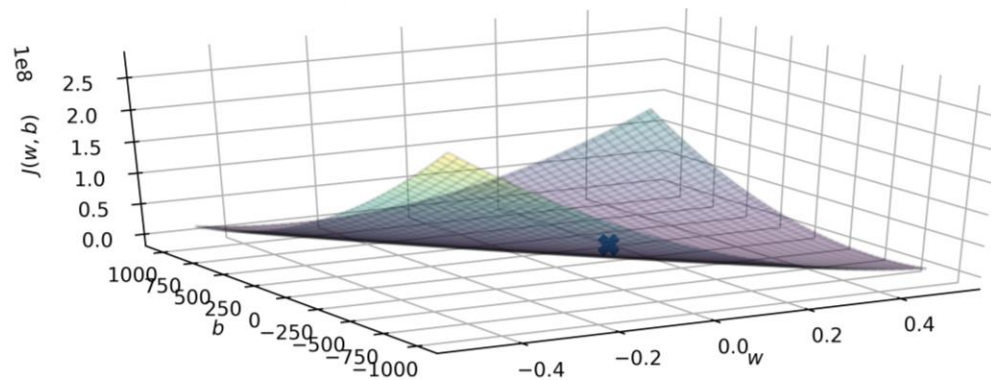
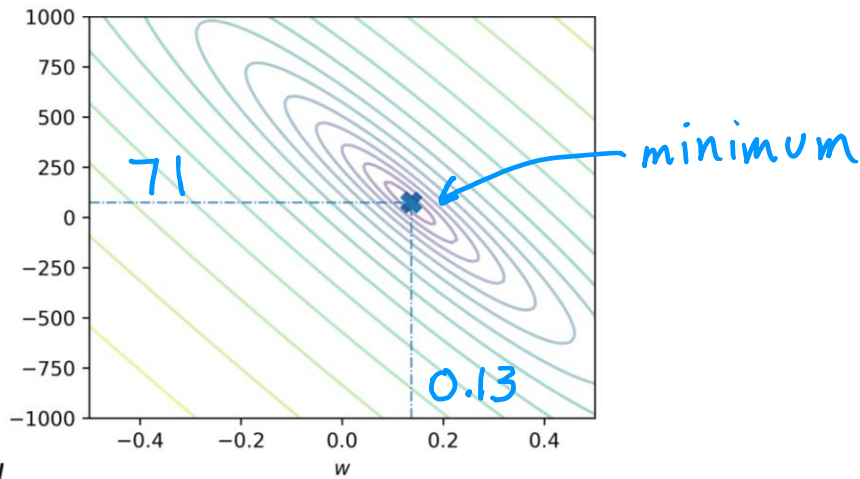
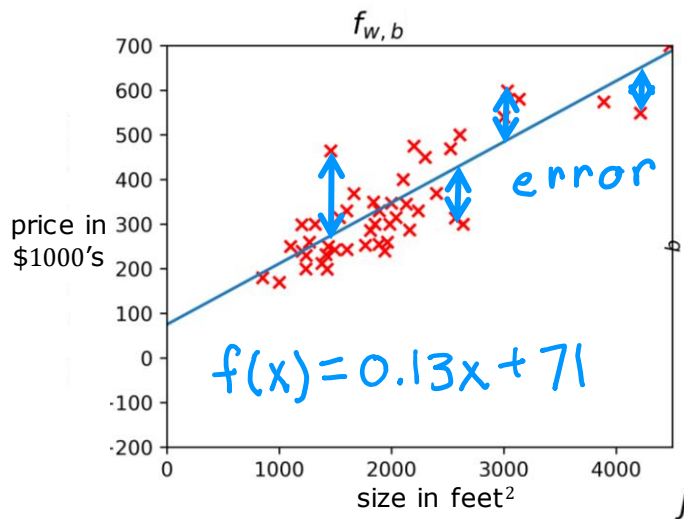




you can rotate this figure







Training Linear Regression

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 $w=0, b=0$)*

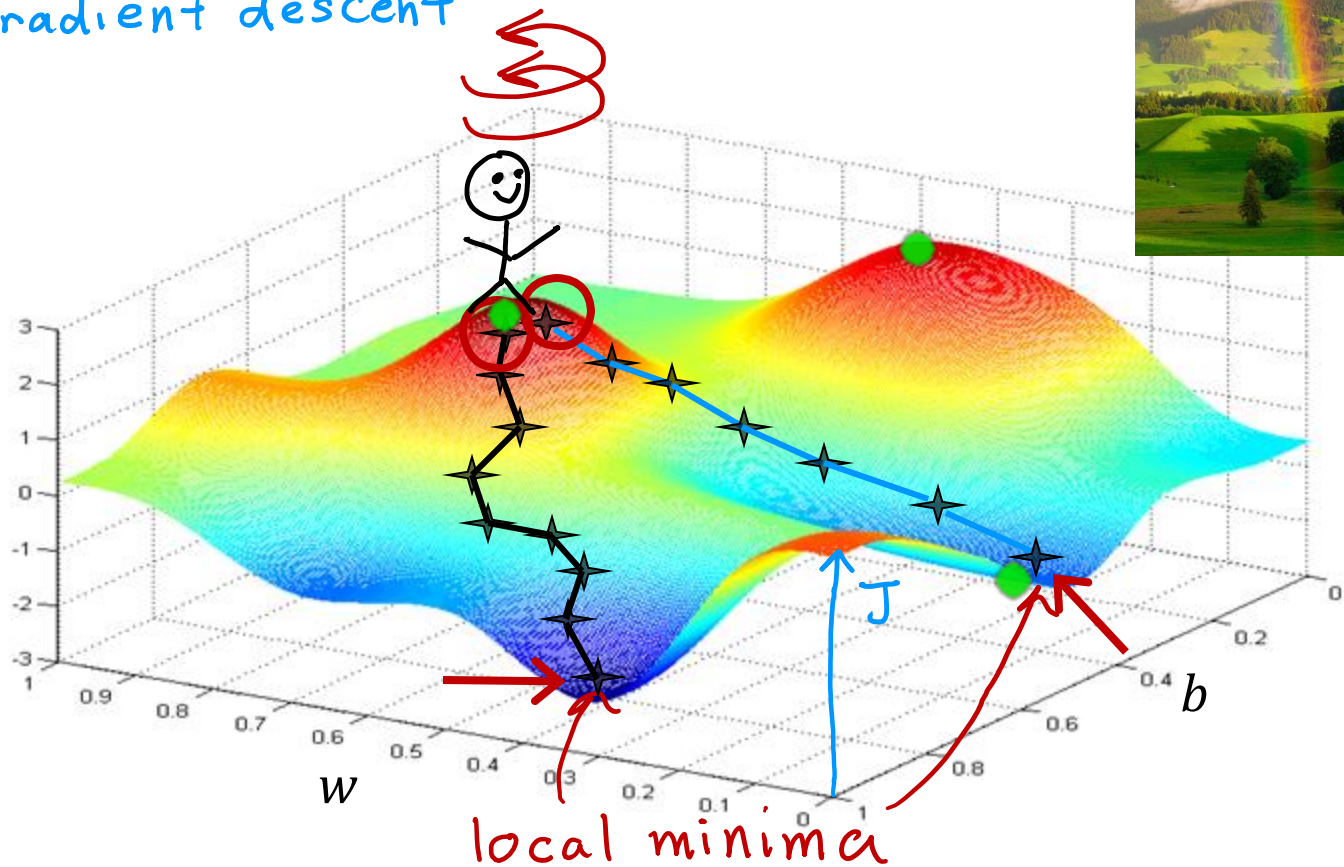
Keep changing w, b to reduce $J(w, b)$

Until we settle at or near a minimum

may have >1 minimum

gradient descent

$J(w, b)$



Training Linear Regression

Implementing
Gradient Descent

Gradient descent algorithm

Repeat until convergence

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

Learning rate
Derivative

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

Simultaneously
update w and b

Assignment

$$a = C$$

$$a = a + 1$$

Correct: Simultaneous update

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

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

$$w = tmp_w$$

$$b = tmp_b$$

Incorrect

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

$$w = tmp_w$$

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

$$b = tmp_b$$

Training Linear Regression

Gradient Descent
Intuition

Gradient descent algorithm

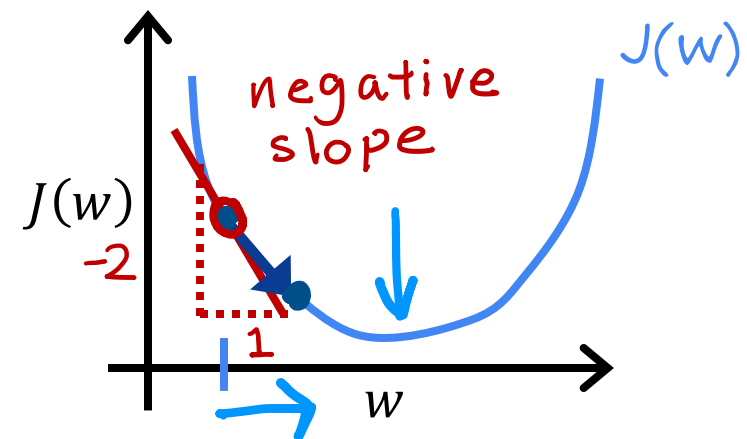
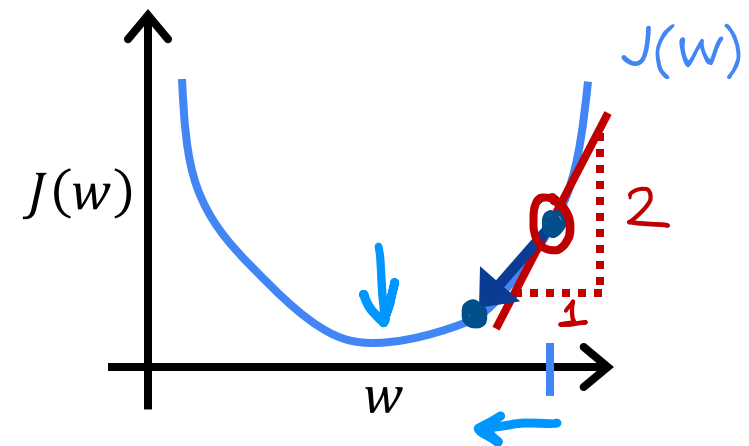
repeat until convergence {

learning rate \rightarrow α $\frac{\partial}{\partial w} J(w, b)$ *derivative*

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

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

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



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

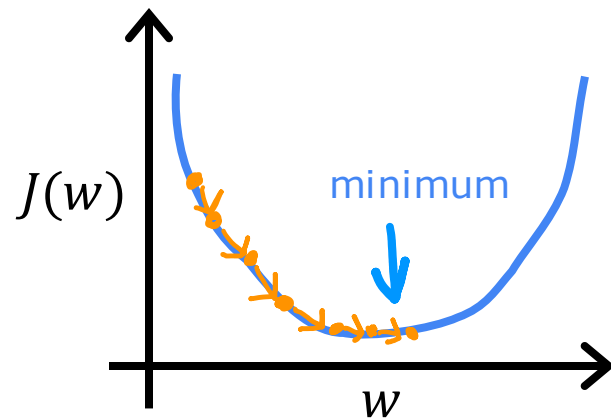
Training Linear Regression

Learning Rate

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

If α is too small...

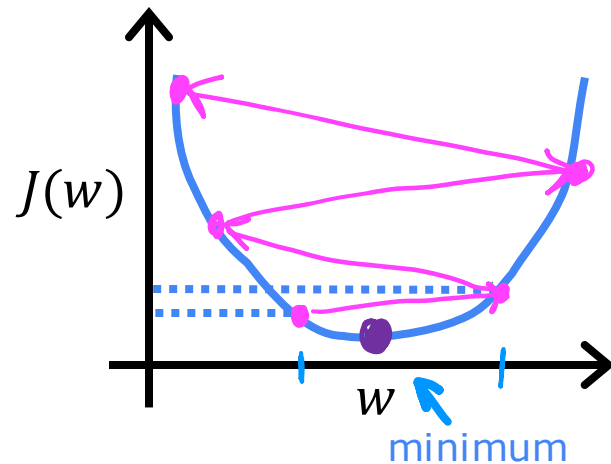
Gradient descent may be slow.

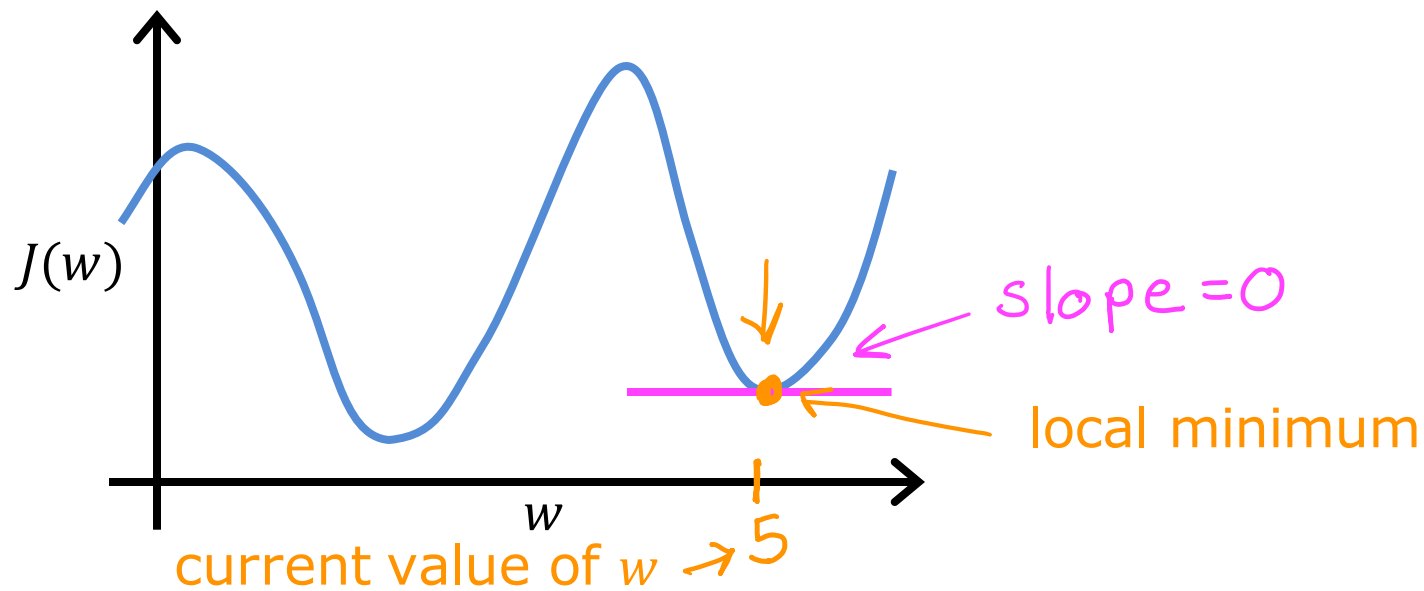


If α is too large...

Gradient descent may:

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





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

$$w = w - \alpha \cdot 0$$

$$w = w$$

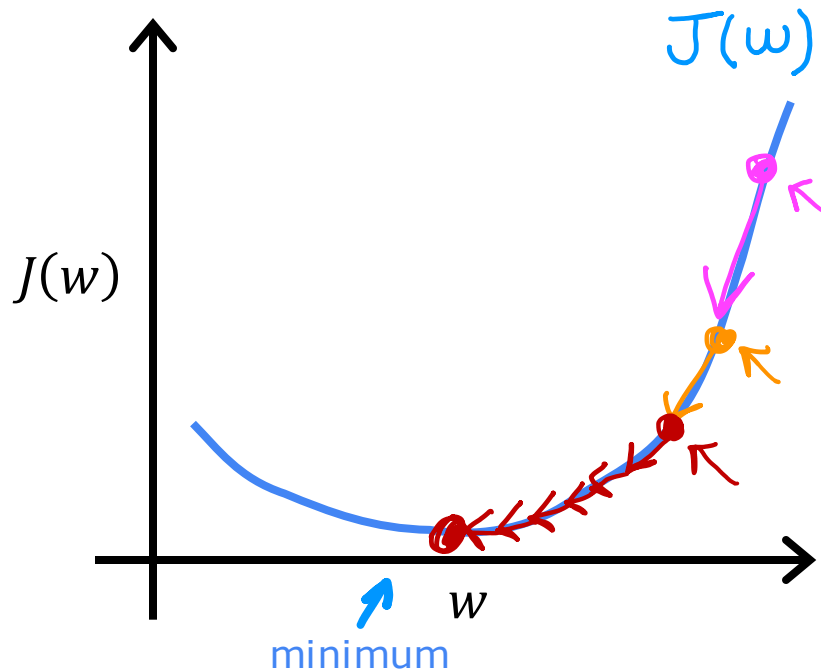
Can reach local minimum with fixed learning rate

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

Near a local minimum,

- Derivative becomes smaller
- Update steps become smaller

Can reach minimum without decreasing learning rate α



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 {

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

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

}

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

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

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

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

Gradient descent algorithm

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

repeat until convergence {

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

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

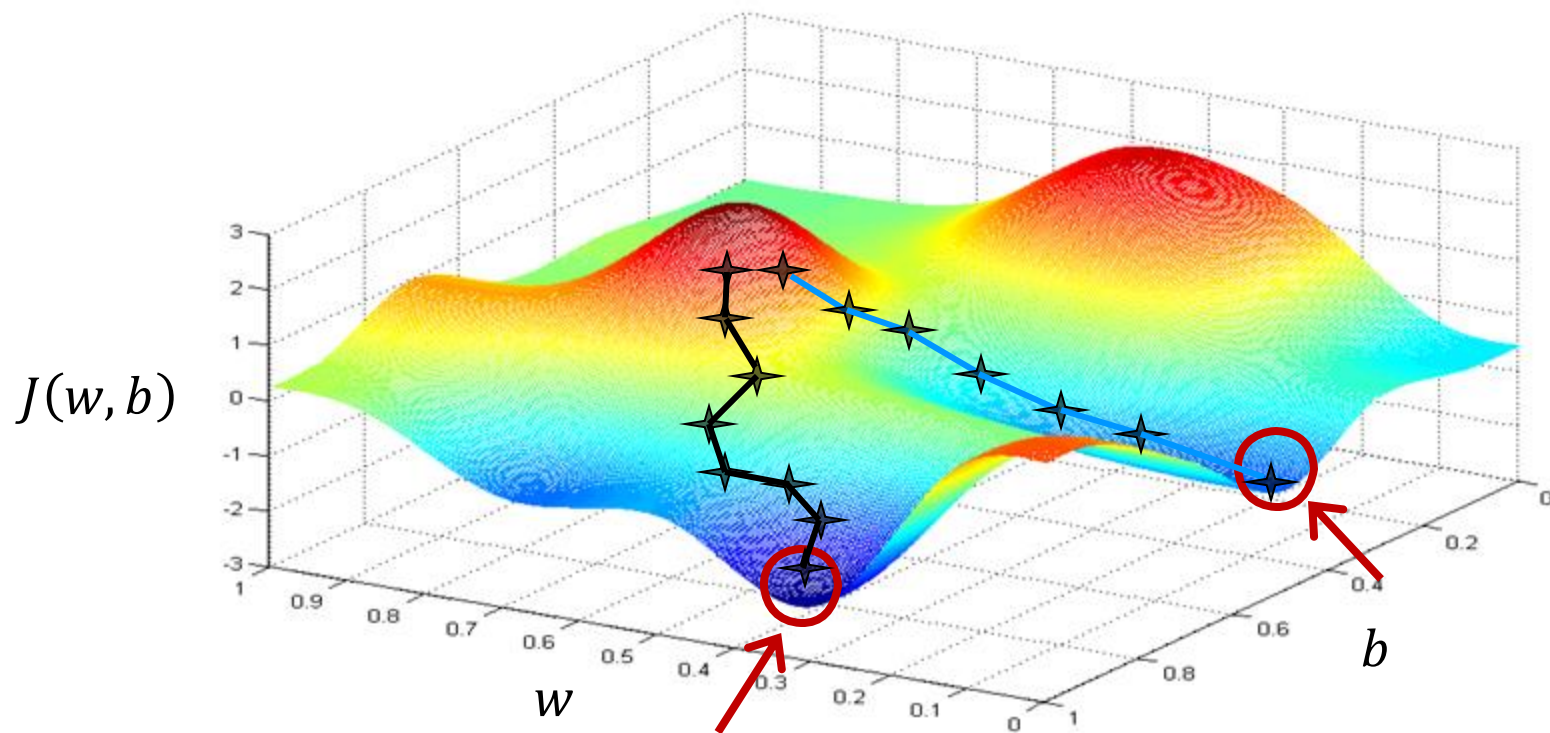
}

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


Update
 w and b
simultaneously

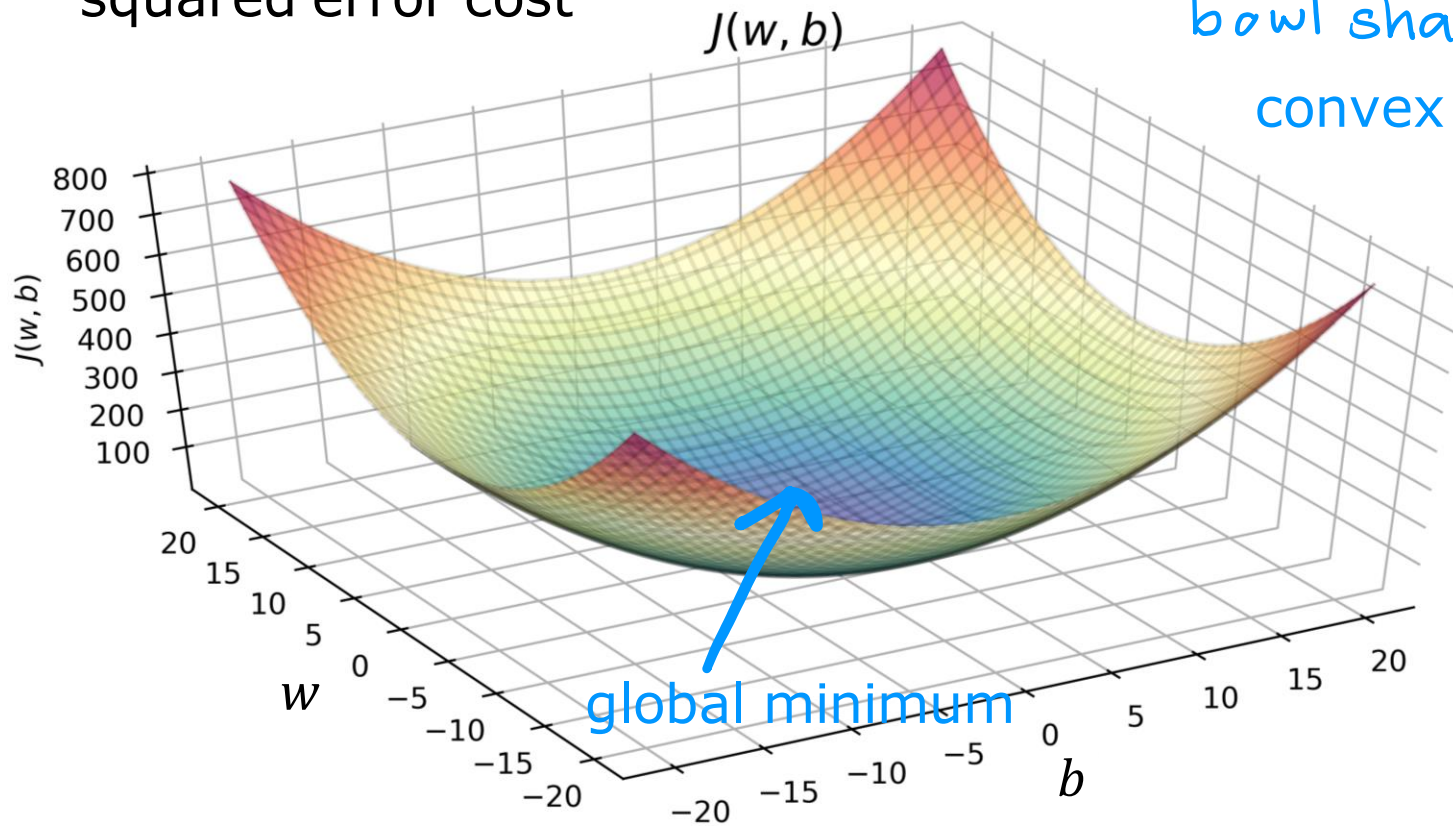
$$f_{w,b}(x^{(i)}) = wx^{(i)} + b$$

More than one local minimum



squared error cost

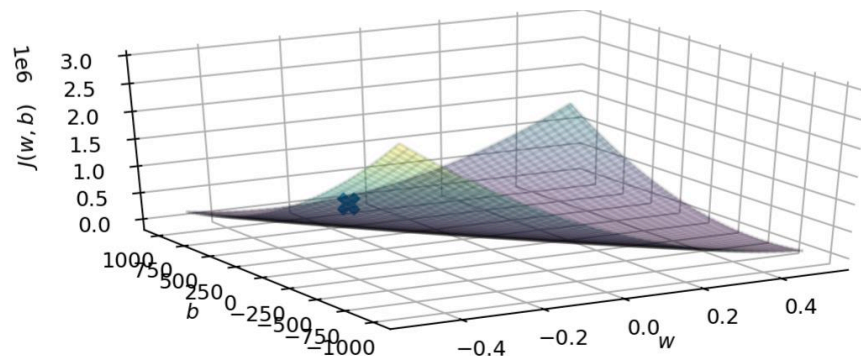
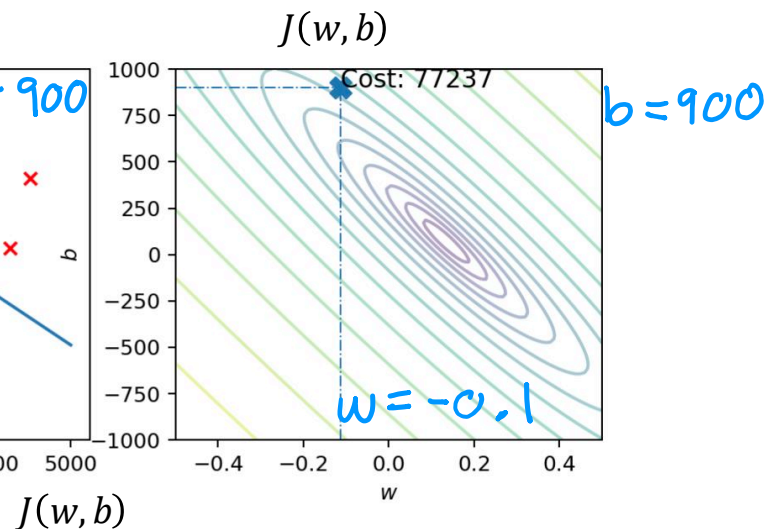
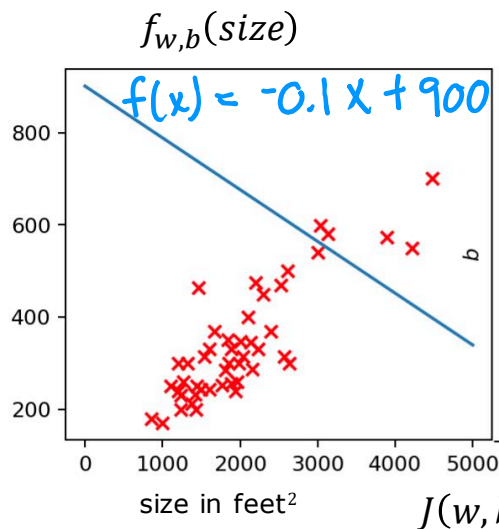
bowl shape 
convex function



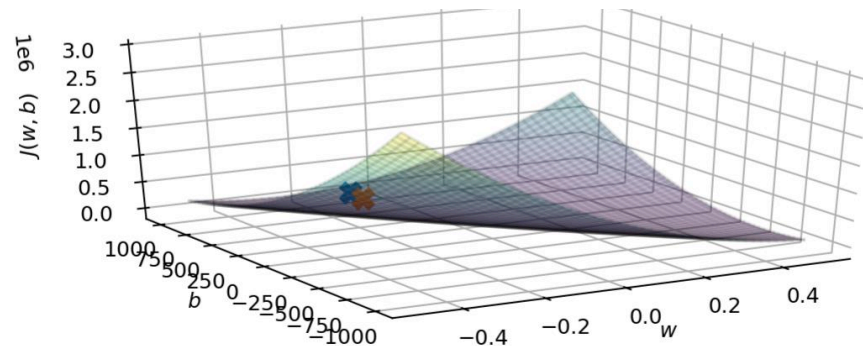
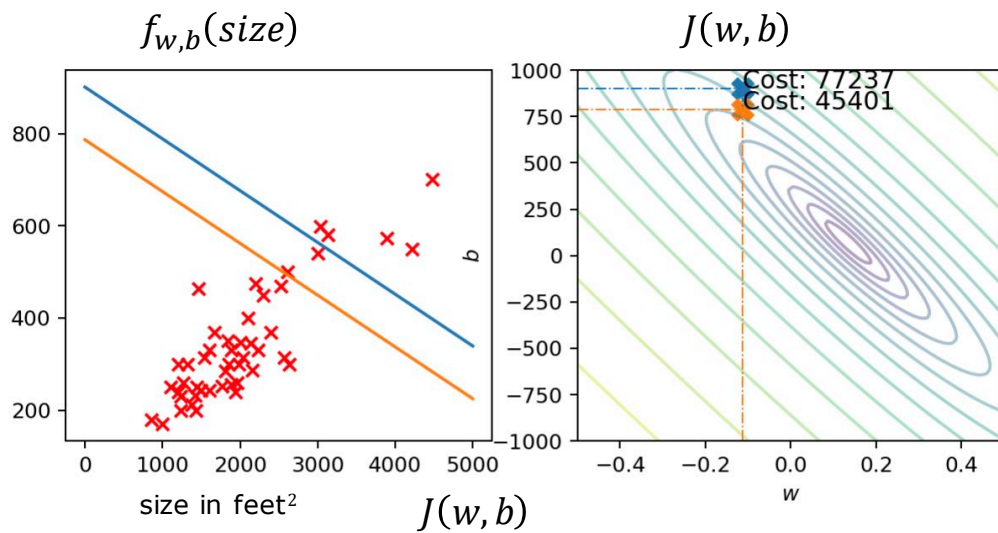
Training Linear Regression

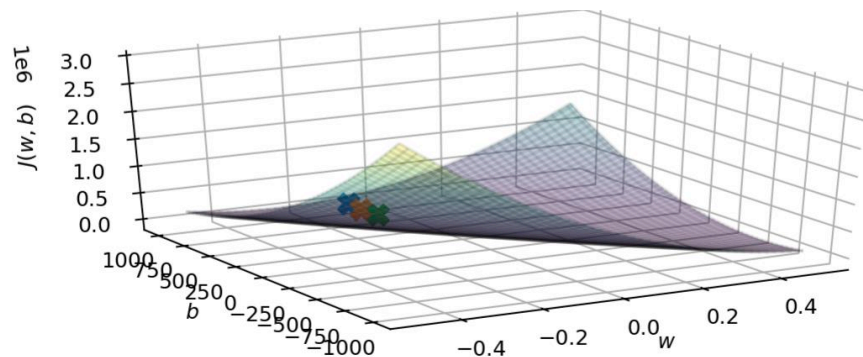
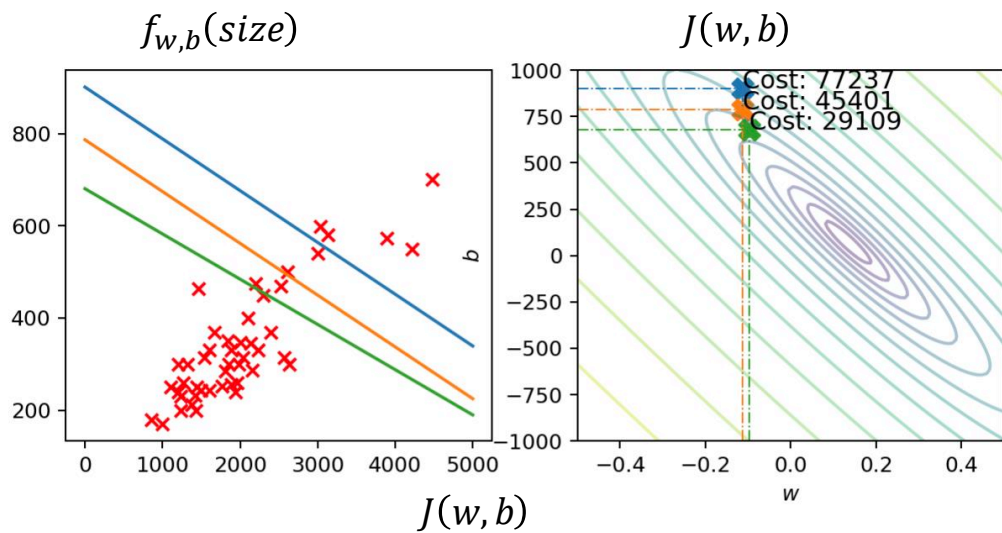
Running
Gradient Descent

price in
\$1000's

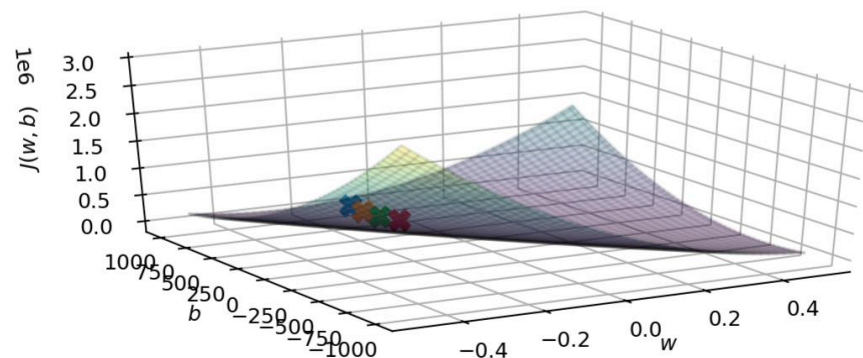
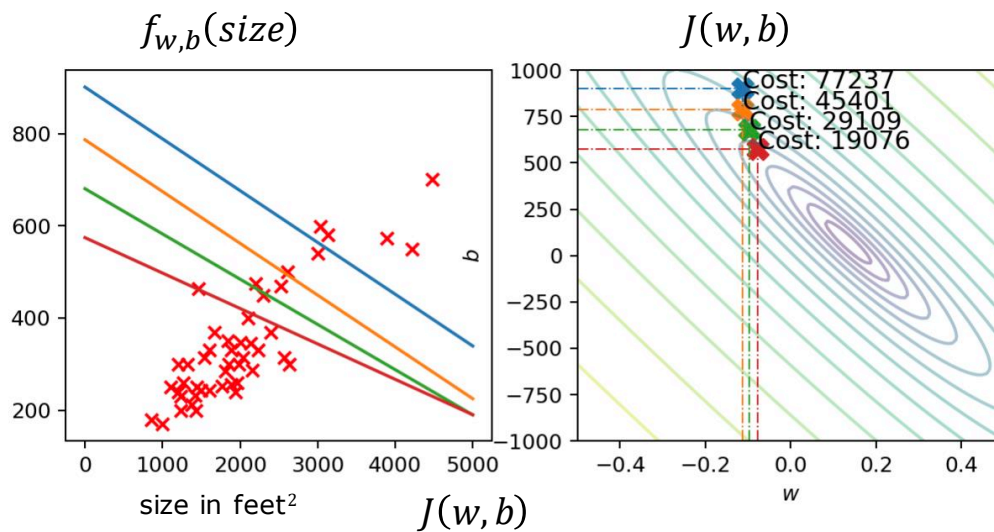


price in
\$1000's

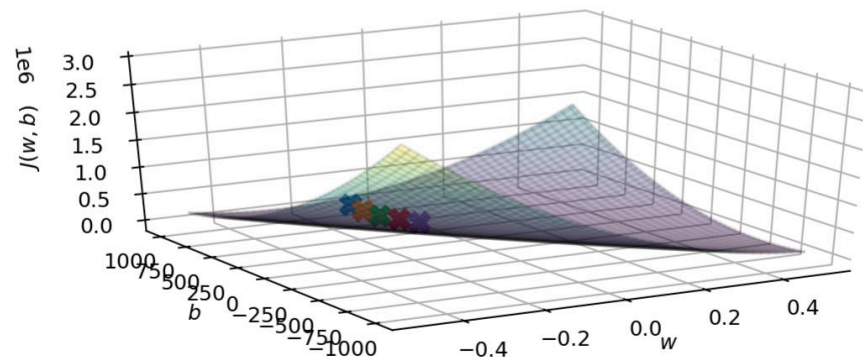
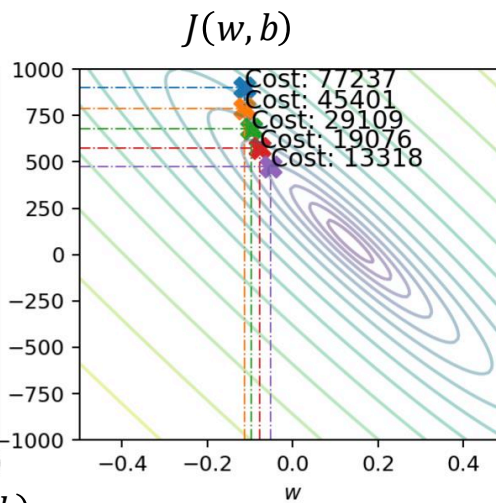
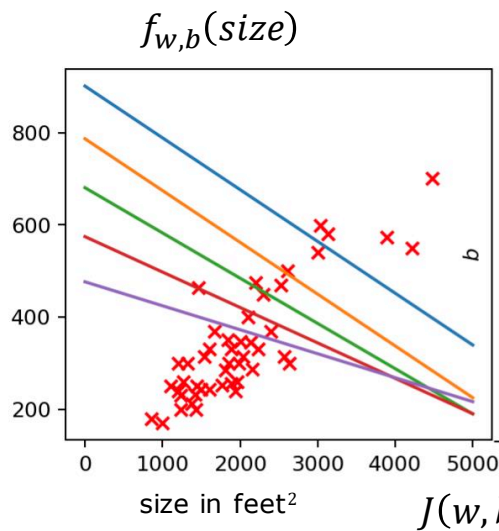




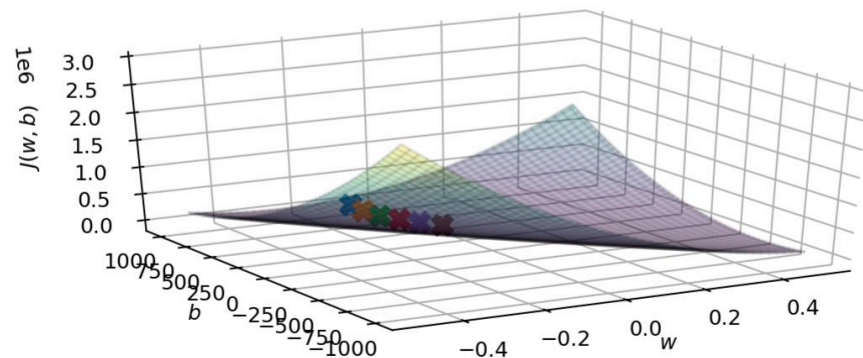
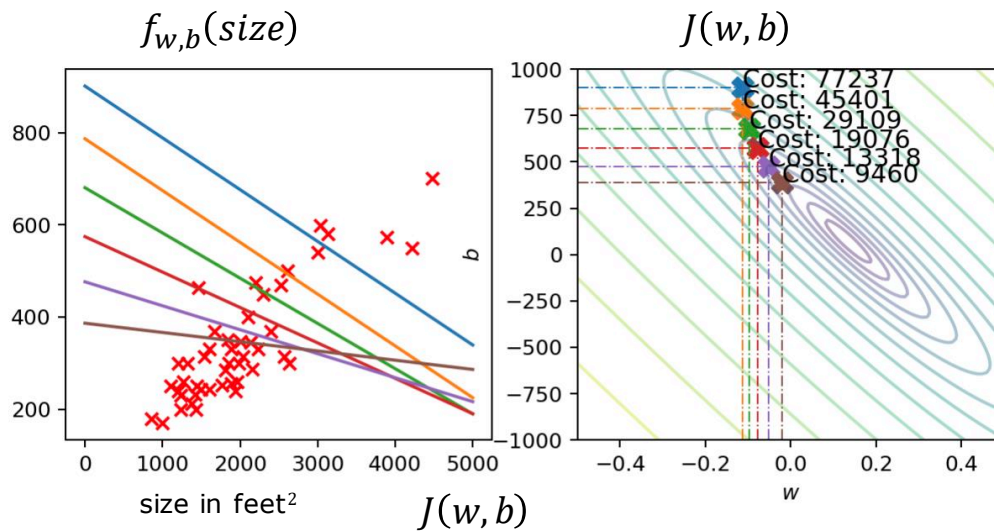
price in
\$1000's



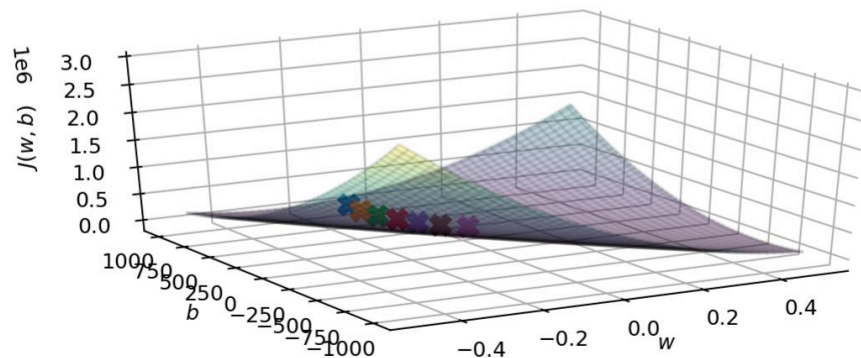
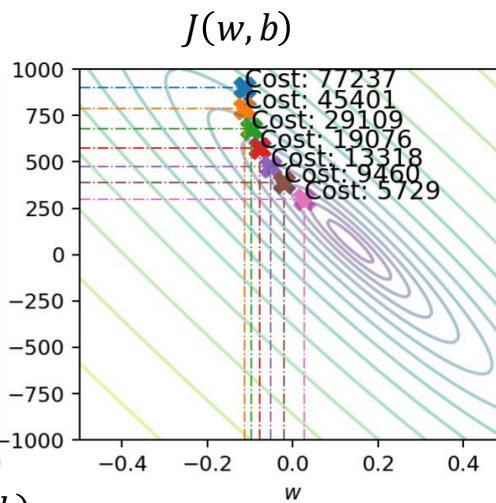
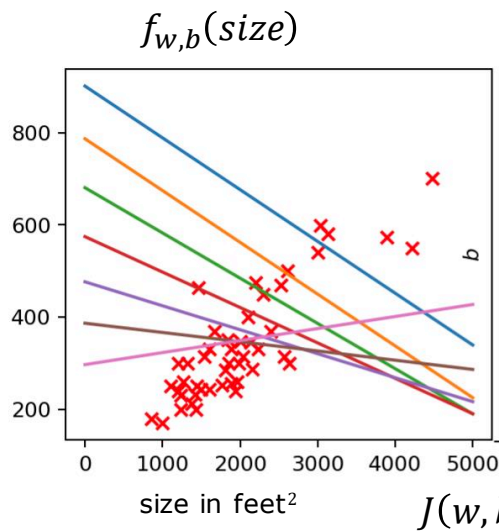
price in
\$1000's



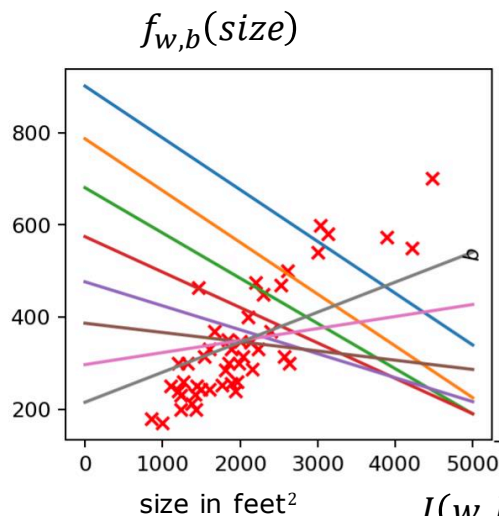
price in
\$1000's



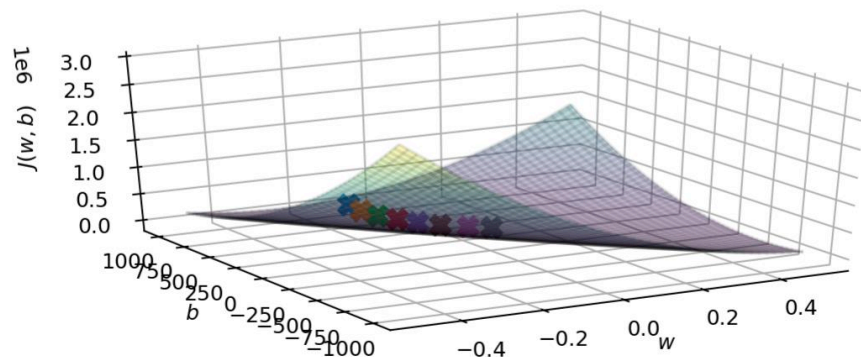
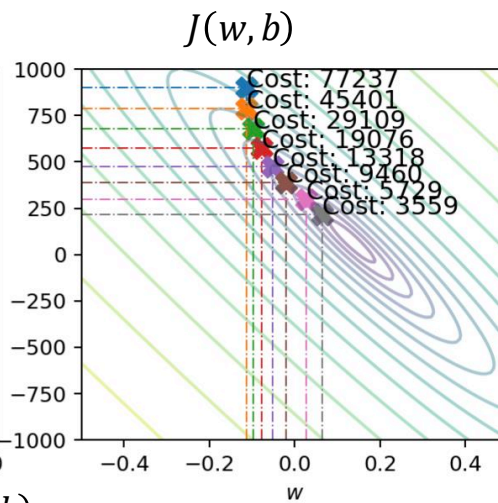
price in
\$1000's



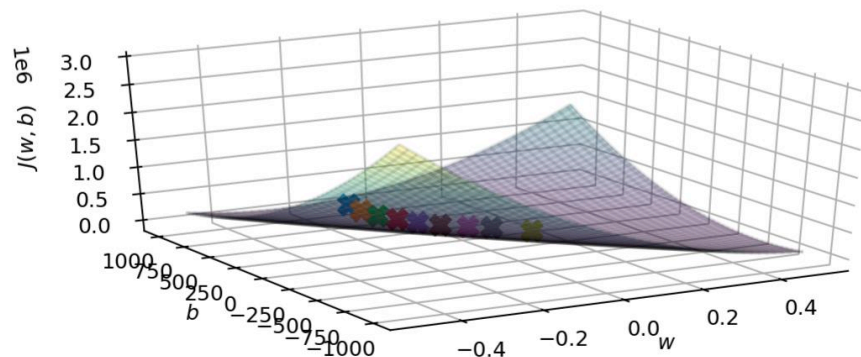
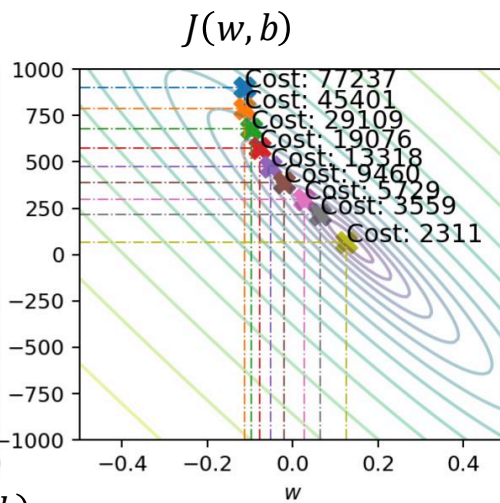
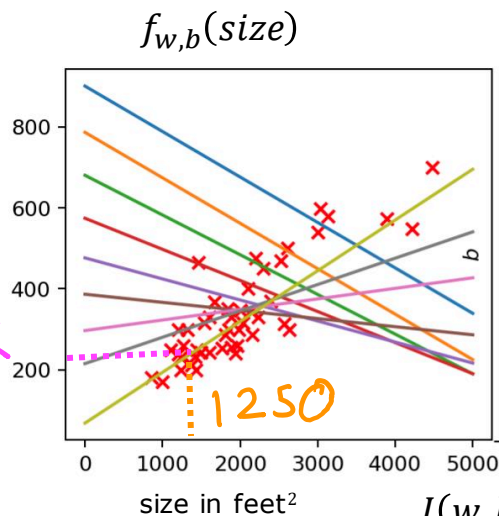
price in
\$1000's



$J(w, b)$



price in
\$1000's
\$250K



“Batch” gradient descent

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

	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

$m = 47$



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

