

Supervised Learning

Algorithms learn

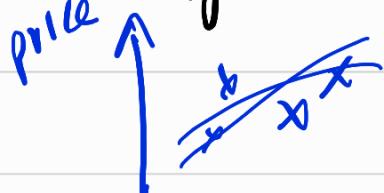
$X \rightarrow Y$
maps? input output label

"learn from being given
"right answers"

<u>Input</u>	<u>Output</u>
email	\rightarrow spam?
audio	\rightarrow text transcript?
ad, user info	\rightarrow click?

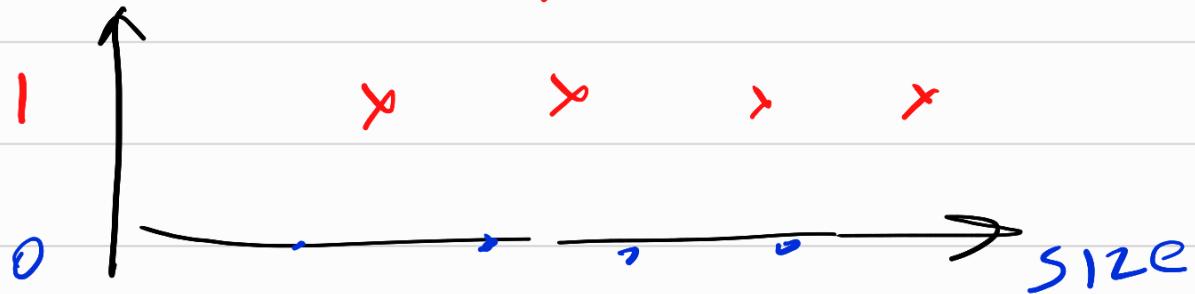
learn from x, y pairs \rightarrow then take new

Regression \rightarrow predicts a number
from many possible outputs (infinite)

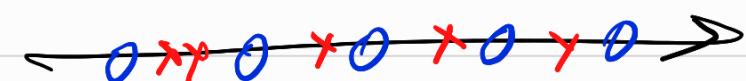


e.g.:

Classify : Breast cancer Detection
malignant or benign



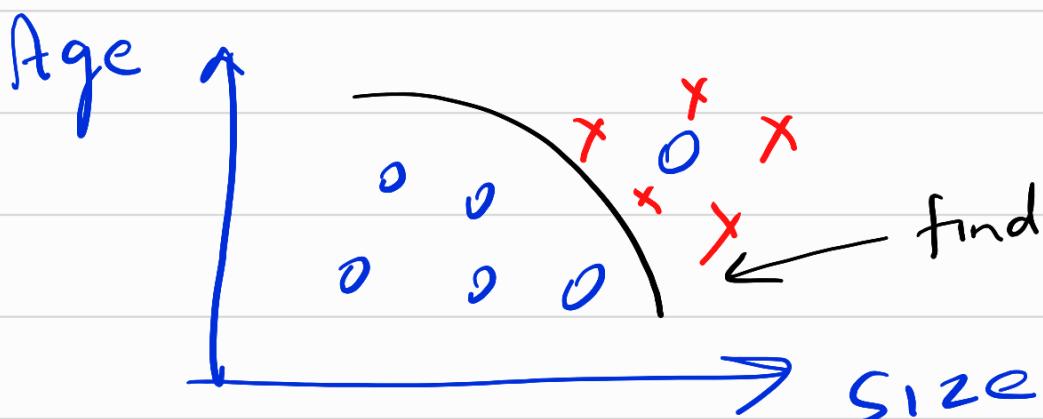
only two possible outputs



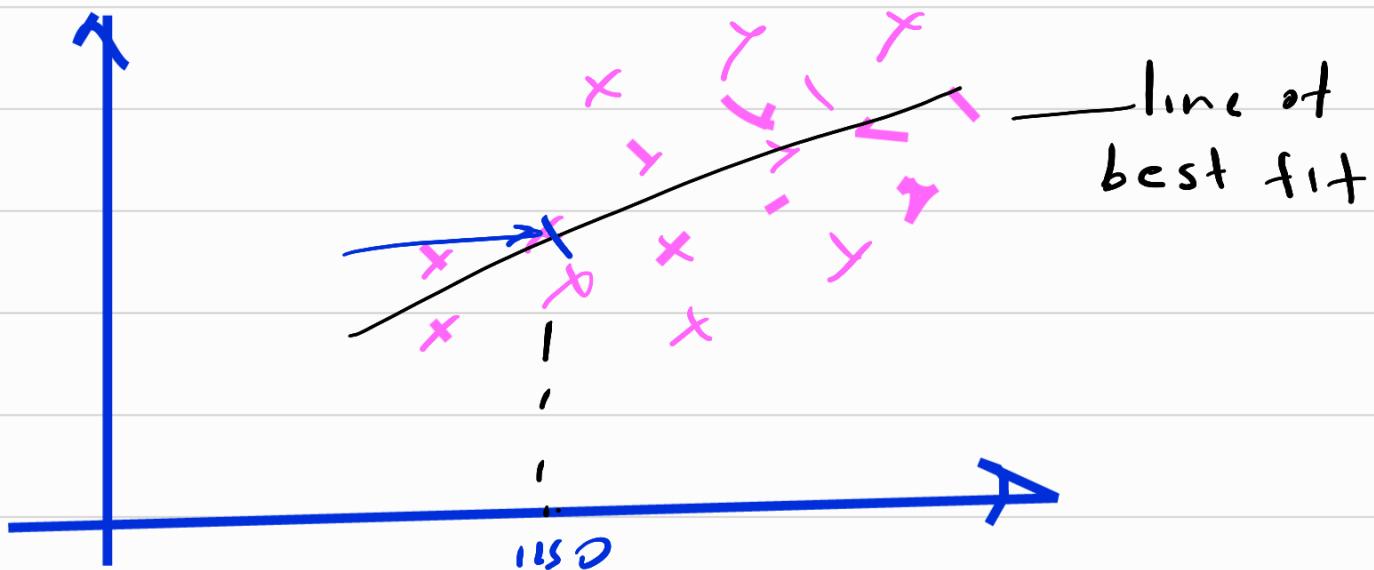
0 Δ Δ - 3 types, ---

Classification \rightarrow predict categories
Algorithms (finite set of outputs)

Can used to 2 inputs too;
or more



Linear Regression



Regression model predicts numbers infinite

Classification model \rightarrow categories
(small no. of outputs)
 \swarrow discrete finite

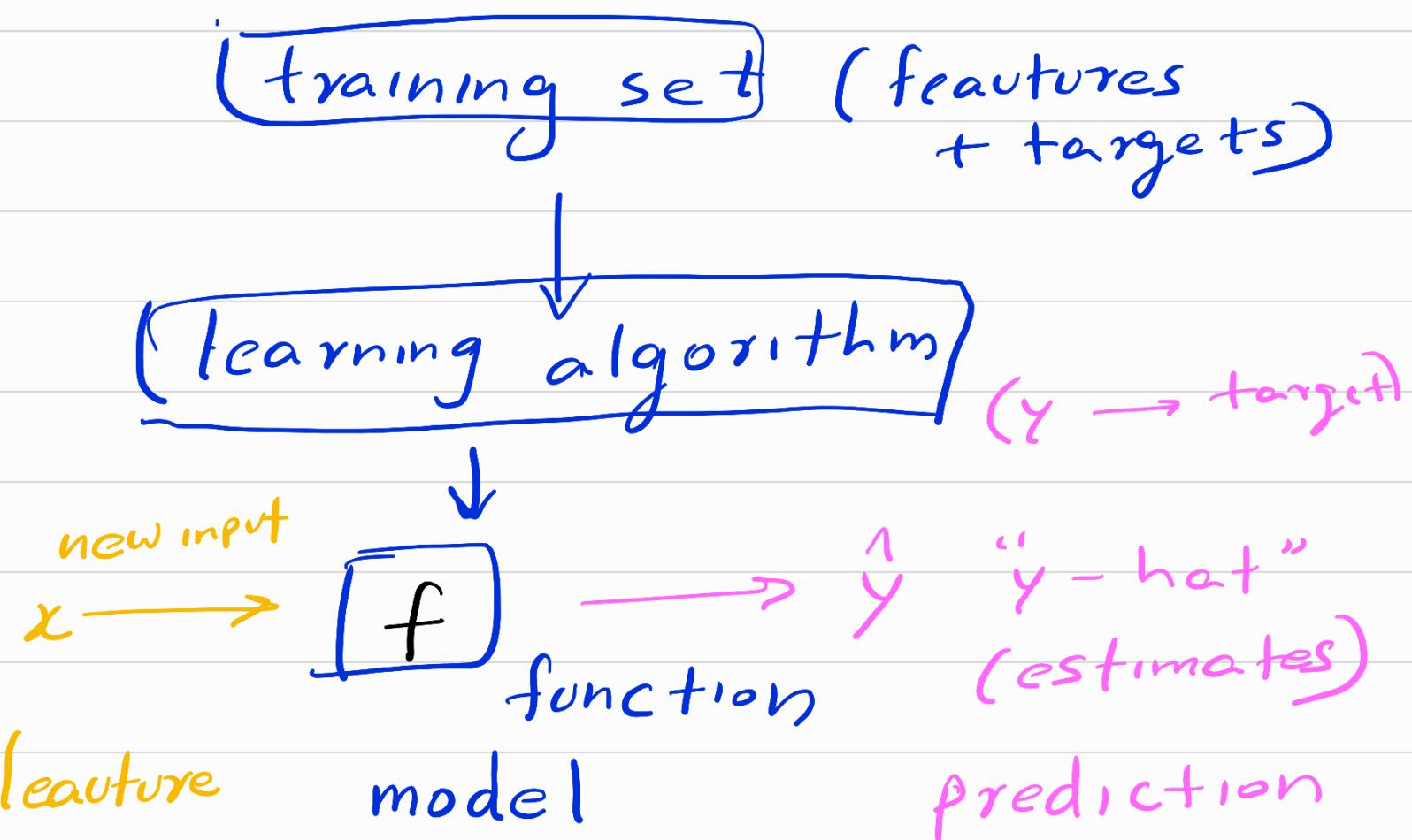
training set \rightarrow data used to train model

x - input variable, feature
 y - output variable, target

m - no. of training examples
(rows)

(x, y) = single training example

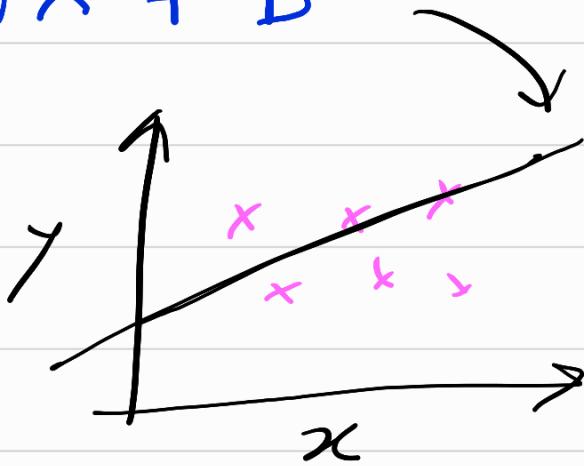
$(x^{(i)}, y^{(i)})$ — i th training example
 $(x^{(1)}, y^{(1)}) \quad i = 1 \rightarrow (2104, 402)$
 . . .



eg 1-
 size \rightarrow **f** \rightarrow price
 (estimated \hat{y})

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

$\underbrace{f(x)}$



Univariate linear regression
(one variable)

$$\frac{y - 300}{x - 1} = \frac{500 - 300}{2 - 1}$$

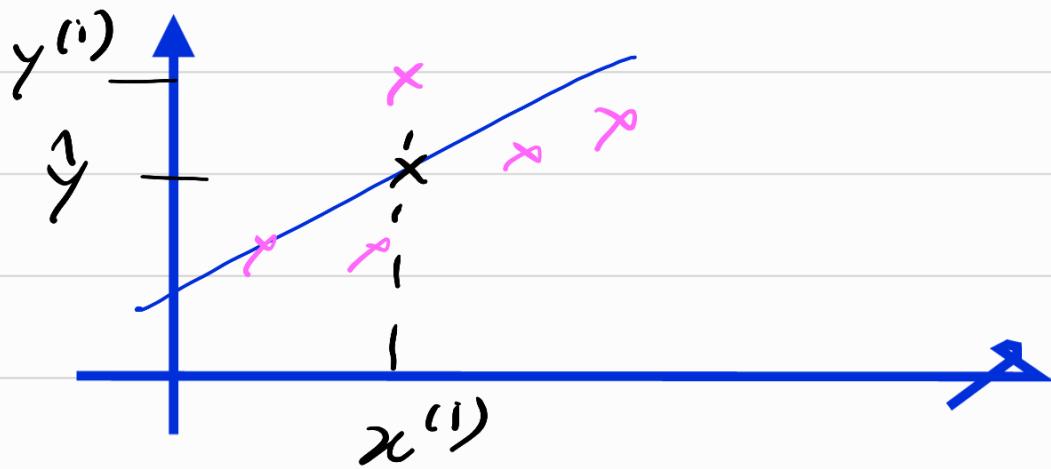
$$\begin{aligned} y &= 200(x - 1) + 300 \\ &= 200x + 100 \end{aligned}$$

Cost Function;

$$f(x) = w_2x + b$$

w, b : parameters

coefficients / weights



$$\frac{\sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2}{2m}$$

↑
(make
reater)
average
for one

error for all
training examples

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

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

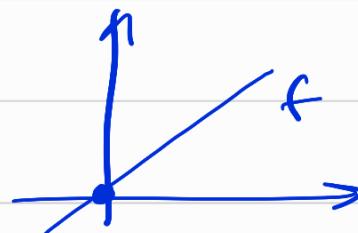


Squared error cost function

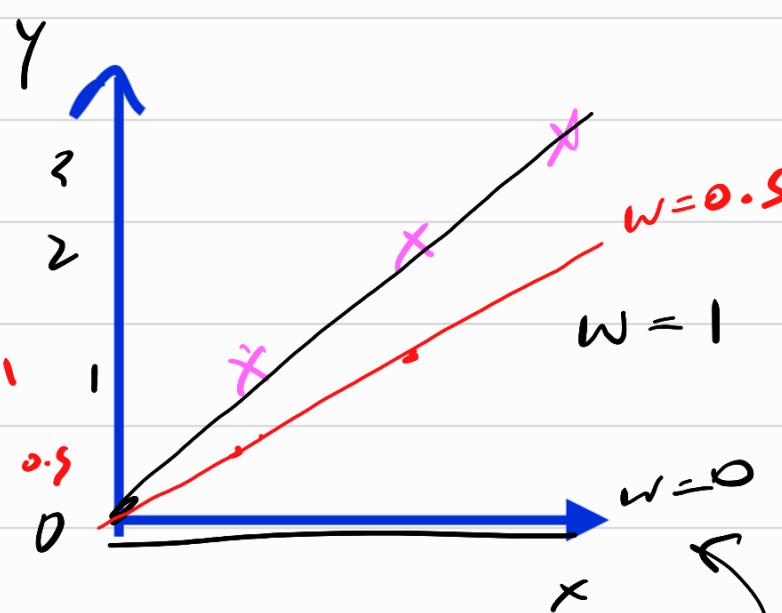
Simplified only for w no b

goal: minimize $J(w, b)$

$$f_w(x) = wx \quad b = \phi$$



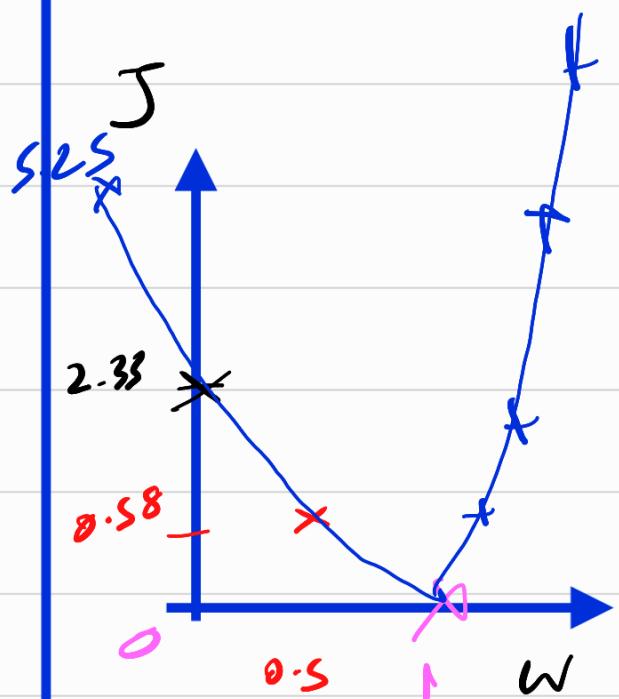
$f_w(x)$
(for fixed w , func of x)
INPUT



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

$0 + 0 + 0$

$J(w)$
(func of w)



$$w=0.5$$

$$\frac{1}{2m} \left\{ 0.5^2 + 1^2 + 1.5^2 \right\} = \frac{3.5}{6} \approx 0.58$$

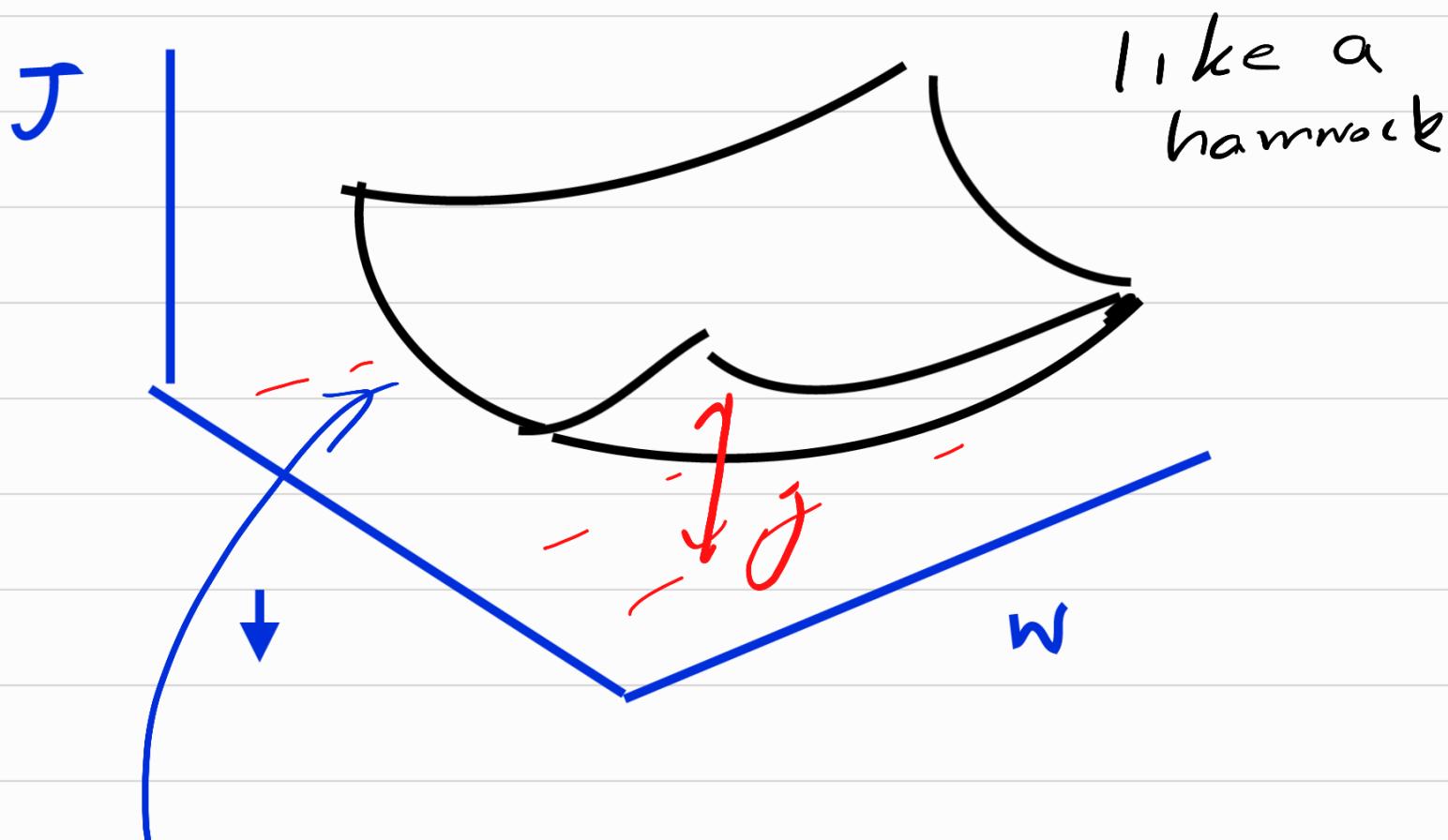
$$\cancel{w=0} \quad \frac{1}{2(3)} (1^2 + 2^2 + 3^2) = \frac{1}{2(3)} \times 14 = 2.33$$

$$w = -0.15 \rightarrow 5.25$$

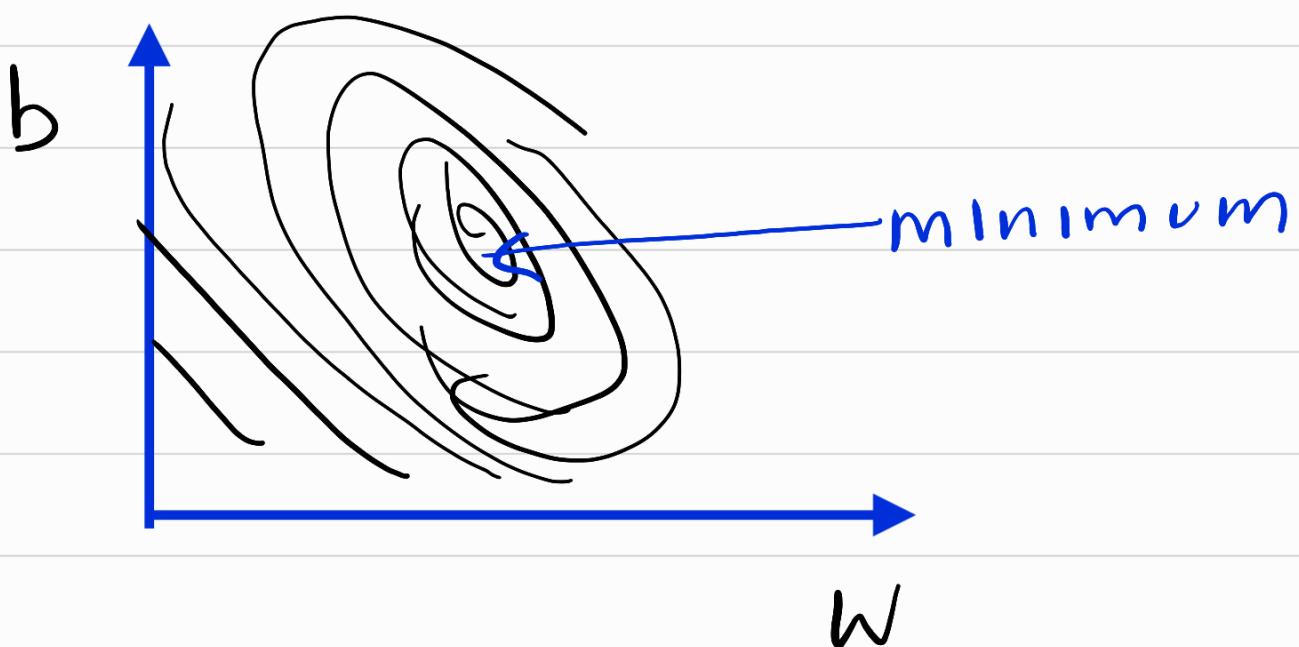
Choose w to minimize $J(w)$

$$\underline{\underline{w = 1}}$$

With both w & b



get horizontal sizes (same heights)
and make a contour plot



Gradient Descent

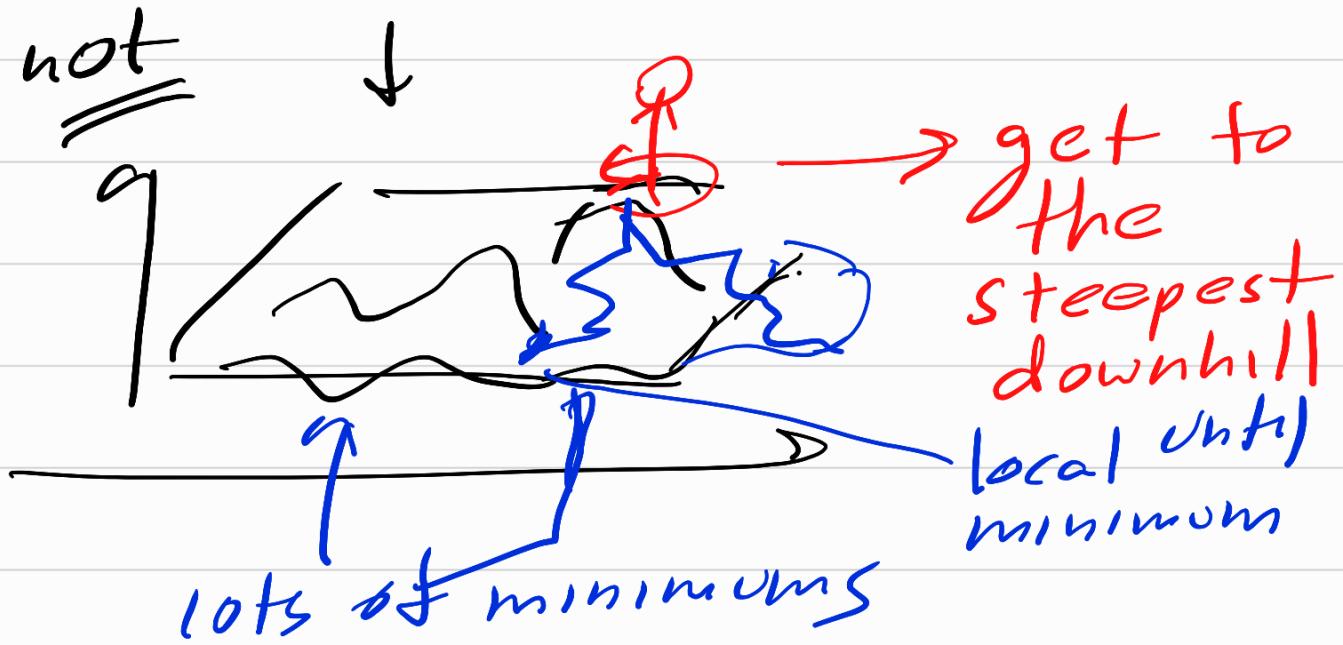
without manually minimizing J
find smallest J

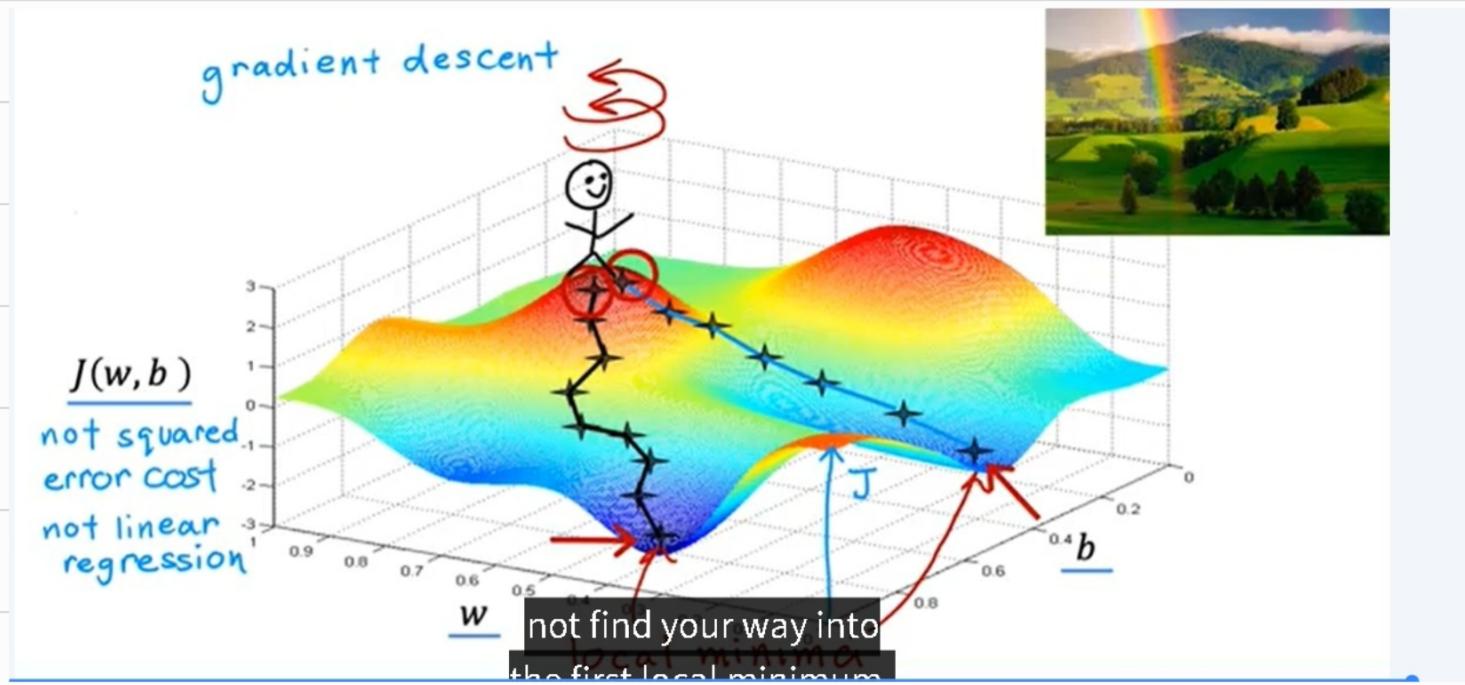
Start with some w, b ($w=0, b=0$)

keep changing w, b to reduce J

Find until we settle at or
near a minimum.

Squared error \rightarrow bow/harrow





Gradient descent algorithm

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

↑
assertion

derivative

α - learning rate
(size of the step)

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

Math

$a = c$

$a = a + 1 - X$

$a = == c$

* Repeat until convergence

- Simultaneously update w & b

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

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

$$w = \text{tmp_}w$$

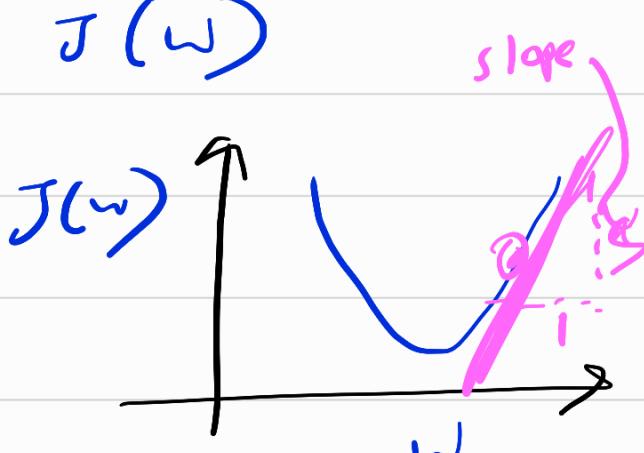
$$b = \text{tmp_}b$$

use
same w
before

$$\cancel{\frac{d w}{d w}}$$

for simplicity

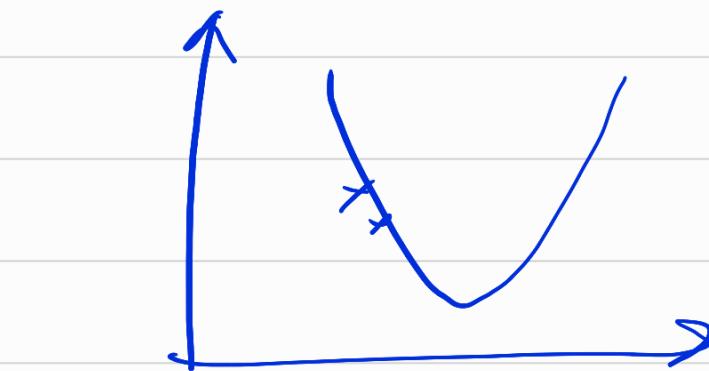
$$w - \alpha \frac{\partial}{\partial w} J(w)$$



$\alpha \geq 0$ always

$\frac{d}{d w} < 0$ if $\frac{d}{d w} > 0$
 w increase $\therefore w$ decreases
 w moves right $\therefore w$ moves left

α



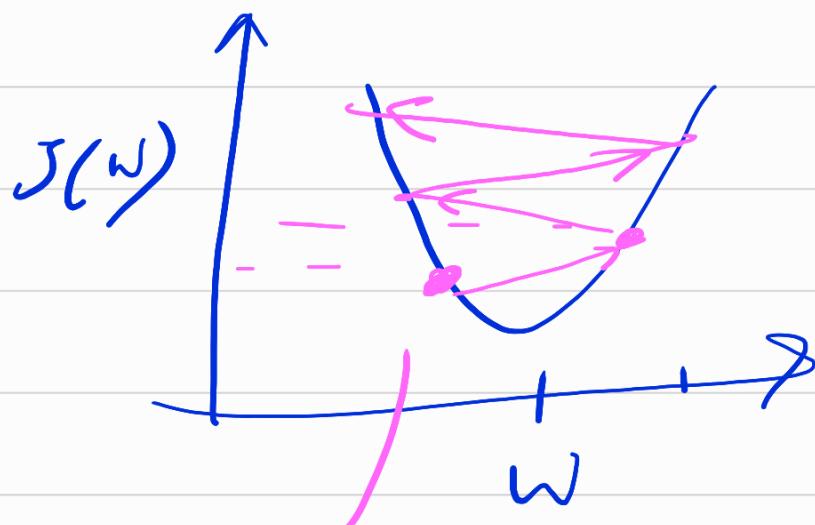
if α is too small

$$w - \alpha \frac{dJ(w)}{dw}$$

↑
(reduce is small)

tiny steps to get to minimum

if α is too large

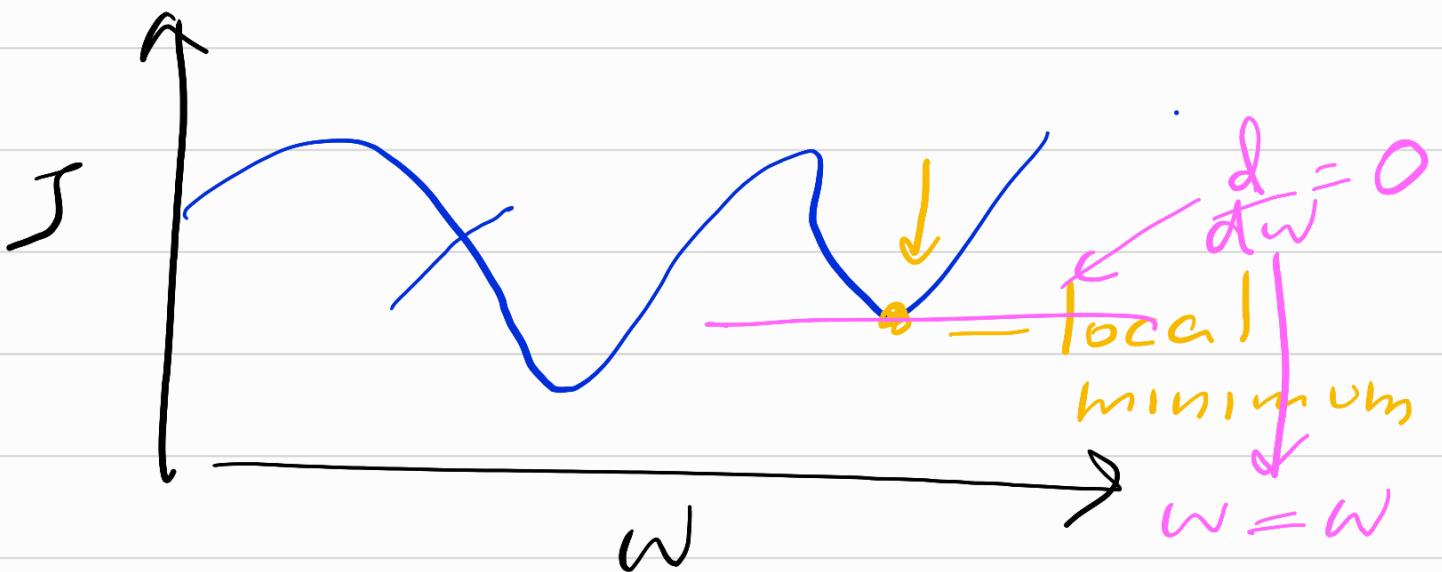


(if close to
minimum)

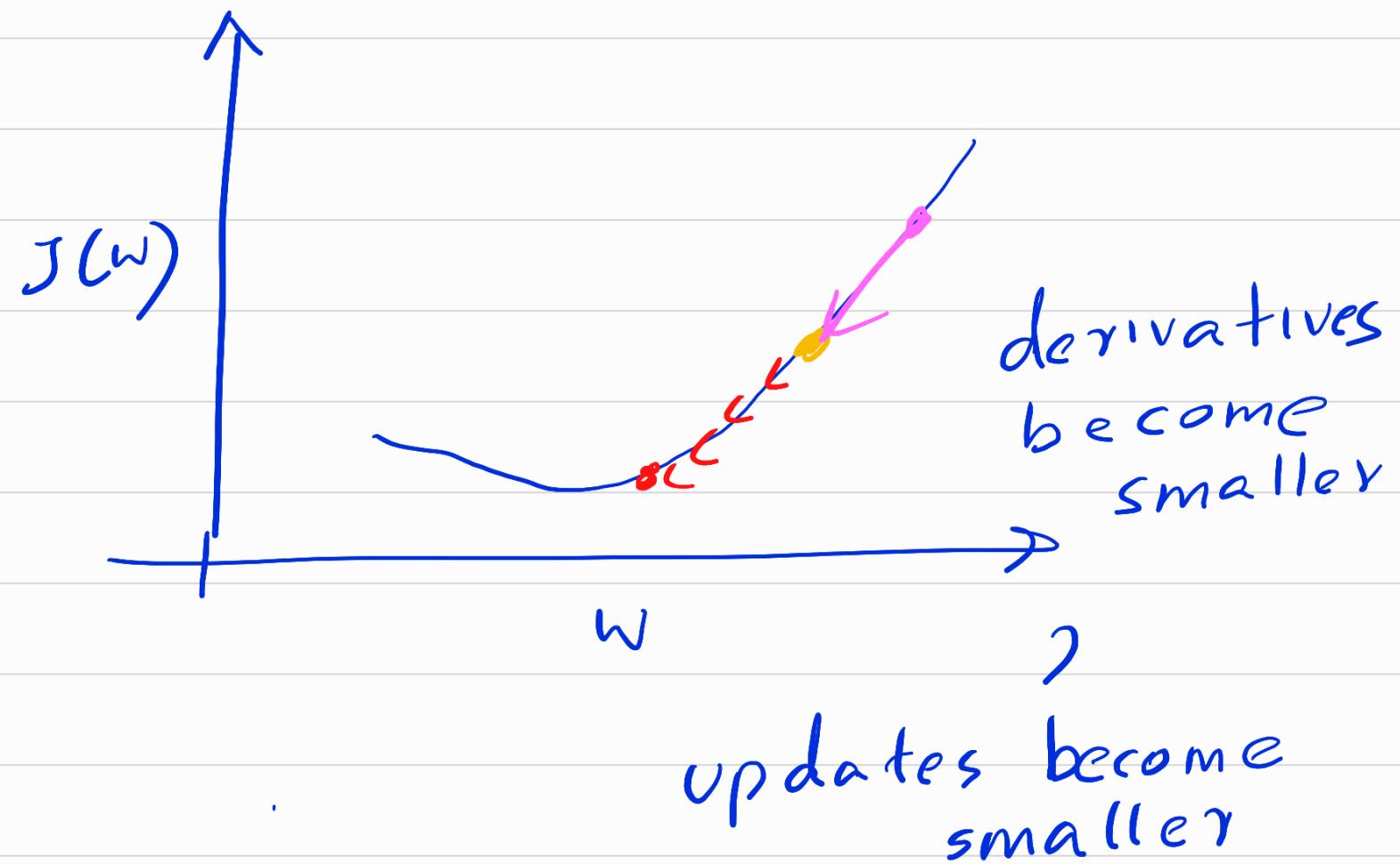
Overshoot,
(fail to
converge)

never reach
minimum

if already in minimum;



(not sq. error → many
minimums)



L.D. Algorithm

(Optional)

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

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

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

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

You can plug them into the
gradient descent algorithm.

with

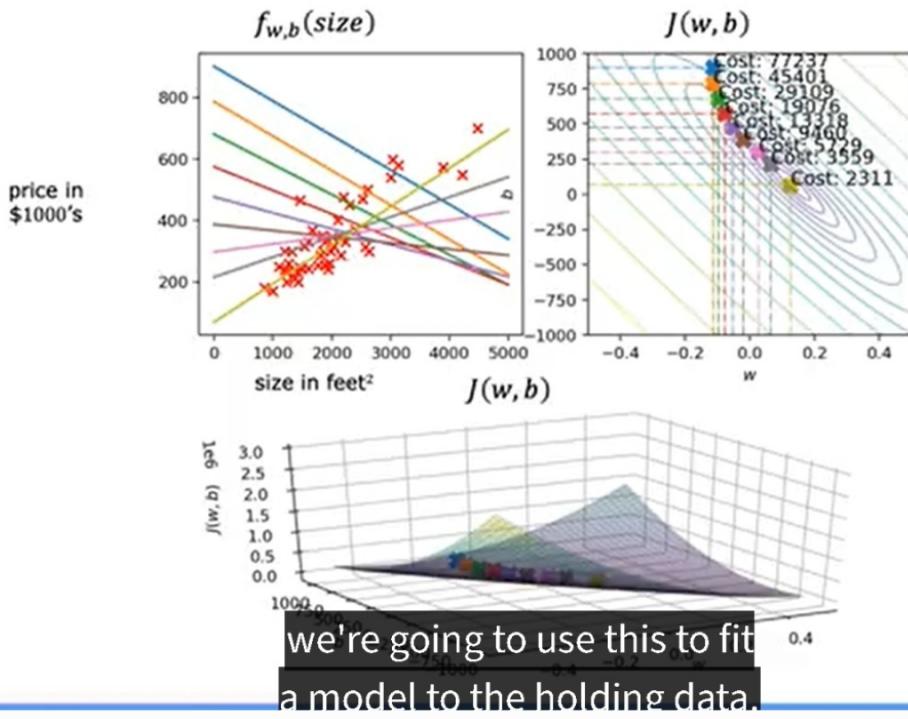
$$w = w - \alpha \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)}$$

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

update simultaneously until convergence.

Convex function (bowl)

→ single global minimum
eg: Squared error cost



Batch gradient descent

each step uses all training examples

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

Multiple features (variables)

x_j - j^{th} feature
 n = number of features
 $\vec{x}^{(i)}$ = features of i^{th} training example

vector
↓

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

$$X_3(2) = 2$$

previously model;

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

Now,

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

e.g.: -

$$0.1x_1 + 4x_2 + 10x_3 - 2x_4 \dots = 180$$

| | | |
size bedrooms years base
price

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

parameters of model

b is a number

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

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

dot product

$$= w_1x_1 + w_2x_2 + \dots + b\theta_nx_n + b$$

multiple linear regression

Vectorization

$$w = \text{np.array}([1.0, 2.5, -3.3])$$
$$b = 4$$

$$x = \text{np.array}([10, 20, 30])$$

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

python

$$f = 0$$

for j in range(0, n):

$$f = f + w[j] * x[j]$$

$$f = f + b \quad \times$$

↓ Vectorization

$$\boxed{f = \text{np.dot}(w, x) + b}$$

Gradient Descent

$$w = \text{np.array}([, \dots])$$
$$d = \text{np.array}([0, \dots])$$

$$w_j = w_j - \underbrace{0.1}_{\alpha} d_j$$

without vector \vec{j}

for j in range(16);
 $w[j] = w[j] - 0.1 * d[j]$

with vector \vec{j}

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

Previous notation

www.coursera.org – To exit full screen, press Esc

Parameters

$$w_1, \dots, w_n$$

$$b$$

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

$$\text{Cost function} \quad J(\underbrace{w_1, \dots, w_n}_J, b)$$

Vector notation

\vec{w} ← vector of length n
 b still a number

$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
 $J(\vec{w}, b)$ dot product

Gradient descent

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

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

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

simultaneously update w, b

}

n features ($n \geq 2$)

```
j=1
     $w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_1^{(i)}$ 
    :
     $\downarrow \frac{\partial}{\partial w_1} J(\vec{w}, b)$ 
j=n
     $w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) x_n^{(i)}$ 
     $b = b - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)})$ 
    simultaneously update
     $w_j$  (for  $j = 1, \dots, n$ ) and  $b$ 
}
```

Alternative to gradient descent

→ Normal equation

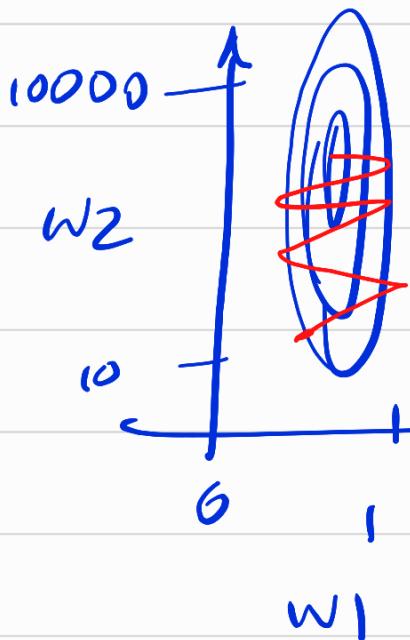
* only for linear regression
solve for w, b without iterations

Doesn't generalize

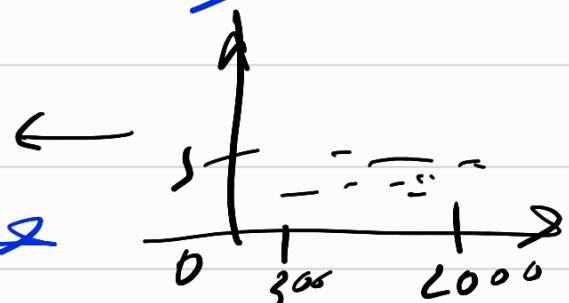
Feature Scaling

feature large relatively
small

$\frac{w}{\text{small}}$
 w
 large



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

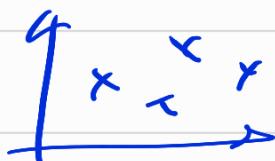


$$300 < x_1 < 2000$$

$$\text{scaled } x_1 = \frac{x_1}{2000}$$

$$0 < x_2 < 5$$

$$\text{scaled } x_2 = \frac{x_2}{5}$$



or Mean Normalization

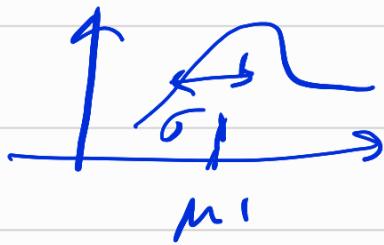
$$x_1 = \frac{x_1 - \mu_1}{s_1}$$

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

$$x_2 = \frac{x_2 - \mu_2}{s_2}$$

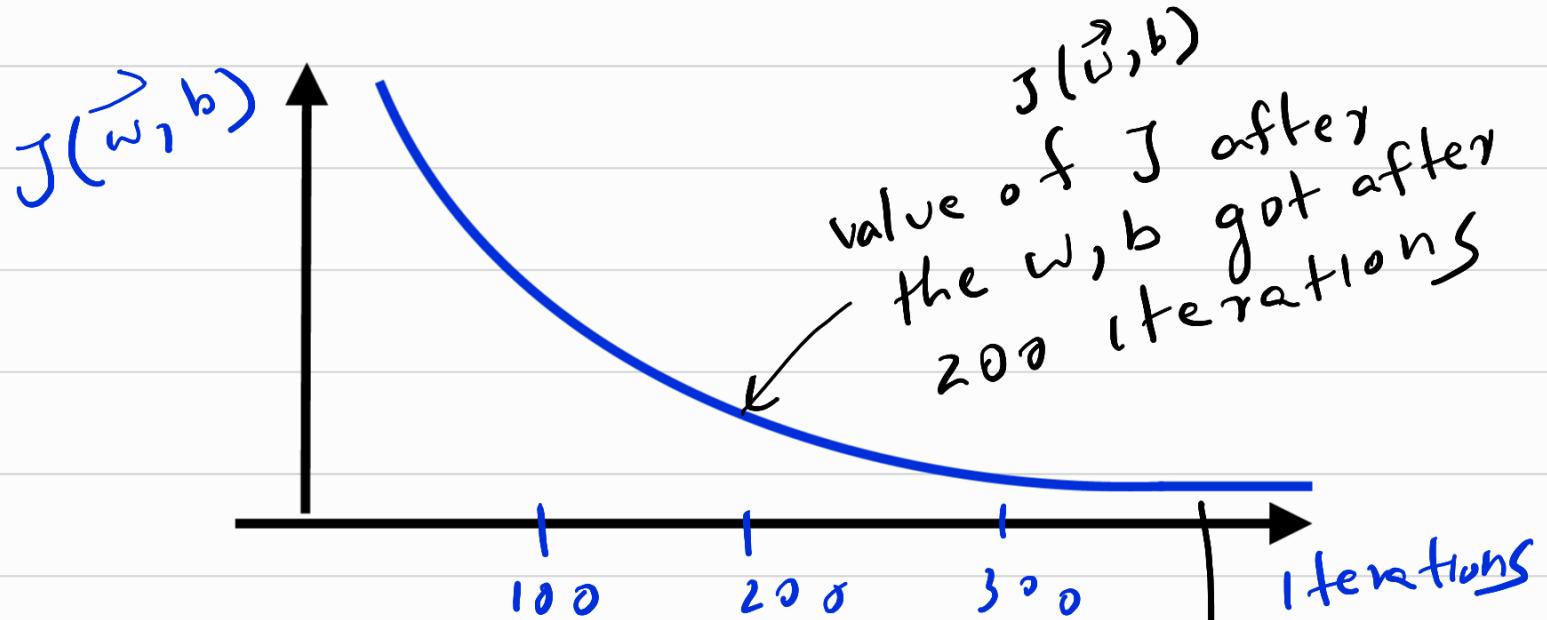
$$-0.46 \leq x_2 \leq 0.54$$

Z-score normalization



$$x_1 = \frac{x_1 - \mu_1}{\sigma_1} \quad x_2 = \frac{x_2 - \mu_2}{\sigma_2}$$

aim about $-1 < x_j < 1$

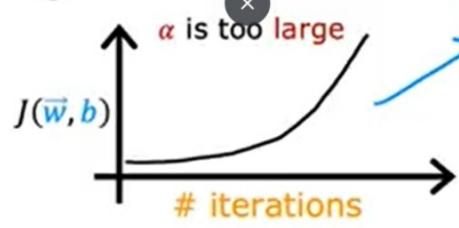
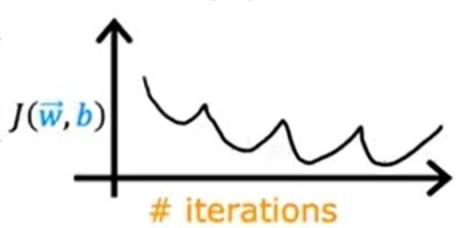


Automatic convergence Test

If $J(\vec{w}, b)$ decreases by $\leq \varepsilon$ in one iteration declare convergence

Choosing Learning Rate

Identify problem with gradient descent



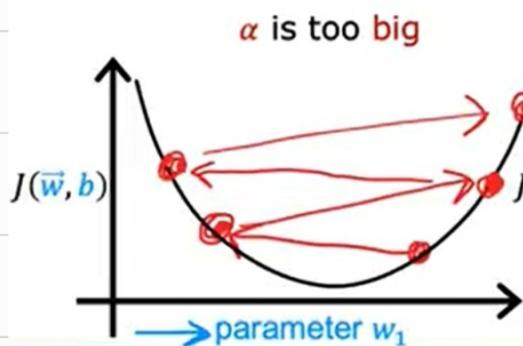
or learning rate is too large

$w_1 = w_1 + \alpha d_1$

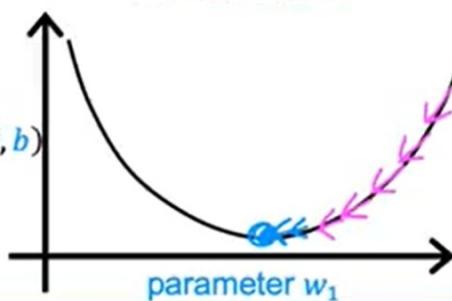
use a minus sign

$w_1 = w_1 - \alpha d_1$

Adjust learning rate



Use smaller α



With a small enough α ,
 $J(\vec{w}, b)$ should decrease
on every iteration

Feature Engineering

transform, combine features
 h, w, l
 \rightarrow *volume

Polynomial Regression



$$w_1 n + w_2 x^2 + w_3 x^3 + b$$

$$\sqrt{x} \quad | \quad w, b(x) = w_1x + w_2\sqrt{x} + b$$