

COURSE 1

What is Machine Learning?

Definition: Field of study that gives computer the ability to learn without being explicitly programmed. [Checkers example]

Major 2 + 1 types of learning:

- 1: Supervised
- 2: Unsupervised
- 3: Reinforcement

Supervised Learning part 1

The key characteristic of supervised learning is that you give your learning algorithms example to learn from.

It learns x to y or input \rightarrow output map[†]
The correct pairs of x and y are given and learning algorithm eventually learns to take just the I/P without O/P label and gives reasonably accurate prediction / guess of O/P

$$x \longrightarrow y$$

I/P O/P label

Learns from being given
"right answers"

Examples:

Email	→ spam (0/1)	- spam filter
audio	→ text trans.	- speech rec
English	→ spanish	- m/c transl
ad, user info	→ click? (0/1)	- O/L adv
image, radar info	→ positi of car	- self driving
image of phone	→ defect (0/1)?	- visual inspec

Supervised Learning part 2

Regression

Predict a **number** from
infinitely many possible O/P

Classification

Predict **category** from small number
of possible O/P.

Unsupervised learning part 1

Find something interesting in **unlabeled** data.
(given I/P x and O/P y , find structure in data)

Part 2

Unsupervised learning

Clustering

group similar points
together

Dimensionality Red

compress data using
fewer numbers

Anomaly det

Find unusual
data points

Linear Regression Model

Linear Regression with one Variable

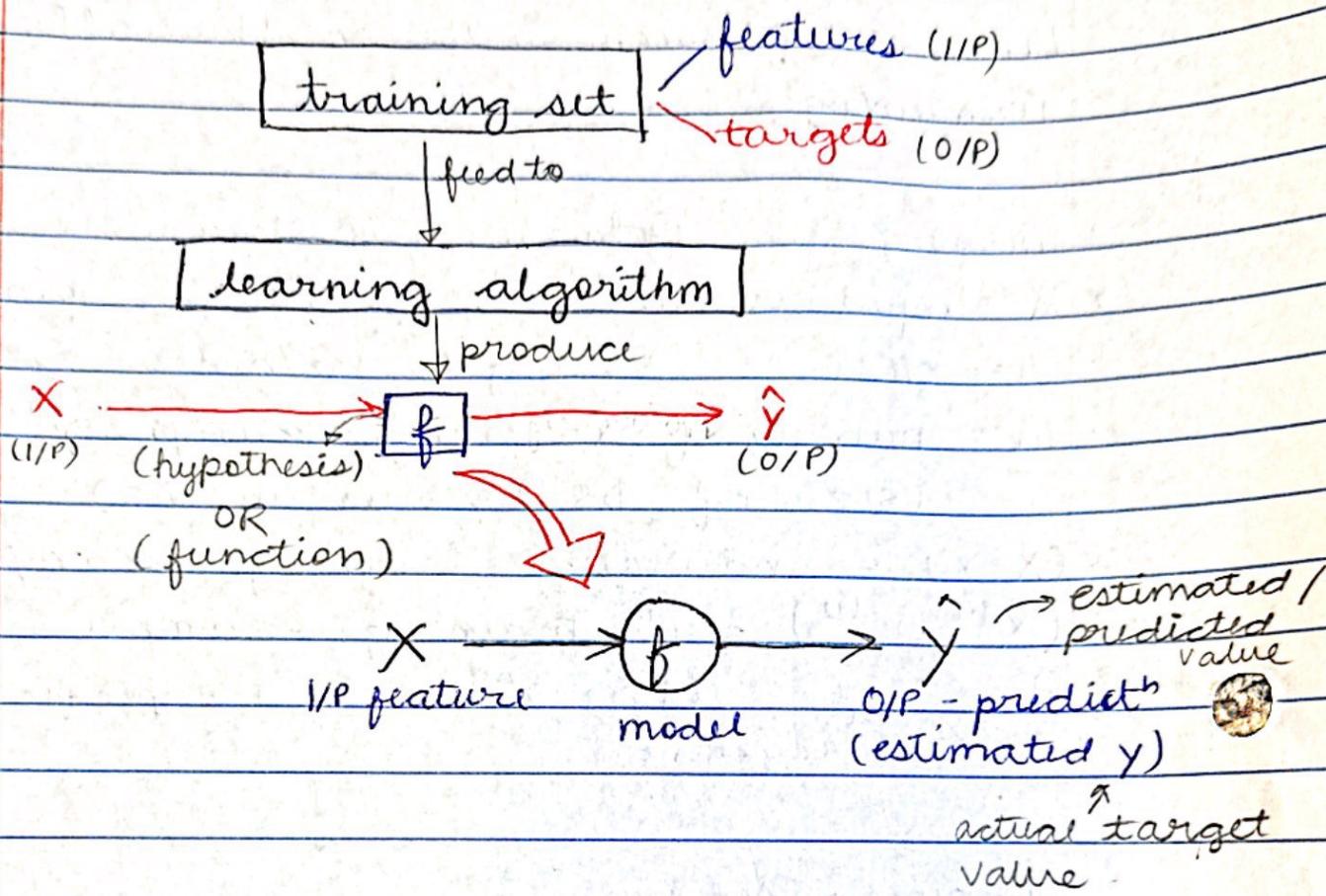
Terminology:

- **Training set**: Data used to train the model
- X = 'input' variable, I/P features, features
- Y = 'O/P', " , target variables
- m = number of training example
(size of DS / rows in DS)
- (X, Y) = individual training example
- $(X^{(i)}, Y^{(i)})$ = i^{th} training example
 $\nearrow \neq X^i$ (1st, 2nd ... i-specific row)
(exponent)

Training set:

	x size in feet ²	y price (\$1000's)	
1)	2104	400	
2)	1416	232	$m = 47$
1)	1534	315	$X^{(1)} = 2104$
1)	'	'	$y^{(1)} = 400$
.	'	'	$X^{(1)}, Y^{(1)} = 2104, 400$
4)	3210	870	

Linear Reg (P2)



How to represent f ?

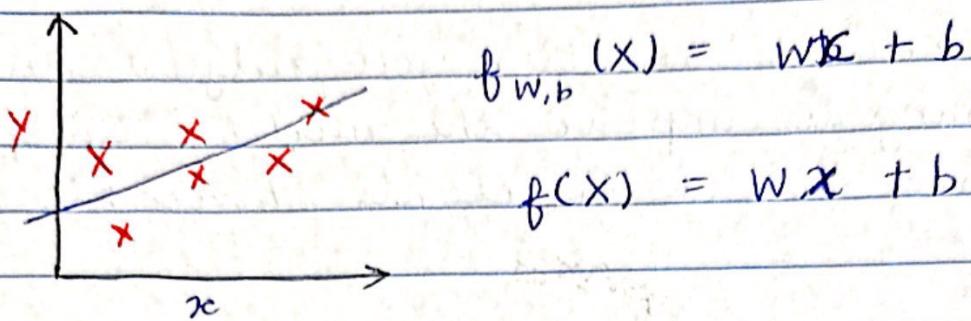
If we assume (for now) that f is a straight line:

$$y = f(x)$$

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

w and b are numbers and predicted \hat{y} will depend on x (1/P feature)

$f_{w,b}(x) = f$ is a fcn, which takes x and depending of w & b it gives y .



This is a linear line : linear regression with one variable / Univariate linear reg.
 ↪ (single feature x)

- - - LAB - - -

Cost Function

Cost func tells us how well a model is doing
 So that we can try to make it better

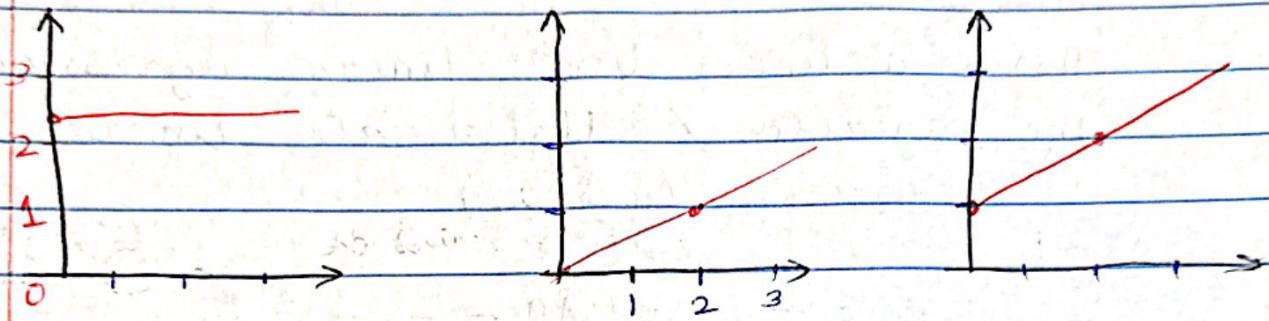
Training set

size in feet ² (x) features	price \$1000's target
2104	460
1416	232
:	:

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

w, b = parameters

Parameters: Variables you adjust during training to improve the model. They are also known as co-efficients / weights.



$$f(x) = 0 \cdot x + 2.5$$

$$\hat{y} = 2.5$$

$$w = 0$$

$$b = 2.5$$

$$f(x) = 0.5 \cdot x + 0$$

$$w = 0.5$$

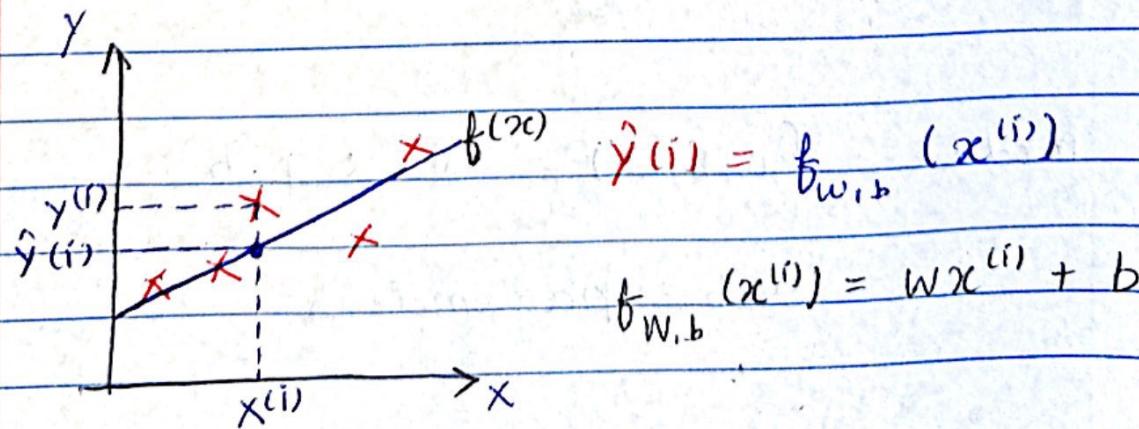
$$b = 0$$

$$f(x) = 0.5x + 1$$

$$w = 0.5$$

$$b = 1$$

Linear Regression: - choose values of params w and b such that straight line you get for $f(x)$ somehow fits the data well passing roughly / somewhere close to almost all train^g ex



Ques: How do you chose values of w and b such that value of $\hat{y}^{(i)}$ is close to $y^{(i)}$

predicted

target

For this, we will see how well a line fits training data. This is measured by cost function.

Cost Function

Squared Error Cost Function

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

of training ex
(2) sq of error terms

average calc
error

because we will have large DS.

Since \hat{y} is predicted value, we can rewrite as

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

Cost function intuition

Recap:

Model

$$f_{w,b}(x) = w_1 x + b$$

Parameters

$$w, b \quad \checkmark \quad t \quad t$$

Cost f(n)

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

Goal

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

Simplified Version:

lets assume that $b = 0$ in our eqn

$$f_w(x) = wx$$

$$b = \emptyset$$

($b = 0 \therefore$ line pass through origin)

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

\nwarrow $wx^{(i)}$ (since $b = 0$)

Visualization:

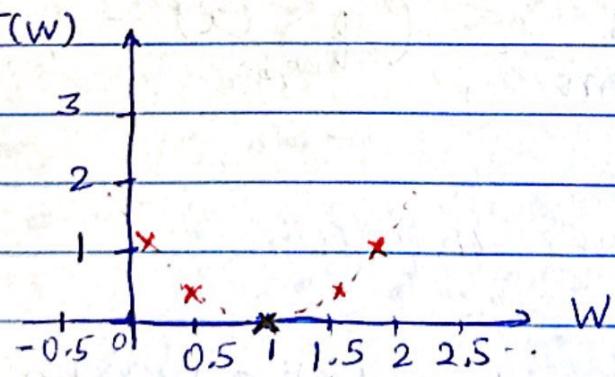
lets suppose we have points as

$(1, 1), (2, 2), (3, 3) \dots$

we start by keeping $w = 0, 0.5, -1, 1 \dots$

and plot graph b/w w and $J(w)$

for cost:



Goal of linear regression:- Minimize $J(w, b)$

Gradient Descent

$$w = w - \alpha \left[\frac{\partial}{\partial w} J(w, b) \right]^*$$

Where

α = learning rate

It controls how big step do we take. If α is large, this means GD is aggressive and we take huge steps & α is too small means very small steps.

* = Derivative term of cost fn J . It tells us about direct of learning step.

α and * together tells about size of steps

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

Correct: Simultaneous update of parameters

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

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

$$w = \text{tmp_w}$$

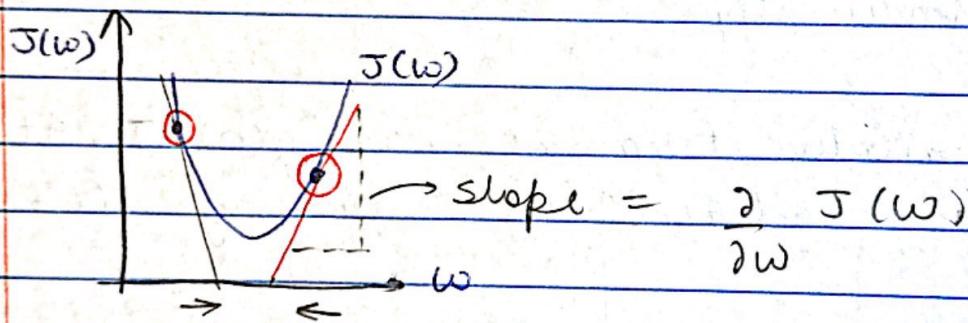
$$b = \text{tmp_b}$$

Intuition :-

Repeat Until Convergence ?

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

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



Learning rate is always +ve no.

So, if slope is +, we decrease w and move left. If slope is -ve negative, we move r/s

If α is :

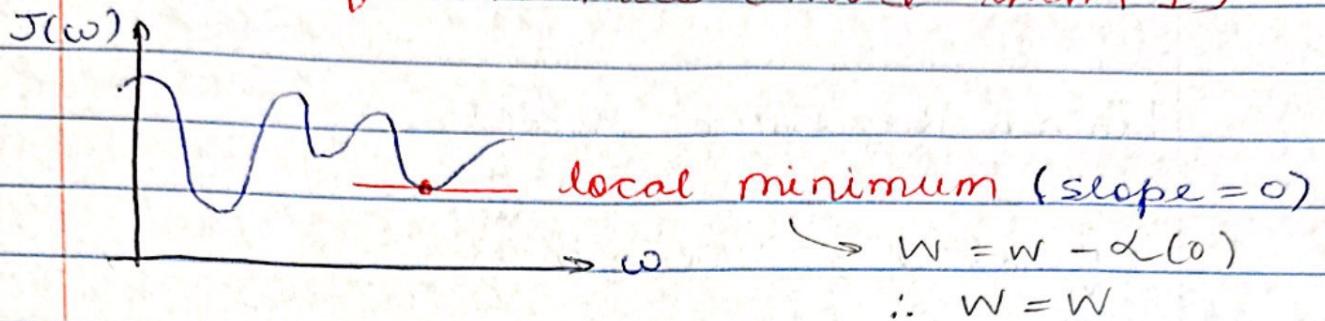
① too small

→ GD will be very slow

② too large

→ may fail to converge
→ overshoot & never reach min

In case of 2 minmas (more than 1)



As we approach minima, slope & i.e derivative becomes +

Near local minima

- Derivative gets smaller
- Update steps become smaller

At local minima GD leave value of w unchanged
Then further GD steps do nothing

Gradient Descent for Linear Algo:

Linear Regression Model

$$f_{w,b}(x) = w_1 x + b$$

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

1

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

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

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

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



Optional [Math Derivation]

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

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

$$\begin{aligned} \frac{d}{dw} (wx + a)^2 & \text{ derivative of } = \frac{1}{2m} \sum_{i=1}^m (\underline{wx^{(i)} + b} - y^{(i)}) \cancel{2x^{(i)}} \\ & = 2(wx + a) \cdot \cancel{wx} \\ & = \frac{1}{m} \sum_{i=1}^m (f_{w,b}(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned}$$

Similarly, we get for $\frac{\partial}{\partial b} J(w, b)$
and put them in GD algo.

Repeat until convergence {

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

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

3

with simultaneous update of w and b

This type of curve is known as convex $f(x)$, which have only one global min & we always converge to it.

(optional)
--- LAB WORK VDO ---

WEEK 2 Multiple Variables

Linear Regr with multiple variables

Price (1000's \$)	Size in feet ²	# bedrooms	# floors	Age of home
y	x_1	x_2	x_3	x_4
①.	②	③.	④	

x_j = j^{th} feature

n = total no. of features (here 4)

row vector $\vec{x}^{(i)}$ = features of i^{th} training example.

$\vec{x}_j^{(i)}$ = value of feature j in i^{th} example
 $\downarrow j^{\text{th}}$ feature of i^{th} sample

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

now, since we have multiple features:

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

eg:

\vec{x}_1	\vec{x}_2	\vec{x}_3	\vec{x}_4	\vec{x}_5
size	#bedrooms	#floors	years	

$$0.1x_1 + 4x_2 + 10x_3 + -2x_4 + 80$$

if house is 1 yr old, its price = price - base

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

$$\vec{w} = [w_1 \ w_2 \ w_3 \ w_4 \dots]$$

b is a number

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

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

multiple linear regression
OR
multivariate " "

Vectorization Part 1

Parameters and features: linear algebra = $\begin{pmatrix} 1 & -1 & 0 \end{pmatrix}$
python = $\begin{pmatrix} 0 & -1 & 1 \end{pmatrix}$

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

b is a number

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

* dot product in np: $f = np \cdot \text{dot}(w, x) + b$

vectorization makes code shorter

Vectorization part 2

W10 Vect.

for j in range(0, 16):

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

With Vect:

$$\text{np. dot}(w, x)$$

$$\text{to } f = w[0] * x[0]$$

$$t_1 \quad f + w[1] * x[1]$$

$$t_2 \quad \vdots$$

$$\text{to } \begin{bmatrix} w[0] \\ w[1] \\ \vdots \\ w[15] \end{bmatrix}$$

$$\text{dot } \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[15] \end{bmatrix}$$

$$\begin{array}{c} t_1 \swarrow \\ \boxed{w[0] * x[0]} + \boxed{} \end{array} \quad \vdots$$

t_{15}

take to to t_{15}

takes to & t_1 only

Vectorization enables parallel processing

c_1, w_2 ↗

$c_0, 1$

Gradient Descent for multiple LR:

n features ($n \geq 2$)

$$(\vec{y} = \vec{w} \vec{x} + b)$$

repeat {

$$\frac{\partial}{\partial w_1} J(\vec{w}, b)$$

$$j=1 \quad w_1 = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f(\vec{w}, b(x^{(i)})) - y^{(i)}) x_1^{(i)}$$

$$j=n \quad w_n = w_n - \alpha \frac{1}{m} \sum_{i=1}^m (f(\vec{w}, b(x^{(i)})) - y^{(i)}) x_n^{(i)}$$

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

Simult. update (w_1, w_2, \dots, w_m) and b

	Previous notation	Vector notation
Parameters	w_1, \dots, w_n b	\vec{w} ↪ vector of length n $w = [w_1 \dots w_n]$ b still a number
Model	$f_{\vec{w}, b}(\vec{x}) = w_1 x_1 + \dots + w_n x_n + b$	$f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$
Cost function	$J(w_1, \dots, w_n, b)$	$J(\vec{w}, b)$ ↪ dot product

Gradient descent

```

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

```

```

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

repeat { One feature

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

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

simultaneously update w, b

n features (*n* ≥ 2)

$j=1$ repeat {
 $\underline{w_1} = w_1 - \alpha \frac{1}{m} \sum_{i=1}^m (f_{\vec{w}, b}(\vec{x}^{(i)}) - y^{(i)}) \underline{x_1^{(i)}}$
 \vdots
 $\bullet = \bullet$ } $\rightarrow \frac{\partial}{\partial w_1} J(\vec{w}, b)$

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

w_i (for $j = 1, \dots, n$) and b

→ P(1)

Alternative to gradient descent

Normal Equation

- only for linear regression
- Solve for w and b w/o iteration

Disadv:

- doesn't solve other (can't generalize)
- Slow if # features is large ($> 10,000$)

Normal eqⁿ method may be used in ML lib
that implement linear regression (in backend)

- LABS - CW2L2

Feature Scaling Part 1

Practical tips for Linear Regression

Print out

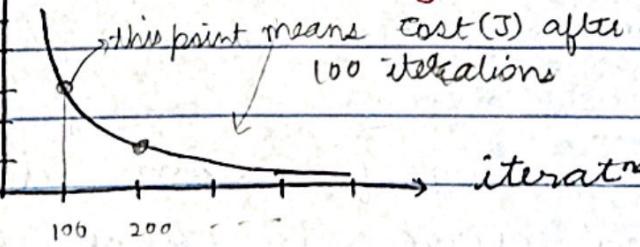
→ P(2)

small change in
 $w \rightarrow$ large change
in price (large
impact)

Part 2: Choosing α → ①

$J(w, b)$

learning curve



→ ②

Auto convergence test

let ϵ be 10^{-3} (0.001)

If $J(\vec{w}, b)$ decreased by
 $\leq \epsilon$ in one iteration,
declare convergence
(found w & b close
to minimal)

better. { gives adv warning
& visualizat^n}

Feature and parameter values

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

x_1 : size (feet²) x_2 : # bedrooms
 range: 300 – 2,000 range: 0 – 5
 ↓ ↓
 size #bedrooms

large small

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

size of the parameters w_1, w_2 ?

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

$$\widehat{\text{price}} = \frac{50 * 2000}{100,000\text{K}} + \frac{0.1 * 5}{0.5\text{K}} + \frac{50}{50\text{K}}$$

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

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

small large

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

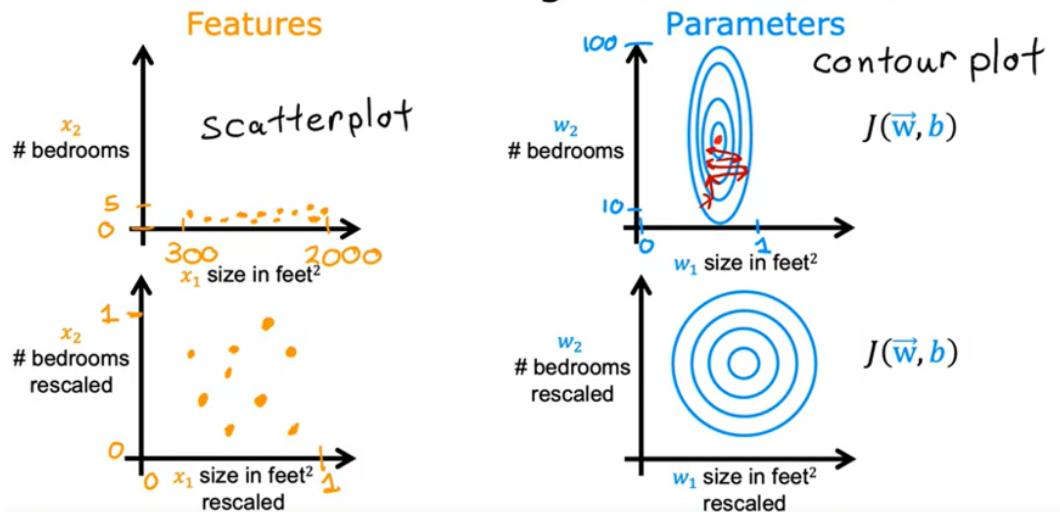
$$\widehat{\text{price}} = \$500\text{k} \quad \text{more reasonable}$$



Feature size and parameter size

	size of feature x_j	size of parameter w_j
size in feet ²	↔	↔
#bedrooms	↔	↔↔

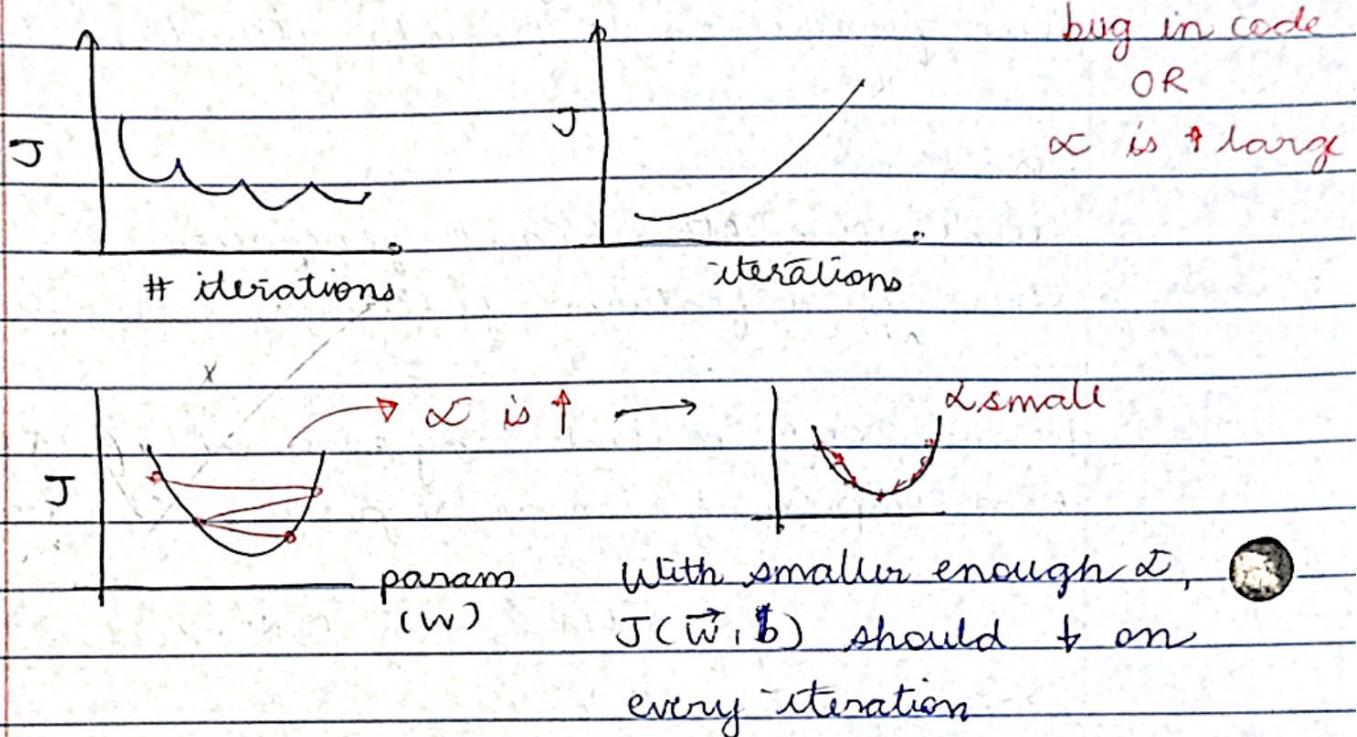
Feature size and gradient descent



$J(w, b)$ i.e cost should always decrease in every iteration

Choosing good learning rate

Problems with GD:



Always do:

Keep α very small & check if $J(\bar{w}, \bar{b})$ [cost] is decreasing in every iteration. If not, there's a bug in code

try: (0.001 0.01 0.1 1)

although these aren't best & may take lots of iteration, but can tell us about buggy code

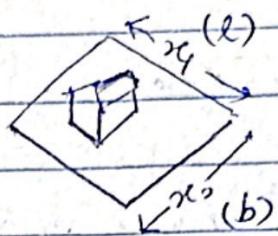
∴ For true α , run for some iterats and plot learning curve to check for improvement

- - IMP - LAB - IMP - -

CIW2 L3: Feature scaling and learning Rate

Feature Engineering

$$f_{\bar{w}b}(x) = w_1 \underbrace{x_1}_{\text{length}} + w_2 \underbrace{x_2}_{\text{width}} + b$$



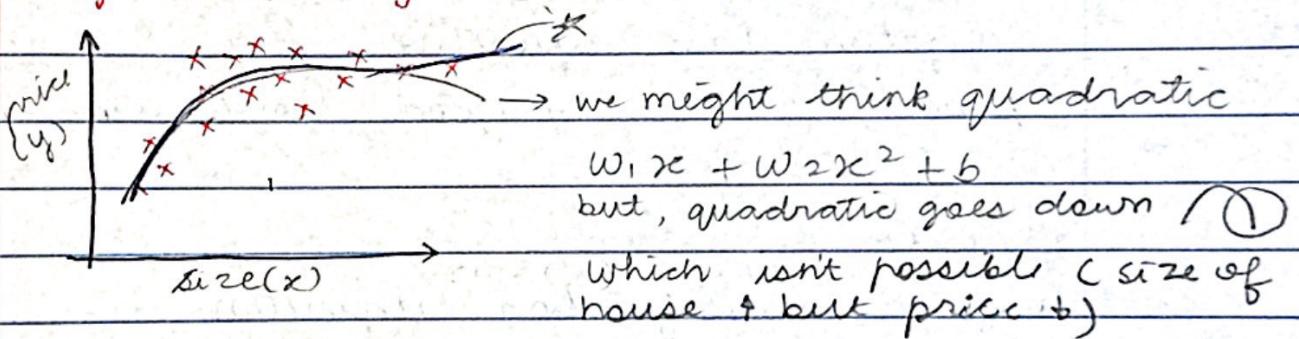
area = $l \times b$

$x_3 = x_1 x_2$ (new feature)

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

defi: Using intuition to design new features by transforming or combining original features

Polynomial Regression



* we might use: $f_{\bar{w}b}(x) = w_1 x + w_2 x^2 + w_3 x^3 + b$

- - - LAB - - -

- - - X - - - WEEK 2 DONE - X -

WEEK 3

Classification

Ques

Ans (y)

Is email spam	no / 0	yes / 1
Is transact' fraudulent	no / 0	yes / 1
Is tumor malignant	no / 0	yes / 1

Here, 'y' can have only 2 values
2 possible classes = categories

{ yes } BINARY
{ no } CLASSIFICATION

False

0

negative class

≠ bad

True

1

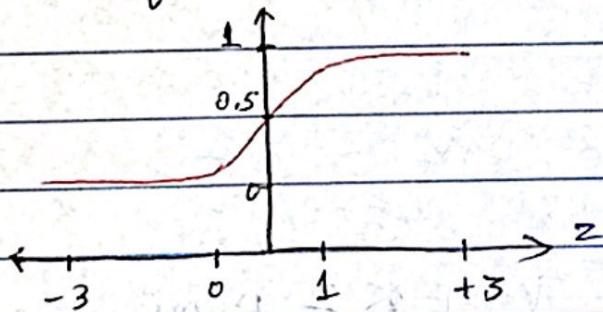
positive class

≠ good (but ans to question)

Logistic Regression

It is a classification algorithm

Sigmoid function



sigmoid fn)

logistic function

O/P b/w 0 and 1

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

$e = 2.7$

If we take $f_{\vec{w}, b}(\vec{x})$

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

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

∴ for logistic regression

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

Interpretation of logistic regression O/P

= Probability that class is one '1'

Example

x is 'tumor size'

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

$f_{\vec{w}, b}(\vec{x}) = 0.7 \rightarrow$ this means 70% chance
that y is 1 (malign. tumor)

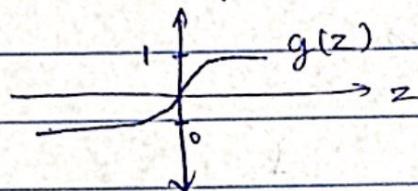
then $f_{\vec{w}, b}[(\vec{x}) \& y=0] = 0.3$
 $\curvearrowright 0.7 + 0.3 = 1$

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

→ Probability that $y=1$ given \vec{x} and parameters \vec{w} and b

Decision Boundary

LR Recap



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

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

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

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

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

$$\begin{aligned} & 0.7 & 1 - 0.7 & = 0.3 \\ & \underbrace{0 \quad 0.5 \quad 1} & & \left. \begin{array}{l} \text{if } f_{\vec{w}b}(\vec{x}) \geq 0.5 \\ \hat{y} = 1 \\ \text{else} \\ \hat{y} = 0 \end{array} \right\} \end{aligned}$$

Now, when is $f_{\vec{w}, b}(\vec{x}) \geq 0.5$

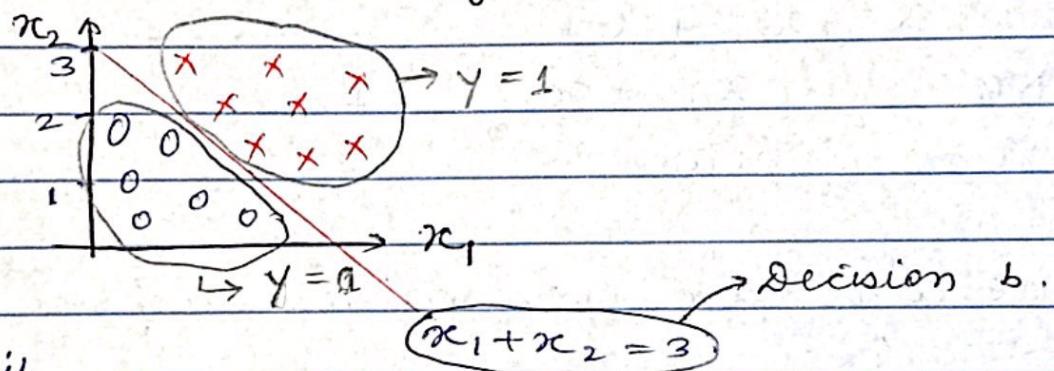
this means when is $g(z) \geq 0.5$
when (by graph) \rightarrow

$$z \geq 0$$

i.e. $\vec{w} \cdot \vec{x} + b \geq 0$ then $\hat{y} = 1$

if $\vec{w} \cdot \vec{x} + b < 0$ then $\hat{y} = 0$

Decision boundary



if

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

let $w_1 = 1$ & $w_2 = 1$ & $b = -3$

$$\therefore z = x_1 + x_2 + b$$

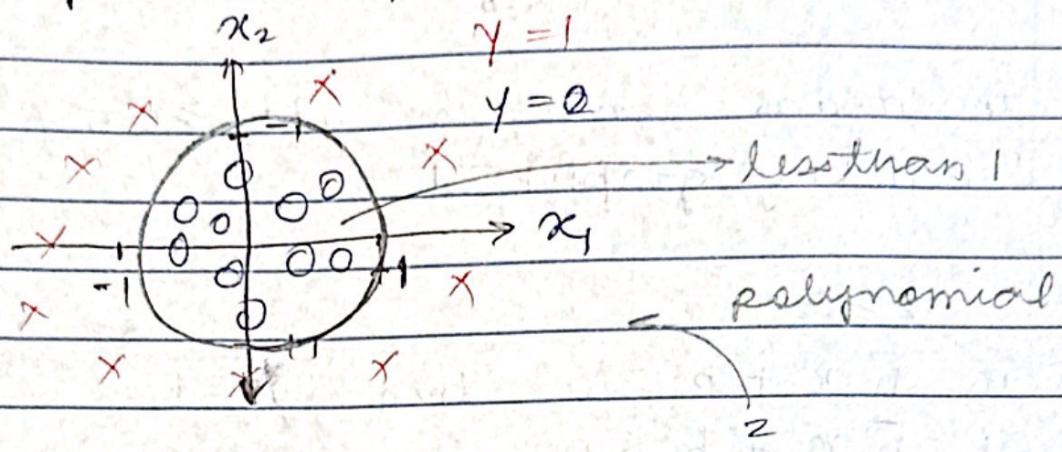
Decision boundary: $z = 0$

$$\therefore x_1 + x_2 - 3 = 0$$

$$x_1 + x_2 = 3$$



Complex example



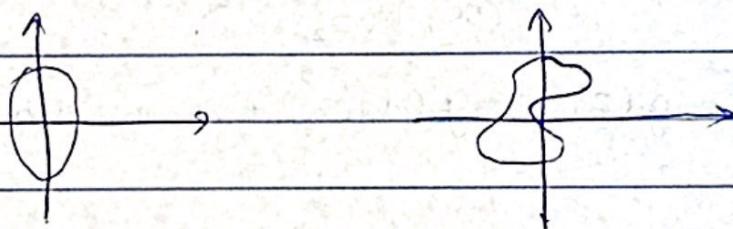
$$f_{\vec{w}, b}(\vec{x}) = g(z) = \underbrace{w_1 x_1^2}_z + \underbrace{w_2 x_2^2}_z + b \quad -1$$

$$z = x_1^2 + x_2^2 + (-1)$$

$$\text{DB: } \underbrace{x_1^2 + x_2^2}_z = +1$$

equat of circle

We can have even more complex DB -



Cost Function (Log. Regression)

Cost function gives us a way to measure how well the parameters fits the training data, giving a way to choose better parameters

Training set

	tumor size	patients age	malignant
	x_1	x_n	y
$i=1$	10	52	1
$i=2$	2	73	0
:	:	:	0
$i=m$:

$i = 1, 2 \dots m$ (training examples)

$j = 1, 2 \dots n$ (features)

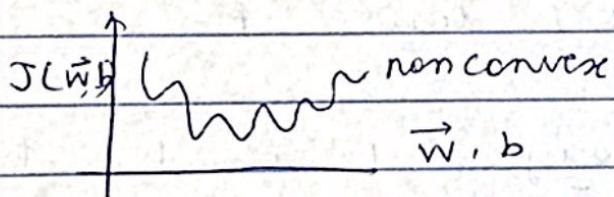
targd y is 0 or 1

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

How to choose w & b ?

$$\rightarrow \underset{\vec{w}, b}{\text{J}}$$

If we plot cost, as in linear reg then we get non-convex curve with multiple local minimas

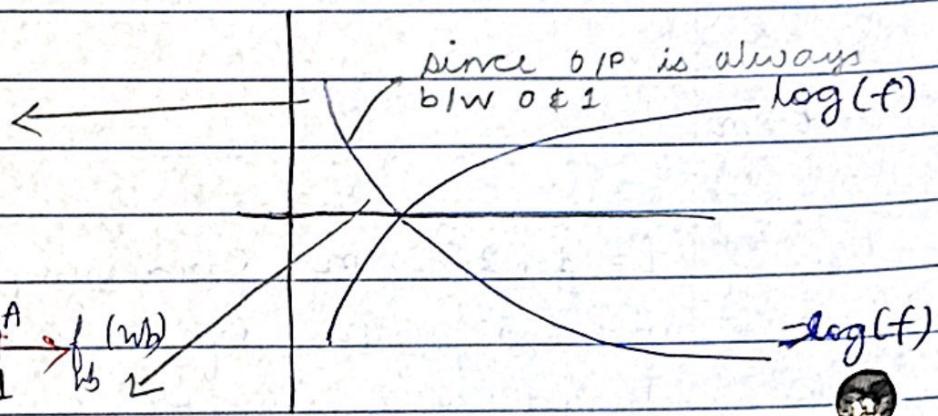
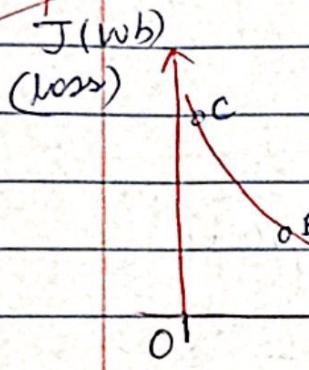


\therefore Squared Error Cost fun
isn't a good choice

Logistic loss function

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

For $y=1$



when $y^{(i)} = 1$

at point A: predicted value = $f_{\bar{w}, b}(\bar{x}^{(i)}) \rightarrow 1$

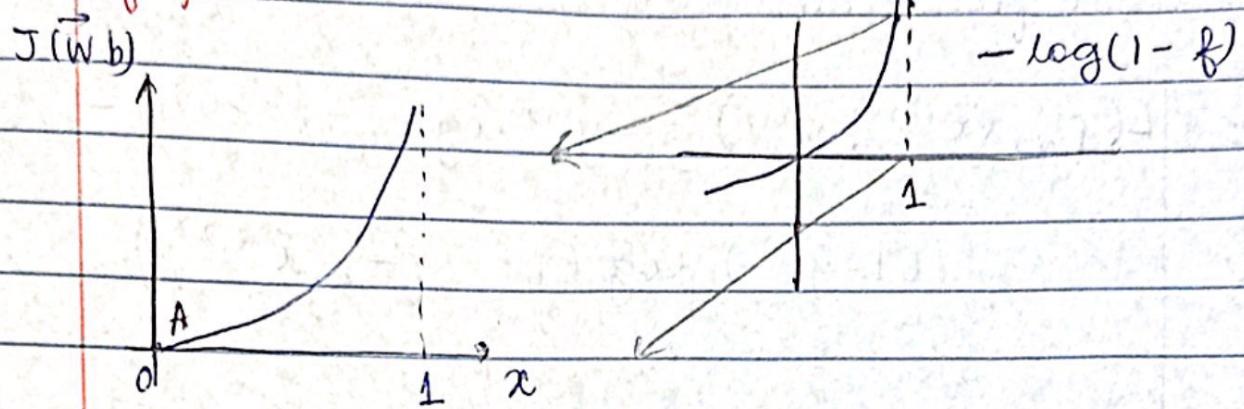
\therefore loss $\rightarrow 0$

at point B: if $f_{\bar{w}, b}(\bar{x}^{(i)}) = 0.5$

ie it says 50% chances, but it is really malignant
then we see at B, loss is ↑

at C, if it says 0.1, ie 10% chance but
actually it is malignant, then we see at
A, loss is really ↑.

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



point A Similar to $y=1$, if $f_{w,b}(\vec{x}) \rightarrow 0$, loss $\rightarrow 0$

Cost =

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

loss

: Find w, b to \downarrow cost

(previous page)

Simplified Cost Function :

$$\begin{aligned} L(f_{w,b}(x^{(i)}), y^{(i)})^2 &= \\ &\quad -y^{(i)} \log(f_{w,b}(\vec{x}^{(i)})) - (1-y^{(i)}) \log(1-f_{w,b}(\vec{x}^{(i)})) \end{aligned}$$

Just a combined form of loss fun
of on previous (-2) page

cost function:

convex only 1 minima

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

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

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

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

Although there are so many cost fns, we choose this because of maximum likelihood

Implementing GD:

Training LR:

① Find \vec{w} and b

② Then give model $\vec{x} = \text{age, size of tm}$
q other param

③ Model predict

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

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

For implementation, we repeat until loss is \downarrow with simultaneous update of w and b .

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

repeat {

$$\text{for } j = 1 \dots n \text{ for all features}$$

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

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

}

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

↑ ^{i} feature
of train example

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

→

Equat^r on previous page might look exactly similar to too linear R. But, funcⁿ f is different in both

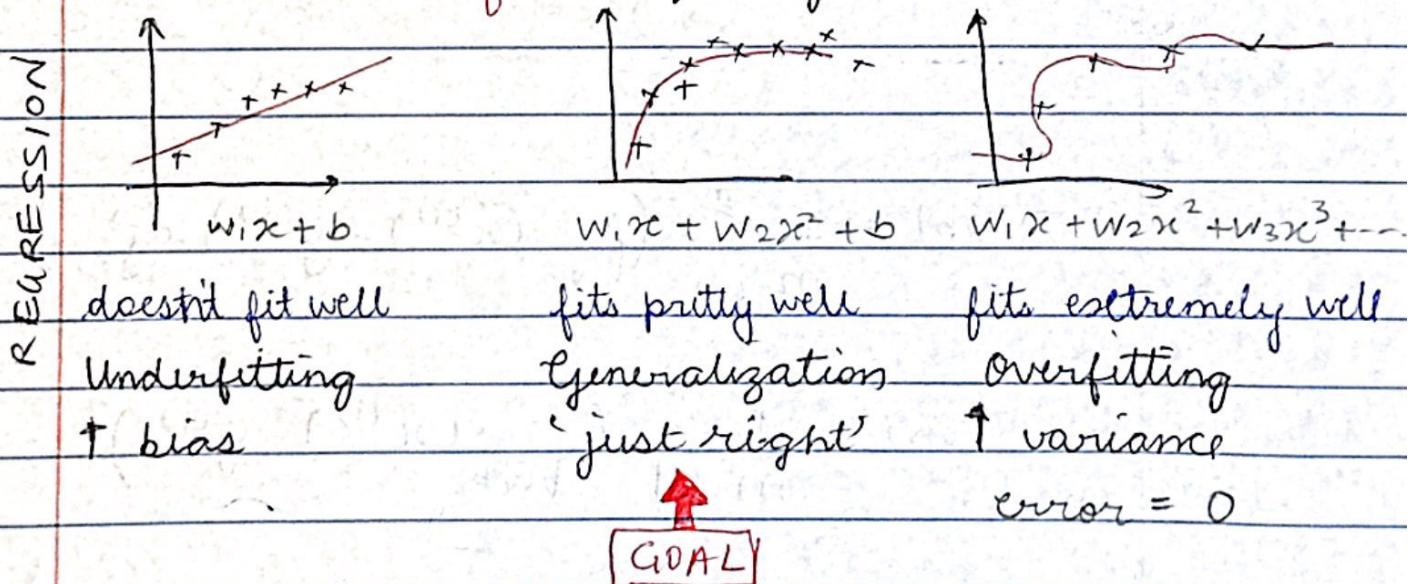
Linear Regression $f_{\vec{w}, b}(\vec{x}) = \vec{w} \cdot \vec{x} + b$

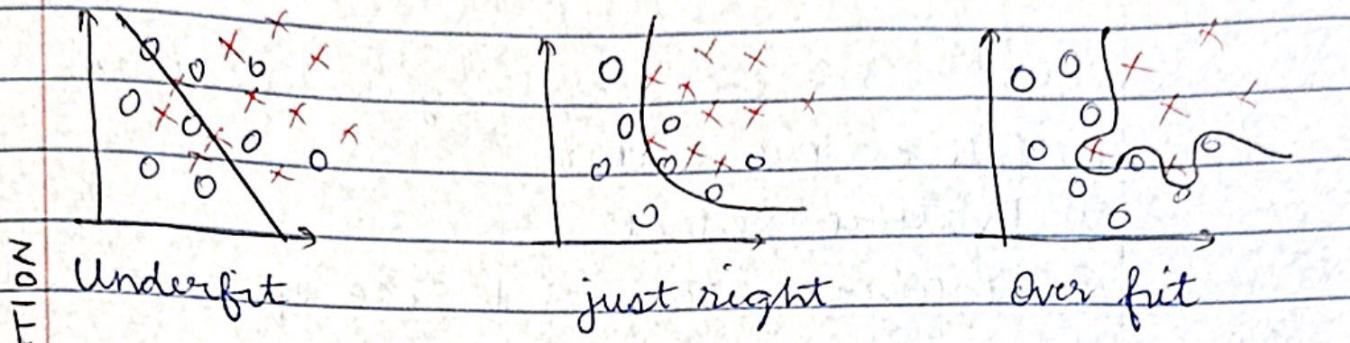
Logistic $f_{\vec{w}, b}(\vec{x}) = \frac{1}{1 + e^{-(\vec{w} \cdot \vec{x} + b)}}$

Same concepts

- Learning curve monitoring
- Vectorized implem
- Feature scaling (to speed up G.D.)

The Problem of Overfitting





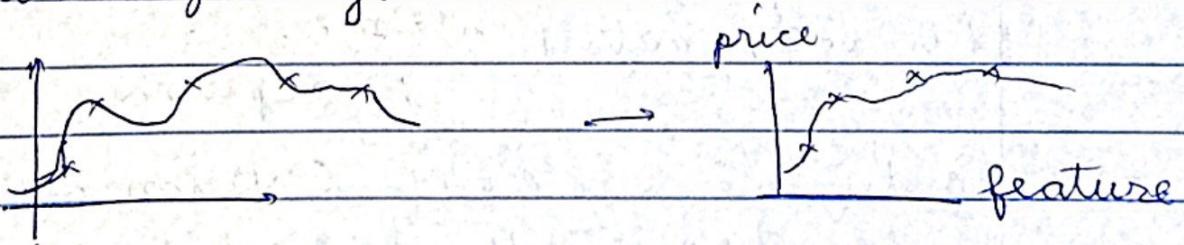
Addressing overfitting

① Collect more training data

② Check if we can use fewer features

A greater # features will + training examples causes overfitting. To avoid, use feature select to select subset of most relevant features

③ Regularization : When we have (suppose 100) features, and all are IMP (cant eliminate any). In this case, instead of just removing (if being harsh) we + the value of parameter. (reduce impact of that feature gently).



$$f(x) = 28x + 385x^2 + \underbrace{39x^3 + 174x^4 + 100}_{\text{elimination}}$$

means setting
 $w=0$ instead of

174 (\uparrow value)

$$f(x) = 128x + 0.23x^2 + 0.00014x^3 +$$

$$0.00001x^4 + 10$$

smaller values of w_j

Recap - Addressing regularization.

- ① Collect more data
- ② Feature Select'
- ③ Regularization (+ size of parameters)

Cost Function with Regularization

If we have w_1, w_2, \dots, w_{100} parameters, we don't know which param is imp & which isn't. so we penalize every param by:

→ gently mild them

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m \left(f_{wb}(x^{(i)}) - y^{(i)} \right)^2 + \lambda \underbrace{\sum_{j=1}^m w_j^2}_{\text{reg. term}}$$

λ = regularization param

only to make sure both three terms are scaled by $1/2m$

Some might have $\frac{1}{2m} \sum w_j^2$
(it doesn't matter)

→ optional

λ defines tradeoff (option of how one balance b/w fitting of data curve to keeping w_j smaller).

if $\lambda = 0$ (if it will overfit)

if λ is $\uparrow\uparrow$ (it will underfit)

if λ is just right (it gives generalizatn)

Regularized Linear Regression

orig cost
fun)

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

reg. term

∴ Now GD algo remains same

but $\frac{\partial J(\vec{w}, b)}{\partial w_j}$ changes and $\frac{\partial J(\vec{w}, b)}{\partial b}$ changes

$\frac{\partial J(\vec{w}, b)}{\partial b}$ — remain same

(bcz we don't reg. b)

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

Print out

Regularized logistic Reg

similar to above for logistic as well

$$J(\vec{w}, b) = \min_{\vec{w}, b} \sum_{i=1}^m \text{cost fun}(\dots) + \frac{\lambda}{2m} \sum_{j=1}^n w_j^2$$

& the change in $\frac{\partial J(\vec{w}, b)}{\partial w_j}$ is

$$\Rightarrow \frac{1}{m} \sum_{i=1}^m [f_{wb}(\vec{x}^{(i)}) - y^{(i)}] x_j^{(i)} + \frac{\lambda}{m} w_j$$

logistic fun)

Urvi Jain