

Gradient descent Algo:

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

update the parameter  $w$  by subtracting with the previous value of  $w$  with  $\alpha \cdot \frac{\partial}{\partial w} J(w, b)$

$\alpha$  = the learning rate (small positive number between 0 and 1)

$\alpha \uparrow$  indicates that the descent is high.

$\alpha \downarrow$  indicates that the descent is with baby steps (low)

$\frac{\partial}{\partial w} J(w, b) =$  derivative of term of loss function  $J$ .

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

continue these till there is no change. ie till it reaches local minima

x Simultaneously update  $w$  and  $b$ .

Correct : Simultaneous way:

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

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

} after this block gets executed  
then do in  $w \leftarrow \text{temp\_}w$  (update)

Δ Do not update immediately.

$$w = \text{temp\_}w$$

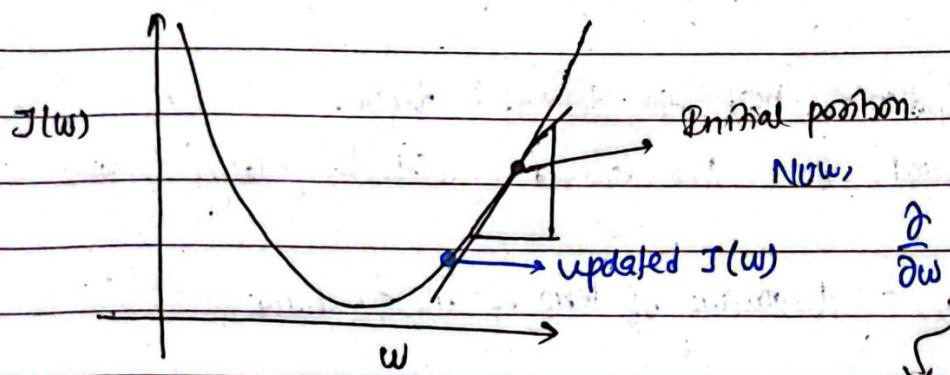
$$b = \text{temp\_}b$$

Now,

Now for a cost function  $J$  of GIA one parameter  $J(w)$

$J(w)$

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



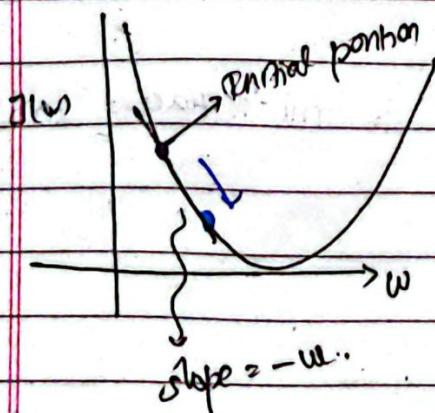
Similarly;

So,

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

i.e.  $w$  decreases.

and updated  $w$  shifts towards left.



$$w = w - \alpha \cdot (-\text{ve slope})$$

= ~~decrease~~  $w + \alpha \cdot \text{number}$ .

=  $w \uparrow$  and shift towards right

Learning rate ( $\alpha$ ):

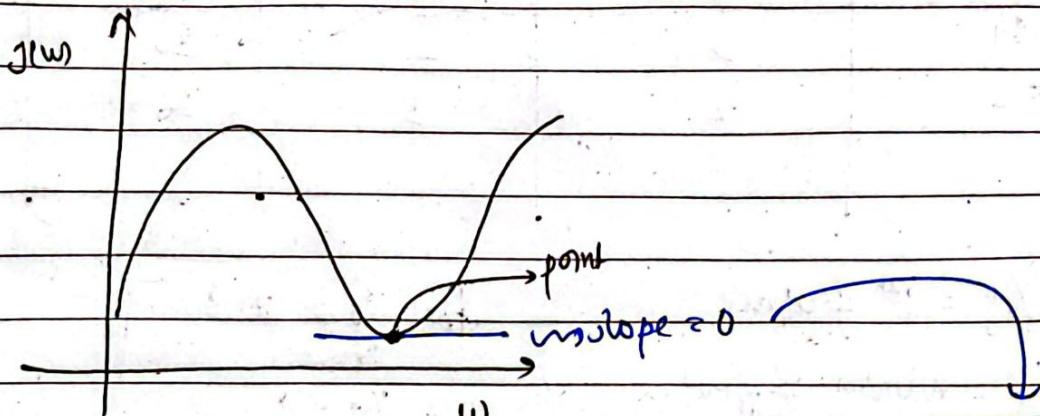
If  $\alpha \rightarrow$  too small

gradient descent may bedow.

If  $\alpha \rightarrow$  too large

may take too huge step and cross the  $\downarrow$  minimum  
which leads to never reaching minimum and may even drusse.  
↳ gradient descent may overburst.

What if the costfunction is already at 'a local minimum'?



$$\text{Now, } w = w - \alpha \cdot \boxed{\frac{d}{dw} J(w)} = 0$$

$$\boxed{w_k} = w$$

therefore no change 😊

## Multiple Linear Regression :

# With multiple variables:

Example:

Size in feet <sup>2</sup>	No of bedrooms	No of floors	Age of home	Price in \$/feet <sup>2</sup>
$x_1$	$x_2$	$x_3$	$x_4$	
2104	5	1	45	460
1416	3	2	40	800
1534	3	2	30	200
852	2	1	36	150

$x_j = j^{\text{th}}$  feature

$n = \text{number of features}$

} for  $x_1 \Rightarrow n=4$

$\begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{matrix}$

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

$$\vec{x}_j^{(i)} \rightarrow x_3^{(2)} = 2$$

Model :

Previously for single feature:

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

Now for multiple features:

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

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

$b$  = a number

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

So it can be written as:

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

$$\vec{w} \cdot \vec{x} = w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_n x_n$$

dot product

aka: multiple linear regression

(Note: Not a multivariate regression)

X. Vectorization:

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

$$= \sum_{j=1}^n w_j x_j + b$$

$$> f = 0$$

for  $j$  in range ( $0, n$ ):

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

$$f = f + b$$

Code quality

( $\therefore$ )  
okay!

But when, vectorization is used;

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

→ import numpy as np

$$> f = np.dot(w, x) + b$$



v-good!

Why vectorized codes run so fast?

Without vectorization

> for j in range (0, 16):

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

$$\text{t}_0 \quad f + w[0] * x[0]$$

$$\text{t}_1 \quad f + w[1] * x[1]$$

$$\vdots$$

$$\text{t}_{15} \quad f + w[15] * x[15]$$

$$\downarrow$$

$$\text{steps} = 16$$

Vectorization:

$$\rightarrow np.dot(w, x)$$

$$\text{to : } [w[0] \mid w[1] \mid \dots \mid w[15]]$$

$$t_1 : [x[0] \mid x[1] \mid \dots \mid x[15]]$$

multiply in parallel ...

and then with specialised hardware  
add together.

$$[w[0]x[0]] + [w[1]x[1]] + \dots + [w[15]x[15]]$$

$$\left\{ \begin{array}{l} \text{operations used} = 2 \\ \text{steps} = 16 \end{array} \right.$$

## Gradient Descent for multiple linear regression

prev. not<sup>n</sup>parameters :  $w_1, w_2, \dots, w_n$  $b$ 

$$\text{Model : } f(\vec{w}, b)(\vec{x}) = w_1x_1 + \dots + w_nx_n + b$$

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

cost fun<sup>n</sup> :  $J(w_1, \dots, w_n, b)$

Vector Not<sup>n</sup>

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

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

gradient descent:

repeat {

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

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

}

repeat {

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

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

}

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

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

 $\Rightarrow$ 

$$\begin{bmatrix} x_0^0 & x_1^0 & \dots & x_{n-1}^0 \\ \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots \\ x_0^{m-1} & x_1^{m-1} & \dots & x_{n-1}^{m-1} \end{bmatrix}$$

## 17 Practical tips for Linear Regression

\* feature scaling

↳ Enables gradient descent to run faster.

Feature & parameter values:

range: 300 - 2000

↑  
(large)

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

$x_1$ : size (feet<sup>2</sup>)  
size  
 $x_2$ : # bedrooms

$a_1$ : size (feet<sup>2</sup>).  
 $a_2$ : # bedrooms

↳ range: 0 - 5  
(small)

Now,

let's take an example case where:

House :  $x_1 = 2000$ ;  $x_2 = 5$ , price = \$500k

So, what would be the parameters  $w_1$  and  $w_2$ ?

Suppose 50 and 0.1 are the numbers but we are not sure which one of them is  $w_1$  and which is  $w_2$

So, there are 2 cases

Case I

$$w_1 = 50; w_2 = 0.1; b = 50$$

$$\begin{aligned} \widehat{\text{price}} &= 50 \times 2000 + 0.1 \times 5 + 50 \\ &= \$100,050.5k \end{aligned}$$

which is way higher

Case II:

$$w_1 = 0.1; w_2 = 50; b = 50$$

$$\begin{aligned} \widehat{\text{price}} &= 0.1 \times 2000 + 50 \times 5 + 50 \\ &= \$500k \end{aligned}$$

and is more reasonable case.

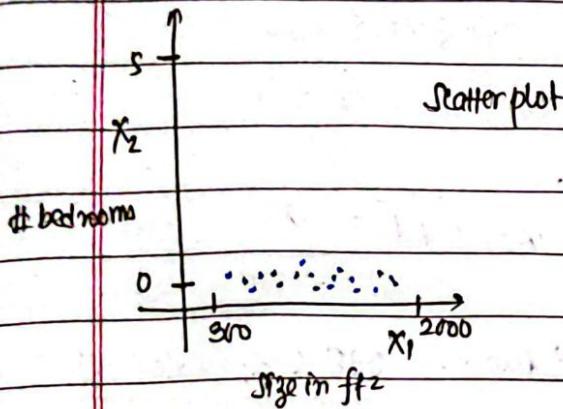
So, the conclusion:

$x_1$  had large range. So, the good model should choose  $w_1$  with less value

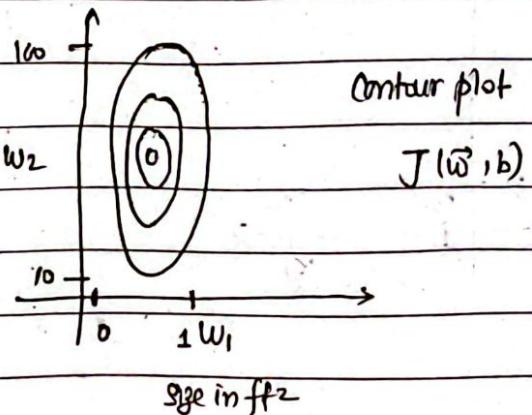
Meanwhile

$x_2$  has smaller range. So, the good model should choose  $w_2$  with greater value.

## Features



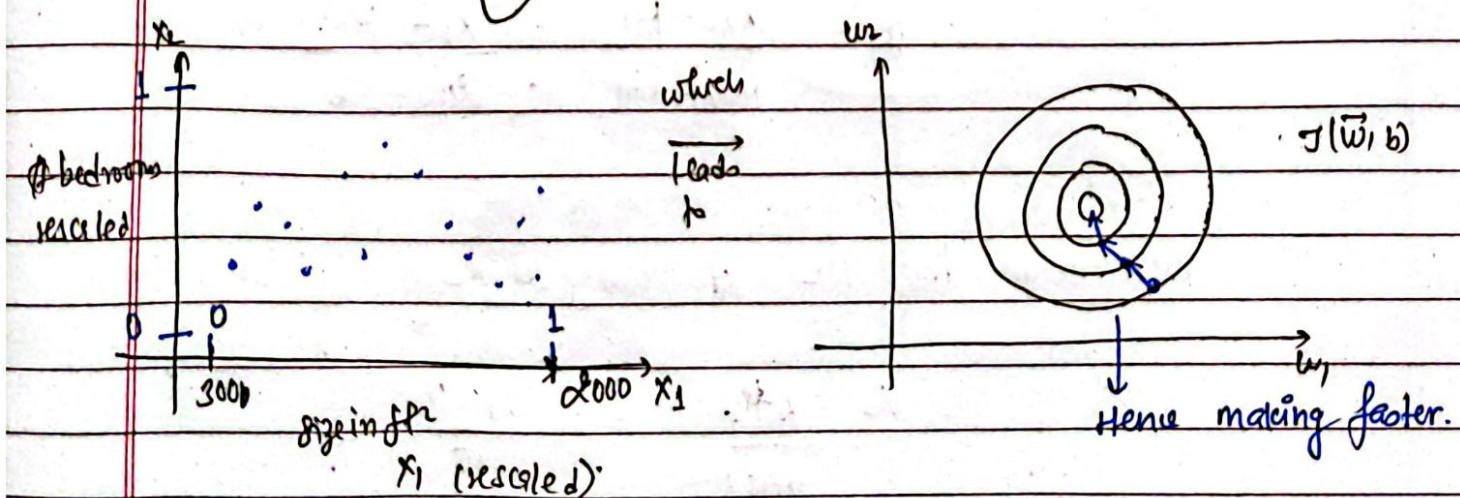
## Parameters



Now, if we were to run the gradient descent with the ~~g~~ training data as it is, since the contour plot is so tall and skinny, gradient descent may end up bouncing back and forth and may take longer to reach the center.

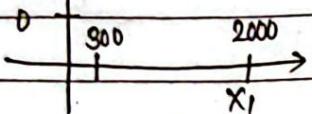
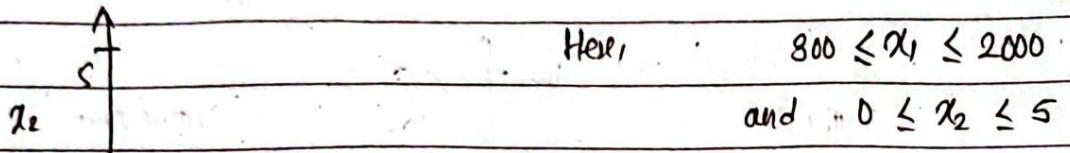


In order to solve this, we simply scale the features leading to:



Q, how do you actually scale features?

- Method 1



$$\text{So, } x_{1,\text{scaled}} = \frac{x_1 - \text{min}}{\text{max} - \text{min}} = \frac{x_1 - 800}{2000 - 800}$$

range:  $0.15 < x_{1,\text{scaled}} < 1$

and

$$x_{2,\text{scaled}} = \frac{x_2 - \text{min}}{\text{max} - \text{min}} = \frac{x_2 - 0}{5}$$

$$\text{range: } 0 < x_{2,\text{scaled}} < 1$$

- Method 2

We can also carry out mean normalization.

- ① first find out the average ( $\mu_1$ ) of training set  $x_1$   
then,

$$x_1 = \frac{x_1 - \mu_1}{\text{max} - \text{min}} = \frac{x_1 - \mu_1}{2000 - 800}$$

- ② then same for  $\mu_2$  of training set  $x_2$

$$x_2 = \frac{x_2 - \mu_2}{\text{max} - \text{min}} = \frac{x_2 - \mu_2}{5}$$

The values range from -ve to +ve.

## -Method 3.

## Z-score normalization

~~Find out ( $\mu$ ) of the data.~~ (standard deviation).  
and  $\mu_1$ .

then,

$$x_1 = \frac{x_1 - \mu_1}{\sigma_1}$$

Note: aim for about  $-1 \leq x_j \leq 1$  for each feature  
 $-3 \leq x_j \leq 3$  } acceptable ranges  
 $-0.3 \leq x_j \leq 0.3$

If  $0 \leq x_1 < 3 \rightarrow$  okay no need to rescale.

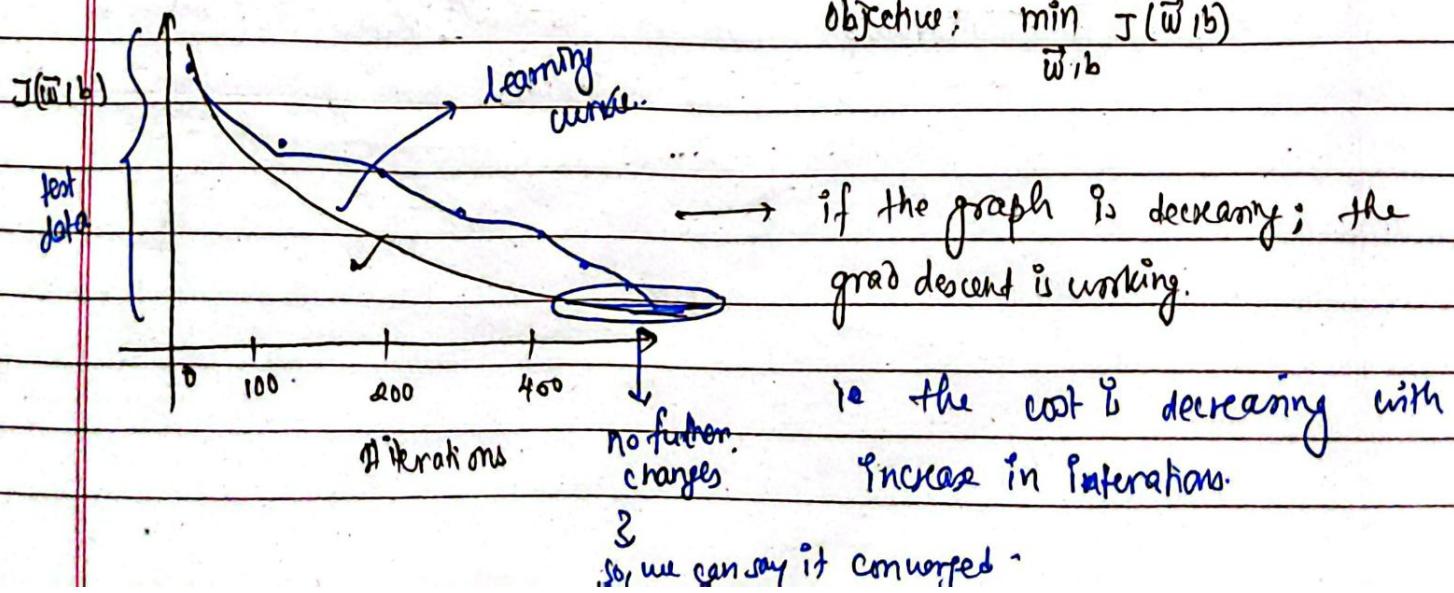
$-2 \leq x_2 \leq 0.05 \rightarrow$  okay no need to rescale

$-100 \leq x_3 \leq 100 \rightarrow$  too large  $\rightarrow$  rescale

$-0.0001 \leq x_4 \leq 0.0001 \rightarrow$  too small  $\rightarrow$  rescale

## # Checking Gradient Descent for Convergence

Objective:  $\min_{\vec{w}, b} J(\vec{w}, b)$



Similarly,

Another way of testing gradient descent:-

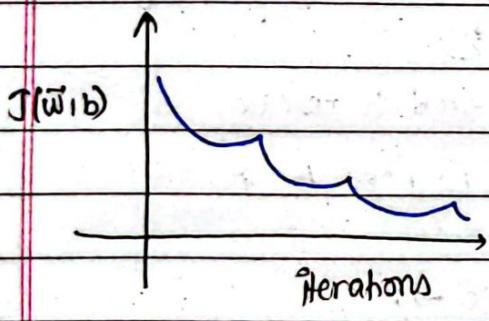
Automatic ~~convergence~~ test

Let "  $\epsilon$  " {epsilon} be  $10^{-3}$ .

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

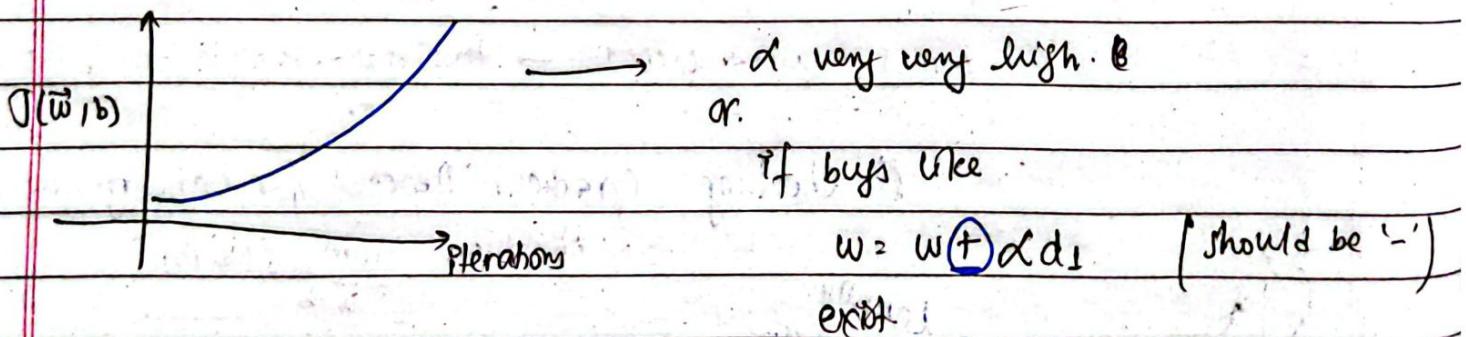
(But graphs could be better in some case)

# Choosing the Learning Rate:



Conclusion: Either there is :

- ① Bug in the code.
- ② the  $\alpha$  is too high  
Causing to overshoot.



Also, if  $\alpha$  is too small, it takes too much of time.

Values of  $\alpha'$  to try:

$$\dots 0.001 \quad 0.01 \quad 0.1 \quad 1 \dots$$

$\curvearrowleft \qquad \curvearrowright$   
 $x10 \qquad x10$

For further tuning; once the rate is found near optimal:

$$\dots 0.001 \quad 0.003 \quad 0.01 \quad 0.03 \quad 0.1 \quad 0.3 \quad 1 \dots$$

$\curvearrowleft \qquad \curvearrowleft \qquad \curvearrowleft \qquad \curvearrowleft \qquad \curvearrowleft$   
 $x3 \qquad x10 \qquad x3 \qquad x3 \qquad x3$

and so on 

\* ~~Polynomial Regression~~

\* Polynomial Regression

→ It is a type of linear regression, but it differs in the way the relationship between independent variable and dependent variable is modeled.

• Model representation:

Linear: Relationship between independent and dependent variable modeled as linear equation.

- assumes straight line relationship

Polynomial: Relationship is modeled as  $n^{\text{th}}$  degree polynomial.

• Eq<sup>n</sup>:

$$\text{Linear} = y = \beta_0 + \beta_1 x$$

$y$  ↓  
dependent       $x$  ↓  
independent

$$\text{Polynomial} = \beta_0 + \beta_1 x + \beta_2 x^2 + \dots + \beta_n x^n$$

Therefore, polynomial regression is the more flexible extension of linear regression that can capture non-linear relationships between variables.

## # Logistic regression Model.

- Single most likely used classification algorithm in the world.

Use cases of linear and logistic regression algo:

### ① Nature of output:

- Linear = Used when the output variable is continuous.
- Logistic = Used when the output variable is binary or categorical.

### ② Model Representation:

• Linear = Output is linear combination of input features

fun to

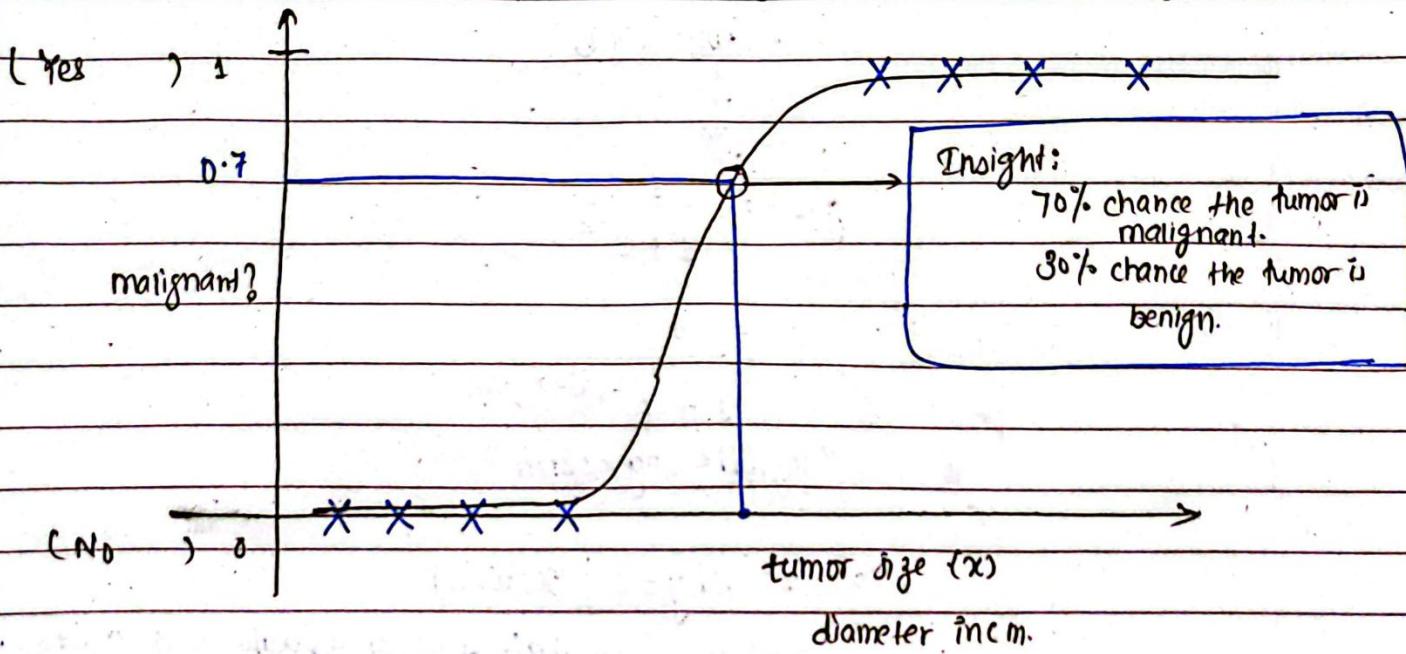
• Logistic = Uses logistic model  $\rightarrow$  probability that given input belongs to certain class.

the output is transformed to fall between 0 and 1.

Example:



Graph for classifying if the tumor is malignant.



# Sigmoid fun (Logistic fun)

output between:  
0 and 1

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

$$(0 < g(z) < 1)$$

if  $z = +ve$  large number:

$e^{-z} \rightarrow$  negligible value

which tends the curve to approach towards '1'

if  $z = -ve$  large number:

$$e^{-z} \Rightarrow e^{-(+2)} \approx e^{+2}$$

and  $\frac{1}{large\ number} \rightarrow$  tends to zero

Likewise for

$$\text{at } z=0 :$$

$$g(z) = \frac{1}{1+1} = \frac{1}{2} = 0.5 \quad \checkmark$$

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

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



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

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

"logistic regression"

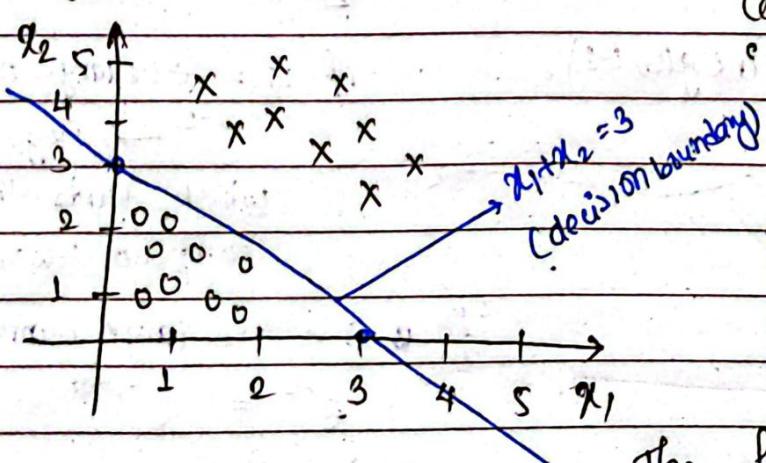
$$= p(y=1 | \vec{x}; \vec{w}, b)$$

read as: probability of  $y$  equals to 1 given  $\vec{x}$  with parameters  $\vec{w}$  and  $b$ .

Decision Boundary

(Learned previously in SVM)

Starting with examples:



Consider we have 2 features  $x_1$  and  $x_2$

$x \Rightarrow y = 1$ (+ve)
$o \Rightarrow y = 0$ (-ve)

The feature vector fun?

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

Then,

$$x_1 + x_2 = 3$$

linear eqn aka.

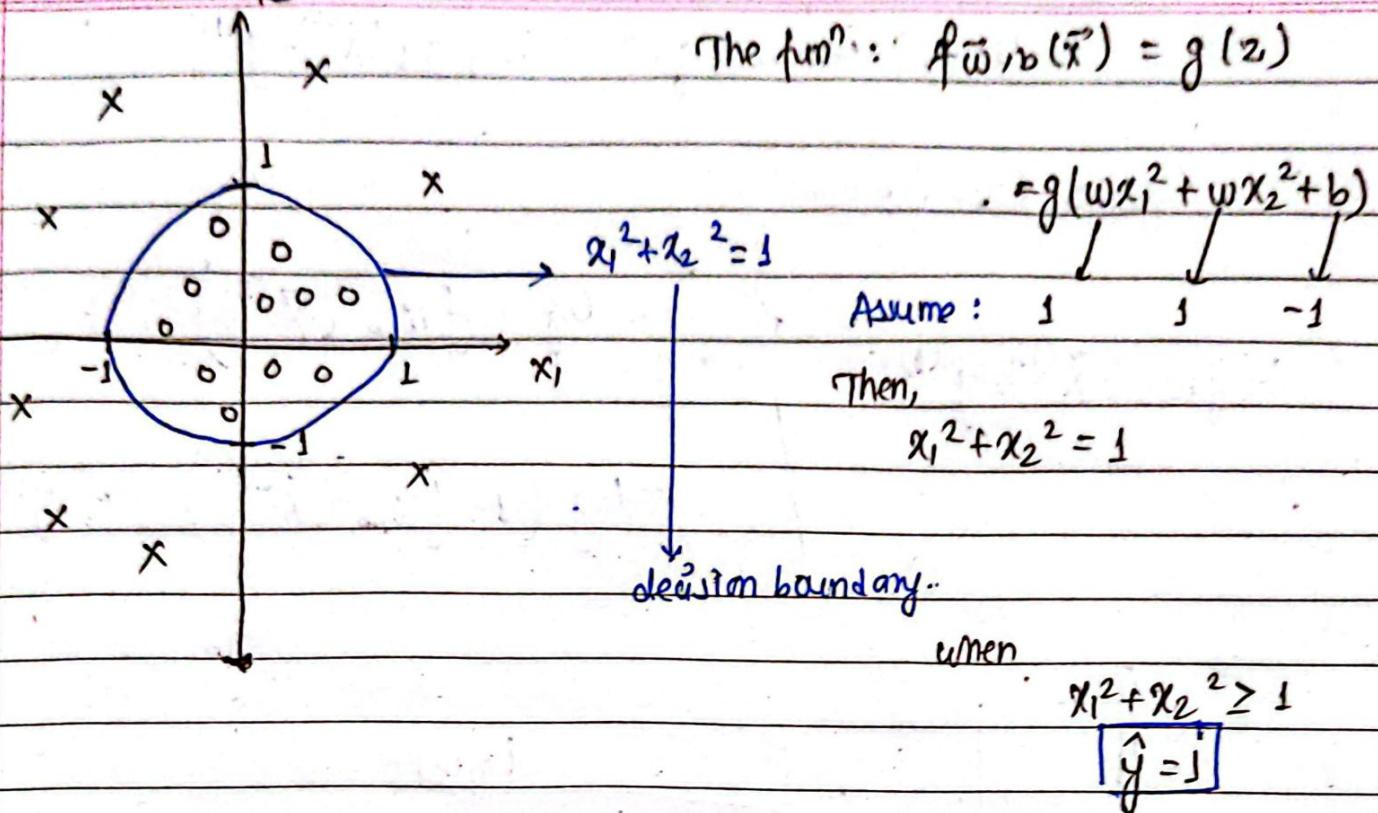
decision boundary

assuming

$$= g(w_1 x_1 + w_2 x_2 + b)$$

+      ⊥      -3

Likewise,  $x_2$



and

$$x_1^2 + x_2^2 < 1$$

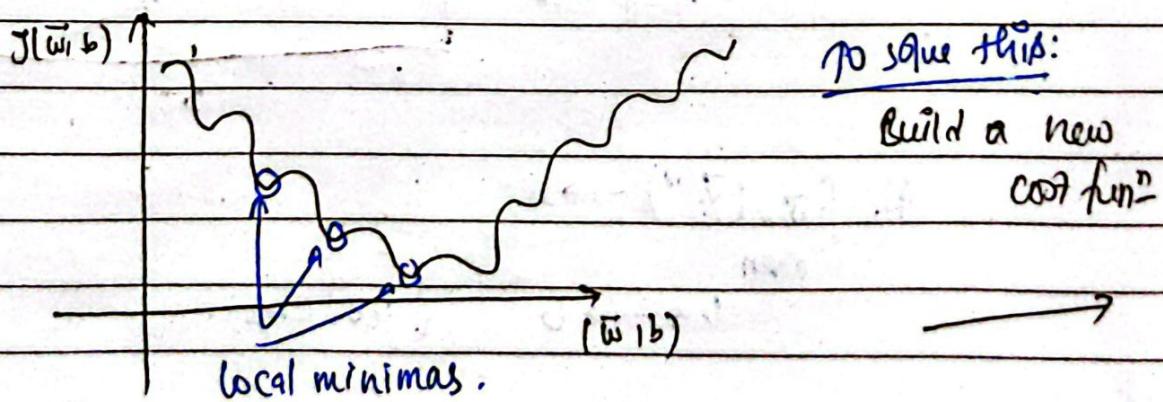
$\hat{y} = 0$

# Cost fun<sup>n</sup> for logistic regression:

→ for logistic reg<sup>n</sup> the squared error cost fun<sup>n</sup> i.e.

$$J(\bar{w}, b) = \frac{1}{m} \sum_{i=0}^m \frac{1}{2} (f_{\bar{w}, b}(\bar{x}^{(i)}) - y^{(i)})^2$$

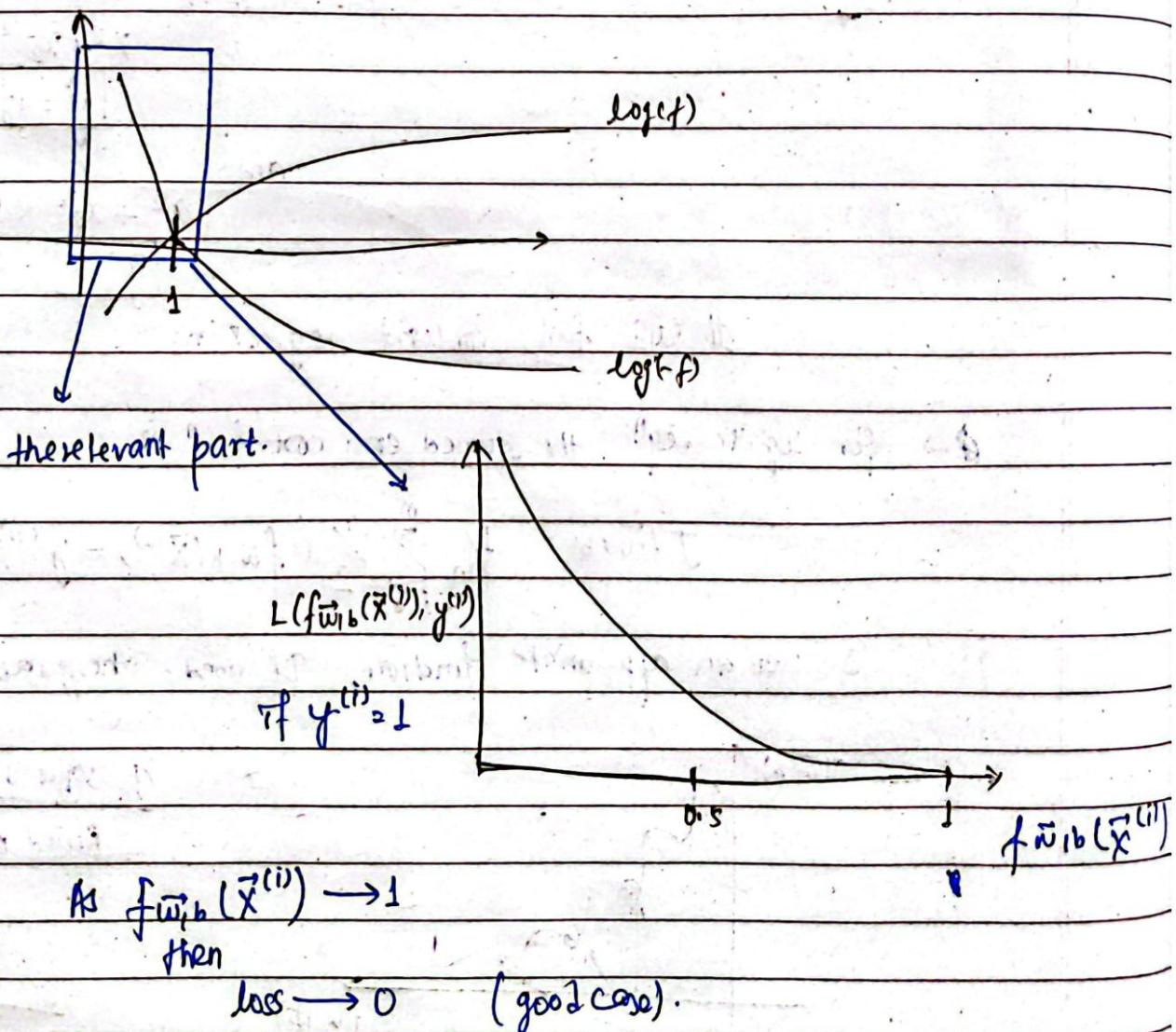
is not an appropriate function. If used the graph becomes like.



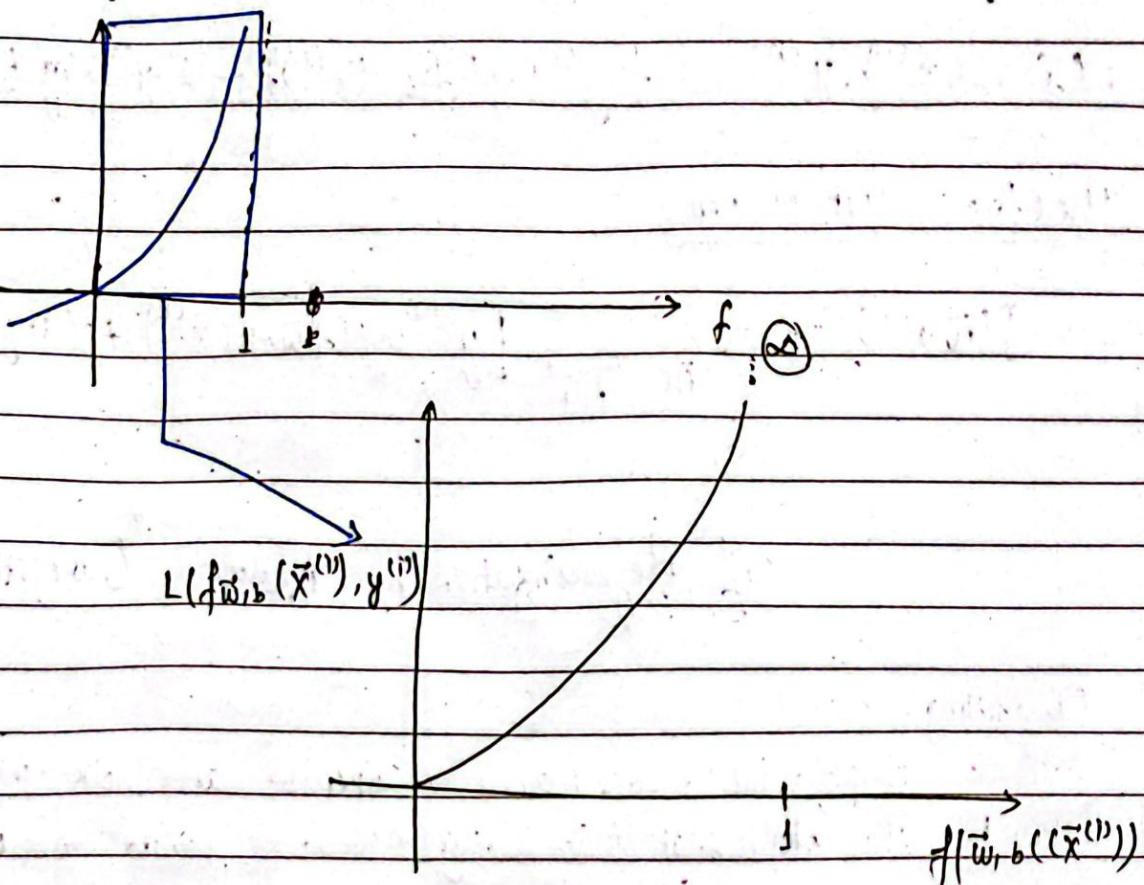
we will use loss fun:  $L(f_{\bar{w}, b}(\vec{x}^{(i)}), y^{(i)})$

→ logistics loss function:

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



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



As  $f(\vec{w}, b)(\vec{x}^{(i)}) \rightarrow \infty$

worst case

$$\text{def } J(\vec{w}, b) = \frac{1}{m} \sum_{i=1}^m L(f(\vec{w}, b)(\vec{x}^{(i)}), y^{(i)})$$

### A Simplified loss function:

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

Therefore: Cost function

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

↓ Regularization to Reduce "Overfitting"

Overfitting:

When ML model learns the details and noise in the training data to the extent that it negatively impacts the performance of the model on new, unseen data.

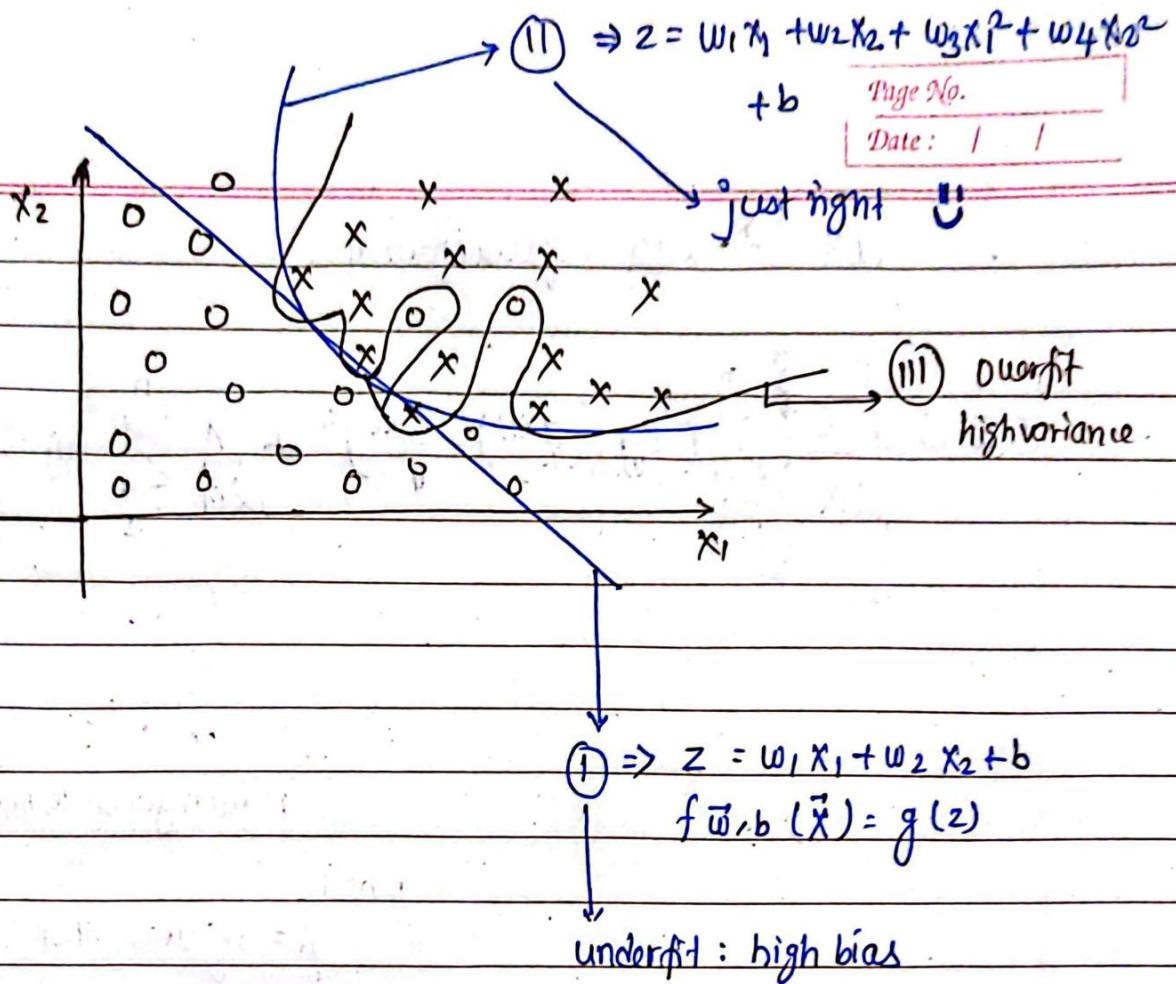
↓  
Can perform well on trained data but fails to "generalize" to new data.

Note:

underfit data  $\Rightarrow$  high bias

overfit data  $\Rightarrow$  high variance





Addressing overfitting:

① Collect more training examples. (when abundant data available)

② Select features to include/exclude.

many features and less training data → leads to overfitting

∴ go for "feature selection"

③ Regularization :

Reduce the size of parameters  $w_j$ .

ie

$$f(x) = 28x - 885x^2 + 89x^3 - 174x^4 + 100$$



$$= 13x - 0.23x^2 + 0.00014x^3 - 0.001x^4 + 10 \quad (\text{Regularized})$$

Example Case

Cost fun' with regularization:

$$J(\vec{w}, b) = \frac{1}{m} \sum_{l=1}^m (f_{\vec{w}, b}(\vec{x}^{(l)}) - y^{(l)})^2 + \frac{\lambda}{m} \sum_{j=1}^n w_j^2 + \frac{\lambda}{2m} b^2$$

regularization term.

where,

$\lambda$  = regularization parameter

If  $\lambda = 0$ : model overfits

Can include or exclude 'b'.

If  $\lambda = \text{very large}$ : model underfit.

if Regularized linear regression:

Gradient descent:

repeat {

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

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

Page No. \_\_\_\_\_  
Date: / /

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

new term.

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



don't have to regularize 'b' !!

$$w_j = w_j - \alpha \left[ \frac{1}{m} \sum_{i=1}^m f(\bar{w}_{i,b}(\bar{x}^{(i)}) - y^{(i)}) x_j^{(i)} + \frac{\lambda}{m} w_j \right] \text{ for } j = 1 \dots n$$

Rearranged to be:

$$w_j = w_j \underbrace{\left( 1 - \alpha \frac{\lambda}{m} \right)}_{\text{New part}} - \alpha \underbrace{\frac{1}{m} \sum_{i=1}^m (f(\bar{w}_{i,b}(\bar{x}^{(i)}) - y^{(i)}) x_j^{(i)})}_{\text{original part.}}$$

This decreases (new part) the value of 'w<sub>j</sub>' each iteration by little bit.

## # Regularized Logistic regression

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



and when applied gradient descent to this cost fun<sup>n</sup>, it looks exactly like that of linear regression.

repeat {

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

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

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

$\Rightarrow$  No change since regularization only applied on  $w$  not  $b$  !!

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