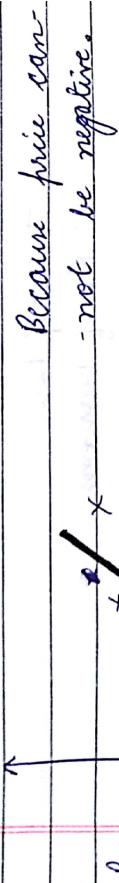


# Deep learning

Date : ...../...../.....

## Housing Price Prediction



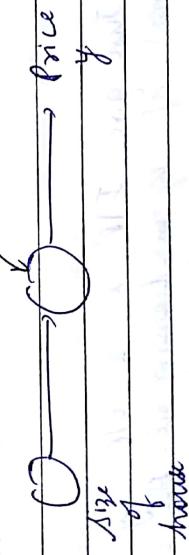
We have to predict price as a function of size of the house.

size of house

4 weeks.

Simplest possible neural net.

Neuron



Neuron  $\rightarrow$  predicts size of house, combines the linear function of the price.

This function is called linear sometimes and sometimes takes off as straight line.

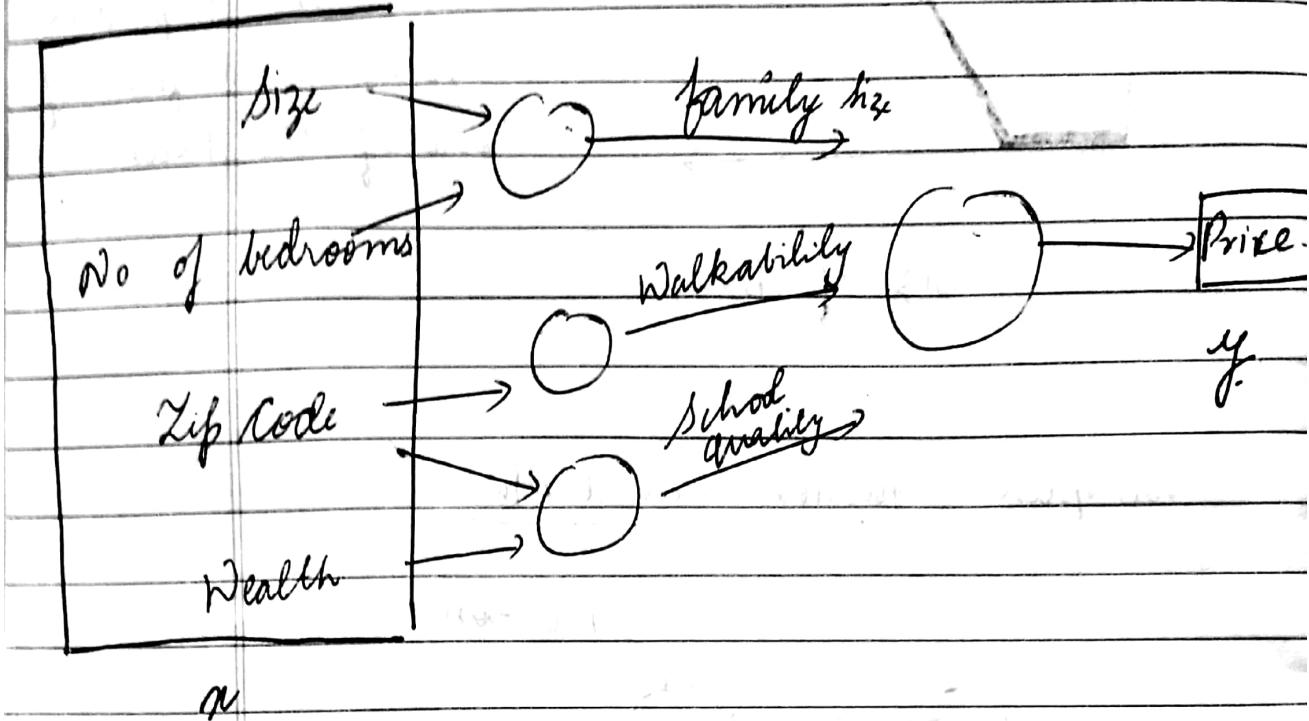
reduces this is called linear unit. Rectify means taking a max of zero.

Predicting price of house based on several features.

→ Size of house

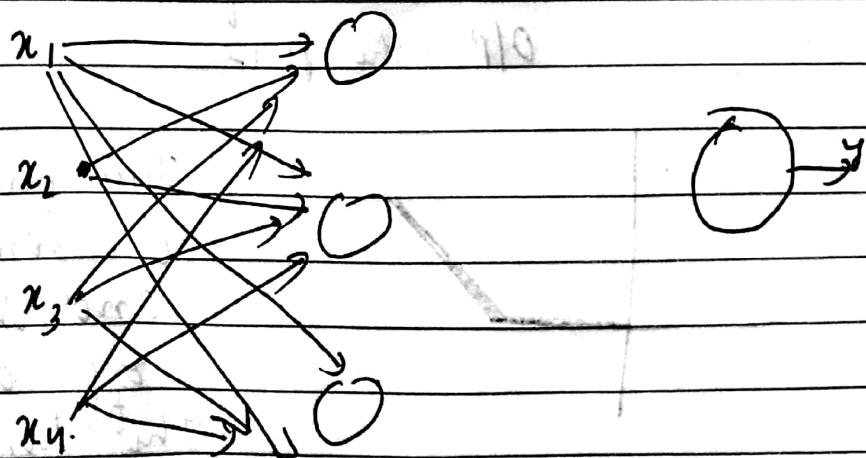
→ No of bedrooms

•



Just give I/P  $x$  & O/P  $y$  for a no of egs in training set. All things in the middle will be figured out on its own

Hidden layer



Each hidden unit of hidden layer takes all 4 I/P features. So eg if first node represent family size & it depends only on I/P  $x_1$  &  $x_2$ . Rather we would say neural node to that you decide whatever you want this node to be & we will give you all 4 I/P features & you can compute whatever you want to.

I/P layer & hidden layer in this case is said to be densely connected bcz every I/P feature is connected to every one of the circles of hidden layer.

If we give large training data of  $x, y$ , the neural net gives good fun<sup>n</sup>s that actually map from  $x$  to  $y$ .

### Supervised learning

In supervised learning, You have an I/P  $x$ . You want a fun<sup>n</sup> that map it to O/P  $y$ .

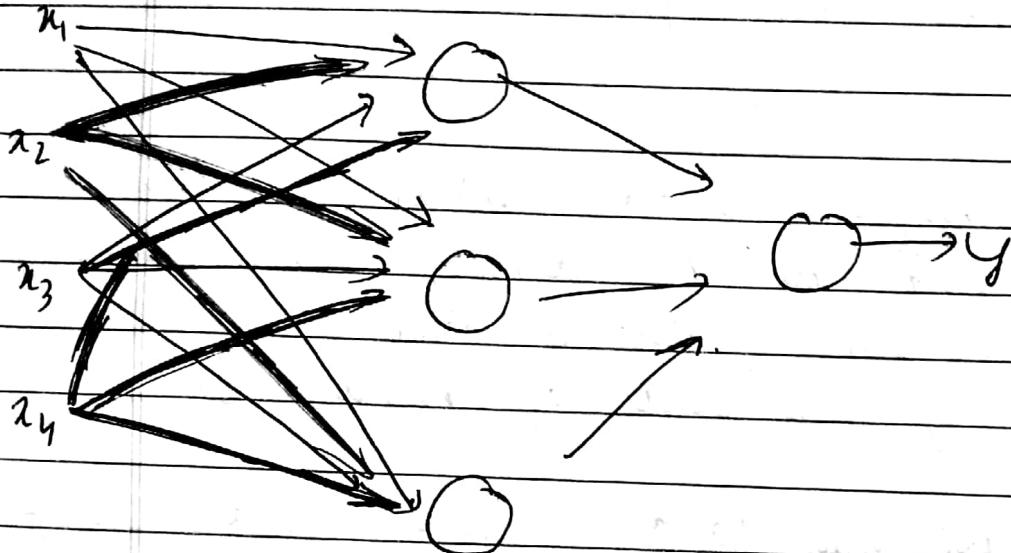
e.g.  $\rightarrow$

I/P ( $x$ )	O/P ( $y$ )	App <sup>n</sup>
→ Home features	Price.	Real Estate
→ Audio	Text transcript	Speech recog
→ English	Chinese	Machine Trans <sup>n</sup>
→ Image, Radar info	Position of other cars.	Autonomous driving

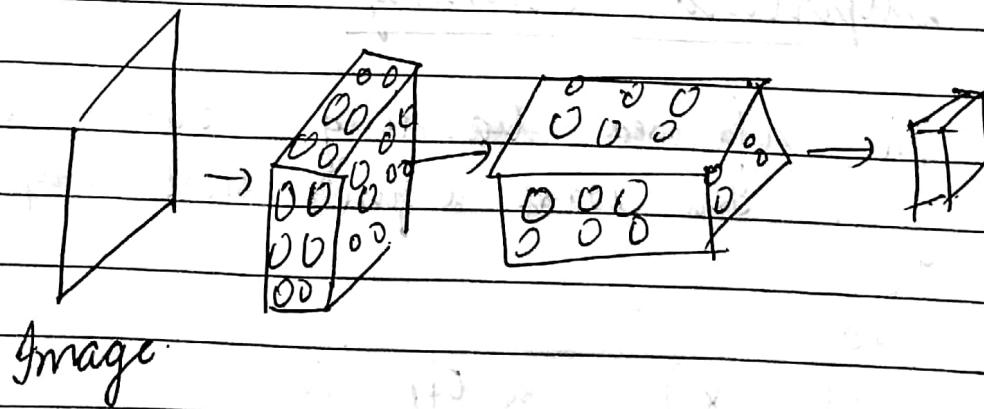
- For image applications we use CNN
- For sequence data like audio, we use RNN
- Languages, complex versions of RNN

## Neural Network examples

### 1. Standard NN



### 2. Convolutional NN (CNN)



It is used mainly for images.

### 3: Recurrent Neural Networks

Structured data

Unstructured data

Means database of data

e.g. → Home price predict<sup>n</sup>

Size	bedrooms	...	Price.
2164	3		400
1600	3		330
:	:		:
3000	4		540

Here Size & no of bedrooms  
are known.

e.g. →

click on ad . There is

info about user such as

age , some info of add

and then labels Y which

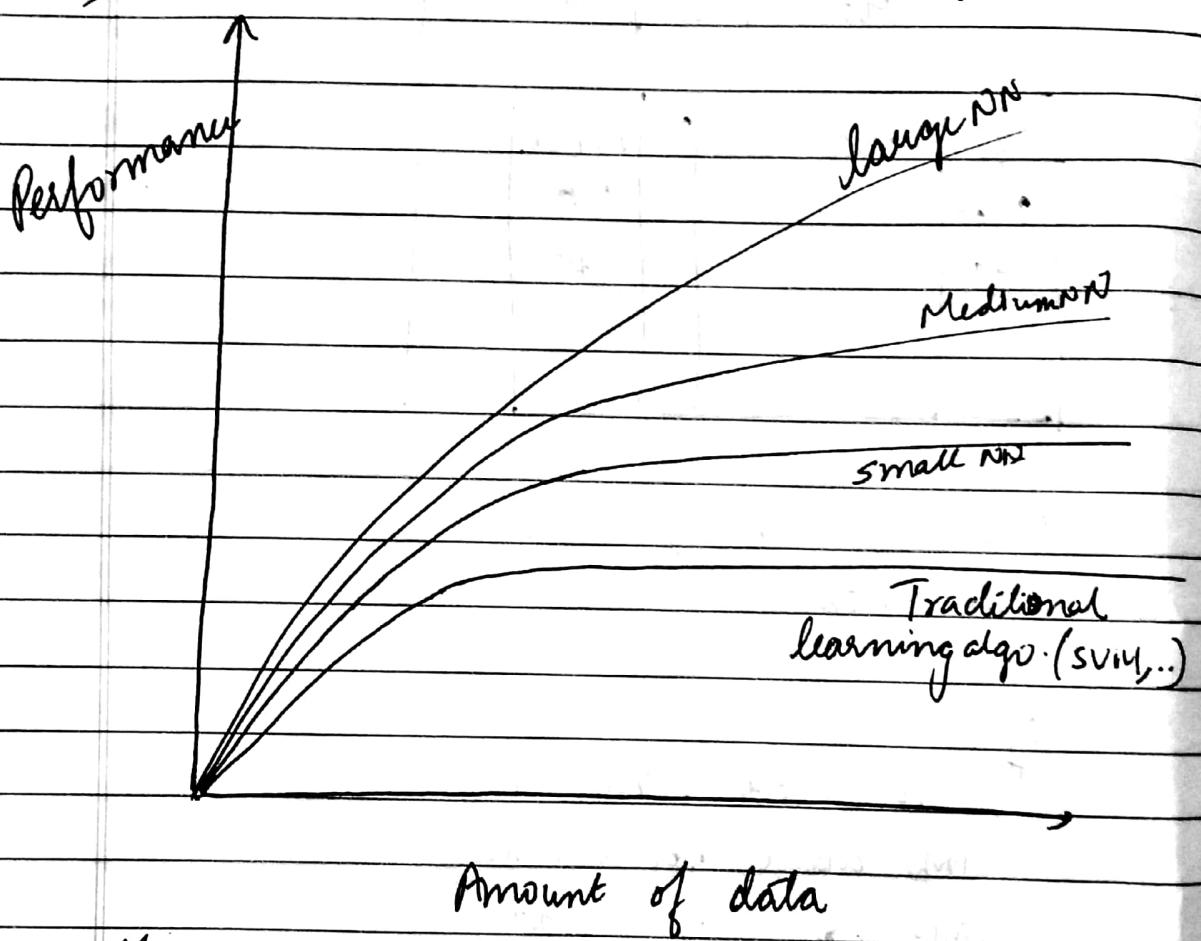
we have to predict

→ In this, each of the  
features like size, no of  
bedrooms, age of user has  
very well defined meaning

→ Refers to things,  
audio, raw  
audio or images  
where you want  
to recognize what's

might be pixel values  
in an image or  
individual words in  
a piece of text.

→ Scale drives deep learning progress



Amount of data

If you want to hit high level of performance then you need 2 things first.

→ Train bigger NN

→ Large amt of data.

# Binary Classification

eg ->

I/P  $\rightarrow$  an image

O/P  $\rightarrow$  Cat (1) or Non-Cat (0)

This is an eg of binary classifi<sup>n</sup> problem.

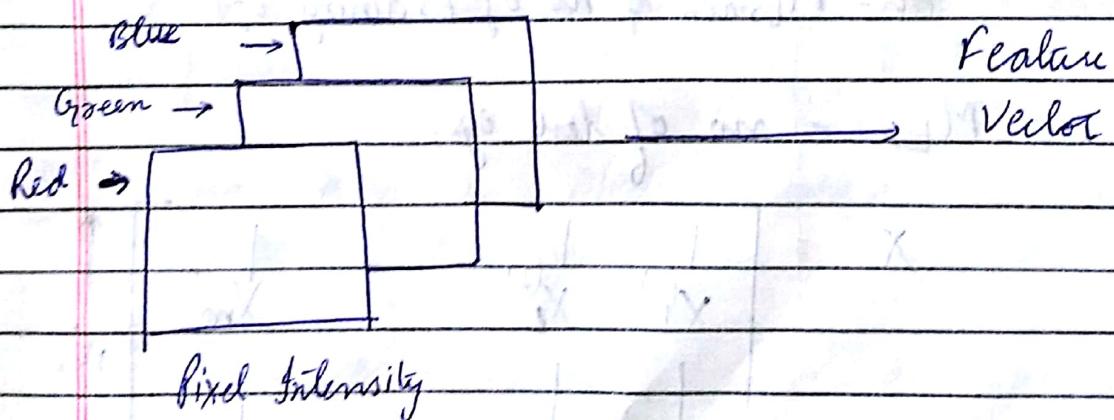
You might have an I/P of an image and want to O/P a label to recognize this image as being either a cat(1) or noncat(0)

Q How an image is stored in comp?

To store an image, comp. separates 3 separate matrices corresponding to red, green and blue color channels of this image.

So if size of image is  $64 \times 64$  pixels.

Then we have 3 matrices of size  $64 \times 64$  corresponding to red, green and blue ~~comp~~ pixel intensities for your image.



Unrolling pixel values to feature vector:

$x = [ \text{red matrix} \quad \text{green matrix} \quad \text{blue matrix} ]^T$  all in column.

Size of feature vector  $x = 64 \times 64 \times 3 = 12288$

$n_x$  represents the dimension of I/P  
feature vector = 12288

⇒ Single training eg is represented by  
 $(x, y)$

$$x \in \mathbb{R}^{N_x}$$

$$y \in \{0, 1\}$$

$m$  training egs

$$\begin{cases} (x_1^{(1)}, y_1^{(1)}) \\ (x_2^{(1)}, y_2^{(1)}) \end{cases} \rightarrow \text{for training eg } 1$$

$$\begin{cases} (x_{m1}^{(m)}, y_{m1}^{(m)}) \\ \vdots \end{cases}$$

$M \approx M_{\text{train}} = \text{no of training egs}$

$M_{\text{test}} = \text{no of test egs.}$

$$X = \left[ \begin{array}{c|c|c|c} & & & \\ \hline x_1^{(1)} & x_2^{(1)} & \cdots & x_m^{(1)} \\ \hline & & & \\ \hline \end{array} \right] \quad \begin{matrix} \uparrow \\ n_x \end{matrix}$$

$\underbrace{\hspace{10em}}$  no of training egs

$n_x \times m$

Date: \_\_\_\_\_  
Sip of feature vector  $x = 14 \times 6 \times 3 = 12288$

mat  $x$  represents the dimension of I/p  
feature vector = 12288

⇒ Single training eg is represented by  
 $(x, y)$

$$x \in \mathbb{R}^{n_x}$$

$$y \in \{0, 1\}$$

m training egs

$$\begin{cases} (x_1^{(1)}, y_1^{(1)}) \\ (x_1^{(2)}, y_1^{(2)}) \end{cases} \rightarrow \text{for training eg } 1$$

$$(x_2^{(1)}, y_2^{(1)}) \rightarrow \text{eg}_2$$

:

$$(x_m^{(1)}, y_m^{(1)})$$

g.

$M \approx M_{train} = \text{no of training egs}$

$M_{test} = \text{no of test egs.}$

$$X = \begin{bmatrix} | & | & | & | \\ x_1^{(1)} & x_1^{(2)} & \cdots & x_1^{(m)} \\ | & | & | & | \\ x_2^{(1)} & x_2^{(2)} & \cdots & x_2^{(m)} \\ | & | & | & | \\ \vdots & \vdots & \vdots & \vdots \\ | & | & | & | \\ x_n^{(1)} & x_n^{(2)} & \cdots & x_n^{(m)} \end{bmatrix} \uparrow n_x$$

↓ no of training egs

UICP  $X \in \mathbb{R}^{N_{train} \times n_x}$

$$y = \begin{bmatrix} y^{(1)} & y^{(2)} & \dots & y^{(m)} \end{bmatrix}$$

$$y \in \mathbb{R}^{1 \times m}$$

$x, y \rightarrow$  All egs (training) are stacked together in columns of a matrix.

## Logistic Regression

It is a learning algorithm used when o/p labels ( $y$ ) is a supervised learning prob. i.e the o/p is 0 or 1 so it is a binary class<sup>n</sup> prob.

Given an I/P feature vector  $x$  corresponding to an image that we want to recognize as cat picture or not a cat picture.

Sol<sup>n</sup> An algo to o/p a pred<sup>n</sup> called as  $\hat{y}$  (y hat)

It is not actual value of o/p.  
It is just the estimate of  $y$ .

$$x \in \mathbb{R}^N$$

Parameters of logistic regression are

$$w \in \mathbb{R}^N, b \in \mathbb{R}$$

So I/P vector  $X$ ,  
parameters  $w, b$  are given.  
& we have to generate O/P  $\hat{y}$ .

One way

$$\hat{y} = w^T x + b$$

Not a good algo for binary class<sup>n</sup>  
bcz you want  $\hat{y}$  to be the chance  
that  $y=1$  ~~or 0~~

So  $\hat{y}$  vary b/w 0 and 1  
but this eq<sup>n</sup>  $\hat{y} = w^T x + b$ . Acc to it  
 $\hat{y} > 1$  or can be -ve which dont  
make sense bcz prob is b/w 0 & 1

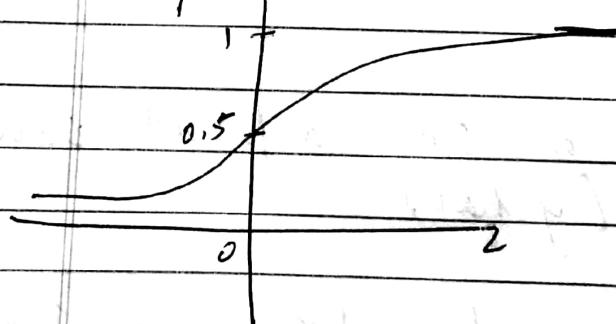
logistic regression

O/P

$$\hat{y} = \sigma(w^T x + b)$$

↓ Sigmoid fun<sup>n</sup>.

$$f = \sigma(z)$$



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

$z$  = real no

If  $z$  = large

$$\text{so } \sigma(z) \approx \frac{1}{1+e^{-\text{large}}} \approx 1$$

ex  
↑

UPCP

If  $z$  = small or large -ve no

$$\sigma(z) = \frac{1}{1+e^{-z}} = \frac{1}{1+e^0} = 0.$$

Logistic regression cost fun<sup>n</sup>:

To train the parameters  $w$  and  $b$  of logistic regression model you need to define cost fun<sup>n</sup>.

$$\hat{y} = \sigma(w^T x + b)$$

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

To learn parameters of model,  
Training set of  $m$  examples is given.

So we want to find parameters  $w$  &  $b$   
so that the predictions are close to actual  
o/p atleast on training data.

So

Given  $\{(x^{(1)}, y^{(1)}) , \dots , (x^{(m)}, y^{(m)})\}$ ,  
 $\hat{y}^{(i)} \approx y^{(i)}$

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b)$$

$$\sigma(z^{(i)}) = \frac{1}{1+e^{-z^{(i)}}} \quad z^{(i)} = w^T x^{(i)} + b$$

Loss / error fun<sup>n</sup> →

Tells how well our algo is.

UPCP	$\text{loss} = L(\hat{y}, y) = \frac{1}{2} (\hat{y} - y)^2$
------	-------------------------------------------------------------

$$L(\hat{y}, y) = - (y \log \hat{y} + (1-y) \log (1-\hat{y}))$$

This loss fun<sup>n</sup> is used in logistic regn.

$$\text{if } y=1 ; L(\hat{y}, y) = -\log \hat{y}$$

so we want  $\log \hat{y}$  to be large.  $\Rightarrow \hat{y}$  to be large.  
 But  $\hat{y}$  is sigmoid fun<sup>n</sup> so  $\hat{y}$  can't be greater than 1  
 so  $\hat{y} = 1$

$$\text{if } y=0; L(\hat{y}, y) = -\log(1-\hat{y})$$

so want  $\log(1-\hat{y})$  to be large.  
 $\Rightarrow 1-\hat{y}$  large  
 $\Rightarrow \hat{y}$  small  
 $\hat{y}$  is b/w 0 &  
 so  $\hat{y} = 0$

### Cost fun<sup>n</sup>

Measures how well you are doing on entire training data.

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

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

\* Cost fun<sup>n</sup> is applied to single training ex.

Cost fun<sup>n</sup> is cost of your parameters

So in training logistic regn model we try to find parameters like weights such that they minimize the overall cost.

- \* Logistic regn can be thought of as small NN.

## Gradient Descent

According to logistic regression algos

$$\hat{y} = \sigma(w^T x + b) \quad \sigma(z) = \frac{1}{1+e^z}$$

Cost fun

$$\begin{aligned} J(w, b) &= \frac{1}{m} \sum_{i=1}^m \ell(\hat{y}^{(i)}, y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m y^{(i)} \log \hat{y}^{(i)} + \frac{(1-y^{(i)}) \log}{(1-\hat{y}^{(i)})} \end{aligned}$$

$\rightarrow$  We want to find  $w, b$  that minimizes  $J(w, b)$

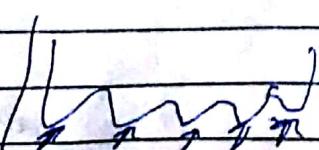
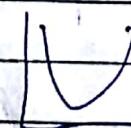
Fun<sup>n</sup> is a convex fun<sup>n</sup> so it

will have minima. This is

the reason we squared the error. Now we will have 1 local minima. This is the reason

we took this fun<sup>n</sup> for

logistic regression.



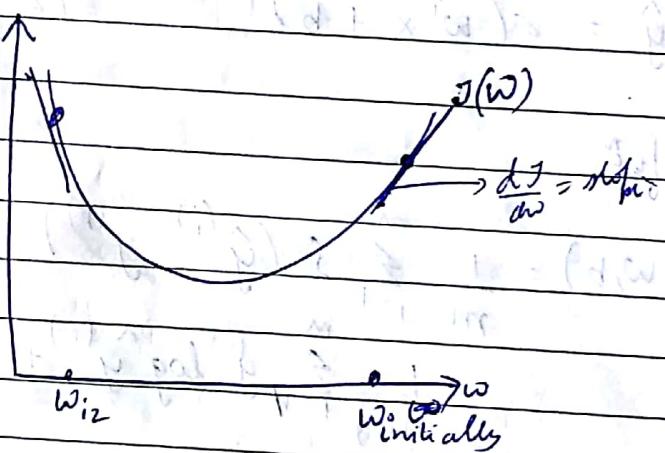
Many local minima in non

→ To find good values for parameters, we first initialize them with some random values.

→ What does Gradient descent do?

It starts at initial point & takes step in direction of deepest downhill.

e.g. Let  $J(w)$  be the fun<sup>n</sup> you want to minimize & you want to



Ignoring  $B$  for now for easiness. Imagine it as a 1-D plot.

Updating  $w$  repeatedly

Repeat {

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

$\alpha \rightarrow$  learning rate & it controls how big each step we take on each iteration.

$\frac{dJ(w)}{dw} \rightarrow$  derivative. This is basically update  
or change you want to do to the  
parameters  $w$ .

In code we will use  $d\omega$  to represent  $\frac{dJ(\omega)}{d\omega}$

$\rightarrow$  If  $w_i$  when we take derivative of  $J$  wrt  $w$   
it is +ve, then  $w := w - \alpha \frac{dJ(w)}{dw}$

So  $w$  is decreased.

& we end up taking a step to the left &  
so this algo slowly decreases the parameter

Case 2: If  $w_i$  is in beginning then slope here  
will be -ve.

$$\text{So } w := w - \alpha \frac{dJ}{dw}$$

So  $w = w + \text{something}$ .

$\rightarrow$  So whether we initialize  $w$  with minor  
max, it will this gradient descent will  
move you towards global minimum.

Since in logistic regression, both  
 $b$  and  $w$  are parameters so we  
do this for both  $w$  &  $b$

$$J(w, b)$$

$$w := w - \alpha \frac{dJ(w, b)}{dw}$$

$$b := b - \alpha \frac{dJ(w, b)}{db}$$

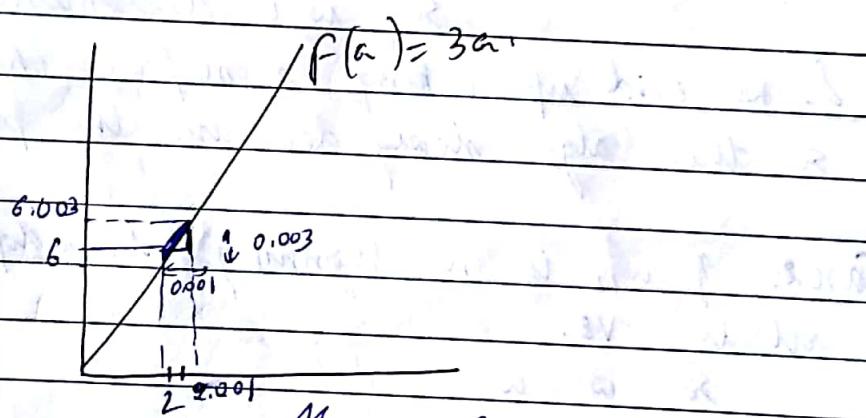
$\frac{dJ(w, b)}{dw}$  is also written as  $\frac{\partial J(w, b)}{\partial w}$

In code we

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

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

## Derivatives



$$\begin{aligned} \text{If } a &= 2 & f(a) &= 6 \\ a &= 2.001 & f(a) &= 6.003 \end{aligned}$$

On moving  $a$  by 0.001 on right,  
 $f(a)$  goes 0.003 on upper side  
 i.e. 3 times  $a$  moves on right

$\Rightarrow$  The slope or derivative of  $f(a)$  at  
 $a = 2$  is 3.

$$\text{slope} = \frac{\text{height}}{\text{width}} = \frac{0.003}{0.001}$$

$$\text{Q} \quad a=5 \quad f(a)=15$$

$$a=5.003 \quad f(a)=15.003$$

So when we bump  $a$  to the right by 0.001,  $f(a)$  goes up by 0.003, i.e.  $a$  moves on right

$\Rightarrow$  slope at  $a=5$  is 3.

$$\frac{d f(a)}{da} = 3$$

Slope of function  $f(a)$  when we change variable by tiny little amt.

$$\Rightarrow \frac{d f(a)}{da} = 3$$

When we nudge change the value of  $a$  by infinitesimally small amount, then  $f(a)$  goes three times the small amt  $a$  has been changed.

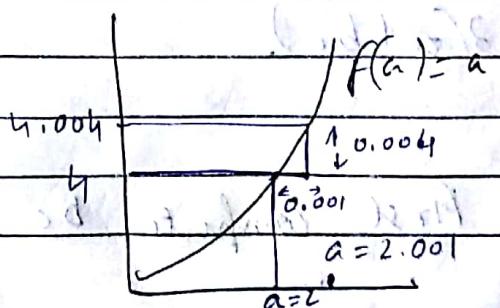
Since for all values of  $a$ , when we change  $a$  by small amt,  $f(a)$  changes by 3 times. Hence slope of fun^n is same everywhere.

Ex

So

$$a=2 \quad f(a)=4$$

$$a=2.001 \quad f(a)=4.004$$



Slope (derivative) of  $f(x)$  at  $a=2$  is 4

$$\frac{d}{da} f(a)=4 \text{ when } a=2$$

In this eg, slope is diff for diff values of  $a$ .

$$a = 5 \quad f(a) = 2.5$$

$$a = 5.001 \quad f(a) = 25.00$$

$$\frac{d}{da} f(a) = 10 \quad \text{when } a = 5.$$

$$\star \frac{d}{da} f(a) = \frac{d}{da} a^2 = 2a.$$

$$\star \frac{d}{da} f(a) \quad f(a) = a^3 \\ \downarrow \\ = 3a^2$$

$$f(a) = \log(a)$$

$$\frac{d}{da} f(a) = \frac{1}{a}$$

## Computation Graph

$$J(a, b, c)$$

$J$  be fun<sup>n</sup> of 3 variables  $a, b, c$

$$J(a, b, c) = 3(a + b + c)$$

3 steps.

$$1.) \quad u = bc \quad \text{First compute } bc$$

$$2.) \quad v = 1/a + u$$

$$3.) \quad J = 3v.$$

a

b

c

$$\text{eg } \rightarrow a = 5$$

$$b = 3 \rightarrow u = 6 \text{ (No )}$$

$$c = 2 \rightarrow v = 11 \text{ (No )}$$

$$v = 11$$

$$J = 33$$

$$V + u \rightarrow J = 3v$$

$$u = bc$$

Derivatives with computation graph.

$$\frac{dJ}{dv} = ? \text{ now } i = 2$$

~~J = 3v~~ with respect of ~~v~~

To calculate the derivative of final variable J wrt v, we have gone backwards 1 step.

$$\frac{dJ}{da}$$

i.e. if we change value of a by small amounts, how does the value of J changes.

$$\text{if } (a = 5) \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

$$\frac{dJ}{da} = 3 \quad \text{since change in J is three times that of a.}$$

So change in a creates change in v which creates change in J

$$\frac{dJ}{da} = \frac{dJ}{dv} \frac{dv}{da}$$

Mostly we calculate derivative of final O/P variable wrt some other variable denoted by dvar.

Final O/P var  $\rightarrow$

other variables  $\rightarrow a, b, c, u, v$  in program

In software what name to give to these variables can be  $\frac{\text{dFinalOutputVar}}{\text{dVar}}$  or  $\frac{\text{dVar}}{\text{dVar}}$   
or dvar

How to compute derivatives to find gradient descent for logistic regression

Using computation graph, we will do it.

logistic regression as follows:

Prediction  $\hat{y}$  is computed as follows.

$$z = w^T x + b$$

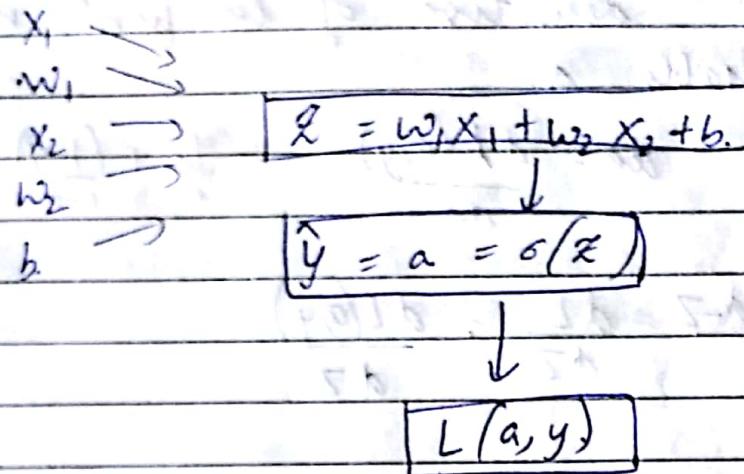
$$\hat{y} = a = \sigma(z)$$

$$l(a, y) = -[y \log(a) + (1-y) \log(1-a)]$$

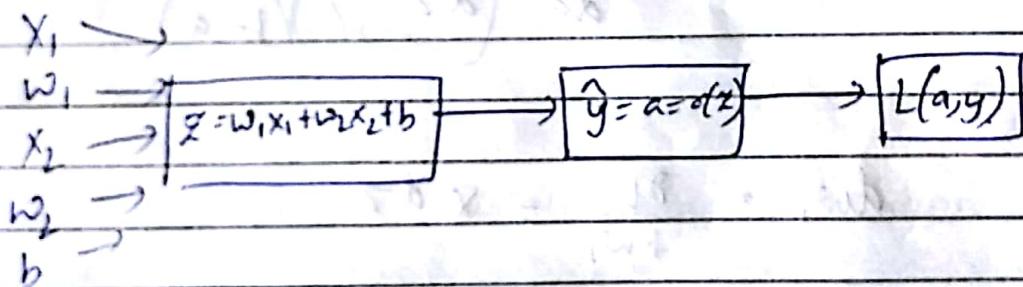
Loss with leg

Consider we have only 2 features  $x_1$  &  $x_2$  for this example.

$$w_1, w_2, b$$



Computing graph



In logistic regression we want to modify the parameters  $w$  and  $b$  in order to reduce/minimize loss.

These are steps of forward prop "on how to compute loss of 1 example training".

Now how to go backwards to compute the derivatives.

What we want to do is to compute derivative wrt this loss.

When going backwards, first thing is to find derivative of loss fun'n wrt variable  $a$ .

$$\frac{da}{d\alpha} = \frac{dL(a, y)}{da} = -\frac{y}{a} + \frac{(1-y)}{1-a}$$

$$\begin{aligned} \frac{dz}{d\alpha} &= \frac{dL}{dz} = \frac{dL(a, y)}{dz} \\ &= a - y \end{aligned}$$

$$= \frac{dL}{da} \times \left( \frac{da}{dz} \right) = (a)(1-a)$$

$$dw_1 = \frac{\partial L}{\partial w_1} = x^T dz$$

$$dw_2 = x_2 dz$$

$$dw_3 = d$$

$$db = dz$$

$$w_1 := w_1 - \alpha dw_1$$

$$w_2 := w_2 - \alpha dw_2$$

$$b := b - \alpha db$$

Gradient descent on m examples.

$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(a^{(i)}, y)$$

$$a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)}) = \sigma(w^T x^{(i)} + b)$$

$$\frac{\partial J(w, b)}{\partial w_i} = \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_i} L(a^{(i)}, y^{(i)})$$

$d w_i^{(i)}$  for each eg.  $\rightarrow (x^{(i)}, y^{(i)})$

Algo

1) Initialize  $J=0$ ,  $dw_1=0$ ;  $dw_2=0$ ;  $db=0$

Use a for loop over training set and compute derivatives wrt each eg & add them up

for  $i=1$  to  $m$

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J + \left[ y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log(1-a^{(i)}) \right]$$

$$d z^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 + = x_1^{(i)} d z^{(i)}$$

$$dw_2 + = x_2^{(i)} d z^{(i)}$$

$$db + = d z^{(i)}$$

] for 2 feature

UPCP

$$\begin{aligned} J &= m \\ dw_1, b &= m \\ db &= m \end{aligned}$$

] for computing avg

So after algo we get

$$\frac{\partial w_1}{\partial w_1} = \frac{\partial J}{\partial w_1}$$

So after this perform 1 step of grad descent.

$$\begin{aligned} w_1 &= w_1 - \alpha \frac{\partial w_1}{\partial w_1}, \\ w_2 & \\ b & \end{aligned}$$

## Vectorization

Date : ..... / ..... / .....

Vectorization is basically for getting rid of explicit for loops in your code.

In Deep learning, large dataset is used for algs to give good results.

But we need to have our codes run faster.

What is vectorization?

$$z = w^T x + b$$

$$w = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad x = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad w \in \mathbb{R}^{N_x} \\ x \in \mathbb{R}^{N_x}$$

If we have non vectorized implement<sup>n</sup>, we do

$$z = 0$$

for i in range( $N_x$ ):

$$z += w[i] * x[i]$$

$$z += b$$

This non vectorized imp<sup>n</sup> will be very slow.

In contrast vectorized implementation just compute  $w^T x$  directly

$$z = np.dot(w, x) + b \quad // \text{compute } w^T x$$

UPCP

Non vectorized implemen takes about 300 times  
longer than vectorized version Date: .....

Eg 1

$$u = A \cdot v$$

vector    matrix    vector

//

// Non-vectorized

$$u_i = \sum_j A_{ij} v_j$$

$$u = np.zeros((n, 1))$$

for i in range(n):

    for j in range(m):

$$u[i] += a[i][j] * v[j]$$

// Vectorized

$$u = np.dot(A, v)$$

Eg 2

You need to apply exponential operation on every element of a matrix/vector.

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

$$u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

This is non vectorized imp.

$$u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = math.exp(v[i])$$

vectorized imp^n

import numpy as np

$$u = \text{np. exp}(v)$$

$\uparrow$   
vector

$\uparrow$   
vector

// Numpy has many vector value func^n

np. log(v)  $\rightarrow$  Computer element wise log

np. abs(v)  $\rightarrow$  absolute

np. maximum(v, 0)  $\rightarrow$  maximum

$v * k^2$   $\rightarrow$  square

$1/v$   $\rightarrow$  invex

Logistic regression derivatives without vectoriz^n

$$J=0, dw_1=0, dw_2=0, db=0$$

for i=1 to m:

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J+ = -[y^{(i)} \log(\hat{y}^{(i)}) + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

$$dz^{(i)} = a^{(i)} / (1 - a^{(i)})$$

$$dw_1+ = x_1^{(i)} dz^{(i)}$$

$$dw_2+ = x_2^{(i)} dz^{(i)}$$

$$db+ = dz^{(i)}$$

$J_{\text{sum}}$

$$J = J/m \quad dw_1 = dw_1 / m \quad dw_2 = dw_2 / m \quad db = db / m$$

But if there were many I/P features then we would have used for loop to find dw<sub>1</sub>, dw<sub>2</sub>, db.

for  $j = 1 \dots n_x$

$$dw_j += \dots$$

So instead of initializing each  $dw_1, dw_2, \dots$  to 0, we will make a vector  $dw$

$$dw = np.zeros((n_x, 1))$$

So now instead of for loop write

$$dw += x^{(1)} dz^{(1)}$$

& later instead of  $dw_j = dw_j/m$  we will write

$$dw = dw/m$$

### Vectorizing logistic regression

We will talk about how to vectorize the implement<sup>n</sup> of logistic regression, so that we can process the entire training set that is implement a single elev<sup>n</sup> of gradient descent wrt an entire training set without even using a single for loop.

Forward prop<sup>n</sup> step of logistic regress<sup>n</sup>. So if there are  $m$  training exs then to make pred<sup>n</sup> on first ex we need

to find:

$$\rightarrow z^{(1)} = w^T x^{(1)} + b.$$

$$\rightarrow a^{(1)} = \sigma(z^{(1)})$$

Now to make pred<sup>n</sup> on 2nd training eg,  
we need to compute:

$$\begin{aligned} z^{(2)} &= w^T x^{(2)} + b \\ a^{(2)} &= \sigma(z^{(2)}) \end{aligned}$$

Similarly for all training egs, predictions  
is done in this manner.

→ So in order to do pred's or carry  
out forward prop<sup>n</sup> steps on our M training  
egs, we can do without explicitly using  
for loop. But how?

Soln

→ Firstly, X is a matrix we defined to be  
our training I/Ps stacked together in  
diff columns

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{bmatrix}$$

X is a  $(n_x, m)$  dimensional matrix.

→ Now computing  $z^{(1)}, z^{(2)}, z^{(3)}$  in 1  
line of code. in 1 step.

→ Construct a  $(1, m)$  matrix which  
is a row vector.

$$\begin{bmatrix} z^{(1)} & z^{(2)} & z^{(3)} & \dots & z^{(m)} \end{bmatrix}$$

It can be expressed as  
 $= w^T X + [b \ b \ b \dots b]$   
 $\downarrow 1 \times M \text{ vector}$   
 $\text{or } m \text{ entries}$

$$= w^T \begin{bmatrix} |^{(1)} & |^{(2)} & \dots & |^{(m)} \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & \dots & | \end{bmatrix}$$

Row vector.

$$= [w^T x^{(1)} + w^T x^{(2)} \quad w^T x^{(3)} \quad w^T x^{(m)}] + [b \quad b \quad \dots]$$

$$= [w^T x^{(1)} + b \quad w^T x^{(2)} + b \quad w^T x^{(3)} + b \quad \dots]$$

$$\begin{matrix} z^{(1)} & z^{(2)} & \dots \end{matrix}$$

So as we obtained  $x$  by taking all training eg. and then stacking them horizontally.

$$z = [z^{(1)} \quad z^{(2)} \quad z^{(3)} \quad \dots \quad z^{(m)}]$$

↓      ↓      ↓      ↓      ↓  
Capital      Small  $z$

Command  $\Rightarrow$

$$z = np.dot(w^T, x) + b$$

$\downarrow$   
 $w$  transpose

In python, it assumes  $x$  to be a real no or  $1 \times 1$  matrix

But when we add the vector  $(w^T x)$  to real no  $b$ , it automatically takes this real no  $b$  and expands it out to this  $(1 \times M)$  row vector.

→ Calculating  $A_1, A_2, \dots, A_M$  all at the same step.

Just as  $X$  i.e. stacking all training egs. in 1 matrix

$$A = [a^{(1)} \ a^{(2)} \ a^{(3)} \ \dots \ a^{(M)}]$$

$$Y = f(Z) = Y - A \cdot S$$

This sigmoid fun<sup>n</sup> inputs  $Z$

### Vectorized logistic Regression's Gradient computation

Here we use vector<sup>n</sup> to perform gradient descent for all  $M$  training egs. all at 1 time.

→ For gradient descent we computed  $d_2$  for first example.

$$d_2^{(1)} = a^{(1)} - y^{(1)} \quad // \text{training eg 1}$$

$$d_2^{(2)} = a^{(2)} - y^{(2)} \quad // \text{training eg 2}$$

and so on

Define new variable  $dZ$

$$dZ = [dZ^{(1)} \ dZ^{(2)} \ dZ^{(3)} \dots \ dZ^{(m)}]$$

$1 \times m$  matrix or  $m^{\text{t}}$  dimension

vector.

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}]$$

$$y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}]$$

$$dZ = A - y = [a^{(1)} - y^{(1)} \ a^{(2)} - y^{(2)} \ \dots \ a^{(m)} - y^{(m)}]$$

Here instead of initializing.

We already removed for loop by  
but we have 2nd for loop  
over training eggs.

we did  $dW = 0$

but we still have to loop over  
the eggs training eggs where

$$dW_{1,t} = x^{(1)} dZ^{(1)}$$

$$dW_{2,t} = x^{(2)} dZ^{(2)}$$

$m$  times

And then  $dW = \frac{dW}{m}$

Similarly for  $db$

$$db = 0$$

$$db + = d z^{(1)}$$

$$db + = d z^{(2)}$$

$$db + = d z^{(m)}$$

So vectorizing it.

$$db = \frac{1}{m} \sum_{i=1}^m d z_i^{(1)} // d z \text{ are in}$$

row vector.

so in python:

$$db = \frac{1}{m} \text{np.sum}(d z)$$

$\sum_{i=1}^m d z_i^{(1)}$  = sum of all elements.

$$dw = \frac{1}{m} X^T d z^T$$

$$= \frac{1}{m} \left[ \begin{array}{c|c} x^{(1)} & | \\ \vdots & | \\ x^{(m)} & | \end{array} \right] \left[ \begin{array}{c} d z^{(1)} \\ \vdots \\ d z^{(m)} \end{array} \right]$$

$$= \frac{1}{m} \left[ \begin{array}{c} x^{(1)} d z^{(1)} \\ \vdots \\ x^{(m)} d z^{(m)} \end{array} \right]$$

In python

$$dw = \frac{1}{m} X^T d z$$

# Implementing logistic regression

## 1. Non vectorized implementation

$$J = 0, \quad dw_1 = 0, \quad dw_2 = 0, \quad db = 0$$

↓  
cost.      ↓ partial derivatives wrt cost.

for  $i = 1$  to  $m$ :

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$j_i = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$\frac{\partial z^{(i)}}{\partial w_1} = a^{(i)} - y^{(i)}$$

$$\frac{\partial J}{\partial w_1} = x_1^{(i)} \frac{\partial z^{(i)}}{\partial w_1}$$

$$\frac{\partial J}{\partial w_2} = x_2^{(i)} \frac{\partial z^{(i)}}{\partial w_2}$$

$$\frac{\partial J}{\partial b} = \frac{\partial z^{(i)}}{\partial b}$$

$$J = \frac{J}{m}, \quad \frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial w_1} / m, \quad \frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial w_2} / m$$

## 2. Vectorized implementation

$$Z = w^T x + b$$

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

$$A = \sigma(Z)$$

$$\frac{\partial Z}{\partial w} = A - y$$

$$w = w - \alpha dw$$

$$b = b - \alpha db$$

Using above code we implemented single iteration of gradient descent for logistic regression.

But for many iterations of implementing gradient descent we need for loop

## Broadcasting in Python

Technique to run codes in python faster.

e.g. In matrix ; no of calories, protein, fats (in 100g) of different food is given.

	Apples	beef	eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
fats	1.8	135.0	99.0	0.9

Goal:

Calculate %age of carb, protein, fats from each of the four foods

In apples,  $(56 + 1.2 + 1.8)$  calories is there.

$$= 59$$

$\frac{56}{59} \times 100\%$  of calories in apples is from carb

for this calculation

We need to add each of columns to get total no of calories in 100g of apple, beg

& then divide throughout the

Can we do this without for loop.

Sol<sup>n</sup>:

→  $A = (3, 4)$  // 3x4 Matrix

→ Using 1 line of python, sum column & we will get 4 elements.

Using it we get total no of cal in 100g of this food.

→ 2nd line of python line, we divide each of 4 columns by these corresponding sum

Code

Import numpy as np

$A = \text{np.array}([[56.0, 0.0, 4.4, 68.0],$

$[1.2, 104.0, 52.0, 8.0]$

$[1.8, 135.0, 99.0, 0.9]])$

np.  
print(A)

$\text{cal} = \text{A.sum (axis=0)}$

$\bar{\text{J}}$  means sum vertically.

`print (cal)`

$[59.0, 239. 155.4 76.9]$

$\text{percentage} = \frac{100}{\text{cal}} / \text{cal.reshape}(1, 4)$

`print (percentage).`

python broadcasting

here A is  $3 \times 4$  matrix

Cal is  $1 \times 4$  matrix

Here cal is already cal is  $1 \times 4$  matrix.

So need for reshaping.

so in python if you are not sure  
of dimensions of matrix, use  
`reshape` command.

`reshape` is  $O(1)$  operation.

$A \rightarrow (3 \times 4)$  Matrix

$\text{cal} \rightarrow (1 \times 4) -$

How can we divide  $(3 \times 4)$  Matrix by  
 $(1 \times 4)$  matrix

## egs of broadcasting

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 \cdot (100) \text{ und } \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}$$

When we add  $1 \times 4$  Matrix to a number, python will take this no & auto expand to  $4 \times 1$  matrix.

So it becomes

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

Here broadcasting occurs

$$2 \cdot \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \end{bmatrix}$$

$2 \times 3$  Matrix       $\downarrow$   
 $(m \times n$  Matrix)       $1 \times n$  matrix  
                                (general)

Python copies the matrix on right  $m$  times to make it  $(m \times n)$  matrix. This is also broadcasting -

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}$$

UPCP

$$= \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 308 \end{bmatrix}$$

$$3. \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}$$

$m \times n$  matrix       $m \times 1$  matrix

Here copy  $m \times 1$  matrix  $n$  times horizontally.

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 \\ 200 & 200 \end{bmatrix}$$

$$= \begin{bmatrix} 110 & 110 \\ 250 & 250 \end{bmatrix}$$

### General principle of broadcasting in python

If you have  $(m, n)$  matrix and you add or subtract or multiply or divide it with a  $(1, n)$  matrix then it will copy  $m$  times into  $(m, n)$  matrix & then apply operation element wise.

If conversely you took  $(m, n)$

— with  $(m, 1)$  —  
 — with  $(1, n)$  —  
 —  $(m, n)$  —

$(m, 1) + R$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 2 & 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 & 102 & 103 \end{bmatrix}$$

III Creating 5 random variables in python

```
import numpy as np
a = np.random.randn(5)
```

print(a)

[ 0.502951 ]

print(a.shape)

Rank 1 array &

neither row nor column vector

print(a[1]) // here ~~trans~~ no change occurs on doing its transpose.

Print(np.dot(a, a[1]))

// gives a no.

Now transpose gives correct result  
In doing this, there will be 2 brackets

$$\begin{bmatrix} [ & [ & ( & ] \\ [ & ] & ] & ] \end{bmatrix}$$

17

However in previous case 1 bracket exists.

### Python / numpy vectors

$\rightarrow a = np.random.randn(5)$   
 $a = a.reshape(5, 1)$

//rank 1 array. It does not give consistent result. Neither row/ column vector.

Rather use this

$\rightarrow a = np.random.randn(5, 1) \rightarrow$  column vector  
 $a = np.random.randn(1, 5) \rightarrow$  row vector

$\rightarrow$  assert ( $a.shape \equiv (5, 1)$ )

$\rightarrow a.reshape((5, 1)) \rightarrow$  if it is ~~not~~ a rank 1 matrix array.

This can eliminate it.

There is a cost fun<sup>n</sup> in logistic regression. Justification of using the cost function

$$\hat{y} = \sigma(w^T x + b)$$

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

We want to interpret  $\hat{y}$  as the prob that  $y=1$  given  $x$ .

So we want our algo to output  $\hat{y}$  as the chance that  $y=1$  given  $x$ .

$$\Rightarrow \text{if } y=1, \text{ then } p(y|x) = \hat{y}$$

$$y=0 \quad \underline{\quad \quad \quad} = 1 - \hat{y}$$

$y=0$  or  $1$  is only possible case as it is binary classification.

On Summarizing.

$$p(y|x) = \hat{y}^y (1-\hat{y})^{1-y}$$

In case 1:

$$y=1$$

$$p(y|x) = \hat{y}^1 (1-\hat{y})^0$$

$$= \hat{y}$$

In case 2:

$$y=0$$

$$(1-\hat{y})^0 (1-\hat{y})^1 = 1 - \hat{y}$$

$$\log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$= y \log \hat{y} + (1-y) \log (1-\hat{y})$$

$$= -L(\hat{y}|y)$$

Cost on m eqs:

$$p(\text{labels in training set}) = \prod_{i=1}^m p(y^{(i)}|x^{(i)})$$

$$\log - \quad \text{log} =$$

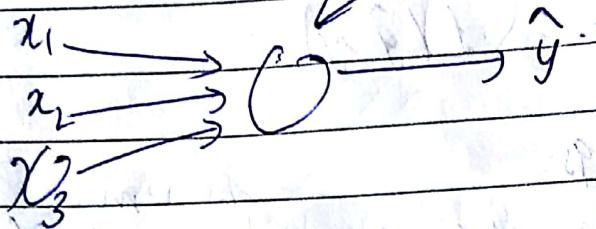
$$(a * b) + (a * c) - (b * c)$$

$$(a * b) + (a * c) - (b * c)$$

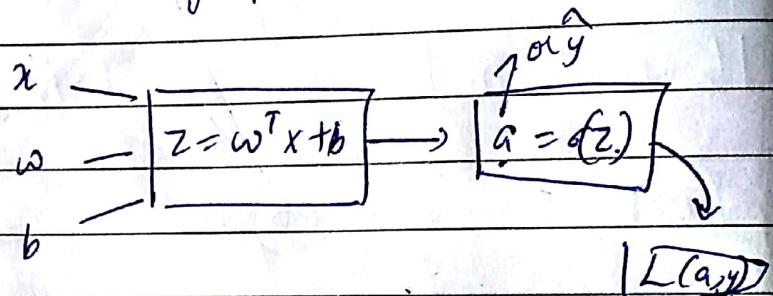
- 1.) np.ufunc(x) → works for any np.array & applies exponential "fun" to every coordinate

# What is Neural Network?

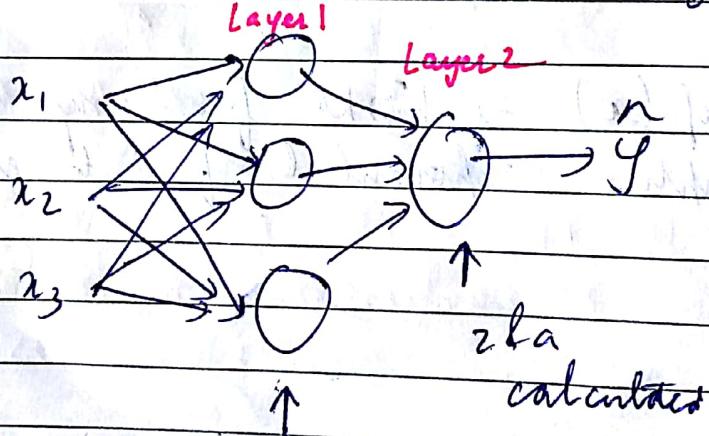
logistic regression



This model corresponds to full computational graph.



Neural Network looks like this



$z \& a$

both

calculated

$$z^{(1)} = w^{(1)}x + b^{(1)}$$

First we I/P  $x$  i.e. features and parameters  $w$  and  $b$ . It helps to

$x^{(1)}$  → round brackets for referring to training eq no  
 $x^{[i]}$  → refers to layer number i.e. referring to Date quantities

calculate  $z^{[1]} = w^{[1]} x + b^{[1]}$  associated with layer [1]

Subscript and square brackets are used to associate with quantities referring to layer 1

[2] → to refer to quantities associated with layer 2

After computing  $z^{[1]}$ , similar to logistic regression, we compute  $a^{[1]}$ .

Then we compute  $z^{[2]}$  using another linear eqn & then compute  $a^{[2]}$ ,  $a^{[2]}$  is the final O/P of NN.

Green - Backward prop

$$\begin{array}{ccccccc} x & \xrightarrow{\text{IP}} & z^{[1]} & = & w^{[1]} x + b^{[1]} & \xrightarrow{a^{[1]} = g(z^{[1]})} & z^{[2]} \\ w^{[1]} & \xleftarrow{\text{IP}} & & & & \xleftarrow{z^{[2]} = w^{[2]} a^{[1]} + b^{[2]}} & a^{[2]} = a(z^{[2]}) \\ b^{[1]} & \xleftarrow{\text{IP}} & & & & & \\ & & & & w^{[2]} & \xleftarrow{\text{IP}} & \\ & & & & b^{[2]} & \xleftarrow{\text{IP}} & \\ & & & & & & L(a^{[2]}, y) \end{array}$$

→ As in logistic regression, we did backward propagation & computed derivatives, gradients. Similarly neural network will end up doing backward calculation.

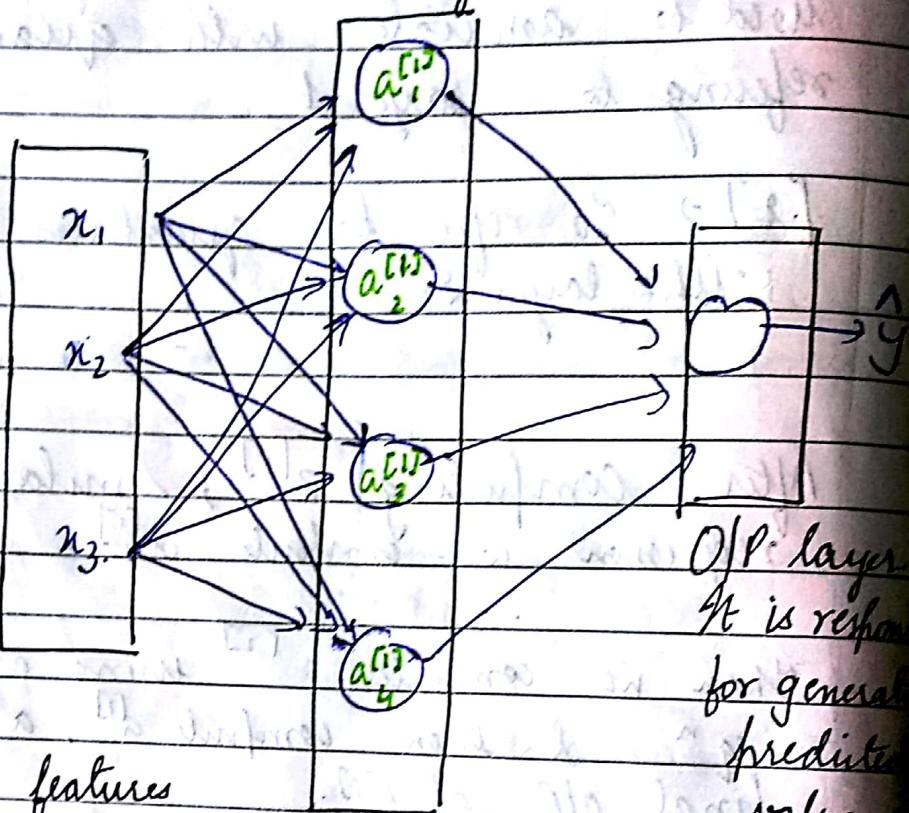
Backward prop occurs from right to left.  $da^{[2]} / dz^{[2]}$  and  $da^{[1]} / dz^{[1]}$

# One layer hidden

Date ..... / ..../ .....

## Neural Network Representation

One hidden layer neural network



O/P layer  
It is responsible for general prediction value

I/P features stacked vertically.  
This is called I/P layer of neural network. i.e it contains I/Ps to NN

Hidden layer of neural network.

In supervised learning ; the training contains values of  $x$  and  $y$  i.e I/P and the target o/p  $y$ .

The things in the hidden layer are not seen in the training set.  
Hidden bcz we don't see in training

$X \rightarrow$  Input feature vector  $\rightarrow a^{[0]}$

↑

alternative representation.

$a$  also stands for activations and it refers to the values that different layers of Neural N/w pass on to the subsequent layers.

So I/P layer passes on values  $x$  to the hidden layer. Hence we call it activations of I/P layer.

The hidden layer will generate activations & we represent that using  ~~$a^{[1]}$~~ .  $a^{[1]}$

Hence the nodes of hidden layer will generate  $a_1^{[1]}, a_2^{[1]}, a_3^{[1]}, a_4^{[1]}$  respectively

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} \rightarrow 4 \times 1 \text{ Matrix or 4dime vector}$$

O/P layer generates value  $a^{[2]}$  which is just a real no.  
so  $y = a^{[2]}$

This N/w is called 2 layer Neural N/w.

bcz when we count layers in N/w, we don't count I/P layer. So, hidden layer is layer 1 & O/P layer is layer 2.

Hidden layer & O/P layer have parameters associated with it parameters  $w$  and  $b$ .

$$\begin{matrix} w^{[1]}, b^{[1]} \\ \downarrow \\ 4 \times 3 \end{matrix}, \quad \begin{matrix} u_{xi} \\ \downarrow \\ 4 \times 1 \end{matrix}$$

parameters of  
Hidden layer

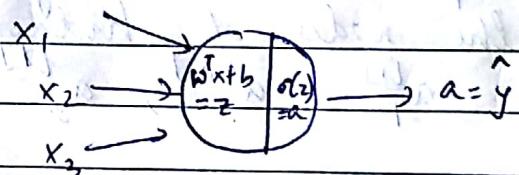
$$\begin{matrix} w^{[2]}, b^{[2]} \\ \downarrow \\ (1 \times 4) \end{matrix}, \quad \begin{matrix} u_{xi} \\ \downarrow \\ (1 \times 1) \end{matrix}$$

parameters of

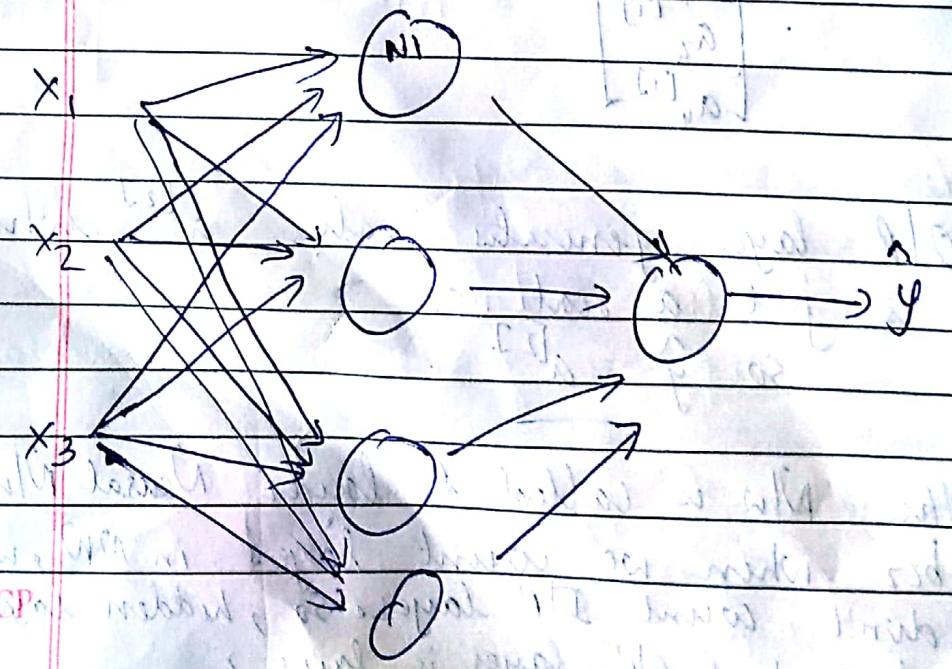
O/P layer

- \* In logistic regression, circle represents 2 steps of computation.  
→ Computing  $Z$

→ Computing "final activation"

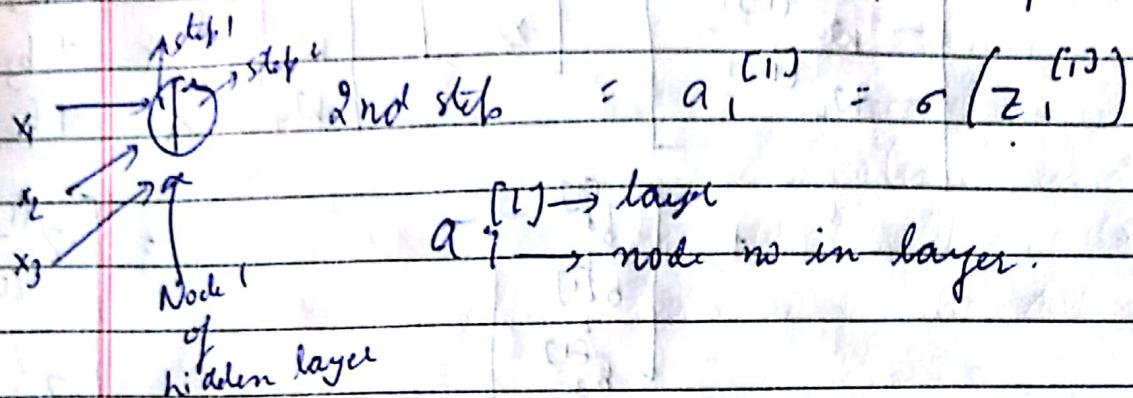


- \* A Neural Net does this many times

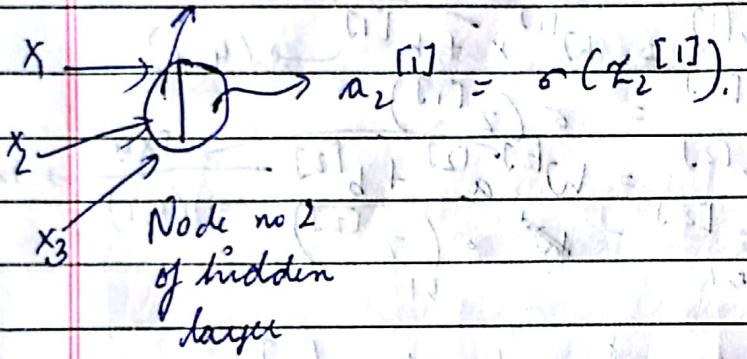


→ Each node of hidden layer does two steps of computation.

$$\text{First step} = w_i^{[1]T} x + b_i^{[1]} = z_i^{[1]}$$



$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}$$



So

$$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}, \quad a_1^{[1]} = \sigma(z_1^{[1]})$$

$$z_2^{[1]} = w_2^{[1]T} x + b_2^{[1]}, \quad a_2^{[1]} = \sigma(z_2^{[1]})$$

$$z_3^{[1]} = w_3^{[1]T} x + b_3^{[1]}, \quad a_3^{[1]} = \sigma(z_3^{[1]})$$

$$z_4^{[1]} = w_4^{[1]T} x + b_4^{[1]}, \quad a_4^{[1]} = \sigma(z_4^{[1]})$$

If we do this wing for loops, it becomes inefficient.

So we vectorize it.

$$\begin{bmatrix}
 -w_1^{[1]T} \\
 -w_2^{[1]T} \\
 -w_3^{[1]T} \\
 -w_4^{[1]T}
 \end{bmatrix}
 \begin{bmatrix}
 x_1 \\
 x_2 \\
 x_3
 \end{bmatrix}
 +
 \begin{bmatrix}
 b_1^{[1]} \\
 b_2^{[1]} \\
 b_3^{[1]} \\
 b_4^{[1]}
 \end{bmatrix}
 = \begin{bmatrix}
 w_1^{[1]T} x + b_1^{[1]} \\
 w_2^{[1]T} x + b_2^{[1]} \\
 w_3^{[1]T} x + b_3^{[1]} \\
 w_4^{[1]T} x + b_4^{[1]}
 \end{bmatrix}
 = \begin{bmatrix}
 z_1^{[1]} \\
 z_2^{[1]} \\
 z_3^{[1]} \\
 z_4^{[1]}
 \end{bmatrix}$$

Given  $\mathbb{R}^4 / P \rightarrow \mathbb{R}^4$

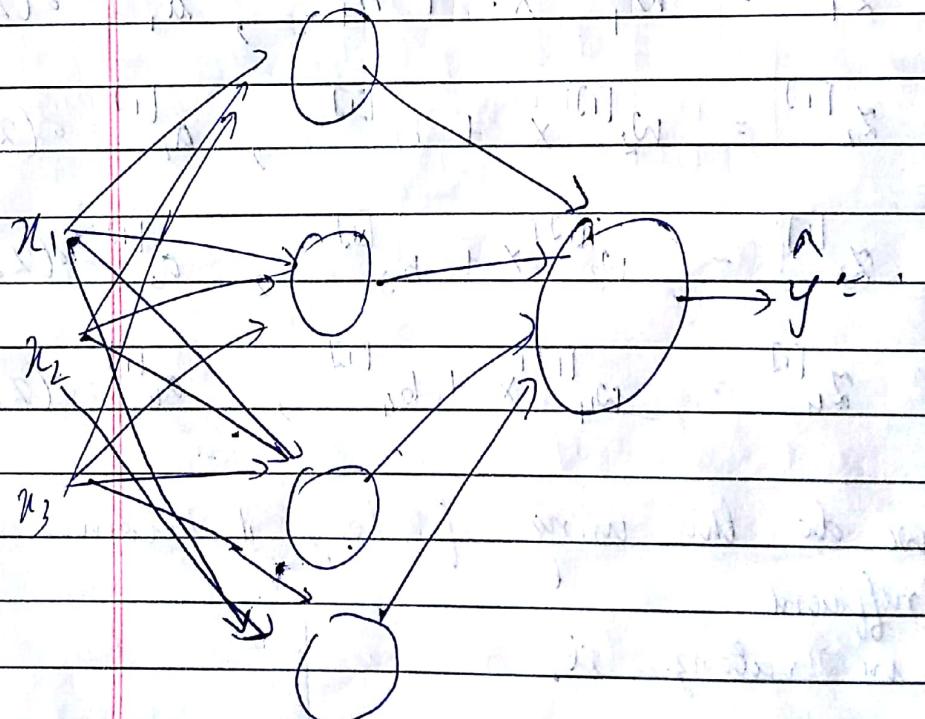
$$\Rightarrow \mathbb{R}^4 \leftarrow z^{[1]} = w^{[1]T} x + b^{[1]} \rightarrow \mathbb{R}^4$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$(1,1) \leftarrow z^{[2]} = w^{[2]T} a^{[1]} + b^{[2]} \rightarrow \mathbb{R}^4$$

$$a^{[2]} = \sigma(z^{[2]})$$

$x$  can be represented by  $a[0]$



$$x \rightarrow a^{[2]} = \hat{y}$$

for training eg 1:

$$x^{[1]} \longrightarrow a^{[2](1)} = y$$

$$\underbrace{x^{(2)}}_{\text{layer 2}} \longrightarrow a^{[2](2)} = y$$

for training eg m

$$x^{[2](m)}$$

layer 2  
training eg i

→ In case of non vectorized implementation,  
if we have m training egs then  
we need to do :

for  $i = 1$  to  $m$ :

$$z^{[1](i)} = w^{[1]} x^{(i)} + b^{[1]}$$

$$a^{[1](i)} = \sigma(z^{[1](i)})$$

$$z^{[2](i)} = w^{[2]} a^{[1](i)} + b^{[2]}$$

$$a^{[2](i)} = \sigma(z^{[2](i)})$$

→ Vectorizing to get rid of for loop

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(m)} \end{bmatrix}$$

↑      ↑      . . .      ↑      ↓

Training eg  
stacked in  
column

Size of matrix X is  $n \times M$

Compute  $z^{[1]}$

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = \sigma(z^{[2]})$$

So we changed from lower case vector  $x$  to capital  $X$  matrix by stacking lower case  $x$  in diff columns.

~~$$z^{[1]} = \begin{bmatrix} | & | & | & | \\ z^{[1](1)} & z^{[1](2)} & \dots & z^{[1](m)} \\ | & | & | & | \end{bmatrix}$$~~

corresponds to 1st hidden unit of this training

~~$$A^{[1]} = \begin{bmatrix} | & | & | & | \\ a^{[1](1)} & a^{[1](2)} & \dots & a^{[1](m)} \\ | & | & | & | \end{bmatrix}$$~~

the horizontal axis, or we index for diff training eqs.

This vertical index corresponds to different nodes in the neural nw.

$x$

# Justific'n for vectorized implement'n

For training eq 1

2

3

$$z^{(1)(1)} = w^{(1)} x^{(1)} + b^{(1)}, z^{(1)(2)} = w^{(1)} x^{(2)} + b^{(1)}$$

$$\text{for eq } z^{(1)(3)} = w^{(1)} x^{(3)} + b^{(1)}$$

Assume  $b = 0$ .

$$w^{(1)} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

$$w^{(1)} x^{(1)} = \begin{bmatrix} ? \\ ? \\ \vdots \\ ? \end{bmatrix}$$

Column  
vector.

$$w^{(1)} x^{(2)} = \begin{bmatrix} ? \\ ? \\ \vdots \\ ? \end{bmatrix}$$

$$w^{(1)} x^{(3)} = \begin{bmatrix} ? \\ ? \\ \vdots \\ ? \end{bmatrix}$$

Some other  
column  
vector.

Other  
column vector.

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & x^{(3)} \\ | & | & | \end{bmatrix}$$

All training egs are stacked  
vertically

$$w^{(1)} x^{(1)} \quad w^{(1)} x^{(2)} \quad w^{(1)} x^{(3)}$$

$$w^{(1)} X = \begin{bmatrix} ? & ? & ? \\ ? & ? & ? \\ ? & ? & ? \end{bmatrix}$$

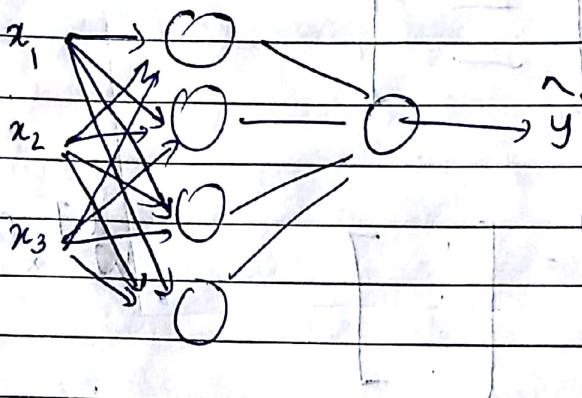
CP

$$= \begin{bmatrix} z^{[1]f(1)} & z^{[2]f(2)} & z^{[3]f(3)} \end{bmatrix} = z^{[0]}$$

## Activation Functions

What activation fun<sup>n</sup> to use in the neural network.

In forward prop<sup>n</sup> step of NN, the steps are:



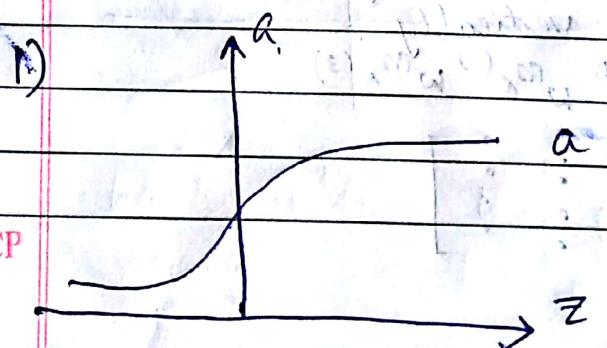
$$z^{[1]} = w^{[1]}x + b^{[1]}$$

$$a^{[1]} = \sigma(z^{[1]})$$

$$z^{[2]} = w^{[2]}a^{[1]} + b^{[2]}$$

$$\rightarrow a^{[2]} = \sigma(z^{[2]}) = x$$

In these two steps, we use the sigmoid functions.



UPCP

In generic can we can use diff activation fun<sup>n</sup>.

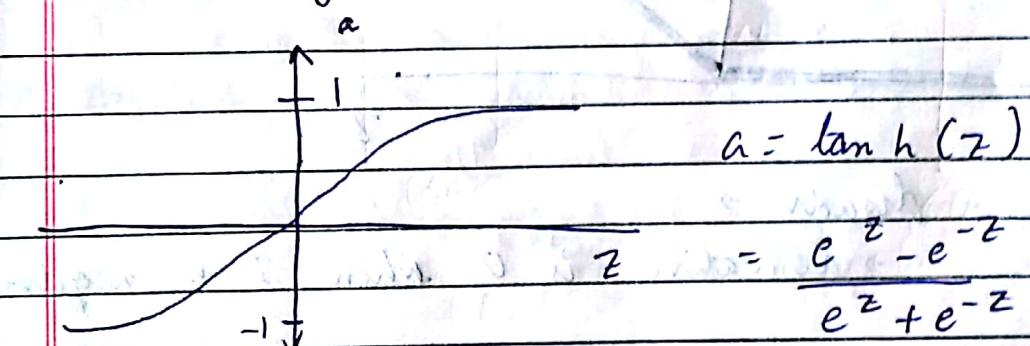
$$a^{[1]} = g(z^{[1]})$$

↓  
generic.

where  $g$  is any non-linear fun<sup>n</sup>  
but not sigmoid function.

Sigmoid ranges b/w 0 & 1.

2. Tanh i.e hyperbolic tangent  
 → is an activa<sup>n</sup> fun<sup>n</sup>. It is works better than sigmoid.  
 → It goes b/w +1 & -1.



If we use tanh ~~as~~ as activation fun<sup>n</sup> in hidden layer than it works better than sigmoid fun<sup>n</sup>.

- In O/P, we prefer using sigmoid fun<sup>n</sup> because we want O/P b/w 0 & 1.

If O/P is b/w -1 and 1, it doesn't make sense. So in case of binary classifi<sup>n</sup> prob we use activa<sup>n</sup> fun<sup>n</sup> of sigmoid.

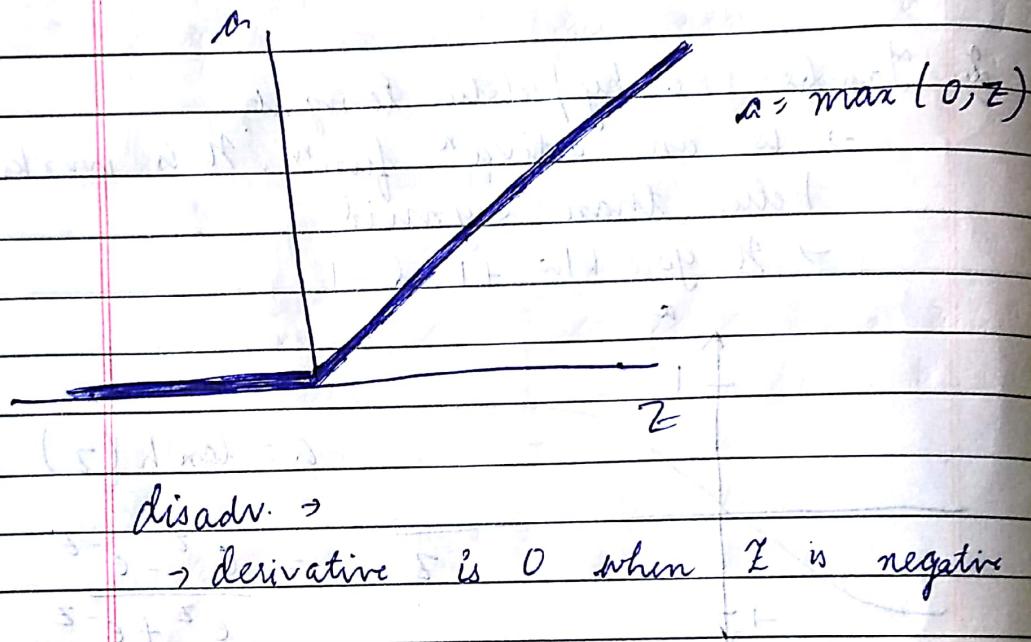
→ We can have diff activ<sup>n</sup> fun<sup>n</sup> for diff layers.

$g^{[1]}(z^{[1]})$  → refers to activation fun<sup>n</sup> for layer 1.

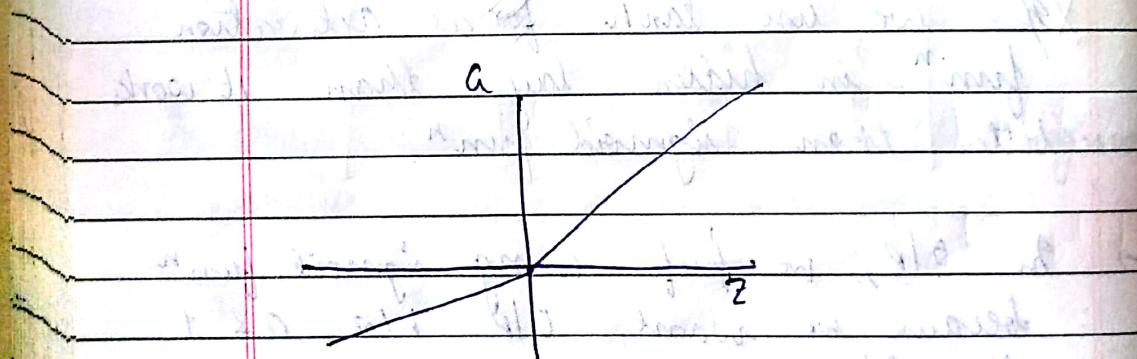
$g^{[2]}(z^{[2]})$

Refers to activ<sup>n</sup> fun<sup>n</sup> for layer 2.

3.) RELU → Rectified linear unit.



4.) Leaky RELU → Another version of REV.



Works better than REV but not used  
mostly.

RELU and leaky RELU has advantage.  
 For most values of  $z$ , the derivative or slope  
 of activation fun<sup>n</sup> is very diff from 0.  
 and so our Neural network will  
 learn very faster than when we use  
 tanh or sigmoid activa<sup>n</sup> fun<sup>n</sup>.

### Pros and cons of different activation fun<sup>n</sup>s

1. Sigmoid:

$$a = \frac{1}{1+e^{-z}}$$

→ Never used.

→ Use only in O/P layer  
 when doing binary  
 classification.

$$\rightarrow \tanh: a = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

→ Better than sigmoid fun<sup>n</sup>

→ RELU

$$a = \max(0, x)$$

→ Leaky RELU.

$$a = \max(0.01z, z)$$

UPCP

Why do we need non-linear activ<sup>n</sup> fun<sup>n</sup>?

⇒ For a Neural Net, the forward prop eqns are:

$$\begin{aligned}z^{[1]} &= w^{[1]}x + b^{[1]} \\a^{[1]} &= g^{[1]}(z^{[1]}) \\z^{[2]} &= w^{[2]}a^{[1]} + b^{[2]} \\a^{[2]} &= g^{[2]}(z^{[2]})\end{aligned}$$

Why can we not get rid of  $a^{[1]} = g^{[1]}(z^{[1]})$   
i.e. if we can remove activ<sup>n</sup> fun<sup>n</sup> &

thus

$$a^{[1]} = z^{[1]}$$

Sometimes also called as linear activ<sup>n</sup> fun<sup>n</sup>. or identity activ<sup>n</sup> fun<sup>n</sup>. bcz it outputs whatever is the IP

$$\text{and } a^{[2]} = z^{[2]}$$

Here  $\hat{y}$  or o/p is just a linear fun<sup>n</sup> of  $x$

$$\Rightarrow a^{[1]} = z^{[1]} = w^{[1]}x + b^{[1]} \quad (1)$$

$$a^{[2]} = z^{[2]} = w^{[2]}a^{[1]} + b^{[2]} \quad (2)$$

from (1) & (2)

$$\begin{aligned}a^{[2]} &= z^{[2]} = w^{[2]}(w^{[1]}x + b^{[1]}) + b^{[2]} \\&= (w^{[2]}w^{[1]})x + w^{[2]}b^{[1]} + b^{[2]} \\&= w^T x + b\end{aligned}$$

If we use linear activ<sup>n</sup> fun<sup>n</sup> then NN  
O/P is a linear fun<sup>n</sup> of I/P.

- Linear activ<sup>n</sup> fun<sup>n</sup> is useless.
- It is used only if we are doing M/c learning on a ~~linear~~ regression prob.

## Derivatives of Activation fun<sup>n</sup>s.

For backpropagation steps, we need to compute derivative / slope of activ<sup>n</sup> fun<sup>n</sup>

Sigmoid activ<sup>n</sup> fun<sup>n</sup>

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

slope of fun<sup>n</sup>

$\frac{d}{dz} g(z) = \text{slope of } g(z) \text{ at } z.$

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

$$= g(z) (1 - g(z))$$

$$z = 10.$$

$$z = -10$$

$$g(z) = 1$$

$$g(z) = 0.$$

$$\frac{d}{dz} g(z) \approx 1/(1-1) \approx 0.$$

$$\frac{d}{dz} g(z) \approx 0(1-0) \approx 0.$$

$$g'(z) = \frac{d}{dz} (g(z))$$

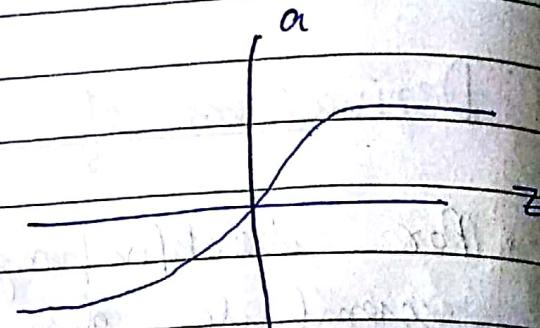
$$a = g(z) = \frac{1}{1+e^{-z}}$$

$$g'(z) = a(1-a).$$

$\rightarrow \tanh$

$$g(z) = \tanh(z)$$

$$\frac{d}{dz} g(z) = \text{Mop of } g(z) \text{ at } z$$



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

$$g'(z) = 1 - (\tanh(z))^2.$$

$$z = 10$$

$$\tanh(z) \approx 1$$

$$g'(z) = 0.$$

$$z = -10$$

$$\tanh(z) \approx -1$$

$$g'(z) = (1 - (-1)^2) \\ = 0.$$

$$z = 0$$

$$\tanh(z) \approx 0$$

$$g'(z) = 1$$

$$a = g(z)$$

$$g'(z) = (1-a)^2$$

Cordless phone  
Pagers  
Mobile phone.

### Few terminologies

- MSC → Mobile switching centre
- BSC → Base station control.

Ericsson company ke derives Mohali BSNL exchange in 1981



Public switch.

telephone n/w

(1)

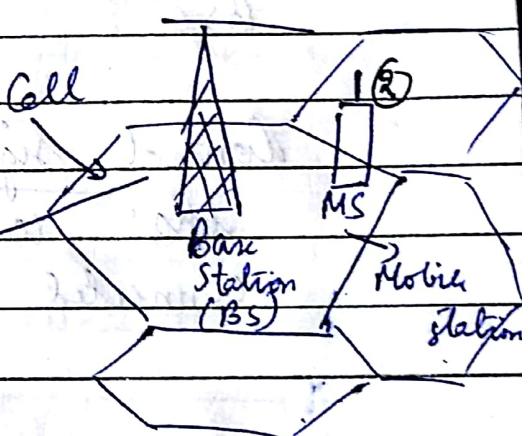
PSTN

(2)

MSC

BSC

BSC



Mobil → Mobile Station

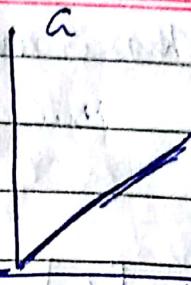
Tower → Base station

Area jisme tower signals dela hai → cell

Sesaciber → jiska ph hai

→ RELU.

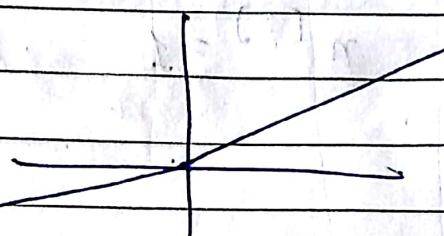
$$g(z) = \max(0, z)$$



$$g'(z) = \begin{cases} 0 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \\ \text{undefined} & \text{if } z = 0 \end{cases}$$

→ Leaky ReLU

$$g(z) = \max(0.01z, z)$$



$$g'(z) = \begin{cases} 0.01 & \text{if } z < 0 \\ 1 & \text{if } z \geq 0 \end{cases}$$

## Gradient descent for neural networks

How to compute grad descent for NN with 1 hidden layer.

NN with 1 hidden layer:

Parameters :  $w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]}$

$n_x = n^{[0]}$        $n^{[1]}$        $n^{[2]}$   
 ↓                  ↓                  ↓  
 no of input features    no of hidden units    no of output units.

Then matrix  $w^{[1]} = n^{[1]} \times n^{[0]}$   
 $n^{[1]} \times n^{[0]}$

$b^{[1]} = n^{[1]} \times n^{[0]} \times 1$

$w^{[2]} = (n^{[2]}, n^{[1]})$

$b^{[2]} = (n^{[2]}, 1)$

$n^{[2]} = 1$  All the examples seen till now.

for NN, loss fun<sup>n</sup>:  $J(w^{[1]}, b^{[1]}, w^{[2]}, b^{[2]})$   
 $= \frac{1}{m} \sum_{i=1}^m l(\hat{y}_i, y_i)$

In case of binary clas<sup>n</sup>, loss fun<sup>n</sup> is same that was in logistic regression.

$L \rightarrow$  is loss when NN predicts  $\hat{y}$ .

To train the parameters of our algz we need to perform gradient descent.

→ for training parameters of NN, it is important to initialize the parameters randomly rather than to all 0's

→ Then each loop of grad descent

would compute predictions  $\hat{y}$ .

→ Then find derivatives of cost fun' w.r.t  
variable  $d w^{[1]}, b^{[1]}, \dots$  & so on.

→ Finally gradient descent update.

$$w^{[1]} = w^{[1]} - \alpha d w^{[1]}$$

→ learning rate

Gradient descent :

→ Repeat 3

Compute predictions ( $\hat{y}^{(i)}$  ... for  $i=1 \text{ to } m$ )  
 $d w^{[1]} = \frac{\partial J}{\partial w^{[1]}}, d b^{[1]} = \frac{\partial J}{\partial b^{[1]}} \dots$

$$w^{[1]} = w^{[1]} - \alpha d w^{[1]}$$

$$b^{[1]} = b^{[1]} - \alpha d b^{[1]}$$

Formulas for computing derivatives.

→ Forward propn.

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = g(z^{[1]})$$

$$z^{[2]} = w^{[2]} A^{[1]} + b^{[2]}$$

$$A^{[2]} = g(z^{[2]}) = \sigma(z^{[2]})$$

↓ for binary  
class n.

→ Backward propagation :

$$d z^{[2]} = A^{[2]} - y$$

$$y = [y^{(1)} \dots y^{(m)}]$$

all rows stacked horizonally

$$d w^{[2]} = \frac{1}{m} d z^{[2]} A^{[1]T}$$

$$d b^{[2]} = \frac{1}{m} \text{np.sum}(d z^{[2]}, \text{axis}=1, \text{keepdims=True})$$

python numpy command for summing across 1-D. in this case summing horizontally.

Keepdims → prevents python from producing Rank 1 matrix & so produce vector.

$$\text{Rank 1} \rightarrow (5,)$$

$$d z^{[1]} = (W^{[2]T} d z^{[2]} * g'(z^{[1]}))$$

$(n^{[1]}, m)$   
 $(n^{[1]}, m)$

derivative of  
 activ'n funn  
 for hidden layer

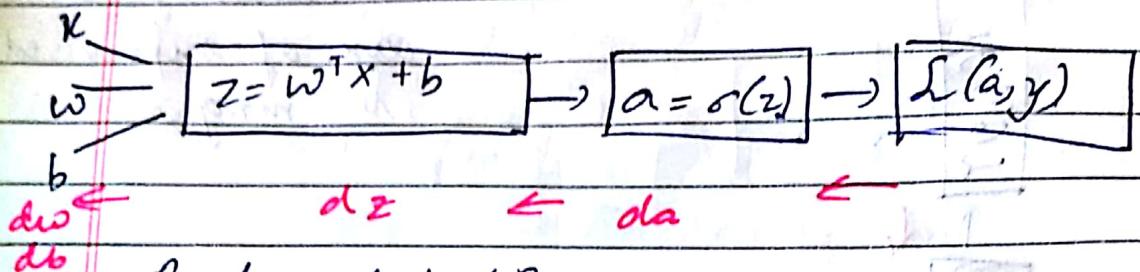
element wise  
product

$$d w^{[1]} = \frac{1}{m} d z^{[1]} x^T$$

$$d b^{[1]} = \frac{1}{m} \text{np.sum}(d z^{[1]}, \text{axis}=1, \text{keepdims=True})$$

## Backpropagation Intuition

In logistic regression,  
The forward prop<sup>n</sup> is  $\rightarrow$   
finding  $\hat{y}$  i.e prediction



Backward prop<sup>n</sup>  
finding derivatives.

$$l(a, y) = \text{loss} = -(y \log a + (1-y) \log(1-a))$$

$$\begin{aligned} da = \frac{d}{da} l(a, y) &= -y \cancel{\log a} - (1-y) \cancel{\log(1-a)} \\ &= -\frac{y}{a} + \frac{(1-y)}{1-a} \end{aligned}$$

$$dz = da \times g'(z) = \frac{d l(a, y)}{da}$$

$$= \frac{d l(a, y)}{dz} \times \frac{da}{dz}$$

$$= da \times g'(z)$$

$$g(z) = \sigma(z)$$

~~$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial x}$$~~

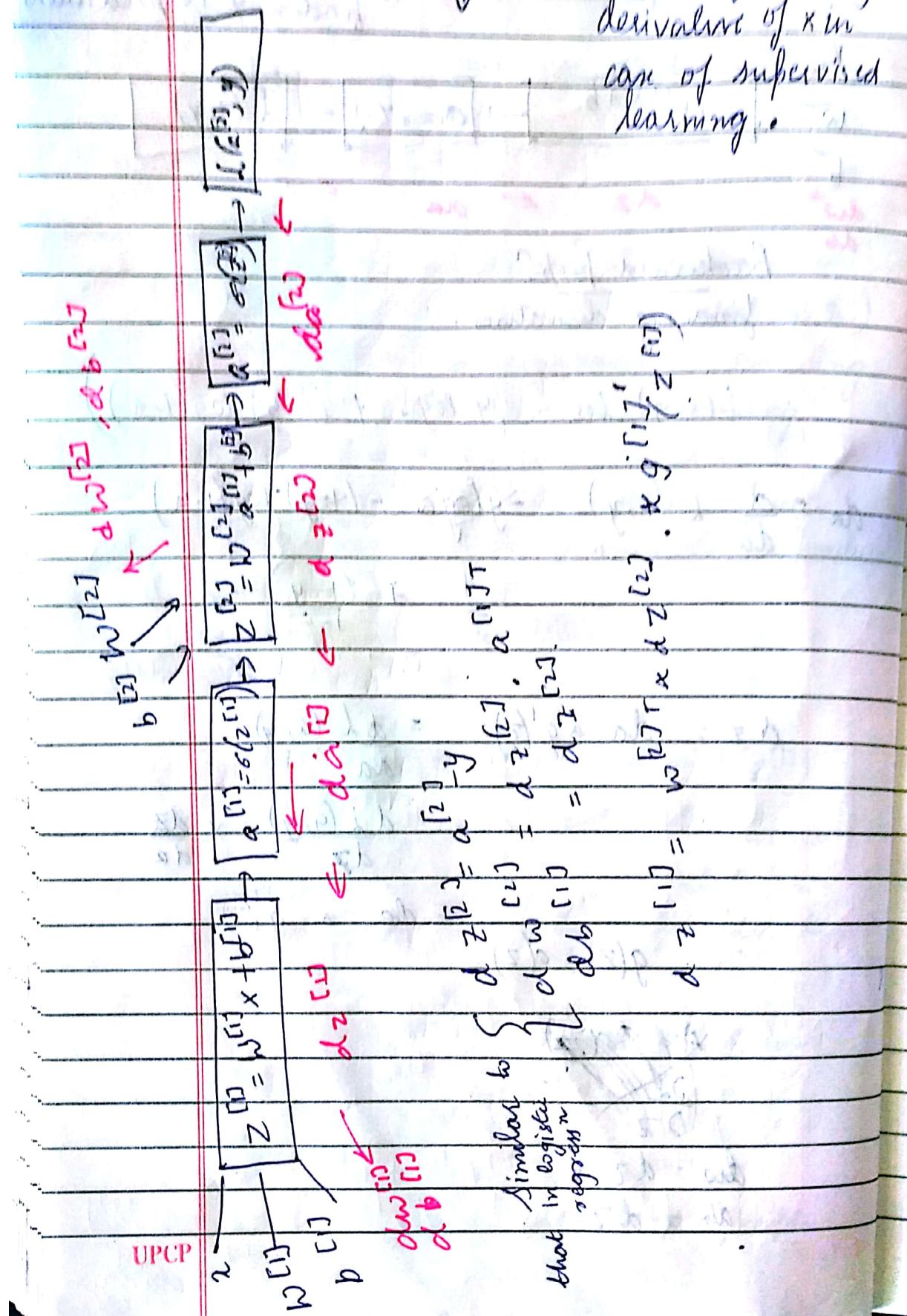
$$dw = dz \cdot x$$

$$db = dz$$

$\Rightarrow$  In NN we do 2. steps b/w 2 layers  
are there.

Steps are  $\downarrow$

No need for taking  
derivative of  $x$  in  
case of supervised  
learning.



Check dimensions

$$\begin{array}{c} x_1 = 0 \\ x_2 = 0 \\ x_3 = 0 \end{array} \rightarrow \begin{array}{c} 0 \\ 0 \\ 0 \end{array} \rightarrow y$$

$$n_x = n^{(0)} \quad n^{(1)} \quad n^{(2)} = 1$$

$$z^{[2]} = \begin{pmatrix} n^{[2]} \\ n^{[1]} \end{pmatrix}$$
$$z^{[2]}, dz^{[2]} \rightarrow (n^{[2]}, 1)$$

$$z^{[1]}, dz^{[1]} \rightarrow (n^{[1]}, 1)$$

$$\begin{aligned} dz^{[1]} &= w^{[2]T} dz^{[2]} * g^{[1]'} z^{[1]} \\ (n^{[1]}, 1) &\quad (n^{[1]}, n^{[2]}) \quad (n^{[2]}, 1) \quad (n^{[1]}, 1) \\ &\quad \downarrow \quad \downarrow \quad \downarrow \\ (n^{[1]}, 1) &\quad (n^{[1]}, 1) \\ &\quad \boxed{(n^{[1]}, 1)} \\ &\quad \downarrow \quad \downarrow \\ &\quad (n^{[1]}, 1) \end{aligned}$$

## Summary of Gradient descent.

$$dz^{[2]} = a^{[2]}$$

$$dw^{[2]} = dz^{[2]} a^{[2]T}$$

$$db^{[2]} = dz^{[2]}$$

$$dz^{[1]} = w^{[2]T} dz^{[2]} * g^{[1]'} (z^{[1]})$$

$$dw^{[1]} = dz^{[1]} x^T$$

$$db^{[1]} = dz^{[1]}$$

# Vectorized Implementation.

$$z^{[1]} = w^{[1]} x + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]})$$

$$z^{[1]} = \begin{bmatrix} z_{[1](1)} & z_{[1](2)} & \dots & z_{[1](m)} \end{bmatrix}$$

Stacking all egs that were there horizontally - i

$$z^{[1]} = w^{[1]} x + b^{[1]}$$

$$A^{[1]} = g^{[1]}(z^{[1]})$$

Summary

$$dz^{[2]} = A^{[2]} - y$$

$$dw^{[2]} = \frac{1}{m} dz^{[2]} A^{[1]\top} \text{ bcz core}$$

$$db^{[2]} = \frac{1}{m} \text{ np.sum}(dz^{[2]}, \text{ axis}=1, \text{ keepdim=True})$$

$$dz^{[1]} = w^{[2]\top} dz^{[2]} * g^{[1]\top} / z^{[1]}$$

$(n^{[1]}, m)$        $(n^{[1]}, m)$        $(n^{[1]}, m)$

element wise product

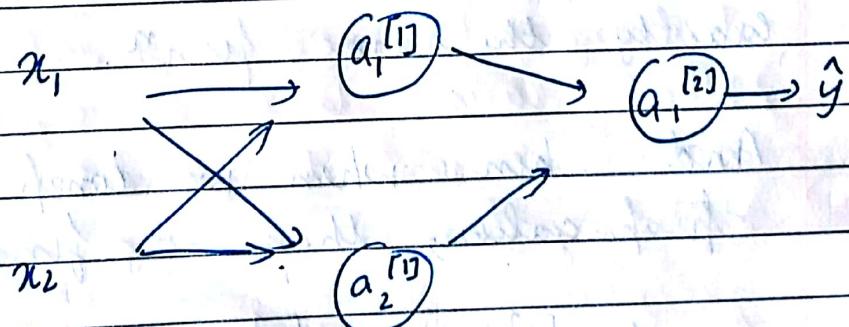
$$\frac{dW^{[1]}}{m} = \frac{1}{m} dZ^{[1]} X^T$$

$$\frac{db^{[1]}}{m} = \frac{1}{m} \text{np.sum}(dZ^{[1]}, \text{axis}=1, \text{keepdims}=\text{True})$$

→ It was ok. to initialize weights to 0 for logistic regression.

But for a NN if we initialize the weights or parameters to all zero and then apply gradient descent then it won't work.

What happens if you initialize the weights to zero?



So if we have 2 I/P features, i.e.  $n_B = 2$  and 2 hidden units so  $n[1] = 2$ .

And so matrix associated with hidden layer,  $W^{[1]} = 2 \times 2$

Let's say you initialize it to all 0s,  
so 0 0 0 0, 2x2 matrix.

$$w^{[1]} = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix}$$

$$b^{[1]} = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

Initializing bias  
to 0 is ok  
but initializing  
weights to 0 is just not ok. It is a  
prob.

Problem with this type of initializ^n is  
that for any eq you give you have:

$a_1^{[1]} = a_2^{[1]}$ . So this activation  
units will be same bcz both of  
these hidden units are computing  
exactly the same fun^n.

And then when we compute back  
propagation, then we find that

$$dZ_1^{[1]} = dZ_2^{[1]}$$

$$w^{[2]} = \begin{bmatrix} 0 & 0 \end{bmatrix}$$

So if we initialize the parameters  
with 0 then for each layer,

all units becomes identical.

$d_w =$

i.e After each iteration of training, the 2 hidden units are computing exactly the same function.

So all hidden units become identical & no matter how long you do the training they remains identical.

### Solution

Initialize your parameters randomly.

$$w[1] = np.random.randn(2, 2) * 0.01$$

generates gaussian random variables multiplied by very small values.

initialize with very small values.

$$b = np.zeros(2)$$

$b$  does not face symmetry prob.

So it is okay to initialize it with zeros.

Why multiplied with 0.01 & not 1 or 2 we have to initialize with very small random values.

Bcz if we are using a tanh or sigmoid activation function, or just at O/P layer.

If weights are too large, then activ<sup>n</sup> fun<sup>n</sup>,  $z^{(1)} = w^{(1)}x + b^{(1)}$   
 $a^{(1)} = g(z^{(1)})$

so if  $w \gg$  very big

then  $a^{(1)}$  applied to  
so in that case we are more likely to end up in the flat part of tanh or sigmoid function - where slope is small & gradient descent = slow & learning = slow.

So if  $w$  large then we are likely to end up even at the start of training with very large values of  $z$ . which causes tanh or sigmoid activation fun<sup>n</sup> to be saturated, thus slowing down learning.

If we don't have any sigmoid or tanh activ<sup>n</sup> function throughout our neural network.

But if we are doing binary class<sup>n</sup> & O/P is a sigmoid fun<sup>n</sup>, then you don't want initial parameters to be too large

So multiplying by 0.01 would be reasonable to try.

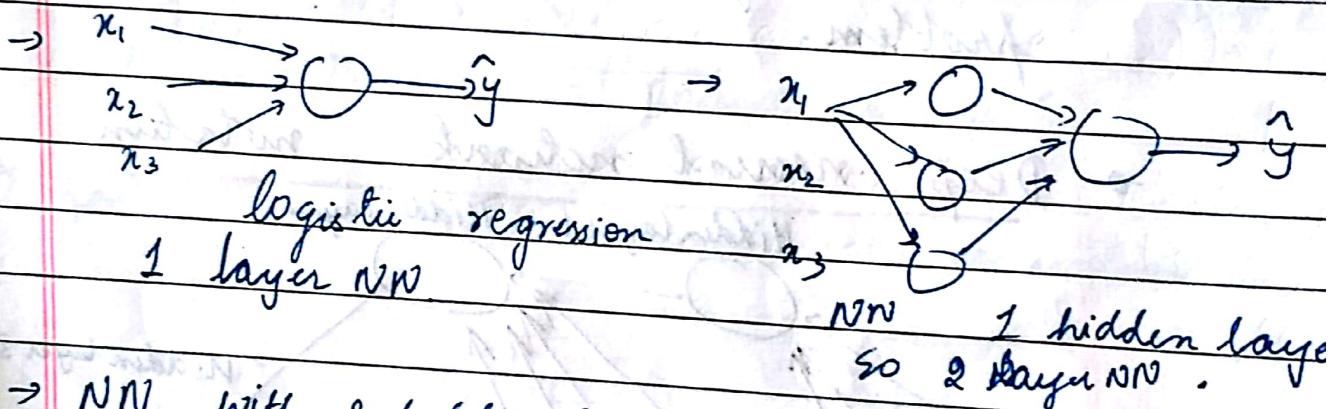
And same for  $w^2$ .

→ When we have small NN then 0.01 will work ok.

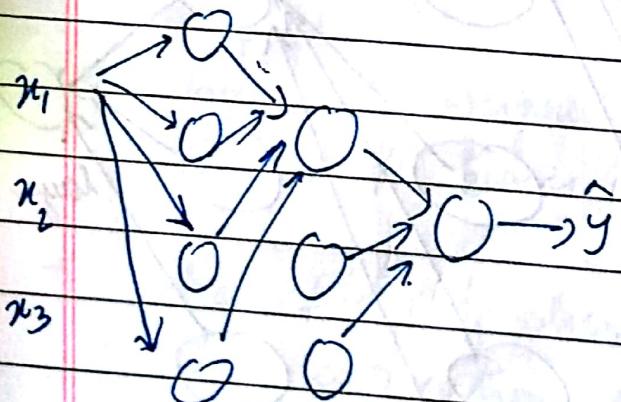
→ But when training a very deep NN, then you might need to pick up a diff constant than 0.01.

## Week 4

What is a deep Neural Network



→ NN with 2 hidden layers



3 layer NN