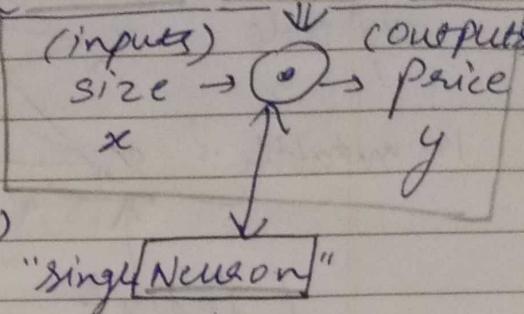
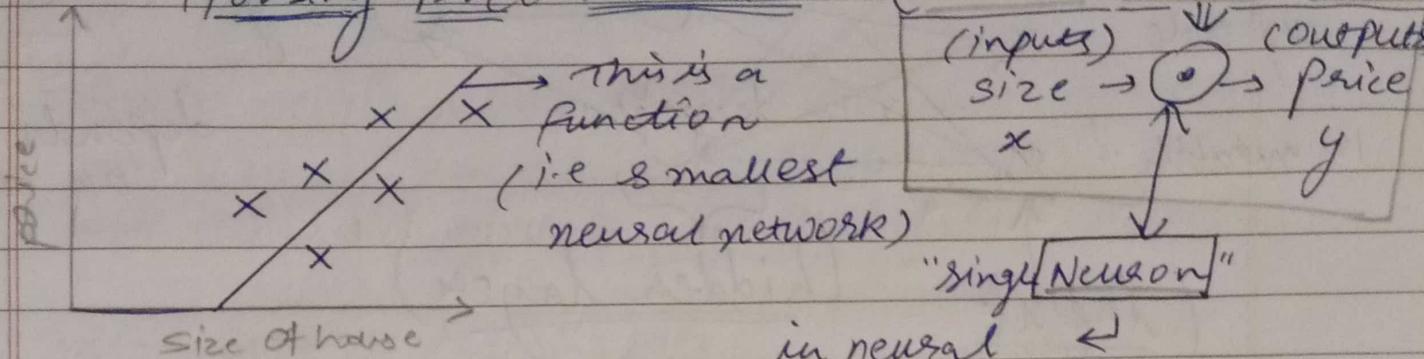


Neural Networks and Deep learning

(Deep learning - AI Specialization - Coursera)

What is neural network?

⇒ Deep learning refers to train Neural Networks
Housing Price Prediction (little neural network)



(as price won't be -ve)
 can continue the line to 0)

in neural network

which implements the func

Neuron → inputs the size, computes the funct, outputs linear estimated price

⇒ function like this tending to 0 & then following straight line is called $y = \text{max}(0)$ (ReLU)

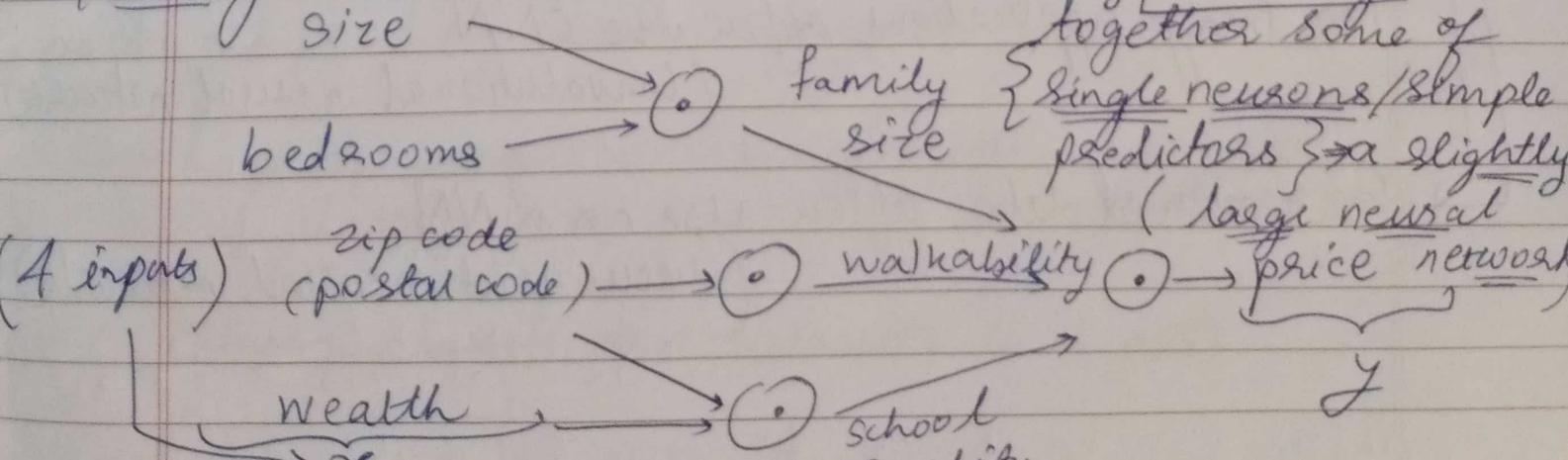
{ "ReLU" → Rectified linear Unit } \Rightarrow max of 0 due to which such shape of func is obtained.

In bigger Neural network many such neurons are stacked together.

Housing Price Prediction -

so by stacking

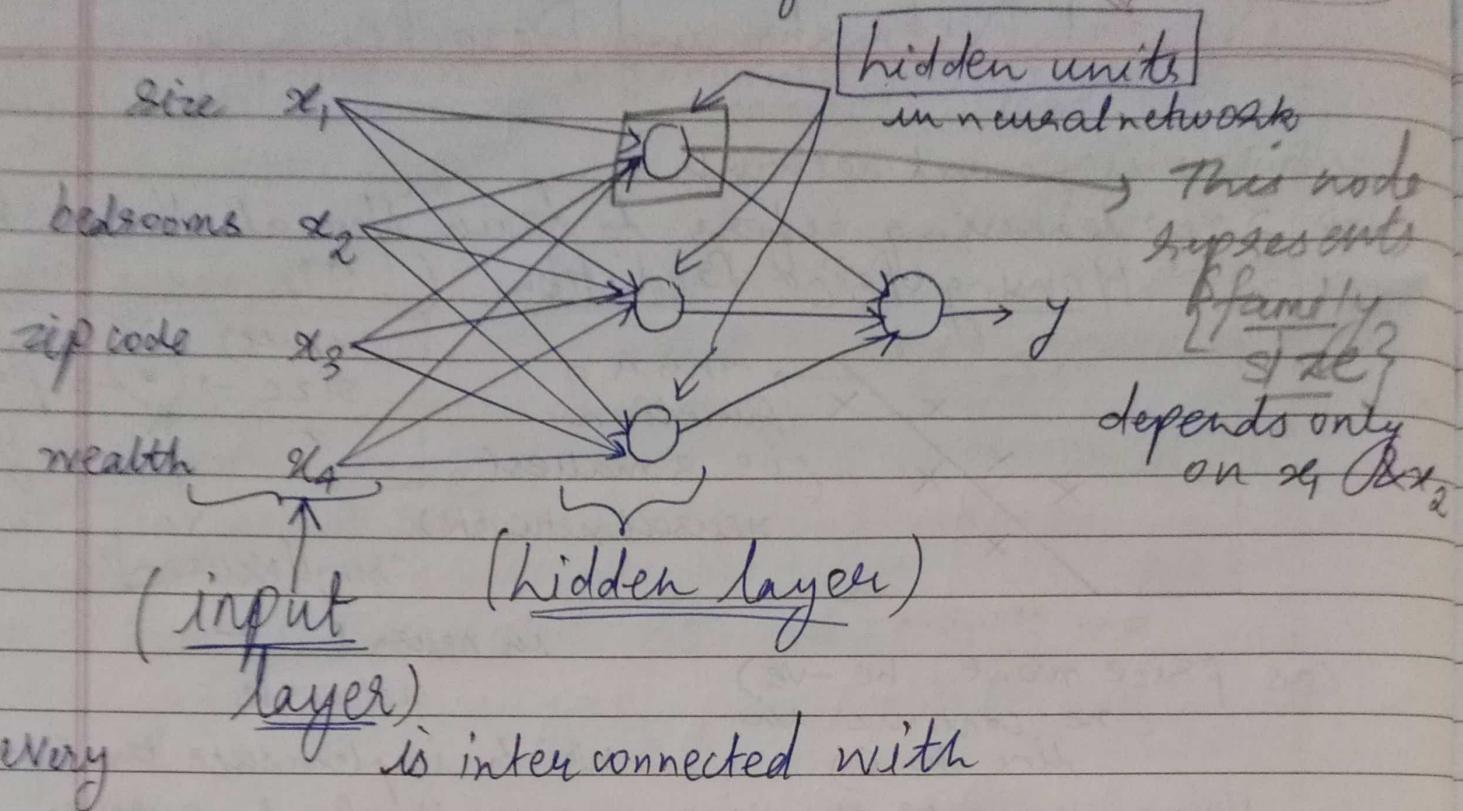
together some of
single neurons / simple predictors \Rightarrow a slightly large neural



⇒ each of these circles may be ReLU or some other slightly non linear func

things occurring in middle will be figured by
the neural networks itself

Date _____
Page _____



Supervised learning with Neural Networks -

Supervised Learning Input (x) Output (y) Application

Audio	Text transcript	Speech recognition
English	Chinese	Machine translation
Ad, user info	click on ad ? (0/1)	Online Advertising

For image applications, often use CNN
(Convolutional neural network)

For sequenced data, often use an RNN
(Recurrent neural network)

Structured Data

Unstructured Data
CLASSMATE

- ⇒ databases of data (difficult) ⇒ audio, images, text
* computers struggle to sense unstructured data compared to structured data *

⇒ Why is Deep learning taking off?

Scale is driving DL

↳ size of neural network

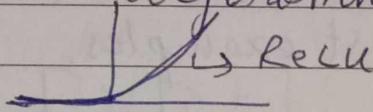
↳ amount of data (labelled data)

{
m → size of training sets }
(number of training examples)}

Sigmoid function to ReLU func

→ gradient descent works fast this

algorithm



Sigmoid

Basics of Neural Network Programming

framework for NN

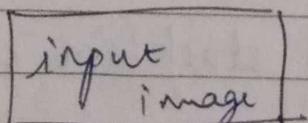
Build deep neural network (CNN)

forward pass / forward propagation step

backward pass / backward propagation step

→ logistic regression is an algorithm for binary classification

Binary classification → (Classify in 2 grp)



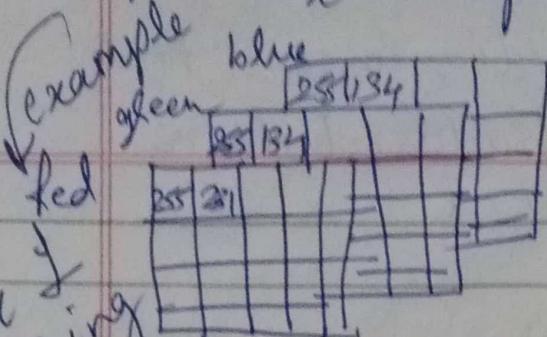
→ 1(cat) vs 0(non-cat)
output

$n = n_x$: represents dimension of input feature vector

classmate

Date _____

Page _____



64x64 pixel value

for each matrix

we will transform 3 into feature vector whose dimension will

corresponding label is (x) then be $n_x = 12288$

Notation:

(x, y)

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

$6 \times$ dimensional feature

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

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \\ 255 \\ 134 \end{bmatrix}$$

determine

m training examples $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)})\}$

$M = M_{\text{train}} \Rightarrow$ number of train examples

$M_{\text{test}} \Rightarrow$ number of test examples

$$X = \left[\begin{array}{c|c|c|c|c} 1 & 1 & 1 & \cdots & 1 \\ \hline x^{(1)} & x^{(2)} & \cdots & x^{(m)} & \\ \hline \end{array} \right] \quad \begin{array}{l} \text{Number of columns} \\ \text{matrix } X \text{ will have} \\ 'm' \text{ where} \\ M \text{ is no. of training} \\ \text{examples} \end{array}$$

↑ column
Training set inputs
↓ m

Height of this matrix is ~~n_x~~ n_x

$$X \in \mathbb{R}^{n_x \times m}$$

{ dimensional matrix

$X.\text{shape} \Rightarrow$ python command

for (finding shape) of
"matrix" that this is a $[n_x, m]$

output labels $y \rightarrow$ its easy to stack in column

$$Y = [y^{(1)} \ y^{(2)} \ \dots \ y^{(m)}] \rightarrow \text{no. of rows} = 1$$

CLASSTIME _____
Date _____
Page _____

$Y \in \mathbb{R}^{1 \times m}$ Y's shape = $(1, m)$

1-dimensional matrix

Logistic Regression

Given x get a prediction $\hat{y} \Rightarrow$ estimate of y

$X \in \mathbb{R}^{n \times x}$ $\hat{y} = P(y=1 | x)$ (y hat)
 X dimensional vector \hookrightarrow probability of chance that y is equal to 1 given input features x

\hat{y} tells what is chance this is a cat pic

Parameters of Logistic Regression:-

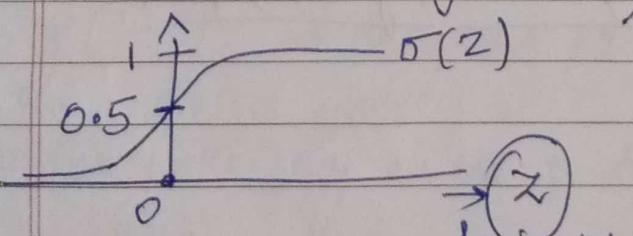
$$W \in \mathbb{R}^{n \times x} \quad b \in \mathbb{R}$$

W is also an X dimensional

Output: $\hat{y} = \underbrace{\{W^T x + b\}}_{0 \leq \hat{y} \leq 1} \xrightarrow{\text{vector}} \text{linear function of } X$

This quantity can even be greater than one or even negative so that's not correct.

Therefore: $\hat{y} = \sigma(W^T x + b)$ sigmoid func



where $\sigma(z) \Rightarrow \sigma(z)$
 here

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

formula
for sigmoid fun

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

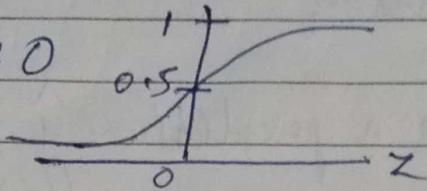
where $z \in \mathbb{R}$

If z is large; $1 + e^{-z} \approx 1$ If z is very small/very large
 $\therefore \sigma(z) \approx 1$ large -ve number.

that relates

when z is very large -ve

then $\sigma(z) \approx 0$



so this is

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

keep the parameters W and b separate

$$X_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$

$$\hat{y} = \sigma(W^T X)$$

$$W = \begin{bmatrix} W_0 \\ W_1 \\ W_2 \\ \vdots \\ W_{n_x} \end{bmatrix} \leftarrow w$$

But don't

follow this approach #

to train w & b of logistic regression parameters

need define a cost function

$$\hat{y}^{(i)} = \sigma(W^T X^{(i)} + b) \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}$$

\Rightarrow given $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$, want

$$\boxed{\hat{y}^{(i)} \approx y^{(i)}}$$

$$Z^{(i)} = \frac{1}{1 + e^{-Z^{(i)}}} = W^T X^{(i)} + b$$

$\left. \begin{array}{l} x^{(i)} \\ y^{(i)} \\ Z^{(i)} \end{array} \right\}$ ith example

Loss(error) function → is restricted to single training measure how well our algorithm is doing. (applied) example

$L(\hat{y}, y)$ may be square error $(\hat{y} - y)^2$
 (loss) true label or $\frac{1}{2}$ of square error $\frac{1}{2}(\hat{y} - y)^2$
 fun algorithm outputs
 ↳ measure how good output \hat{y} is when true label is y
 required to define to
 (need)

But loss function used in logistic regression is as follows:

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

Why this is preferred?

If $y=1$ $L(\hat{y}, y) = \boxed{-\log \hat{y}}$ ← we want loss func as small as possible
 $\log \hat{y} \rightarrow$ as large as possible

$\hat{y} \rightarrow$ want large $\hat{y} = \sigma(z)$
 it can never be bigger than one

so for $y=1$ \hat{y} will be close to 1

If $y=0$ $L(\hat{y}, y) = \boxed{-\log(1-\hat{y})}$ ← as small as possible

$\log(1-\hat{y})$ as large as possible . . . want \hat{y} as small as possible
 $\therefore \hat{y}$ lies between 0 & 1

so for $y=0$ \hat{y} as close to zero as possible

In general $y=1$ \hat{y} is as large as possible
 $y=0$, \hat{y} is as small as possible
 true labels

Cost function -
of parameter W and b

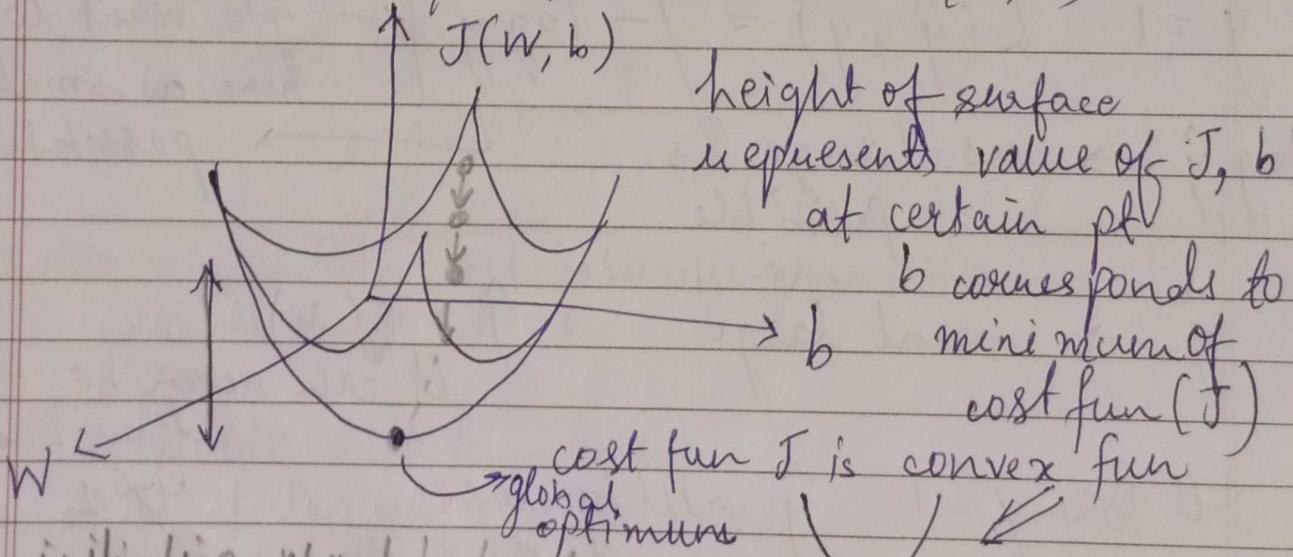
$$J(W, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \log \hat{y}^{(i)} + (1-y^{(i)}) \log(1-\hat{y}^{(i)})]$$

Gradient Descent -

↳ algorithm to train parameter W on training set

Cost function measures how well are the parameters W and b are doing on training set

want to find W, b that minimize $J(W, b)$

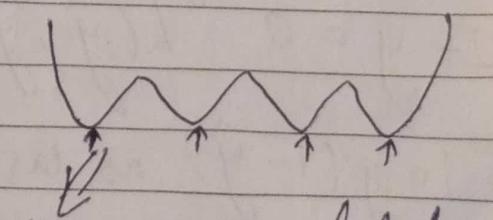


initialize W and b by some initial value

for logistic regression
any initialization works

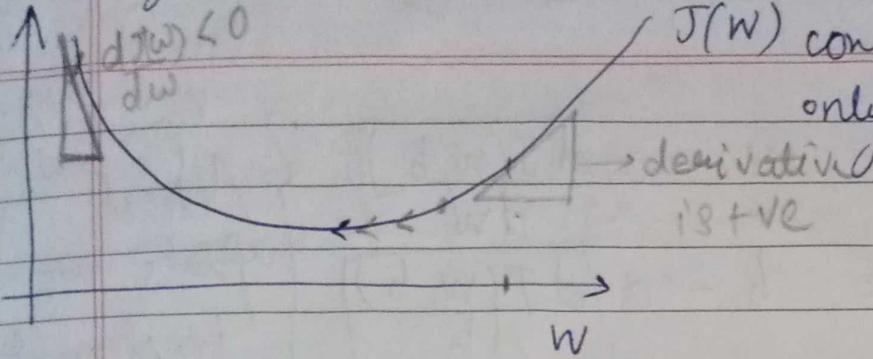
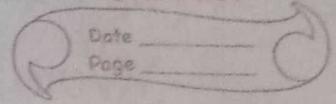
gradient descent starts at that initial pt and takes step in the steepest downhill direction.

(this is 1 iteration of gradient descent)



non-convex lots of different local optima

This iterations continue until it converges to global optimum or close to global optimum.



Gradient descent algorithm repeats and this gets repeated until algorithm converges.

Repeat { learning rate

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

learning rate: controls how big a step we take on each iteration in gradient descent algorithm.

: update of the change
Want to make to parameters w

to represent derivative term in code $\Rightarrow 'd w'$

$$w := w - \alpha d w$$

gradient descent update

\therefore derivative is +ve

$$w := w - \alpha d w$$

subtracting from w

so w will move towards left

\Rightarrow gradient descent algorithm slowly decrease parameter

-ve slope/derivative

so thereby w will increase and move towards right

So basically from wherever u start you will try to move towards centre(global minimum)

In previous examples considering cost function
Where w is only parameter

But in reality its

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

$$J(w, b)$$

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

Actual partial derivative update we implement.

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

$\frac{dw}{(update w)}$

If J is a func of 1 variable
we use ' d '

Since J is a func of 2 parameters

This basically
represents
how much
 $J(w, b)$ slopes
in b direction

$$\frac{db}{\partial b}$$

(update b)

for a straight line, slope of the line
remains same (derivative of function is
constant i.e same)

1 derivative of a func just means slope of fun
slope can be diff at diff pts on function.

Computation Graph -
(computes output)

eg - $J(a, b, c) = 3(a + bc)$

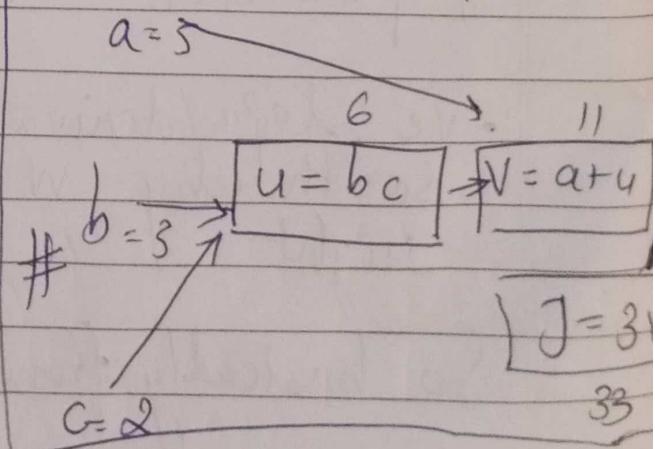
$$u = bc$$

$$v = a + u$$

$$J = 3v$$

for computing
derivatives, we
need to go left pass
will use right

Computation graph :-



(left to right pass)

forward pass / forward propagation step to compute output

backward pass / backward propagation step to compute
Computing derivatives (Right to left pass) (gradient of output)

$$\frac{dJ}{dv} = ? \quad J = 3v$$

$$v = 11 \rightarrow 11.001$$

Derivatives with

$$= 3 \quad J = 33 \xrightarrow{\text{compute}} \text{Final } 33.003 \quad \text{Computation Graph}$$

If we want derivative of output variable

we are one step of backpropagation \rightarrow one step
doing

$$a = 5 \rightarrow 5.001$$

$$v = 11 \rightarrow 11.001$$

$$J = 33 \rightarrow 33.003$$

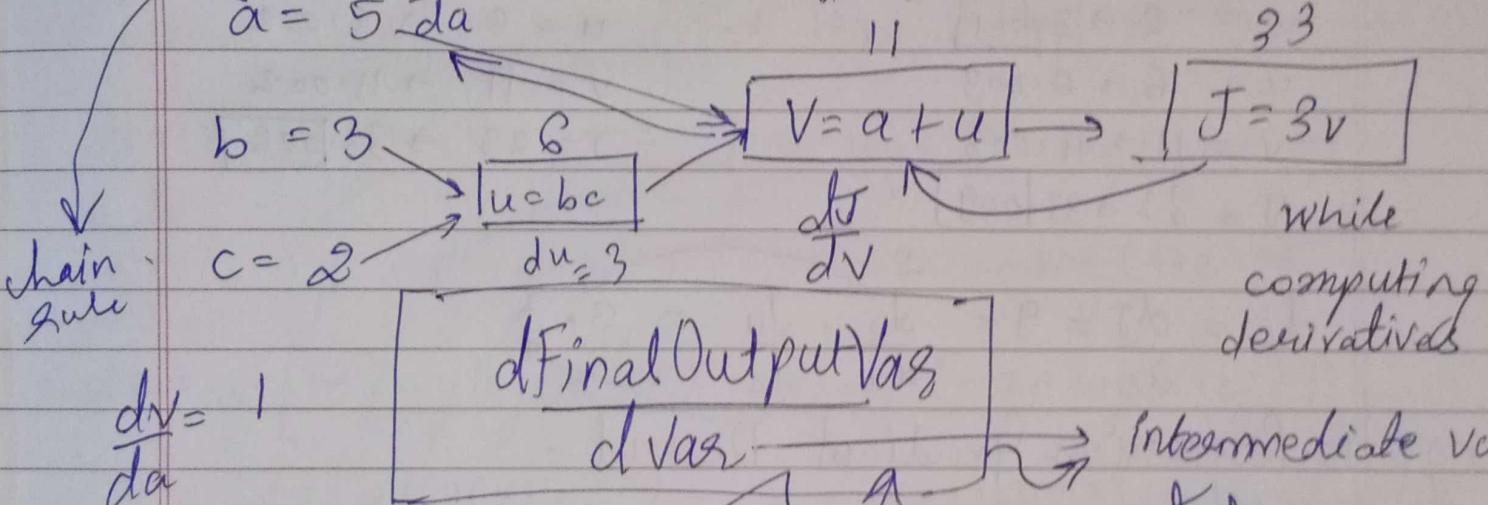
backwards
in graph

$$\frac{dJ}{da} = 3 = \frac{dJ}{dv} \cdot \frac{dv}{da}$$

$$(3 \cdot 1) \frac{dJ}{da} \\ a = 5 \cdot \frac{dJ}{da}$$

$$dv = 3$$

$$33$$



while
computing
derivatives

$\frac{dJ}{dv} = 3$ \downarrow representation
 $\frac{dJ}{dv} = 3$ \downarrow format

$\frac{dJ}{dvar}$ $\frac{dJ}{dvar}$ $\frac{dJ}{dvar}$

$$dv = \left(\frac{dJ}{dv} \right)$$

$$\frac{dJ}{da} = 3 \quad \frac{dJ}{da} = 3 \quad \frac{dJ}{du} = 3$$

$$u = 6 \rightarrow 6.001 \\ v = 11 \rightarrow 11.001 \\ J = 33 \rightarrow 33.003$$

$$\frac{dJ}{db} = \frac{dJ}{du} \cdot \frac{du}{db} \\ \frac{dJ}{db} = 3 \cdot 6 = 18$$

$$b = 3 \rightarrow 3.001 \\ u = 6 \rightarrow 6.002$$

$$\frac{dJ}{du} = 3 \cdot \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 3 \cdot 3 = 27$$

dvar: derivative of final gp variable
wrt to intermediate variables.

$$\begin{aligned}
 a &= 5 \\
 da &= 3 \quad (\frac{dJ}{da} = 3) \\
 b &= 3 \\
 db &= dJ = 6 \quad (\frac{dJ}{db} = 6) \\
 c &= 2 \\
 dc &= dJ = 9 \quad (\frac{dJ}{dc} = 9)
 \end{aligned}$$

$v = u + a \quad J = 3v$
 $dv = 3 = \frac{dJ}{dv}$

$u = bc \quad J = 33$
 $du = \frac{dJ}{du} = 3 \quad dv = 3 = \frac{dJ}{dv}$

$$\begin{aligned}
 \frac{dJ}{dc} &= 9 & u &= 6 \rightarrow 6.001 \\
 3 &= \frac{dJ}{du} = \frac{dJ}{dv} \cdot \frac{dv}{du} = 3 \cdot 1 & v &= 11 \rightarrow 11.001 \\
 \frac{dJ}{db} &= 6 = \frac{dJ}{du} \cdot \frac{du}{db} = 3 \cdot 2 & J &= 33 \rightarrow 33.003 \\
 c &= 2 \rightarrow 2.001 & b &= 3 \rightarrow 3.001 \\
 u &= 6 \rightarrow 6.003 & u &= 6 \rightarrow 6.002 \\
 v &= 11 \rightarrow 11.003 & v &= 11 \rightarrow 11.002 \\
 J &= 33 \rightarrow 33.009 & J &= 33 \rightarrow 33.006
 \end{aligned}$$

$$\frac{dc}{dc} = \frac{dJ}{dc} = 9 = \frac{dJ}{du} \cdot \frac{du}{dc} = 3 \cdot 3$$

Logistic Regression Gradient Descent -

for logistic regression $z = w^T x + b$

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

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

Computational Graph -

$$\hat{y} = a \quad (\text{denoted})$$

$$\begin{array}{c}
 x_1 \xrightarrow{} \\
 \xrightarrow{} x_2 \xrightarrow{} \\
 \xrightarrow{} w_1 \xrightarrow{} \\
 \xrightarrow{} w_2 \xrightarrow{} \\
 \xrightarrow{} b \xrightarrow{} \\
 \boxed{z = w_1 x_1 + w_2 x_2 + b} \xrightarrow{} \\
 \boxed{\hat{y} = \sigma(z)} \xrightarrow{} \\
 \boxed{L(\hat{y}, y)} \xrightarrow{} \\
 \boxed{L(a, y)} \xrightarrow{} \\
 \downarrow \frac{da}{dL(a, y)}
 \end{array}$$

In logistic regression, we have to modify parameters w and b in order to minimize the loss(L)

This is entirely for one training example

Now computing for derivatives (backward propagation step)

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

$$\frac{da}{da} = \frac{dL(a, y)}{da} = -\left[\frac{y}{a} + \frac{(1-y)(-1)}{1-a}\right] = -\frac{y}{a} + \frac{(1-y)}{1-a}$$

$$d_z = \frac{dL(a, y)}{dz} = d[-(y \log a + (1-y) \log(1-a))]$$

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

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

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

$$dz = \frac{d}{dz} [-\{ -y \log(1+e^{-z}) + (1-y)(-z - \log(1+e^{-z})) \}]$$

$$= \frac{d}{dz} [y \log(1+e^{-z}) + z + \log(1+e^{-z}) - yz - y \log(1+e^{-z})]$$

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

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

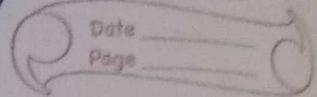
$$dz = a - y = \frac{dL}{dz}$$

$$\frac{dL}{dw_1} = dw_1 = x_1 \cdot dz$$

$$\frac{dL}{dw_2} = dw_2 = x_2 \cdot dz$$

$$\left. \begin{array}{l} w_1 := w_1 - \alpha dw_1 \\ w_2 := w_2 - \alpha dw_2 \\ b := b - \alpha db \end{array} \right\} \text{updating values} \quad \alpha : \text{learning rate}$$

Now, we'll check for m training examples how this works-



Gradient Descent on m Examples - (Training eg)

$$J(w, b) = \frac{1}{m} \sum_{i=0}^m L(a^{(i)}, y^{(i)}) \quad (x^{(i)}, y^{(i)})$$

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

$$dw_1^{(i)}, dw_2^{(i)}, db^{(i)}$$

derivative w.r.t to w_1 , of overall cost function is also average of derivatives w.r.t to w_1 , of individual loss terms.

$$\frac{\partial}{\partial w_1} J(w, b) = \frac{1}{m} \sum_{i=0}^m \frac{\partial}{\partial w_1} L(a^{(i)}, y^{(i)}) \times [dw_1^{(i)} - x^{(i)} y^{(i)}] \times$$

\Rightarrow algorithm for m training examples -

$$J = 0; \{ dw_1 = 0; dw_2 = 0; db = 0 \};$$

for $i = 1$ to m —

$$\begin{cases} z^{(i)} = w^T x^{(i)} + b \\ a^{(i)} = \sigma(z^{(i)}) \end{cases}$$

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

for $j = 1$ to n

$$\begin{cases} dz^{(i)} = a^{(i)} - y^{(i)} \\ dw_1 += x_1^{(i)} dz^{(i)} \\ dw_2 += x_2^{(i)} dz^{(i)} \\ db += dz^{(i)} \end{cases} \quad \text{considering } n=2 \text{ (two features)} \quad \text{end of loop}$$

$dw_1/m = \frac{1}{m} \sum_{i=0}^m x_1^{(i)} dz^{(i)}$

$dw_2/m = \frac{1}{m} \sum_{i=0}^m x_2^{(i)} dz^{(i)}$

$db/m = \frac{1}{m} \sum_{i=0}^m dz^{(i)}$

$dw_1/m = \frac{1}{m} \sum_{i=0}^m x_1^{(i)} (a^{(i)} - y^{(i)})$

$dw_2/m = \frac{1}{m} \sum_{i=0}^m x_2^{(i)} (a^{(i)} - y^{(i)})$

$db/m = \frac{1}{m} \sum_{i=0}^m (a^{(i)} - y^{(i)})$

$dw_1/m = \frac{1}{m} \sum_{i=0}^m x_1^{(i)} (\sigma(w^T x^{(i)} + b) - y^{(i)})$

$dw_2/m = \frac{1}{m} \sum_{i=0}^m x_2^{(i)} (\sigma(w^T x^{(i)} + b) - y^{(i)})$

$db/m = \frac{1}{m} \sum_{i=0}^m (\sigma(w^T x^{(i)} + b) - y^{(i)})$

$d w_1, d w_2, db$ are used as accumulators so they do not have a superscript i over it.

whereas $d z^{(i)}$ is for the respective training set example i.e calculated.

$$d w_i = \frac{\partial J}{\partial w_i}$$

also updated values -

$$w_1 := w_1 - \alpha d w_1$$

$$w_2 := w_2 - \alpha d w_2$$

$$b := b - \alpha d b$$

This is over a single step of gradient descent

You have to repeat this process over multiple times for gradient descent algorithm to implement.

We have to write 2 for loops for this :

(gradient descent for single step)

→ 1st for loop ←

2nd for loop ← for all the n features over here

(we have considered 2 features here)
($n = 2$ & $n_x = 2$)

Now basically writing such for loops makes the code inefficient for bigger datasets.

In Deep-learning algorithm (getting rid of for loops)
Therefore, Vectorization Techniques are used

PYTHON & VECTORIZATION

Vectorization art of getting rid of explicit loops in your code.

What is vectorization -

$$Z = \underbrace{W^T X}_{} + b$$

non-vectorized implementation -

$$Z = 0$$

for i in range(n - x)

$$Z += W(i) * X(i)$$

$$Z += b$$

$$W = \begin{bmatrix} : \\ : \end{bmatrix} \quad X = \begin{bmatrix} : \\ : \end{bmatrix}$$

$W \in R^{n_x}$ dimensional
 $X \in R^{n_x}$ vectors

Vectorized (in python or numpy)

$$Z = np.dot(W, X) + b$$

$$W^T, X$$

So for deep learning we prefer to use GPU rather than CPU but using vectorised approach we can even perform them on CPU

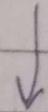
These computations involve SIMD instructions (single instructions multiple data)

so as GPU's are good at SIMD instructions than CPU's. (but CPU's are not too bad)

$$u = A v \rightarrow \text{matrix } A$$

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

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



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

for i . . . \leftarrow 1st for loop

for j . . . \leftarrow 2nd for loop

$$u_i += A[i][j] * V[j]$$

Vectorised

non
vectorised

Vectors & matrix valued functions -

Say you need to apply exponential operations on every element of a matrix/vector.

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

non-vectorised

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

$\rightarrow \text{for } i \text{ in range}(n) :$

$\rightarrow u[i] = \text{math.exp}(v[i])$

This makes the vector u set to zero

Vectorised

$\rightarrow u = np.exp(v)$

numpy consists of vector valued functions -

$np.log(v)$ element wise

$np.abs(v)$ log

element wise

absolute

value

$np.maximum(v, 0)$ take max

$v**2$ square of every element

$1/v \rightarrow v^{-1}$ inverse

Vectorising Logistic Regression -

Propagation step -

$$\begin{array}{l} X \ni \text{Training inputs} = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \quad (n_x, m) \\ \text{matrix} \qquad \qquad \qquad \mathbb{R}^{n_x \times m} \end{array}$$

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = w^T X + [b \ b \ \dots \ b]$$

M-dimensional row vector (m : training eg.)

W^T is a row vector

$$Z = [z^{(1)} z^{(2)} \dots z^{(m)}] \xrightarrow{W^T_{1 \times n_x} \text{ matrix}} \begin{bmatrix} & & & & & \\ & x^{(1)} & x^{(2)} & \dots & x^{(m)} & \\ & | & | & & | & \\ \xrightarrow{\text{W}^T x^{(1)} + b} & z_1 & & & & \\ \xrightarrow{\text{W}^T x^{(2)} + b} & z_2 & & & & \\ & \vdots & & & & \\ & & & & & \xrightarrow{\text{W}^T x^{(m)} + b} \\ & & & & & z_m \end{bmatrix}_{n \times m}$$

Capital Z consists of all z (lowercase) in horizontal stacking

Capital X consists of all x (lowercase) in horizontal stacking

* confusion :-

$$\textcircled{1} \quad Z = \underbrace{\text{np.dot}(W^T \xrightarrow{\text{Transpose}} X)}_{\text{vector}} + \underbrace{b}_{\text{real no.}}$$

b: real no. / $(1 \times 1 \text{ matrix})$
but when we add this vector to this real no.
python automatically takes this real no b & expands it out into $1 \times m$ row vector.

This is called BroadCasting in Python.

$$A = [a^{(1)} a^{(2)} \dots a^{(m)}] = \sigma(Z) \xrightarrow{\text{for sigmoid fun}} \text{use np.exp}$$

capital A consists of all a (lower case) in horizontal stacking.

↳ ⁶ This is for m training examples ↴

Vectorising Logistic Regression's Gradient

Output -

$$dz^{(1)} = a^{(1)} - y^{(1)}$$

$$dz^{(2)} = a^{(2)} - y^{(2)}$$

$$dZ = [dz^{(1)} dz^{(2)} \dots dz^{(m)}]_{1 \times m \text{ rowvector}}$$

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

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

$dw = 0 \quad db = 0$

to

$\left. \begin{array}{l} dw^+ = X^{(1)} dz^{(1)} \\ dw^+ = X^{(2)} dz^{(2)} \\ \vdots \\ dw^+ = X^{(m)} dz^{(m)} \end{array} \right\}$

$\left. \begin{array}{l} db^+ = dz^{(1)} \\ db^+ = dz^{(2)} \\ \vdots \\ db^+ = dz^{(m)} \end{array} \right\}$

eliminate for loop required for this purpose

$dw = dw/m \quad db = m$

$\{dw = m\}$

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$

four step propagation

$$db = \frac{1}{m} np.sum(dz)$$

$$dw = \frac{1}{m} X dz$$

(3)

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

(4)

$$w := w - \alpha dw$$

$$= \frac{1}{m} \begin{bmatrix} | & | & | \\ X^{(1)} & \dots & X^{(m)} \\ | & | & | \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

$$= \frac{1}{m} \begin{bmatrix} X^{(1)} dz^{(1)} + X^{(2)} dz^{(2)} + \dots + X^{(m)} dz^{(m)} \end{bmatrix}$$

$n \times 1$ vector

Getting rid of main for loop -

multiple

for iterations of gradient descent we need a for loop for that

$$Z = w^T X + b = np.dot(w.T, X) + b$$

$$A = \sigma(Z)$$

$$dz = A - Y$$

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

$$db = \frac{1}{m} np.sum(dz)$$

single iteration of gradient descent

descent

$w := w - \alpha dw$ } update
 $b := b - \alpha db$

Broadcasting in Python -

Calories from Carbs, Proteins, fats in 100g of different foods-

	Apples	Beef	Eggs	Potatoes	
Carbs	56.0	0.0	4.4	68.0	
Proteins	1.2	104.0	52.0	8.0	
Fats	1.8	185.0	99.0	0.9	

calories from

Calculate % of carbs, proteins, fats { without using for loop }

To create a matrix, we use numpy library

eg - $A = \text{np.array}([[2, 1, 3, 4], [5, 6, 7, 8]])$

O/P: $\begin{bmatrix} 2 & 1 & 3 & 4 \\ 5 & 6 & 7 & 8 \end{bmatrix}$

$\text{cal} = A \cdot \text{sum}(\text{axis} = 0)$ # To sum vertically

print(cal)

O/P : $\begin{bmatrix} 7 & 7 & 10 & 12 \end{bmatrix}$ ↓ axis 0 vertically → axis 1 {sum horizontally}

percentage = $100 * A / \text{cal} \cdot \text{reshape}(1, 4)$

print(percentage)

redundant

reshape: while writing

O/P $\begin{bmatrix} 28.55 & 14.28 & 30 & 33.33 \\ 71.42 & 85.55 & 70 & 66.66 \end{bmatrix}$ python codes if not sure

what are dimensions of matrix.

eg of Python broadcasting -

take a matrix A (2x4) & divide it by

(1x4) matrix

Operation of (2x4) matrix divide by (1x4) matrix

Broadcasting -

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

considered as

Broadcasting works with

both row vectors & column vectors

Broadcasting done

$$\textcircled{2} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix}_{1 \times 3} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}_{2 \times 3}$$

$$\textcircled{3} \quad \begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}_{2 \times 3} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}_{2 \times 1} \begin{bmatrix} 100 & 100 \\ 200 & 200 \end{bmatrix}_{2 \times 2} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}_{2 \times 3}$$

General Matrix :-

Principle Broadcasting

$$(m, n) \xrightarrow{\neq} (1, n) \rightsquigarrow (m, n)$$

→ matrix

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

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

$a = np.random.randn(5)$

$a.shape = (5,)$

"This is a rank1 array"

} Don't use

this in your code

It is neither a row vector nor a column vector

✓ $a = np.random.randn(5, 1) \rightarrow a.shape = (5, 1)$
 \rightarrow column

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

Logistic regression cost function -

$$\text{If } y=1 : p(y|x) = \hat{y} \quad ①$$

$$y=0 : p(y|x) = 1 - \hat{y} \quad ②$$

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

both ① & ②

Since logarithmic function is strictly monotonically increasing function; & we have to maximize $\log p(y|x)$

$$\begin{aligned}\log [p(y|x)] &= y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ &= -\{L(\hat{y}, y)\}\downarrow\end{aligned}$$

We minimize loss function

Statistics Principle :- Maximum Likelihood Estimation
product of m

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

To maximise this;

$$\log p(\text{labels . . .}) = \log \prod_{i=1}^m p(y^{(i)}|x^{(i)})$$

\log of products terms is \log individual terms
equal to sum of \log individual terms

$$\log p(\cdot \cdot \cdot) = \sum_{i=1}^m \underbrace{\log p(y^{(i)}|x^{(i)})}_{\text{individual terms}}$$

$$\underset{\text{(maximise)}}{=} \sum_{i=1}^m -L(\hat{y}^{(i)}, y^{(i)})$$

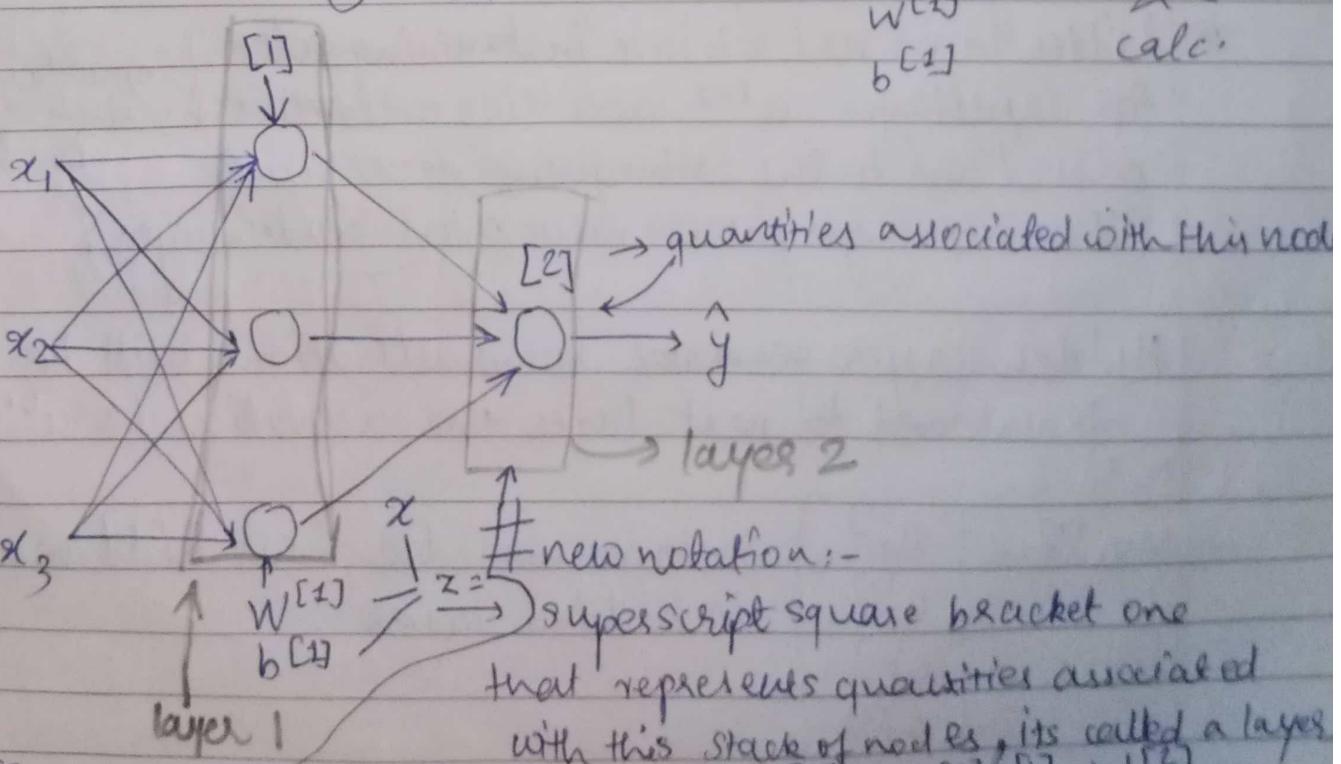
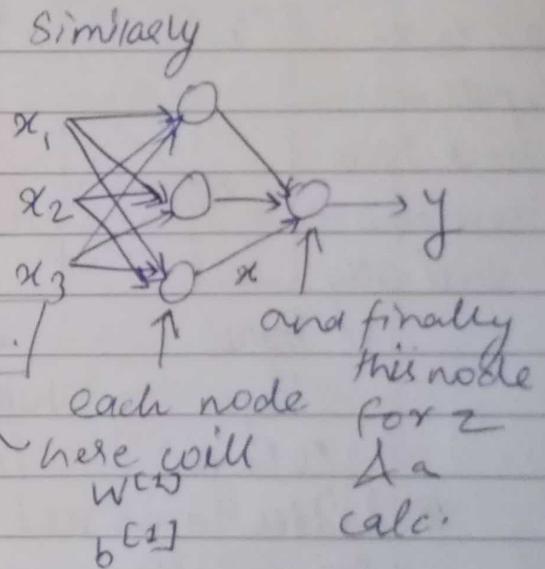
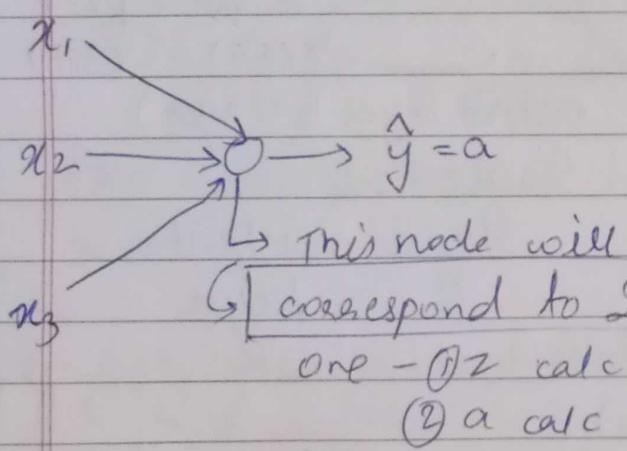
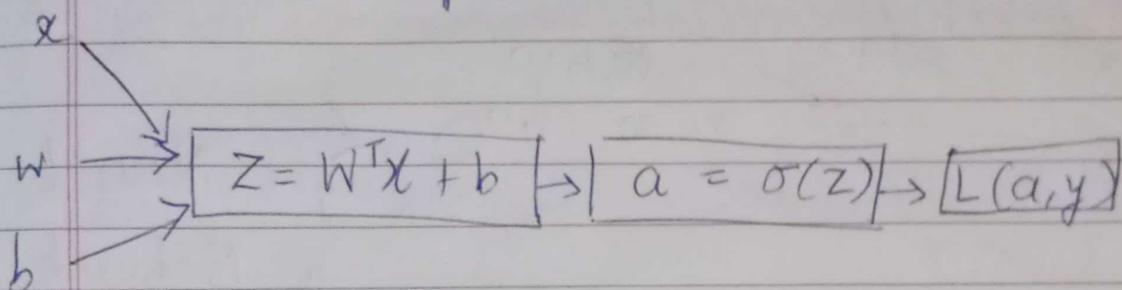
$$= - \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)})$$

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

(minimise)

why we minimise cost function?

What is a Neural Network?
and how to implement



$$\rightarrow a^{[2]} = \sigma(z^{[2]}) \longrightarrow L(a^{[2]}, y)$$

'z' calculation followed by 'a' calculation.

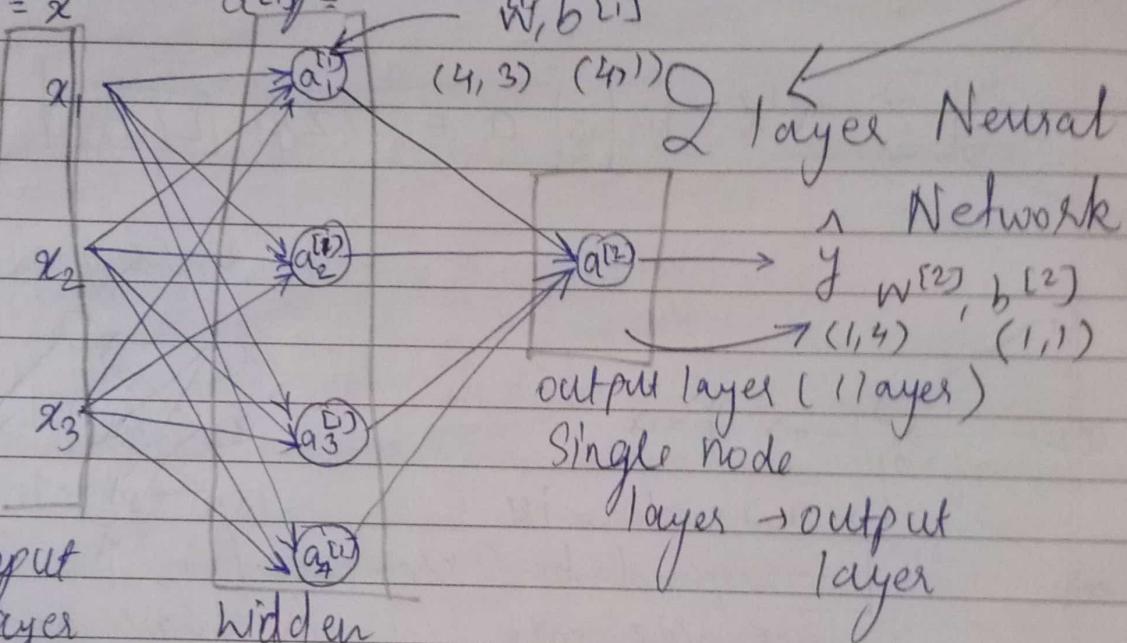
To compute derivatives - to get backward calculations in a similar manner

Neural Network Representation -

In neural network,

x_1, x_2, x_3
inputs are
stacked over
each other
vertically

input
(0 layers) layer (1 layer)



* hidden layer not visible in training sets.
for input layer $a^{[0]}$ acts like activations of input layer which gets passed on to the subsequent next layer
(Here 0 layer passes value x to hidden layer)

Hidden layer generates some activations will be transferred to next layer the activation is $a^{[1]}$

$$a^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix}$$

So these are the four hidden units in hidden layer.

Output layer will generate some activation $a^{[2]}$ which will be passed on \hat{y} so $\hat{y} = a^{[2]}$

So, this network is referred as to "2 layer neural network". We don't consider input layer only considering hidden & output layer. (1 hidden + 1 o/p)

We mention input layer as a zero & count the rest layers (hidden + output) as per notation convention.

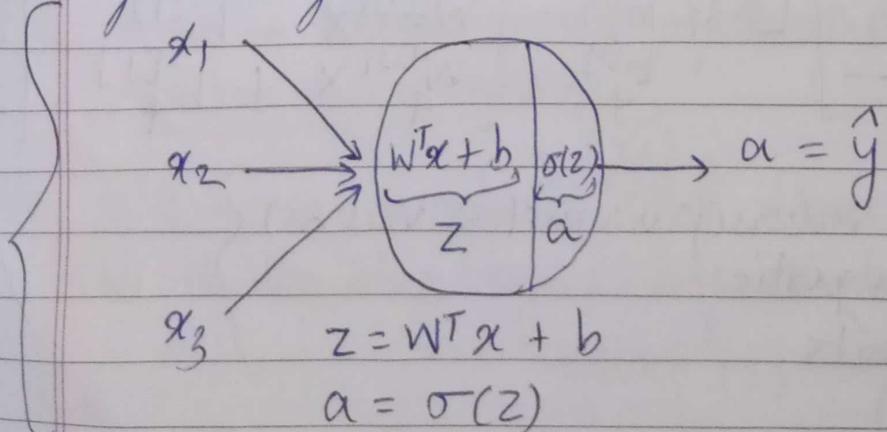
Input hidden layers & output layers will have parameters associated with it.

Hidden layer has parameters w & b associated with it.

Also some o/p layers has parameters w & b associated with it.

Computing a Neural Network's Output -

In Logistic Regression -



This 2 step computation of z & a is repeated many times in neural network.

Let's consider the first hidden node.

classmate
page

subscript 1 : indicates 1st node in that corresponding layer

superscript [1] : indicates 1st layer in neural network

$z_1^{[1]} = w_1^{[1]T} x + b_1^{[1]}$

$a_1^{[1]} = \sigma(z_1^{[1]})$

i.e. $a_i^{[l]}$ → layer $[l]$
 $i \rightarrow$ node in layer $[l]$

#

one logistic regression unit

Second node in hidden layer

Similarly, we write z & a calc. for all nodes

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

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

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

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

Vectorising this

to calculate z & a as activations

consider;

$Z^{[1]} = \begin{bmatrix} -w_1^{[1]T} - \\ -w_2^{[1]T} - \\ -w_3^{[1]T} - \\ -w_4^{[1]T} - \end{bmatrix}$

\downarrow

$\uparrow (4 \times 3)$

$(3,1) \rightarrow$ input features (n_x)

$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} + \begin{bmatrix} b_1^{[1]} \\ b_2^{[1]} \\ b_3^{[1]} \\ b_4^{[1]} \end{bmatrix}$

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

We will run all these values stacked vertically i.e. we have 4 logistic regression units

Stacking 4 w vectors (parameters vectors)

we stacked all z values for each node of that layer

$$\underline{\underline{a}}^{[1]} = \begin{bmatrix} a_1^{[1]} \\ a_2^{[1]} \\ a_3^{[1]} \\ a_4^{[1]} \end{bmatrix} = \sigma(\underline{\underline{z}}^{[1]})$$

for hidden layer;

$$x = a^{[e]}$$

$$\underline{\underline{z}}^{[1]} = W^{[1]} \underline{\underline{x}} + b^{[1]}$$

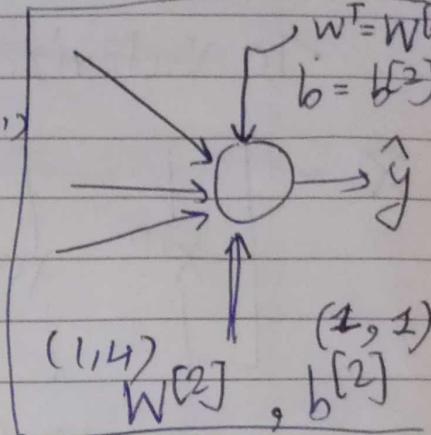
for output layer, we have parameters W & b as well

here;

$$\underline{\underline{z}}^{[2]} = W^{[2]} \underline{\underline{a}}^{[1]} + b^{[2]}$$

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

O/P layer



$$\underline{\underline{z}}^{[1]} = W^{[1]} \underline{\underline{a}}^{[0]} + b^{[0]} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{hidden layer}$$

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

for 2 layer neural network similarly

Vectorising across multiple training examples for neural networks -

$$x \longrightarrow a^{[2]} = \hat{y} \quad (\text{over 1 training eg})$$

for m training egs; repeat above step over m egs.

$$x^{(1)} \longrightarrow a^{[2](1)} = \hat{y}^{(1)} \quad (i) \rightarrow \text{refers to training eg } i$$

$$x^{(2)} \longrightarrow a^{2} = \hat{y}^{(2)}$$

$$\vdots$$

$$x^{(m)} \longrightarrow a^{[2](m)} = \hat{y}^{(m)} \quad \text{layers 2}$$

#

training examples

↓ different units (i.e. hidden units of layer)

for $i = 1$ to m

$$\cancel{X \in \mathbb{R}^{n \times m}} \Rightarrow \left\{ \begin{array}{l} z^{[1](i)} = W^{[1]} X^{(i)} + b^{[1]} \\ a^{[1](i)} = \sigma(z^{[1](i)}) \\ z^{[2](i)} = W^{[2]} a^{[1](i)} + b^{[2]} \\ a^{[2](i)} = \sigma(z^{[2](i)}) \end{array} \right\}$$

To vectorize it.

$$X = \begin{bmatrix} | & | & & | \\ X^{(1)} & X^{(2)} & \cdots & X^{(m)} \\ | & | & & | \end{bmatrix} \quad \text{for } m \text{ training egs}$$

explained (n_x, m)

or $W^{[1]} A^{[0]} + b^{[1]}$ rather (column vector stacking)

$$\begin{aligned} Z^{[1]} &= W^{[1]} X + b^{[1]} \\ A^{[1]} &= \sigma(Z^{[1]}) \\ Z^{[2]} &= W^{[2]} A^{[1]} + b^{[2]} \\ A^{[2]} &= \sigma(Z^{[2]}) \end{aligned}$$

So $X^{(i)}$ elements/examples are stacked vertically manner in a column manner one after another to get

Similarly doing this for $Z^{[1]}$, and $A^{[1]}$, $Z^{[2]}$ & $A^{[2]}$

$$Z^{[1]} = \begin{bmatrix} | & | & & | \\ Z^{1} & Z^{[1](2)} & \cdots & Z^{[1](m)} \\ | & | & & | \end{bmatrix}$$

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

Scanning from left to right refers to the scanning of all training examples

scanning from top to bottom consists of various nodes present in that layer (corresponding layer nodes) different units in layer

Why we use this activation func (hyperbolic tangent func)
 → It causes the mean of our data to come close to zero as it ranges between -1 & 1 (rather than 0.5)
 → Due to this centering of data takes place

Activation functions

Till now we were using sigmoid activation function, but there are still different choices available of activation function.

If a function is a non-linear function & not a sigmoid function.

Activation function working better than sigmoid func is Tangent function / Hyperbolic tangent function.

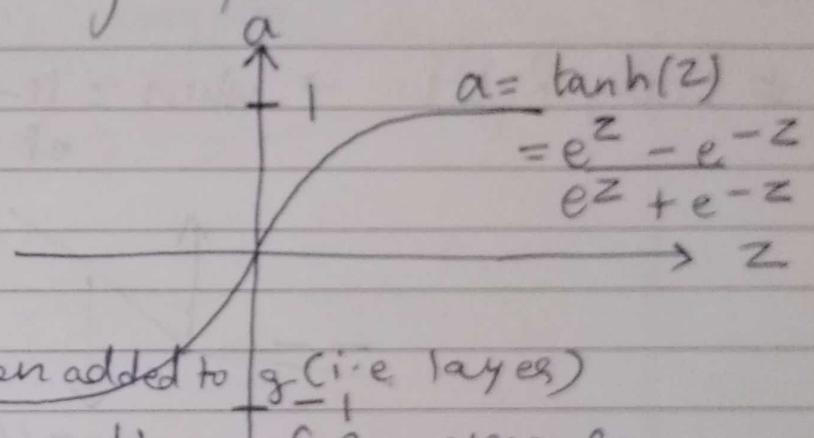
$$a = \tanh(z)$$

(goes between -1 & $+1$)

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

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

↑ superscript even added to g (i.e. layers)



But there is one exception for output layer we have to use the sigmoid activation function as it ranges b/w 0 & 1 the value of \hat{y} easy to use: $\hat{y} \in (0,1)$, $\therefore \hat{y} \in G(0,1)$ for binary classification

$$\therefore \text{Here } g(z^{[2]}) = \sigma(z^{[2]})$$

used as a output layer

If z is very large or is very small, the slope of function is close to zero (very small)

This

→ (slow downs gradient descent)

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

$$\begin{aligned} \text{slope} &= (e^z + e^{-z})(e^z + e^{-z}) - \\ &\quad (e^z - e^{-z})(e^z - e^{-z}) \\ &= (e^{2z} + 2 + e^{-2z}) - (e^{2z} + e^{-2z} - 2) \\ &= \frac{(e^z + e^{-z})^2 - 4}{(e^z + e^{-z})^2} \end{aligned}$$

Sol for this \rightarrow Rectified Linear Unit (ReLU function)

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

if z is +ve

slope(derivative) is one
i.e $da = 1$

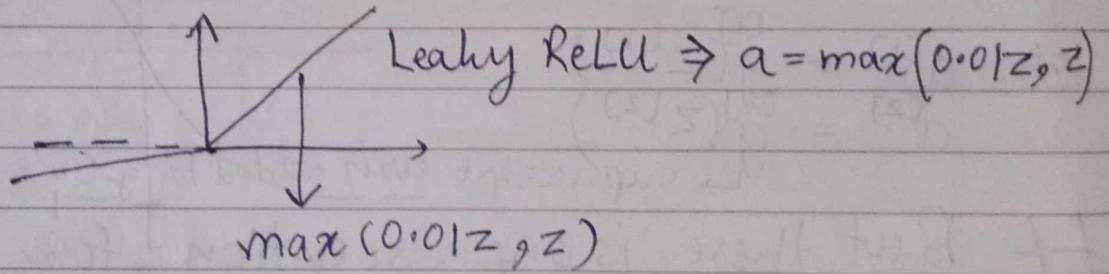
if z is -ve

slope(derivative) is zero
i.e $da = 0$

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

ReLU activation function is the default choice now-a-days.

Leaky ReLU \rightarrow when z is -ve it's a slight slope instead of zero.



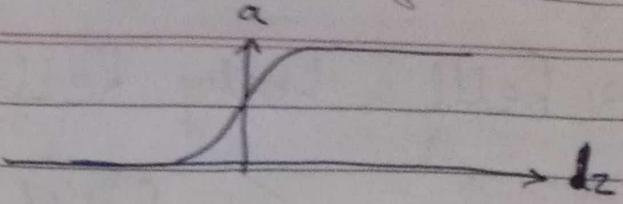
→ Choose Activation Function that
Bests Suit Your Applications

If we use linear activation function (i.e $a^{(l)} = z^{(l)}$)
the presence of hidden layer then it behaves like
logistic regression. It is basically calculating y_p to
give us opp using linear activation function and there
is no factor of non-linearity in that case

Derivatives of Activation Functions -

Sigmoid activation function -

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



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

$$\begin{aligned} &= \frac{-1}{(1+e^{-z})^2} (-e^{-z}) = \left[\frac{(1+e^{-z}) - 1}{(1+e^{-z})^2} \right] \\ &= \frac{1}{(1+e^{-z})} \left[1 - \frac{1}{(1+e^{-z})} \right] \end{aligned}$$

$$\begin{aligned} &= g(z) (1 - g(z)) \\ z = 10, \quad g(z) &\approx 1 \\ \frac{d}{dz} g(z) &\approx 1 (1 - 1) \approx 0 \end{aligned}$$

we know
 $a = g(z)$

$$z = -10 \quad g(z) \approx 0$$

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

$$z = 0 \quad g(z) \approx 1/2$$

$$\frac{d}{dz} g(z) = \frac{1}{2} \left(\frac{1-1}{2} \right) = \frac{1}{4}$$

$$g'(z) = g(z)[1 - g(z)] = a(1-a)$$

Tanh activation function -

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

$$\begin{aligned} \frac{dg(z)}{dz} = g'(z) &= \text{slope of } g(z) \text{ at } z \\ &= \frac{(e^z + e^{-z})^2 - (e^z - e^{-z})^2}{(e^z + e^{-z})^2} \\ &\quad \{ \text{previous page} \} \end{aligned}$$

$$z = 10 \quad \tanh(z) \approx 1$$

$$= 1 - [\tanh(z)]^2$$

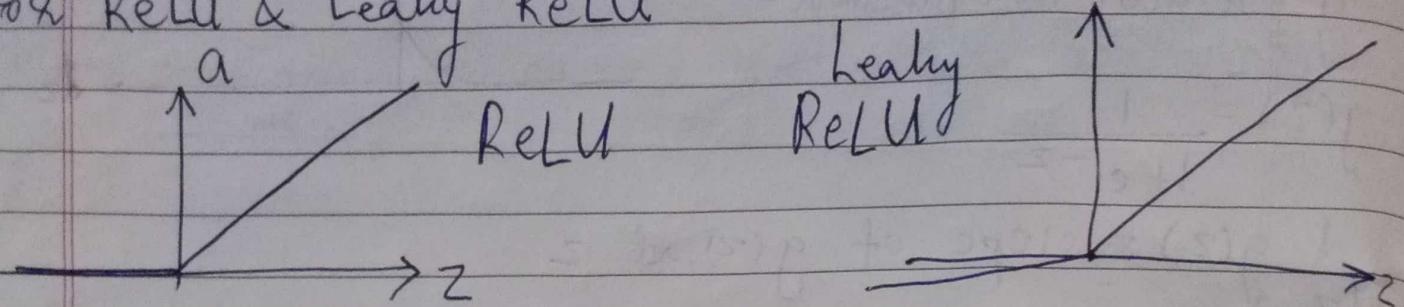
$$g'(z) \approx 0$$

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

$$z = -10 \quad \tanh(z) \approx -1$$

$$\begin{aligned} z = 0 \quad \tanh(z) &= 0 \\ g'(z) &\approx 1 \end{aligned}$$

for ReLU & Leaky ReLU



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

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

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

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

~~†~~ technically
this will be treated as
 $0.00000\ldots$

Gradient descents for neural networks -
(1 hidden layer)

Parameters: $W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]}$
 $(n^{[1]}, n^{[0]})^T, (n^{[1]}, 1), (n^{[2]}, n^{[1]}), (n^{[2]}, 1)$
 $n_x: n^{[0]}, n^{[1]}, (n^{[2]}=1)$
 input units, hidden units, output units

Cost function: $J(W^{[1]}, b^{[1]}, W^{[2]}, b^{[2]})$

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

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

Gradient Descent:
Repeat {

np.sum → Python numpy command for summing across
1 dimension of matrix

role of keepdims : It prevents Python from outputting 1 array
(rank one array)

Compute predictions ($\hat{y}^{(i)}$, $i=1, \dots, m$)
compute $dW^{[1]} = \frac{\partial J}{\partial W^{[1]}}$, $db^{[1]} = \frac{\partial J}{\partial b^{[1]}}$,
derivatives
Similarly $dW^{[2]}$ & $db^{[2]}$

update values
 $W^{[1]} := W^{[1]} - \alpha dW^{[1]}$
 $b^{[1]} := b^{[1]} - \alpha db^{[1]}$
 $W^{[2]} := W^{[2]} - \alpha dW^{[2]}$
 $b^{[2]} := b^{[2]} - \alpha db^{[2]}$

This is over one gradient descent iteration and continues until parameters are converging.

forward propagation :-

$$Z^{[1]} = W^{[1]} X + b^{[1]}$$

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

$$Z^{[2]} = W^{[2]} A^{[1]} + b^{[2]}$$

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

for o/p only

back propagation :-

$$dZ^{[2]} = A^{[2]} - Y$$

where :-

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

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1] T}$$

horizontally

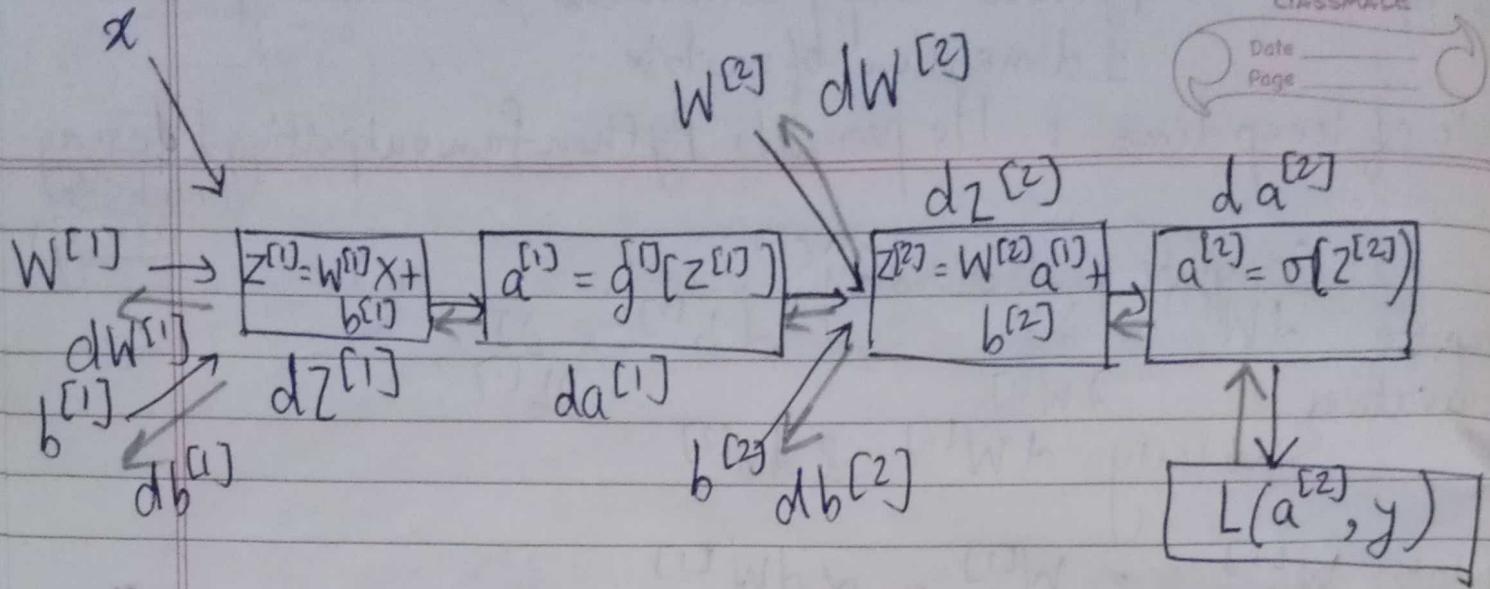
$$db^{[2]} = \frac{1}{m} np.sum(dZ^{[2]}, axis=1, keepdims=True)$$

~~neural network layers~~
neural network layers
(with hidden layers) $dZ^{[1]} = W^{[2] T} dZ^{[2]} *$
 $(n^{[1]}, m) \ g^{[1]}, (Z^{[1]})$
 $T \ n^{[2]}, m$

Here, element wise product of
2 matrices

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

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



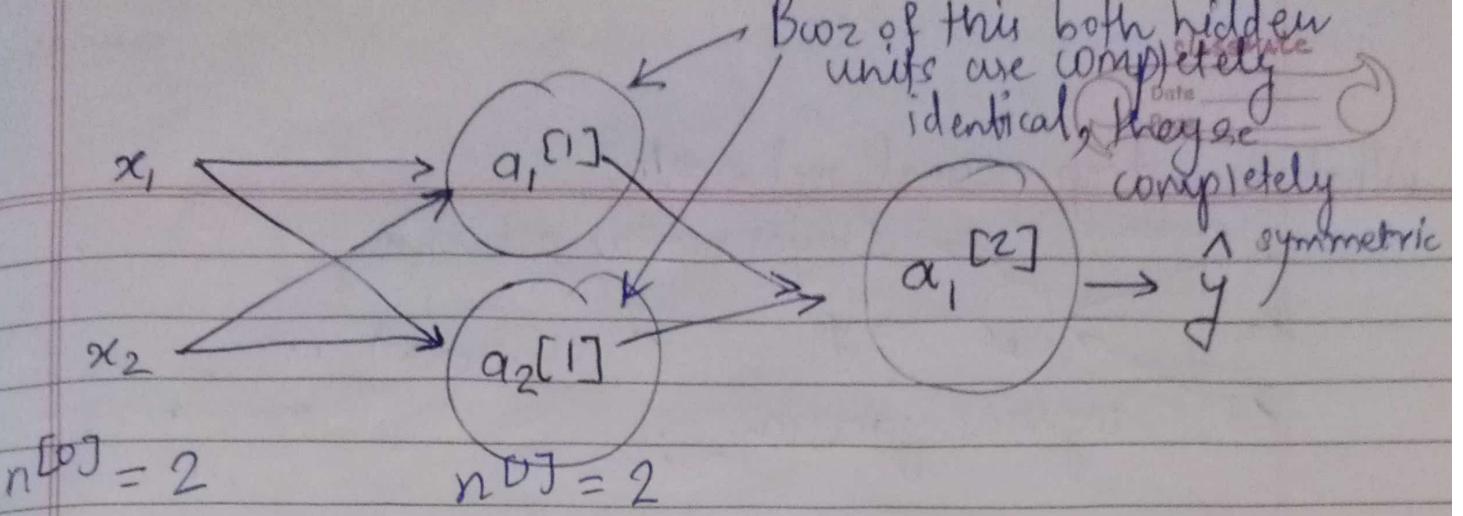
for W & dW , z & dz should have same dimensions
 i.e Variable & its derivative should have same dimension.

* : element-wise product

Random Initialization -

What happens if you initialize weights to zero?

for logistic regression, we initialized weights with zero which was fine but for a neural network initializing weights to parameters (W) to all zero and applying gradient descent won't work



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

↓
Problem

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

↳ this is fine

$$W^{[2]} = [0 \ 0]$$

$$a_1^{[1]} = a_2^{[1]} \rightarrow \text{It would be same}$$

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

how

So many time you run neural network, both hidden

units are still computing the same function.

So to this :- $W^{[1]} = np.random.randn(2, 2) * 0.01$

initializing to very small values

$$\text{np for } b^{[1]} = np.zeros((2, 1)) \quad ? \text{ why only } 0.01$$

(no problem)

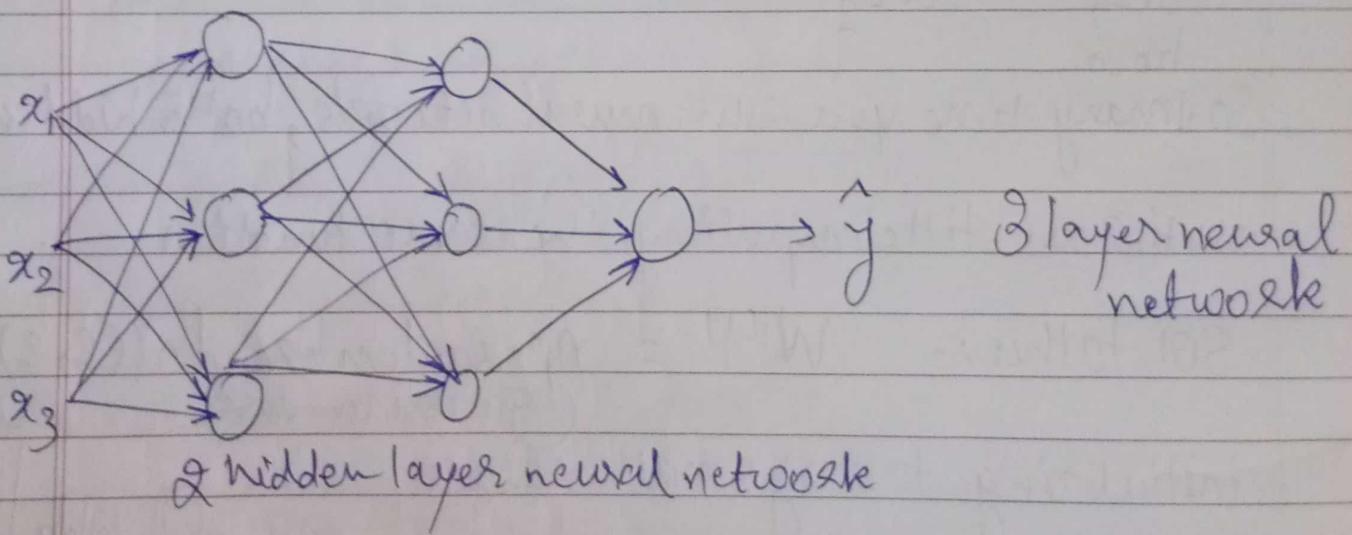
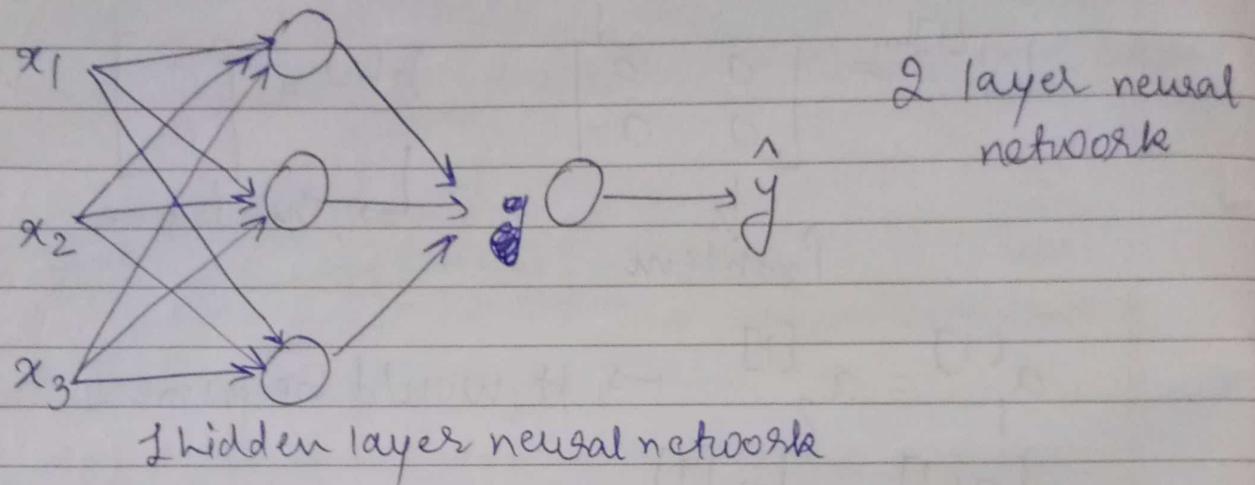
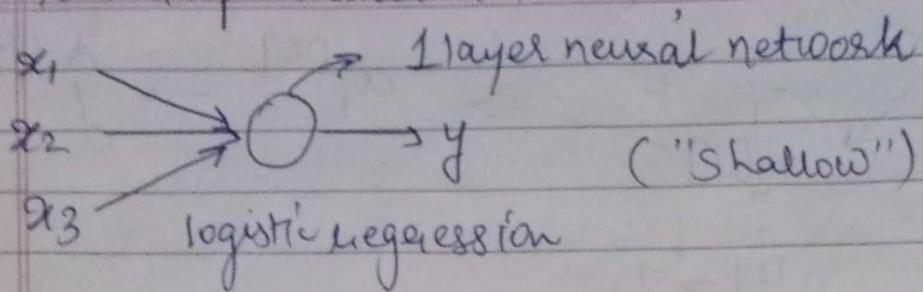
$$W^{[2]} = \dots$$

$$b^{[2]} = 0$$

if you take / initialize with larger values $Z^{[1]}$ will be very large, i.e. these very small

will be the parts of the activation function where slope / gradient is very small so gradient descent is very slow

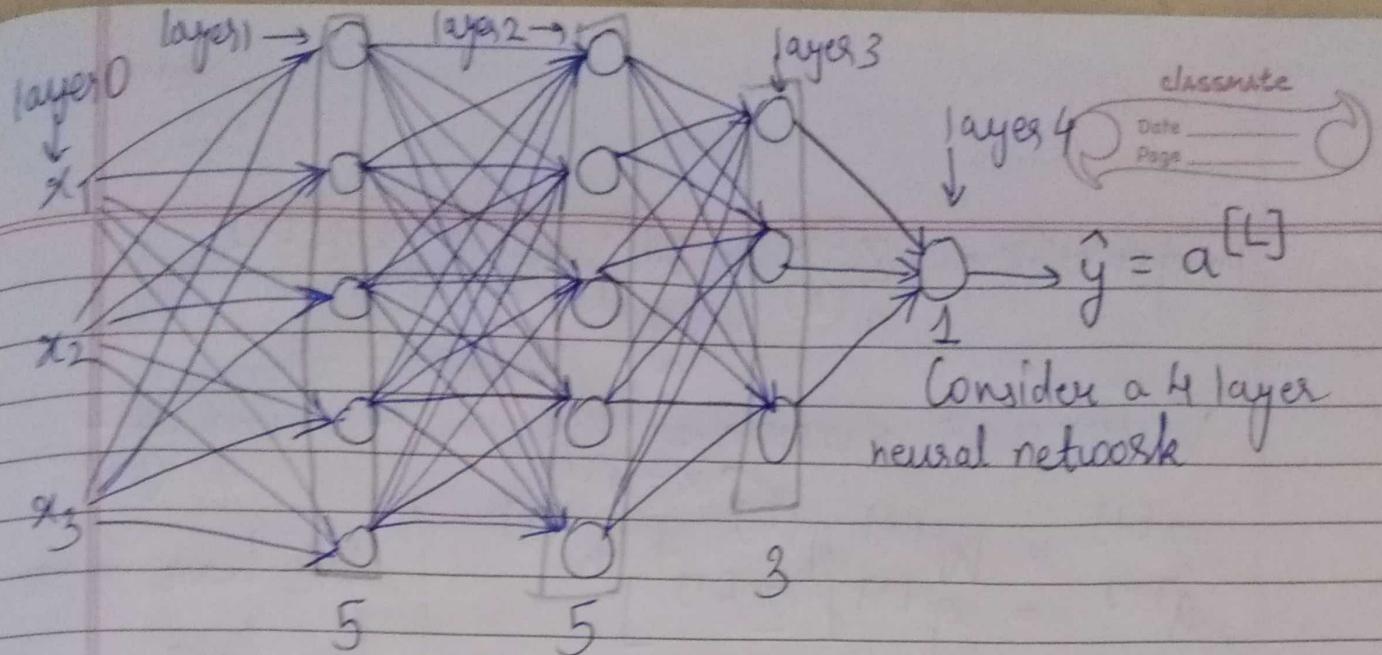
What is a deep neural network?



5 layer neural network - ("deeper")

The degree depends upon the depth of neural network

deep neural network notation -



Consider a 4 layer neural network

h : Number of layers in the network $h = 4 (\# \text{ layers})$

$n^{[l]}$: Number of nodes/units present in layer l

$a^{[l]}$: activations in layer l

$$n^{[0]} = 3 = n_x$$

$$n^{[1]} = 5$$

$$n^{[2]} = 5$$

$$n^{[3]} = 3$$

$$n^{[4]} = n^{[L]} = 1$$

input features : X

X is called activations of layer zero $a[0] = X$

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

$z^{[l]}$ \Rightarrow we require $W^{[l]}$ (weights)
& $b^{[l]}$ of parameters

$\hat{y} = a^{[L]}$: predicted output

for a single training example -

$$x: z^{[0]} = W^{[0]} X + b^{[0]} = W^{[0]} a^{[0]} + b^{[0]} = a^{[0]}$$

$$a^{[1]} = g^{[1]}(z^{[0]})$$

$$z^{[1]} = W^{[1]} a^{[0]} + b^{[1]}$$

$$a^{[2]} = g^{[2]}(z^{[1]})$$

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

$$a^{[3]} = g^{[3]}(z^{[2]})$$

$$z^{[3]} = W^{[3]} a^{[2]} + b^{[3]}$$

$$a^{[4]} = g^{[4]}(z^{[3]})$$

In general : $Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$
 $a^{[l]} = g^{[l]}(Z^{[l]})$

Vectorised version : For m training eg -

$$\begin{aligned} Z^{[1]} &= W^{[1]} A^{[0]} \\ A^{[1]} &= g^{[1]}[Z^{[1]}] \end{aligned}$$

$$\begin{aligned} Z^{[2]} &= W^{[2]} A^{[1]} + b^{[1]} \\ A^{[2]} &= g^{[2]}[Z^{[2]}] \end{aligned}$$

$$Y = g^{[4]}[Z^{[4]}] = A^{[4]}$$

This (for loop) needs to get implemented

We have to use
for loop
for $l = 1, \dots, L$
(Here $L = 4$)

This for loop can't
be avoided

dimensions of matrix $\Rightarrow W$

$$W^{[l]} \Rightarrow (n^{[l]}, n^{[l-1]})$$

dimensions of matrix $\Rightarrow b$

$$b^{[l]} \Rightarrow (n^{[l]}, 1)$$

dimension of dW & W should be same.
Similarly for db & b .

dimensions of X, Z, A depends on vectorization
(number of training examples m)

But $Z^{[l]}, a^{[l]} \Rightarrow (n^{[l]}, 1)$

$Z^{[l]}$ $\underset{l=0}{\bullet} A^{[l]} \Rightarrow (n^{[l]}, m)$
 $A^{[0]} = X = (n^{[0]}, m)$

for a layer in deep neural network:-

layer l : $W^{[l]}, b^{[l]}$

forward propagation : $a^{[l-1]}$, output $a^{[l]}$

$$Z^{[l]} = W^{[l]} a^{[l-1]} + b^{[l]}$$

$$a^{[l]} = g^{[l]}(Z^{[l]})$$

cache $Z^{[l]}$ is stored
for back propagation
(storing)

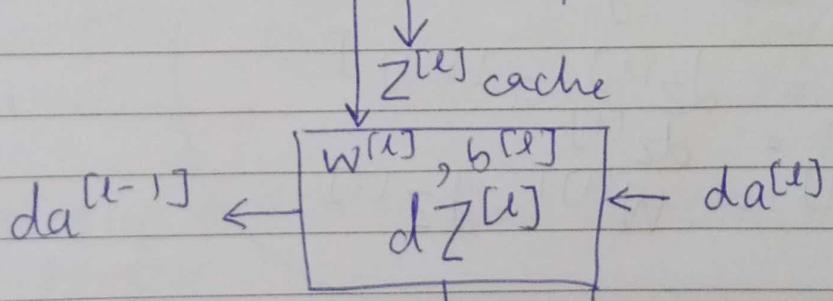
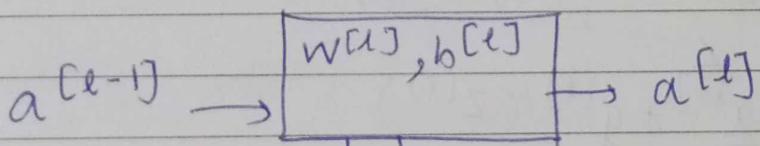
backward : Input $da^{[l]}$,
cache ($Z^{[l]}$)

Outputs $da^{[l-1]}$

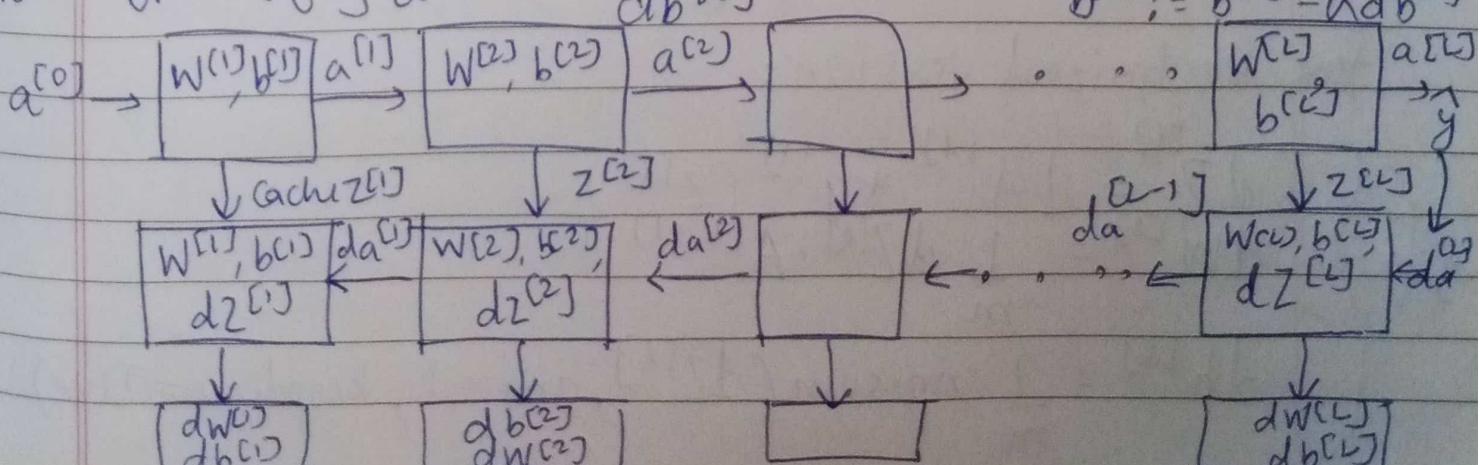
$$\frac{dW^{[l]}}{dZ^{[l]}}$$

$$\frac{db^{[l]}}{dZ^{[l]}}$$

layer l



$W^{[l]}$ is over one iteration one of gradient descent



forward propagation for layer l

input $a^{[l-1]}$
 Output $a^{[l]}$, cache ($Z^{[l]}$) $\rightarrow W^{[l]}, b^{[l]}$

$$\begin{aligned} Z^{[l]} &= W^{[l]} \cdot a^{[l-1]} + b^{[l]} \\ a^{[l]} &= g^{[l]}(Z^{[l]}) \end{aligned}$$

Vectorized: -

~~$$\begin{aligned} Z^{[l]} &= W^{[l]} \cdot A^{[l-1]} + b^{[l]} \\ A^{[l]} &= g^{[l]}(Z^{[l]}) \end{aligned}$$~~

In general

backward propagation for layer l

input data - d_{out}

output - $-d_{\text{out}}^{[l-1]}, dW^{[l]}, db^{[l]}$

$$\left\{ \begin{array}{l} dZ^{[l]} = \underline{da^{[l]}} * g^{[l]}'(Z^{[l]}) \\ dW^{[l]} = \frac{\partial Z^{[l]}}{\partial Z^{[l]}} \cdot a^{[l-1]T} \\ db^{[l]} = \frac{\partial Z^{[l]}}{\partial b^{[l]}} \\ da^{[l-1]} = W^{[l]T} \cdot dZ^{[l]} \\ dZ^{[l]} = W^{[l-1]T} \cdot dZ^{[l-1]} * g^{[l-1]}'(Z^{[l-1]}) \end{array} \right.$$

for vectorized version:

$$\left\{ \begin{array}{l} dZ^{[l]} = dA^{[l]} * g^{[l]}'(Z^{[l]}) \\ dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T} \\ db^{[l]} = \frac{1}{m} np.sum(dZ^{[l]}), axis=1, keepdims=True \end{array} \right.$$

$$dA^{(l-1)} = W^{(l)T} \cdot dZ^{(l)}$$

$$da^{(l)} = -\frac{y}{a} + \frac{(1-y)}{(1-a)}$$

$$dA^{(l)} = \left\{ \begin{array}{l} -\frac{y^{(1)}}{a^{(1)}} + \frac{(1-y^{(1)})}{(1-a^{(1)})}, \dots, -\frac{y^{(m)}}{a^{(m)}} + \frac{(1-y^{(m)})}{(1-a^{(m)})} \end{array} \right.$$

Parameters V/s Hyperparameters

Parameters: $W^{(l)}, b^{(l)}, W^{(2)}, b^{(2)}, \dots$

Hyperparameters: learning rate α
 no of iterations
 no of hidden layers
 no of hidden units
 choice of activation function

} These the parameters
 that control the
 parameters W & b
 ultimately. Hence
 they are called
 hyperparameters
 whose values need
 to be set

Hyperparameters influence the
 parameters

a lot of setting up needs to be done to better implement that
 neural network

Its like -

* idea \rightarrow code \rightarrow experiment *

Best to be identified directly \rightarrow research is ongoing
 hyperparameters