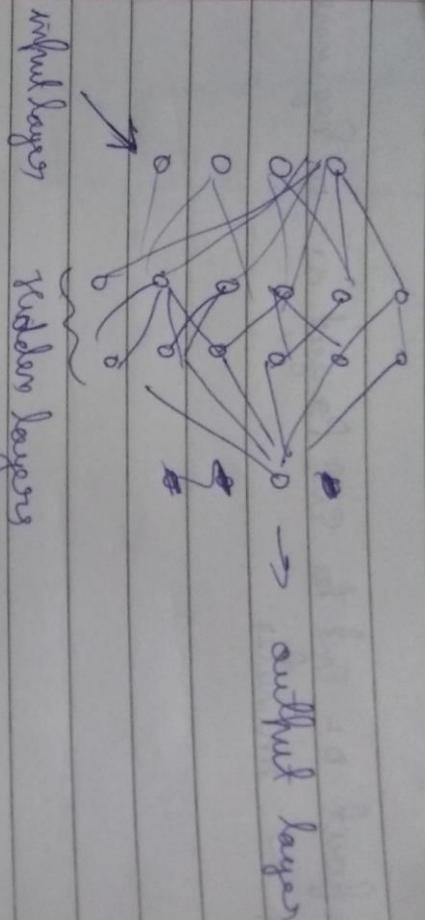


CNN
└ Convolution → summing
→ Pooling

Neural Network



Input layer - accepts input in different forms .

Hidden layer → performs calculation on these inputs

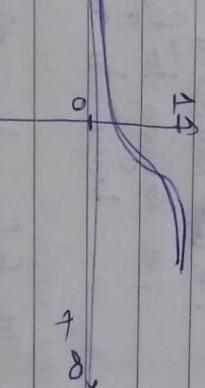
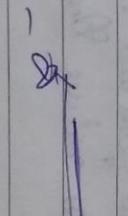
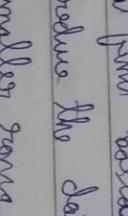
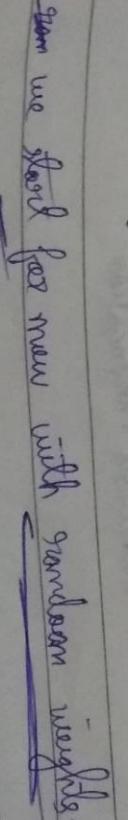
Output layer → then delivering the outcome of the calculations & extraction

- Each of these layers ~~are~~ contains neurons that are connected with neurons of previous layers .

- Each neuron has its own weight .

In \rightarrow hidden layer \rightarrow output function

Neural Network

- Artificial neuron \rightarrow "Perceptron"
- Hidden layer is a layer responsible for learning.
- Input layer \rightarrow receiver layer
- Activation functions
 - more than 50-60 AF.
 - Sigmoid ($0-1$)

$$\text{Sigmoid} = \frac{1}{1 + e^{-x}}$$
 - ReLU \rightarrow 
 - Ridge \rightarrow data Normalization.
 $\text{Normalization} = \frac{x - \text{mean}}{\text{std}}$
- What you basically helps changing the scale to reduce the data into a smaller range.

- Normalizing the data \rightarrow change the scale. By having the data to a similar "func", doesn't matter which type of activation function we use.
- weights \rightarrow defines a new "set" input taken & input that goes into the next layer.

- bias we start from new with random weights.

Backpropagation

PAGE NO.	
DATE	

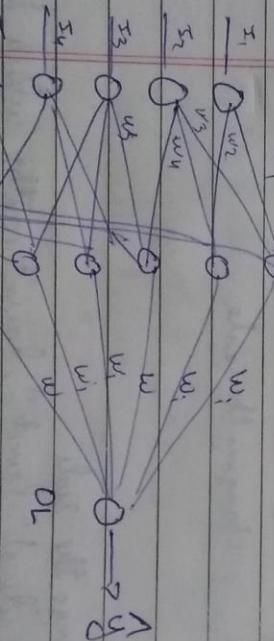
- less \rightarrow difference b/w y & \hat{y} ($y - \hat{y}$) \rightarrow error
- there are different types of loss functions

- we pass the loss/ cost to the optimizer

- in forward pass we don't learn anything
- optimizes tries to perform backward propagation
 - i.e. learning
 - i.e. trying to find out new values of weights to reduce the loss

We keep on doing this until we find minimum loss

\rightarrow forward propagation



all PNC \leftarrow backward propagation

w_i \rightarrow weights

$Z_i = \sum_j (I_j w_{ij} + b)$ this goes into hidden layer
In hidden layer we have activation layer which will pass

Learning \rightarrow finding pattern

PAGE NO.

DATE

in Normalizing the data
eg: sigmoid function

$$\frac{1}{1+e^{-x}}$$

will result in values

$$w_0^1 \ 0 \ 2 \ 1$$

Initially random weights are assigned

From hidden layers it goes to output layer where we use Probabilistic Function like sigmoid to give the Probability

$$y = (y - \hat{y}) \text{ we send to loss function}$$

We have optimizers which are used to reduce the loss function by back propagation and changing of weights

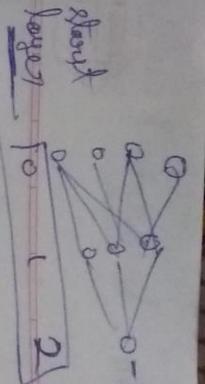
This process is done until we get the lower loss function

Biases can be introduced anywhere (in hidden layers or output layer)
it periodically helps in better learning.

new wt \rightarrow older wt \rightarrow etc \rightarrow decreasing rate \rightarrow lowest

$$w_n = w_{n-1} - \eta \Delta E_w$$

$$w_n = w_{n-1} - \eta \left(\frac{\partial E}{\partial w} \right)$$



$$w_n = w_{n-1} - \eta \left(\frac{dE}{dw} \right)$$

$$\boxed{\frac{dc}{dw} = \frac{\partial c}{\partial z} \cdot \frac{\partial z}{\partial w}} \quad \begin{matrix} \rightarrow \text{cost function} \\ w \rightarrow \text{weights} \end{matrix}$$

$$z = \sum_{i=1}^N (x_i w_i + b)$$

↑ the input which going inside hidden layer / output layer

forward

- $L \rightarrow$ No of layers in Network
- N^l - Dimensionality of the layer $\leftarrow \{0, \dots, L\}$
- $\left\{ \begin{array}{l} N^0 \rightarrow \text{dimensionality of the input} \\ N^L \rightarrow \text{dimensionality of the output} \end{array} \right.$

→ hidden layer

- $w^m \in R^{N^m \times N^{m-1}}$ - weight matrix from layer $m \in \{1, \dots, L\}$.
- w_j^m is the wt b/w the j^{th} unit in layer m & j^{th} unit in layer $m-1$.

- $b^m \in R^{N^m}$ - Bias vector for layer m .
- σ^m - Non linear activation function of the units in layer m , applied elementwise.
- $\underline{x}^m \in R^{N^m}$ - Linear sum of the inputs to layer m , computed by $\underline{x}^m = w^m a^{m-1} + b^m$.

$\sigma^m(h^m)$

Running same value in the output function

Page No.	
Date	

- $a^m \in R^{n^m}$ \rightarrow activation of units in layer m - computed by
 $a^m = \sigma^m(h^m) = \sigma^m(w^m a^{m-1} + b^m)$. a^l is the output of the network \leftarrow Network. We define a final use a^o as the output of the network

- y_{ER^N} - Target output of the network

- $c \rightarrow$ cost "error function" of the network which is a formula of a^l (~~the network output~~) & y (target as ~~answ~~)

$$\text{from this } Z = (x_1 w_{00} + x_2 w_{10} + b_0)$$

summed from others

x_1 \circ w_{00} \rightarrow output of the activation

a^m \rightarrow output of the activation

x_2 \circ w_{10} \rightarrow output of the activation

a^l \rightarrow output of the final network

$c = (y - \hat{y})$

x_3 \circ w_{20}

w_{21}

w_{22}

w_{23}

w_{24}

w_{25}

w_{26}

w_{27}

w_{28}

w_{29}

w_{210}

w_{211}

w_{212}

w_{213}

w_{214}

w_{215}

w_{216}

w_{217}

w_{218}

w_{219}

w_{220}

w_{221}

w_{222}

w_{223}

w_{224}

w_{225}

w_{226}

w_{227}

w_{228}

w_{229}

w_{230}

w_{231}

w_{232}

w_{233}

w_{234}

w_{235}

w_{236}

w_{237}

w_{238}

w_{239}

w_{240}

w_{241}

w_{242}

w_{243}

w_{244}

w_{245}

w_{246}

w_{247}

w_{248}

w_{249}

w_{250}

w_{251}

w_{252}

w_{253}

w_{254}

w_{255}

w_{256}

w_{257}

w_{258}

w_{259}

w_{260}

w_{261}

w_{262}

w_{263}

w_{264}

w_{265}

w_{266}

w_{267}

w_{268}

w_{269}

w_{270}

w_{271}

w_{272}

w_{273}

w_{274}

w_{275}

w_{276}

w_{277}

w_{278}

w_{279}

w_{280}

w_{281}

w_{282}

w_{283}

w_{284}

w_{285}

w_{286}

w_{287}

w_{288}

w_{289}

w_{290}

w_{291}

w_{292}

w_{293}

w_{294}

w_{295}

w_{296}

w_{297}

w_{298}

w_{299}

w_{300}

w_{301}

w_{302}

w_{303}

w_{304}

w_{305}

w_{306}

w_{307}

w_{308}

w_{309}

w_{310}

w_{311}

w_{312}

w_{313}

w_{314}

w_{315}

w_{316}

w_{317}

w_{318}

w_{319}

w_{320}

w_{321}

w_{322}

w_{323}

w_{324}

w_{325}

w_{326}

w_{327}

w_{328}

w_{329}

w_{330}

w_{331}

w_{332}

w_{333}

w_{334}

w_{335}

w_{336}

w_{337}

w_{338}

w_{339}

w_{340}

w_{341}

w_{342}

w_{343}

w_{344}

w_{345}

w_{346}

w_{347}

w_{348}

w_{349}

w_{350}

w_{351}

w_{352}

w_{353}

w_{354}

w_{355}

w_{356}

w_{357}

w_{358}

w_{359}

w_{360}

w_{361}

w_{362}

w_{363}

w_{364}

w_{365}

w_{366}

w_{367}

w_{368}

w_{369}

w_{370}

w_{371}

w_{372}

w_{373}

w_{374}

w_{375}

w_{376}

w_{377}

w_{378}

w_{379}

w_{380}

w_{381}

w_{382}

w_{383}

w_{384}

w_{385}

w_{386}

w_{387}

w_{388}

w_{389}

w_{390}

w_{391}

w_{392}

w_{393}

w_{394}

w_{395}

w_{396}

w_{397}

w_{398}

w_{399}

w_{400}

w_{401}

w_{402}

w_{403}

w_{404}

w_{405}

w_{406}

w_{407}

w_{408}

w_{409}

w_{410}

w_{411}

w_{412}

w_{413}

w_{414}

w_{415}

w_{416}

w_{417}

w_{418}

w_{419}

w_{420}

w_{421}

w_{422}

w_{423}

w_{424}

w_{425}

w_{426}

w_{427}

w_{428}

w_{429}

w_{430}

w_{431}

w_{432}

w_{433}

w_{434}

w_{435}

w_{436}

w_{437}

w_{438}

w_{439}

w_{440}

w_{441}

w_{442}

w_{443}

w_{444}

w_{445}

w_{446}

w_{447}

w_{448}

w_{449}

Now do we change wts?

$$w_m = w_{m-1} - \eta \left(\frac{dc}{dw} \right) \quad \text{row wise}$$

where w_k is old wt

$\frac{\partial C}{\partial a^L} \rightarrow$ rate of change of cost function depends on the output of the Network (a^L)

$\frac{\partial C}{\partial z^m} \rightarrow$ rate of change of O/P of the network depends on the input (z^m), found (a^m) depends on the input (z^m), which it takes in that particular layer.

Our requirement \rightarrow To find out "rate of change of error w.r.t weights

Now we can write the chain rule as:-

$$\frac{\partial C}{\partial w_i^m} = \sum_{k=1}^{N^m} \frac{\partial C}{\partial z_k^m} \cdot \frac{\partial z_k^m}{\partial w_i^m} \quad (1)$$

$m \Rightarrow$ layers

$a^L \rightarrow$ the output of any layer

$z^m \rightarrow$ net back error of m layer

102a

$w_1, w_2, \dots, w_n \rightarrow$ input of this layer

$$= (w_1 n_1 + w_2 n_2 + b)$$

layer $\rightarrow m-1 \rightarrow m$

$[Q]$

$[K]$

$w_{ki}^m \rightarrow$ wt back error of m layer i.e. w_i

$a_i^{m-1} \rightarrow$ output of all of $m-1$ layers n_i

$$\frac{\partial C}{\partial w_{ij}^m} = \sum_{k=1}^{N^m} \frac{\partial C}{\partial z_k^m} \left(\frac{\partial z_k^m}{\partial w_{ij}^m} \right)$$

Replacing $\frac{\partial z_k^m}{\partial w_{ij}^m}$ in eq (1) Finding derivative of z_k^m w.r.t w

$$\frac{\partial z_k^m}{\partial w_{ij}^m} = \frac{\partial}{\partial w_{ij}^m} \left(\sum_{x=1}^{N^m} w_{ix}^m a_x^{m-1} + b_i^m \right)$$

derivative of
constant = 0

$$= a_j^{m-1}$$

$$a_1 \quad O_1 \quad w_{10}$$

$$a_2 \quad O_2 \quad w_{11}$$

$$a_3 \quad O_3 \quad w_{12}$$

$$a_k \quad O_k \quad w_{1k}$$

$$j \quad O_j \quad w_{ij}$$

$$i \quad O_i$$

$$R$$

$$\frac{dy}{dx} = x_1 w_{10}$$

$$\frac{dy}{dx} = w_{10}$$

$$a_1 \quad O_1 \quad w_{10}$$

$$a_2 \quad O_2 \quad w_{11}$$

$$a_3 \quad O_3 \quad w_{12}$$

$$a_k \quad O_k \quad w_{1k}$$

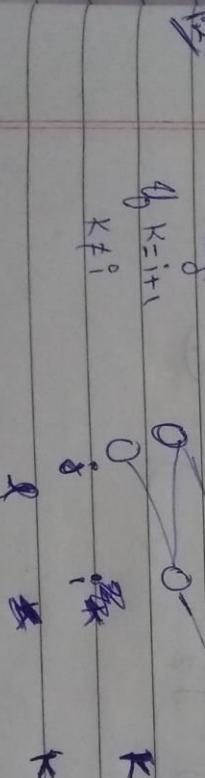
$$j \quad O_j \quad w_{ij}$$

$$i \quad O_i$$

$$R$$

If $k \neq i$ then there is no path "back" layer ① & ②
so we can't find the derivative $\therefore 0$.

- eq:
 $y = x_1 w_{10}$
 $\frac{dy}{dx} = w_{10}$



If $k = i$, the input of layer i (i.e. j) will be
 $x_1 w_{10}$ for a specific neuron
 and if we derivatives it w.r.t w we get x_1 , i.e. input
 of previous layer
 i.e. a_{i-1}^{m-1}

$$\frac{\partial z_k^m}{\partial w_{ij}^m} = \begin{cases} 0 & k \neq i \\ a_j^{m-1} & k = i \text{ (some layer)} \end{cases}$$

Now, the full-fledged $\frac{\partial c}{\partial a_k^m}$ is 0 unless $k = i$ causes the summation to collapse, giving

$$\frac{\partial c}{\partial w_{ij}^m} = \frac{\partial c}{\partial z_i^m} a_j^{m-1} \quad \text{--- (2)}$$

$a_j^{m-1} \rightarrow 1/p$ from the previous layer.

in vector form

$$\frac{\partial c}{\partial w^m} = \frac{\partial c}{\partial z^m} a^{m-1}$$

$$\boxed{\text{in simple form } \frac{dc}{dw} = \frac{dc}{dz} \cdot a^{m-1}} \quad \text{--- (3)}$$

$$b_k = b_{k-1} - \eta \left(\frac{\partial C}{\partial b} \right)$$

* Similarly for the bias variable b^m , we have :

$$\frac{\partial C}{\partial b_i^m} = \sum_{k=1}^{N^m} \frac{\partial C}{\partial z_k^m} \cdot \frac{\partial z_k^m}{\partial b_i^m}$$

finding $\frac{\partial z_k^m}{\partial b_i^m}$

$$\frac{\partial z_k^m}{\partial b_i^m} = \frac{\partial}{\partial b_i^m} \left(\sum_{l=1}^{N^m} w_{il}^m a_l^{m-1} + b_i^m \right)$$

$$= 1$$

$$\frac{\partial z_k^m}{\partial b_i^m} = \begin{cases} 0 & k \neq i \\ 1 & k = i \text{ (same layer)} \end{cases}$$

The summed " collapses to give

$$\frac{\partial C}{\partial b_i^m} = \frac{\partial C}{\partial z_i^m} \quad \text{--- (3)}$$

in vector form.

$$\frac{\partial C}{\partial b^m} = \frac{\partial C}{\partial z^m} \quad \text{in simple form}$$

$$\frac{dz}{db} = 0 + 1$$

$$\boxed{\frac{dz}{db} = 1}$$

If we are able to find $\frac{\partial C}{\partial z}$ which is same in Eqⁿ

② & Eqⁿ ③ we can find out $\frac{\partial C}{\partial w} \frac{\partial C}{\partial b}$

or values of b_n & w_n

$$b_n = b_{n-1} - \eta \left(\frac{\partial C}{\partial b} \right)$$

$$w_n = w_{n-1} - \eta \left(\frac{\partial C}{\partial w} \right)$$

* Now, we must compute $\frac{\partial c}{\partial z^m}$. For the final layer ($m=L$), this term is straightforward to compute using the chain Rule

$$\boxed{\frac{\partial c}{\partial z^L} = \frac{\partial c}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L}} - (F_1)$$

* $\frac{\partial c}{\partial a^k}$ rate of change of C.F w.r.t

final outcome of the network a^L

e.g. if ~~a^L~~ y is close to y then C.F loss F^L would be less

in vector form

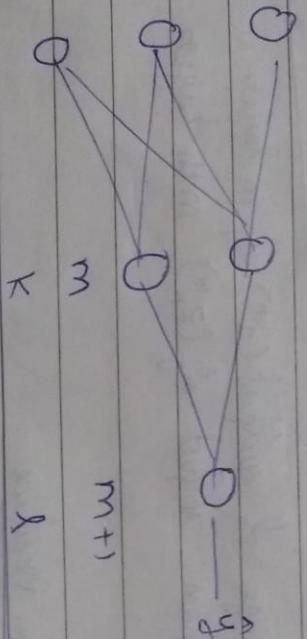
$$\boxed{\frac{\partial c}{\partial z^L} = \frac{\partial c}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L}}$$

for m th layer we can write the eq.

$$\frac{\partial c}{\partial z^m} = \frac{\partial c}{\partial a^m} \cdot \frac{\partial a^m}{\partial z^m}$$

$$= \left(\sum_{k=1}^{N^{m+1}} \frac{\partial c}{\partial z_k^{m+1}} \frac{\partial a_k^m}{\partial a_k^m} \right) \cdot \frac{\partial a_k^m}{\partial z^m}$$

applying chain rule



$\frac{\partial c}{\partial z^{m+1}}$ → rate of change of c w.r.t. deflections when input (z) in $(m+1)$ th layer.

$\frac{\partial c}{\partial a_k^m}$ → rate of change of c w.r.t. deflections when input (a_k^m) of m th layer.

$$= \left(\sum_{k=1}^{N^{m+1}} \frac{\partial c}{\partial z_k^{m+1}} \frac{\partial a_k^m}{\partial a_k^m} \left(\sum_{k=1}^{N^m} w_{kk}^{m+1} a_k^m + b_k^{m+1} \right) \right) \frac{\partial a_k^m}{\partial z^m}$$

$$\boxed{\frac{\partial c}{\partial z^m} = \left(\sum_{k=1}^{N^{m+1}} w_{kk}^{m+1} \cdot \frac{\partial c}{\partial z_k^{m+1}} \right) \frac{\partial a_k^m}{\partial z^m} - 0} \quad (4)$$

Simple form

$$\boxed{\frac{dc}{dz} = \left(w \cdot \frac{dc}{da} \right) \left(\frac{da}{dz} \right)} \quad - (5)$$

1.46

$$S_1 = \frac{\partial C}{\partial w^m} = \left(\frac{\partial C}{\partial z^m} \right) \left(a^{m-1} \right)$$

$$S_2 = \frac{\partial C}{\partial b^m} = \frac{\partial C}{\partial z^m} \cdot 1$$

$$S_3 = \frac{\partial C}{\partial z^m} = \left(w^m \frac{\partial C}{\partial z^{m+1}} \right) \left(\frac{\partial a}{\partial z^m} \right) \quad - \text{final Eq.}$$

$\frac{\partial a}{\partial z^m} \rightarrow$ rate of change of output (a^m) of m layer w.r.t
 z_k^m depends upon input of (z_k^m) in k th layer.

Backpropagation in Practice

Nonlinearity $a^m = \sigma^m(z^m)$ Note
 sigmoid $\frac{e^{z^m} - e^{-z^m}}{e^{z^m} + e^{-z^m}}$ $1 - (\sigma^m(z^m))^2 = 1 - (a^m)^2$ Equivalent of \tan^{-1}
 - squashes any
 input to the range [0,1]

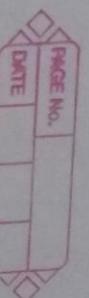
$$\text{Tanh} \quad \frac{e^{z^m} - e^{-z^m}}{e^{z^m} + e^{-z^m}}$$

ReLU $\max(0, z^m)$ $0, z^m < 0; 1, z^m \geq 0$ Commonly used in NN with
 many layers.

Now, if we use sigmoid function as our activation layer, the
 Eq. will become

$$\frac{\partial C}{\partial z^m} = \left(w^{m+1} \frac{\partial C}{\partial z^{m+1}} \right) \cdot \left(a^m(1-a^m) \right)$$

$$RMSE = \frac{1}{2} (y - \hat{y})^2$$



with different sigmoid function and Δa value changes.

^{imp} Do we use "softmax" function in Backward Propagation?
→ yes, we do

Cost function C.F $\frac{\partial C}{\partial a^l}$ Notes

Squared Error $\frac{1}{2} (y - a^l)^2$ commonly used when the value of θ is not constrained to a specific graph.

Entropy $(y - 1) \log(1 - a^l)$ $\frac{a^l - y}{a^l(1 - a^l)}$ commonly used for binary classification tasks, can yield faster convergence.

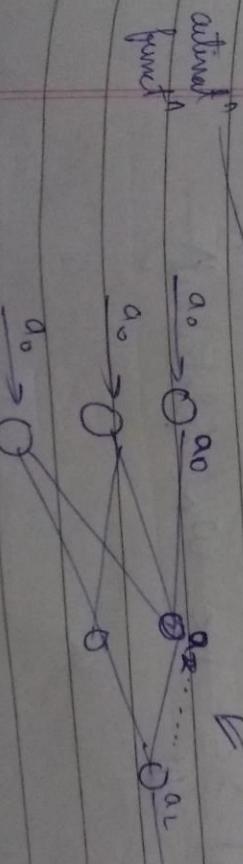
* In practice, backpropagation proceeds in the following manner for each training sample:

1. Forward Pass : Given the network input a^0 , compute

a^m , recursively by

$$a^l = \sigma'(w^l a^{l-1} + b^l) \dots a^0 = \sigma'(w^0 a^{-1} + b^0)$$

→ many layers in between



2. Backward Pass : Compute $\frac{\partial C}{\partial z^L}$ ^{cost func'}

$$\left[\text{Imp Eq} \right] \rightarrow \frac{\partial C}{\partial z^L} = \frac{\partial C}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L} \rightarrow \text{activation func'}$$

for the final layer based on the tables in previous page; then recursively compute.

$$\frac{\partial C}{\partial z^m} = \left(w^{m+1T} \frac{\partial C}{\partial z^{m+1}} \right) \cdot \frac{\partial a^m}{\partial z^m} - s_3.$$

for all other layers. Plug these values into

$$\frac{\partial C}{\partial w^m} = \frac{\partial C}{\partial z^m} \cdot a^{m-1T} - s_1$$

and

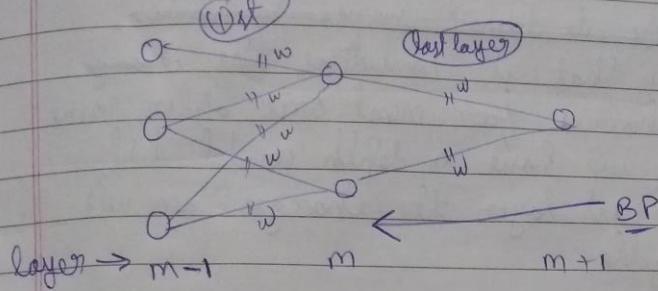
$$\frac{\partial C}{\partial b^m} = \frac{\partial C}{\partial z^m} - s_3$$

to obtain the updates.

By substituting s_3 in $\frac{\partial C}{\partial b^m}$ we get $\frac{\partial C}{\partial b^m}$ i.e b_n

and in s_3 we get $\frac{\partial C}{\partial w^m}$ i.e w_n (updated wt)

Consider eqⁿ S₃ :-



of m th layer.

$$\frac{\delta C}{\delta z^m} \rightarrow \text{rate of change of C.F., w.r.t. input of } m\text{th layer}$$

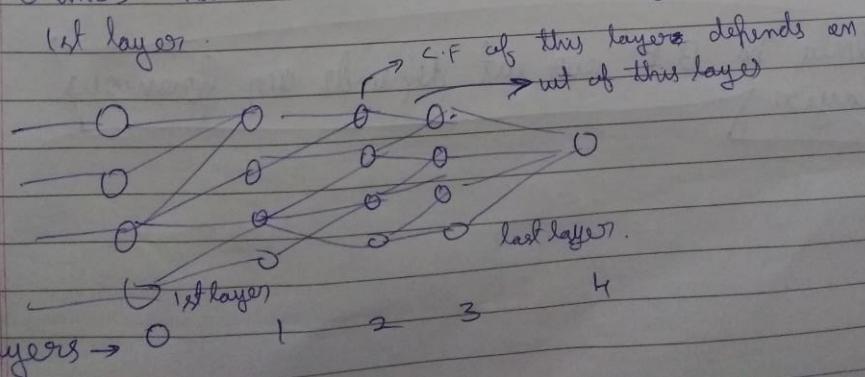
$\frac{\delta C}{\delta z^m} \quad (z^m) \text{ defn}$

~~depends upon weight of previous layer (w^{m+1})~~

In Backward
Propagation η
(BP)

Therefore we say that if we have many layers in between & as C.F. is dependent on wts of previous layer so if the model is deep, there are very high chances that there will be no change in wts of

1st layer.



Similarly C-F depends on wts of previous layer
in B.P.

Therefore if the wts doesn't change by much,
Therefore we say that it considers the ~~change~~
Previous wt changes from next one after some
time i.e if we have a depth model till
we reach the 1st layer the changes in wts
will be negligible.

Therefore we avoid making model with ~~deep~~ many
layers.

∴ There will be no learning if wts doesn't
change

$$\text{eg } y = x^2 \quad \frac{dy}{dx} = 2x \quad , \quad x=3$$

$$q \quad \frac{d^2y}{dx^2} = 2$$

$$\frac{d^3y}{dx^3} = 0$$

similar things happens in case
of wts if we keep taking
derivative of it.

Only in B.P our wts depends on previous
layer

Example :-

Sigmoid Network with cross-entropy loss using gradient descent

$$\text{Step 1} \quad \frac{\partial C}{\partial z^L} = \frac{\partial C}{\partial a^L} \cdot \frac{\partial a^L}{\partial z^L} \quad \xrightarrow{\text{actual "funtal"}}$$

$Ty \rightarrow \text{actual value}$

↳ L.F.

$$= \left(\frac{a^L - y}{a^L(1-a^L)} \right) a^L (1-a^L) \quad \text{from tables}$$

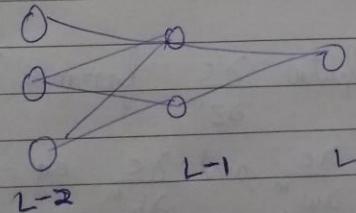
$$\frac{\partial C}{\partial z^L} = a^L - y \quad \begin{array}{l} \rightarrow \text{as we have got the value of } \frac{\partial C}{\partial z} \\ \text{we can find out } s_3 \end{array}$$

$$\frac{\partial C}{\partial b^L} = \eta(a^L - y)$$

Step 2. Forward here, we can compute

$$\therefore b_n = b_{n-1} = \eta(a_n - y)$$

Bias final value



$$\begin{aligned} \frac{\partial C}{\partial z^{L-1}} &= \left(w^{LT} \frac{\partial C}{\partial z^L} \right) \cdot \frac{\partial a^{L-1}}{\partial z^{L-1}} \\ &= w^{LT} (a^L - y) \cdot a^{L-1} (1 - a^{L-1}) \end{aligned} \quad \begin{array}{l} \text{finding } s_3 \text{ for} \\ m = L-1 \end{array}$$

$$\begin{aligned} \frac{\partial C}{\partial z^{L-2}} &= \left(w^{L-1T} \frac{\partial C}{\partial z^{L-1}} \right) \cdot \frac{\partial a^{L-2}}{\partial z^{L-2}} \\ &= w^{L-1T} (w^{LT} (a^L - y) \cdot a^{L-1} (1 - a^{L-1})) \cdot a^{L-2} (1 - a^{L-2}) \end{aligned}$$

~~We can also find now~~ Eq_3

$$w_n = w_{n-1} - \eta (a^t - y) \quad \text{Updating the weight.}$$

for $L-1$ layer Similarly we can find for all layers.
 $w_{L-1} = \dots$ as shown.

Important Note:-

If our cost Function $= 0$ \therefore in this case $y - a^L = 0$

$$\text{then } \boxed{\text{Imp Eq} = 0} \quad \therefore \frac{\partial c}{\partial z^m} = 0$$

$\because \frac{\partial c}{\partial b^m} = 0$ therefore there will be no
updates of bias & weights

and so on, until we have computed $\frac{\partial c}{\partial z^m}$ for $m \in \{1 \dots L\}$

This allows us to compute $\frac{\partial c}{\partial w_{ij}}$ and $\frac{\partial c}{\partial b_i}$ eg.

$$\frac{\partial c}{\partial w^l} = \frac{\partial c}{\partial z^l} a^{l-1}$$

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

$$\frac{\partial c}{\partial w^{l-1}} = \frac{\partial c}{\partial z^{l-1}} a^{l-2}$$

$$= w^L (a^l - y) \cdot a^{l-1} (1 - a^{l-1}) a^{l-2}$$

So standard gradient descent then updates each parameter as follows:

$$\left. \begin{array}{l} w^m = w^n - \lambda \frac{\partial C}{\partial w^n} \\ b^m = b^n - \lambda \frac{\partial C}{\partial b^m} \end{array} \right\}$$

this was our
aim i.e
new weights &
biases get updated

where λ is the learning rate. This process is repeated until some stopping criteria is met.