

# Neural Network :-

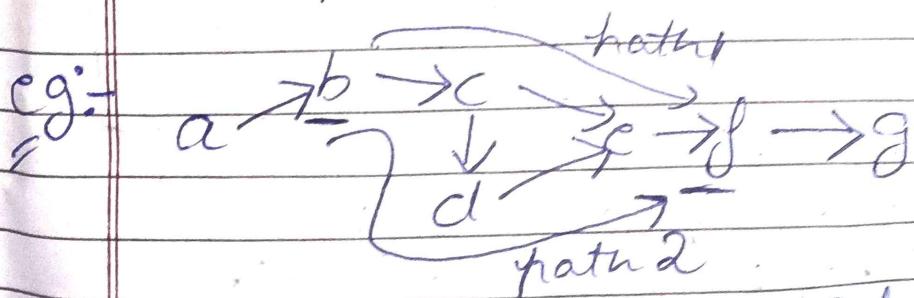
Prerequisite:-

# Chain of derivation:-

$x \rightarrow z \Rightarrow z$  is dependent on  $x$   
 $y \rightarrow z \quad z$  is dependent on  $y$   
i.e.  $z = f(x, y)$   
 $x, y$  both appear in  $z$

$x \rightarrow y \rightarrow z \Rightarrow y$  dependent on  $x$   
and  $z$  dependent on  
(only  $y$  appear in  $z$ )

$\rightarrow$  Derivative of a var w.r.t another var  
is sum of product of derivative along  
each path



$$\frac{\partial f}{\partial a} = \frac{\partial f}{\partial e} \frac{\partial e}{\partial c} \frac{\partial c}{\partial b} + \frac{\partial f}{\partial d} \frac{\partial d}{\partial e} \frac{\partial e}{\partial a}$$

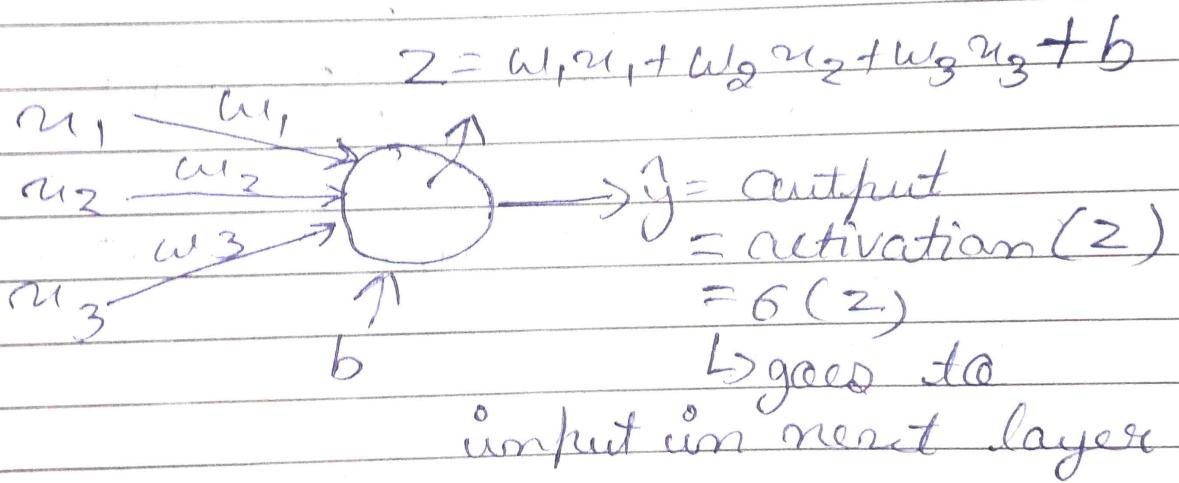
↓                      ↓  
path 1                  path 2

eg  $t \rightarrow u = f(t) \rightarrow z = h(u, y)$   
 $\rightarrow y = g(t)$

$$\frac{\partial z}{\partial t} = \frac{\partial z}{\partial u} \frac{\partial u}{\partial t} + \frac{\partial z}{\partial y} \frac{\partial y}{\partial t}$$

$\hookrightarrow$   $z, y$  are only dependent on  $t$ .

# Neuron:



Matrix form:

$$w = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \quad x = \begin{bmatrix} u_1 \\ u_2 \\ u_3 \end{bmatrix}$$

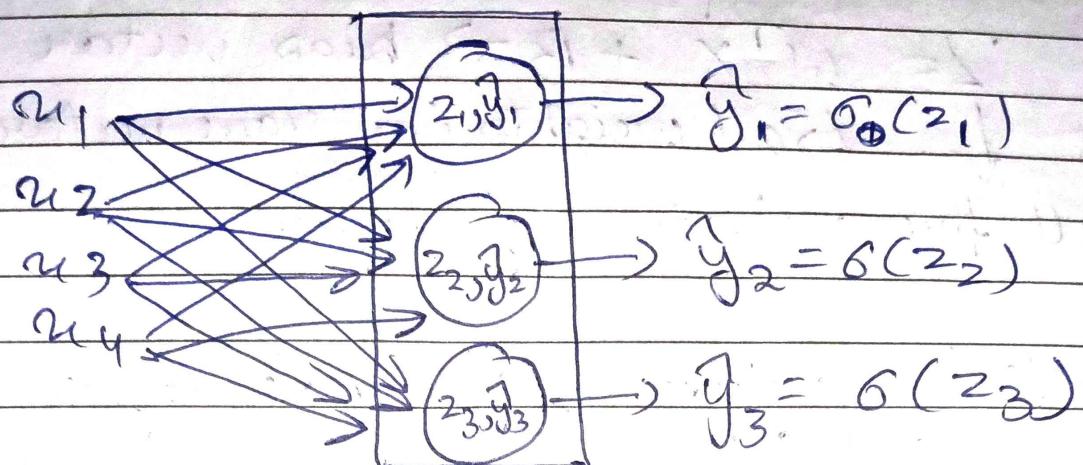
$$z = \vec{w}^T \vec{u} + b$$

$\hookrightarrow$  dot product ( $\vec{w} \cdot \vec{u}$ )

$\hookrightarrow$  helps in vectorization

$\rightarrow$  output calculated by passing to  $\sigma(z)$ .

# layer: Multiple neurons in a layer are fed same input.



$b_i$  = bias of  $i$ th neuron layer

$W_i$  = weight matrix of  $i$ th neuron

$z_i$  = pre activation of  $i$ th neuron

$y_i$  = prediction

$X$  = input vector (given to all)

→ activation function  $\sigma$  is same for every neuron in layer

Matrix form:

$$z_1 = w_1^T X + b_1$$

$$z_2 = w_2^T X + b_2$$

$$z_3 = w_3^T X + b_3$$

$$Z = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} = \begin{bmatrix} \leftarrow w_1^T \rightarrow \\ \leftarrow w_2^T \rightarrow \\ \leftarrow w_3^T \rightarrow \end{bmatrix} X + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \end{bmatrix}$$

$$= \begin{bmatrix} \overset{\uparrow}{w_1} & \overset{\uparrow}{w_2} & \overset{\uparrow}{w_3} & \overset{T}{\downarrow} \end{bmatrix} X + B$$

$$= w^T X + B$$

→ assume weight of neurons are  
# column vectors

$\rightarrow$  input

$$Z = w^T x + \beta \rightarrow \text{bias vector}$$

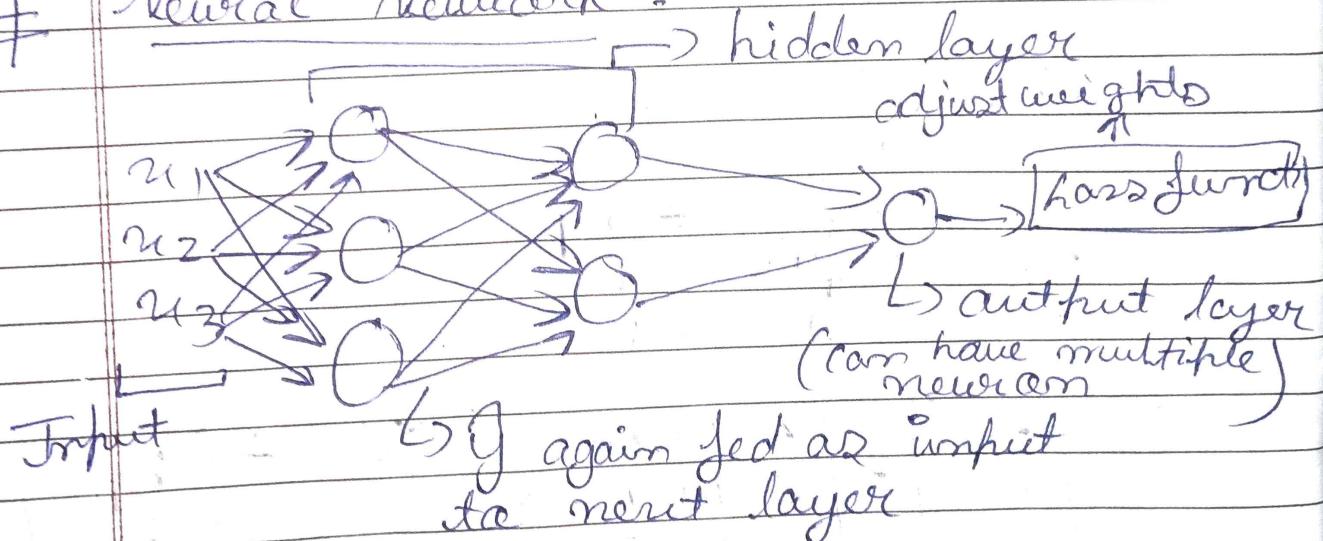
$\downarrow$  all weights in same vector form  
all pre activation of layer

$\sigma(z) = \text{activation of whole layer}$

$$= \begin{bmatrix} \sigma(z_1) \\ \sigma(z_2) \\ \sigma(z_3) \end{bmatrix}$$

$\rightarrow$  Different layer can have different activation function

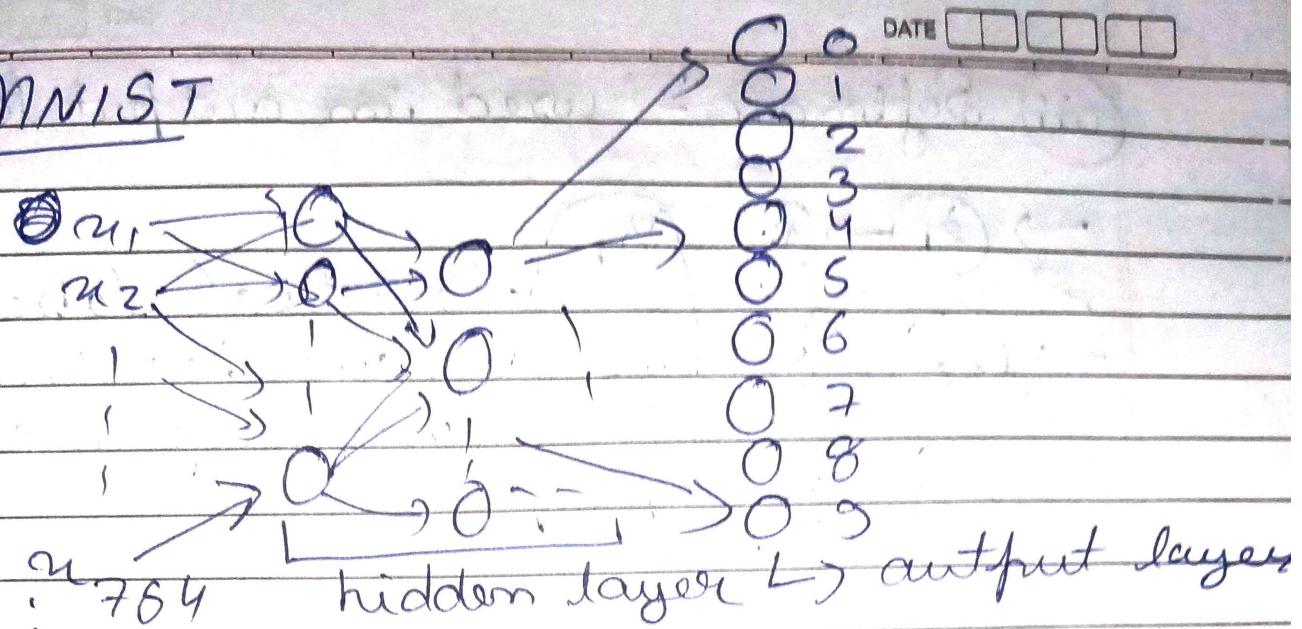
Neural Network :



$\rightarrow$  each layer can have different activation functions

$\rightarrow$  output can be many or one

## MNIST



↳ grayScale value of MNIST image

→ output layer uses Softmax

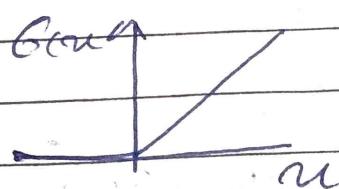
↳ output probability

→ neuron with max probability is chosen as prediction

## Activation function

→ Takes  $z$  and give output  $\sigma(z)$

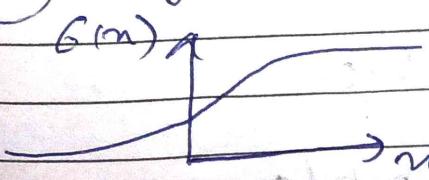
(i) ReLU



$$\sigma(u) = \max(0, u)$$

$$\text{Derivative} = \sigma'(u) = \begin{cases} 1 & u > 0 \\ 0 & u \leq 0 \end{cases}$$

(ii) Sigmoid

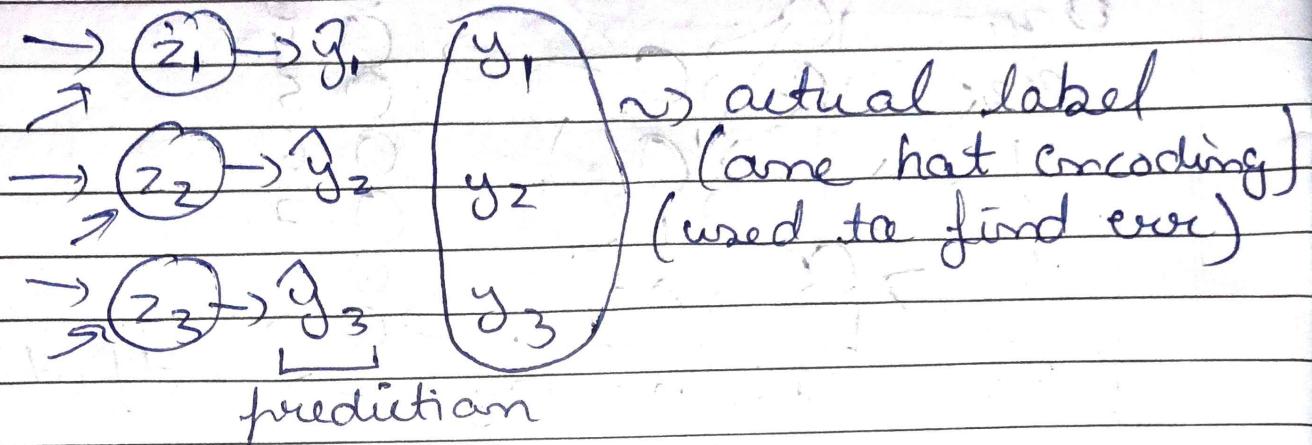


$$\sigma(u) = \frac{1}{1 + e^{-u}}$$

$$\text{Derivative} = \sigma(u)(1 - \sigma(u))$$

$$\sigma'(u) = \sigma(u)(1 - \sigma(u))$$

iii) Softmax :- used in output layer



$$y_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$y_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$y_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\rightarrow y_1 + y_2 + y_3 = 1$$

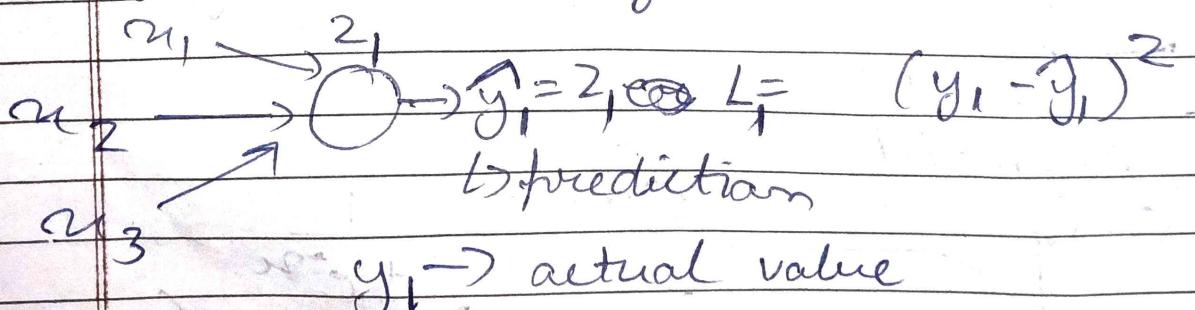
Calculating  $z_i$  requires  $z_s$  of all layers

Derivative

$$\frac{\partial y_i}{\partial z_j} = -y_i y_j \quad k \neq j$$

# Error function :-

i) SSE (Sum of Square) :-  
in linear regression



↳ for a single sample

for multiple

$$L = L_1 + L_2 + \dots$$

$$= \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$\hookrightarrow$  optimal averaging

### (iii) Cross Entropy loss -

In classification (MNIST)

$$L = -\sum y_i \log \hat{y}_i \quad [\text{Single Sample}]$$

$y_i$  = probability of  $i$ th class

e.g.  $O \rightarrow \hat{y}_1$   $\oplus$  In actual the class

$O \rightarrow \hat{y}_2$  label is 2

$O \rightarrow \hat{y}_3$  Hence  $y = [0, 1, 0, 0, 0]$

$O \rightarrow \hat{y}_4$   $\hat{y}_1 \hat{y}_2 \hat{y}_3 \hat{y}_4 \hat{y}_5$

$O \rightarrow \hat{y}_5$   $\hookrightarrow$  one hot encoding

prediction for class 2

$$\text{Loss} = -[y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3]$$

$$+ y_4 \log \hat{y}_4 + y_5 \log \hat{y}_5]$$

$$= -[0 + \log \hat{y}_2 + 0 + 0 + 0]$$

$$= -\log \hat{y}_2$$

$\hookrightarrow$  probability for class 2

given by model

→ for multiple reward

$$L = \frac{L_1 + L_2 + \dots + L_n}{n}$$

$$= \frac{\sum_{i=1}^n \log y_i}{n} + \frac{\sum_{i=1}^n \log y_{\bar{i}}}{n}$$

## # Optimization :-

a function  $z = f(w_1, w_2, w_3)$

decreases in direction of  $\left\langle -\frac{\partial z}{\partial w_1}, -\frac{\partial z}{\partial w_2}, -\frac{\partial z}{\partial w_3} \right\rangle$

i.e. if we are at point  $(2, 5, 3)$   
and then we calculate  $\left\langle -\frac{\partial z}{\partial w_1}, -\frac{\partial z}{\partial w_2}, -\frac{\partial z}{\partial w_3} \right\rangle$

we can make a small nudge in  
that direction

$$\vec{R}_{\text{next}} = \vec{R}_{\text{current}} - \alpha \nabla f$$

As we move,  $\nabla f$  also change

so we recalculate  $\nabla f$  & keep

moving until minimal  $z$  is achieved

→ only gives local minima/maxima

## # Gradient descent:-

Loss function is dependent on  
weight & biases

$$L = f(w_1, w_2, w_3, b_1, b_2)$$

We calculate  $\nabla f$ , then move in negative direction.

$$\nabla f = \left\langle \frac{\partial f}{\partial w_1}, \frac{\partial f}{\partial w_2}, \frac{\partial f}{\partial w_3}, \frac{\partial f}{\partial b_1}, \frac{\partial f}{\partial b_2} \right\rangle$$

Then update weight & bias

$$w_1' = w_1 - \alpha \frac{\partial f}{\partial w_1} \quad w_3' = w_3 - \alpha \frac{\partial f}{\partial w_3}$$

$$w_2' = w_2 - \alpha \frac{\partial f}{\partial w_2} \quad b_1' = b_1 - \alpha \frac{\partial f}{\partial b_1}$$

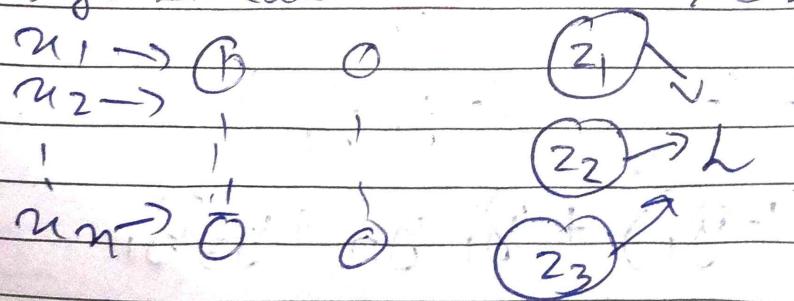
$$b_2' = b_2 - \alpha \frac{\partial f}{\partial b_2}$$

# Calculating  $\frac{\partial L}{\partial w_i}$  &  $\frac{\partial L}{\partial b_j}$

- We use chains of derivative
- for MNIST classification outer layer

# For Single Sample

→ first calculate  $\frac{\partial L}{\partial z_i}$  for each layer



Für die am outer layer first werden wir  
Dependency graph benutzen.

$$\begin{aligned} z_1 &\rightarrow \hat{y}_1 = \frac{e^{z_1}}{S} \\ z_2 &\rightarrow \hat{y}_2 = \frac{e^{z_2}}{S} \quad \rightarrow L = -\sum y_i \log \hat{y}_i \\ z_3 &\rightarrow \hat{y}_3 = \frac{e^{z_3}}{S} \quad \boxed{\hat{y}_1, \hat{y}_2, \hat{y}_3} \end{aligned}$$

$\rightarrow y_i$  abhängt von all  $z_i$

$$S = e^{z_1} + e^{z_2} + e^{z_3}$$

$$L = -\sum y_i \log \hat{y}_i = -[y_1 \log \hat{y}_1 + y_2 \log \hat{y}_2 + y_3 \log \hat{y}_3]$$

$\Rightarrow \frac{\partial L}{\partial z_1} =$  all paths from  $L$  to  $z_1$ ,

$\hookrightarrow z_1 = \text{logit of first class}$

$$\Rightarrow \frac{\partial L}{\partial z_1} = \frac{\partial L}{\partial \hat{y}_1} \frac{\partial \hat{y}_1}{\partial z_1} + \frac{\partial L}{\partial \hat{y}_2} \frac{\partial \hat{y}_2}{\partial z_1} + \frac{\partial L}{\partial \hat{y}_3} \frac{\partial \hat{y}_3}{\partial z_1}$$

$$= \cancel{\left( \frac{\partial}{\partial \hat{y}_1} \right)} y_1 (1 - \hat{y}_1) + \left( \frac{\partial}{\partial \hat{y}_2} \right) (-\hat{y}_2 \hat{y}_1)$$

$$+ \left( \frac{\partial}{\partial \hat{y}_3} \right) (-\hat{y}_3 \hat{y}_1)$$

$$= -y_1 (1 - \hat{y}_1) + \hat{y}_1 y_2 + y_3 \hat{y}_1$$

$$= -y_1 + y_1 \hat{y}_1 + \hat{y}_1 y_2 + y_3 \hat{y}_1$$

$$= -y_1 + \bar{y}_1(y_1 + y_2 + y_3)$$

$$= -y_1 + \bar{y}_1 \quad y_1 + y_2 + y_3 = 1 \quad (\text{Same hat})$$

$$= \bar{y}_1 - y_1$$

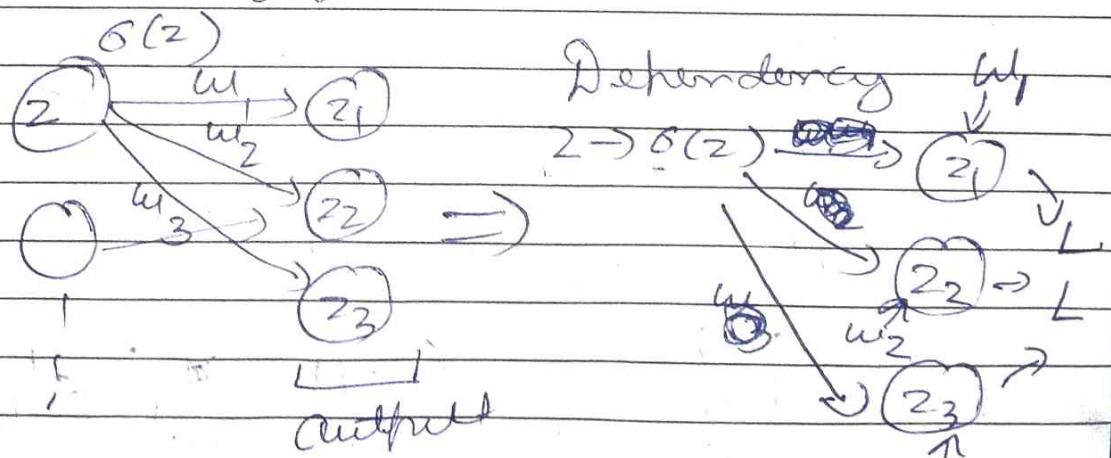
$$\boxed{\frac{dL}{dz_i} = \bar{y}_i - y_i}$$

→ In General

→ We calculate  $\frac{\partial L}{\partial z_1}, \frac{\partial L}{\partial z_2}, \frac{\partial L}{\partial z_3} \delta s_a$

an for last layer.

# To calculate  $\frac{\partial L}{\partial z_i}$  for previous layer



$$\begin{aligned} \frac{\partial L}{\partial z} &= \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial \sigma(z)} \frac{\partial \sigma(z)}{\partial z} + \frac{\partial L}{\partial z_2} \frac{\partial z_2}{\partial \sigma(z)} \frac{\partial \sigma(z)}{\partial z} \\ &\quad + \frac{\partial L}{\partial z_3} \frac{\partial z_3}{\partial \sigma(z)} \end{aligned}$$

$$\begin{aligned} &= \frac{\partial L}{\partial z_1} \cdot w_1 \cdot \sigma'(z) + \frac{\partial L}{\partial z_2} \cdot w_2 \cdot \sigma'(z) \\ &\quad + \frac{\partial L}{\partial z_3} \cdot w_3 \cdot \sigma'(z) \end{aligned}$$

Hence

$$\frac{\partial L}{\partial z} = \left( \sum_{j=1}^n \frac{\partial L}{\partial z_j} w_j \right) \sigma'(z)$$

In matrix form

$$\frac{\partial L}{\partial z} = (w^T) \sigma'(z)$$

↳ derivative of activation w.r.t  $z$

↳ weight vector of connected edge

$$S = \begin{bmatrix} \frac{\partial L}{\partial z_1} \\ \frac{\partial L}{\partial z_2} \\ \frac{\partial L}{\partial z_3} \end{bmatrix} = \text{delta vector}$$

For whole layer

→ To find  $\frac{\partial L}{\partial z}$

$$z_{11} \quad z_{21} \quad z_{31} \\ z_{12} \quad z_{22} \quad z_{32} \\ z_{13} \quad z_{23} \quad z_{33}$$

(1)                          (2)  $\rightarrow \frac{\partial L}{\partial z}$  formula

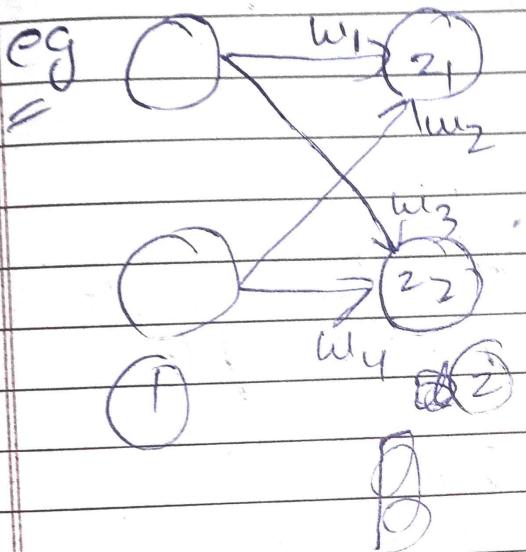
$$z = \begin{bmatrix} z_{11} \\ z_{12} \\ z_{13} \\ z_{21} \\ z_{22} \\ z_{23} \\ z_{31} \\ z_{32} \\ z_{33} \end{bmatrix}$$

$$S = \begin{bmatrix} \frac{\partial L}{\partial z_{11}} \\ \frac{\partial L}{\partial z_{12}} \\ \frac{\partial L}{\partial z_{13}} \end{bmatrix} = (w^T S) \sigma'(z)$$

( $w$  = weight of layer (2))

$S$  = delta of layer 2

$\sigma'(z)$  = activation derivative of layer 1



$$W_2 = \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix}$$

$\hookrightarrow$  weight matrix of layer 2

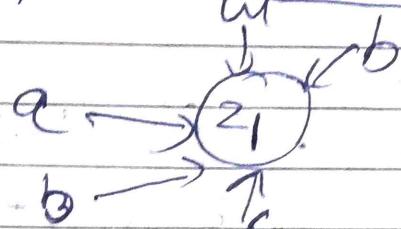
$$\begin{aligned} W_2 S &= \begin{bmatrix} w_{11} & w_{12} \\ w_{21} & w_{22} \end{bmatrix} \begin{bmatrix} \frac{\partial L}{\partial z_{21}} \\ \frac{\partial L}{\partial z_{22}} \end{bmatrix} \\ &= \left[ \begin{array}{l} (\frac{\partial L}{\partial z_{21}}) w_{11} + (\frac{\partial L}{\partial z_{22}}) w_{12} \\ (\frac{\partial L}{\partial z_{21}}) w_{21} + (\frac{\partial L}{\partial z_{22}}) w_{22} \end{array} \right] \end{aligned}$$

$\hookrightarrow$  matches with chain of derivative

$\rightarrow$  May need to use  $W^T S$  if weights of neurons are row vector.

→ Name all other layer  $\partial S$  can be calculated with only using next layer  $S$  without calculating full path derivative.

# To calculate  $\frac{\partial L}{\partial w_i}$ :



$a$  = activation from previous layer or input

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial w_i} = \left( \frac{\partial L}{\partial z_1} \right) \cdot a$$

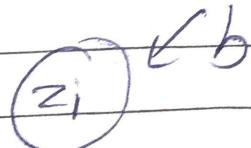
activation of weight  $w_i$

$$z_1 = w_i a + w_i b + w_2 c + b$$

$$\frac{\partial z_1}{\partial w_i} = a$$

$$\boxed{\frac{\partial L}{\partial w_i} = \left( \frac{\partial L}{\partial z_i} \right) a} \rightarrow \text{for single sample}$$

# To calculate  $\frac{\partial L}{\partial b}$ :

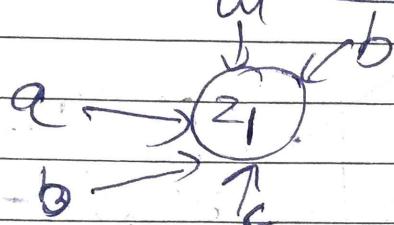


$$\frac{\partial b}{\partial L} = \frac{\partial L}{\partial z_1} \frac{\partial z_1}{\partial b} = \frac{\partial L}{\partial z_1}$$

$$\frac{\partial z_1}{\partial b} = 1$$

→ Now all other layer  $\theta_8$  can be calculated with only using next layer  $\theta_8$  without exploring full path derivative

# To calculate  $\frac{\partial L}{\partial w_i}$ :



$a$  = activation from previous layer or input

$$\frac{\partial L}{\partial w_i} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_i} = \left( \frac{\partial L}{\partial z_1} \right) \cdot a$$

activation of weight  $w_i$

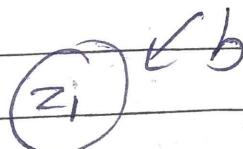
$$z_1 = w_1 a + w_2 b + w_3 c + b$$

$$\frac{\partial z_1}{\partial w_i} = a$$

$$\frac{\partial L}{\partial w_i} = \left( \frac{\partial L}{\partial z_1} \right) a$$

for single Sample

# To calculate  $\frac{\partial L}{\partial b}$ :



$$\frac{\partial L}{\partial b} = \frac{\partial L}{\partial z_1} \cdot \frac{\partial z_1}{\partial b} = \frac{\partial L}{\partial z_1}$$

$$\frac{\partial z_1}{\partial b} = 1$$

$$\frac{\partial L}{\partial b_i} = \frac{\partial L}{\partial z_i}$$

Name after calculating  $\frac{\partial L}{\partial w_i}$ ,  $\frac{\partial L}{\partial b_i}$   
we update weight and bias

$$w_i' = w_i - \alpha \frac{\partial L}{\partial w_i}$$

$$b_i' = b_i - \alpha \frac{\partial L}{\partial b_i}$$

In matrix form

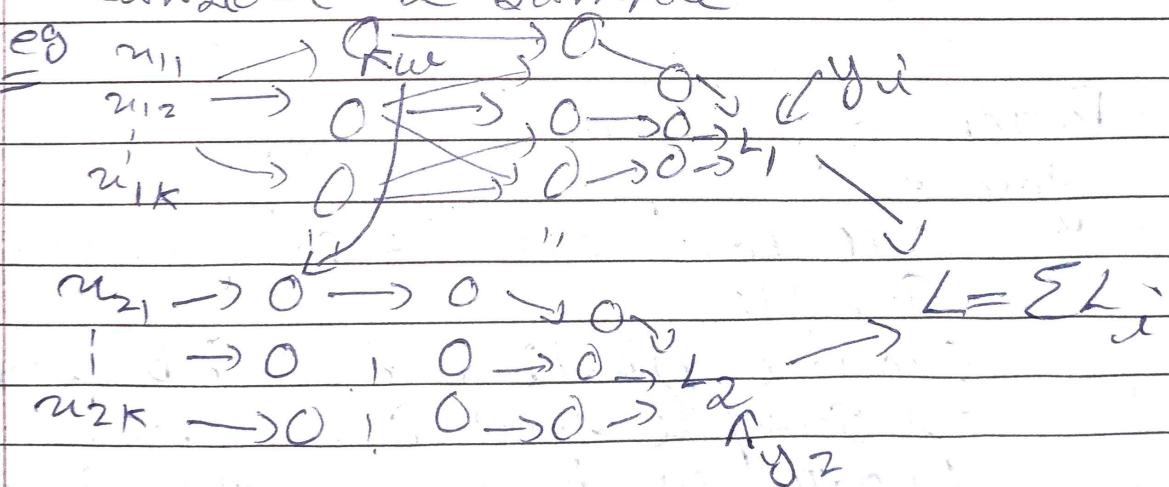
$$W = W - \alpha \nabla W$$

$$B = B - \alpha \nabla B$$

## # Batch gradient descent:

In Batch GD Loss of whole batch  
is considered

Consider 2 Sample



Consider the weight  $w_i$  of a neuron.

In dependency graph the weight ( $w_i$ ) affect both  $L_1$  &  $L_2$

(pre activation changes per sample hence  $\frac{\partial L}{\partial z_i}$  changes)

$\rightarrow$  weight has path through all samples

Hence

$$w = w - \alpha \left( \frac{\partial L}{\partial w} \right)$$

() Sample sum of gradient

Similarly for  $b$   $b = b - \alpha \left( \sum \frac{\partial L}{\partial b} \right)$

However  $\frac{\partial L}{\partial z}$  are not connected

to other graph  $\frac{\partial L}{\partial z}$  is different for different sample  
Hence  $\frac{\partial L}{\partial z}$  is not added

Batch algorithm -

for  $e$  in epochs (100),

calculate all  $\frac{\partial L}{\partial z}$  for all samples.

Calculate  $\frac{\partial L}{\partial w}$  &  $\frac{\partial L}{\partial b}$  for all samples and sum them

Update  $w, b$  using the sum

DATE

Stochastic  
for each (100):  
 $\theta$

for each Sample  $S$ :  
calc. all  $\frac{\partial L}{\partial z}, \frac{\partial L}{\partial w}, \frac{\partial L}{\partial b}$  for  $S$   
update  $w, b$

Adam Optimizer & Momentum  
based learning rates: