

Perception Loss Function

Perception: $L(y_i, f(x_i)) = \max(0, -y_i f(x_i))$

Hinge loss: $L(y_i, f(x_i)) = \max(0, 1 - y_i f(x_i))$

SGD → optimization

$$E(w, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda R(w)$$

where L is a loss function that measures model (mis) fit & R is regularization that penalizes model complexity; λ is a non-negative hyperparameter that controls the regularization strength.

$$L(w_1, w_2, b) = \frac{1}{n} \sum_{i=1}^n L(y_i, f(x_i)) + \lambda R(w_1, w_2) \quad \text{Regularization}$$

$$L(y_i, f(x_i)) = \max(0, 1 - y_i f(x_i))$$

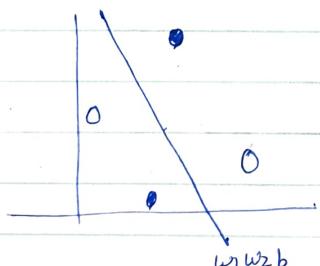
$f(x_i) = w_1 x_1 + w_2 x_2 + b$

$$L = \underset{w, b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, 1 - y_i f(x_i))$$

minimum
Gradient descent

where $f(x_i) = w_1 x_1 + w_2 x_2 + b$

E.g



$$\begin{array}{|c|c|c|} \hline & x_1 & x_2 & y \\ \hline 1 & x_{11} & x_{12} & y_1 \\ \hline 2 & x_{21} & x_{22} & y_2 \\ \hline \end{array}$$

$$f(x_i) = w_1 x_{1i} + w_2 x_{2i} + b$$

cpa	i	row	placed
7	8	1	0
6	8	-1	0
4	2	1	0
1	1	-1	0

y _i	y _i
1	1
-1	-1
1	-1
-1	1

Let $L = \frac{1}{2} [\max(0, -y_1 f(x_1)) + \max(0, -y_2 f(x_2))] \text{ for 2 points}$

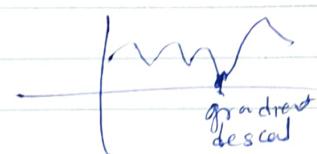
$$w_1 x_{1i} + w_2 x_{2i} + b$$

$$\max(0, -y_i f(x_i)) = \begin{cases} -y_i f(x_i) & -y_i f(x_i) \geq 0 \\ 0 & -y_i f(x_i) < 0 \end{cases}$$

y_i	\hat{y}_i	$\max(0, -y_i f(x_i))$	$w_1 x_1 + w_2 x_2 + b \geq 0$	$\text{ie } \max(0, -ve) = 0$
1	1	ve	$w_1 x_1 + w_2 x_2 + b < 0$	$\text{ie } \max(0, -ve) = 0$
-1	-1	$y_i = -ve$	$f(x_i) = -ve$	$\text{ie } \max(0, -ve) = 0$
1	-1	$y_i = +ve$	$f(x_i) = -ve$	$\text{ie } \max(0, +ve) = +ve$
-1	1	$y_i = -ve$	$f(x_i) = +ve$	$\text{ie } \max(0, +ve) = +ve$

$$L = \underset{w_1, w_2, b}{\operatorname{argmin}} \frac{1}{n} \sum_{i=1}^n \max(0, -y_i f(x_i))$$

where $f(x_i) = w_1 x_{i1} + w_2 x_{i2} + b$



for i in epochs: $\eta = 0.01$

$$w_1 = w_1 + \eta \frac{\partial L}{\partial w_1}$$

$$w_2 = w_2 + \eta \frac{\partial L}{\partial w_2}$$

$$b = b + \eta \frac{\partial L}{\partial b}$$

$$\frac{\partial L}{\partial w_1} = \frac{\partial L}{\partial f(x_i)} \times \frac{\partial f(x_i)}{\partial w_1}$$

$$\frac{\partial L}{\partial f(x_i)} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial f(x_i)}{\partial w_1} = x_{i1}$$

$$\frac{\partial L}{\partial w_1} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i1} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial L}{\partial w_2} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i x_{i2} & \text{if } y_i f(x_i) < 0 \end{cases}$$

$$\frac{\partial L}{\partial b} = \begin{cases} 0 & \text{if } y_i f(x_i) \geq 0 \\ -y_i & \text{if } y_i f(x_i) < 0 \end{cases}$$

Def Loss Function: Binary Cross Entropy

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

⇒ Perception

Step function

use as a perception

⇒ Perception == Logistic Regression if

Activation funⁿ = Sigmoid
Loss funⁿ = Binary cross entropy

⇒ Activation function = Softmax
Loss function = Categorical cross entropy

$f = \frac{e^{z_i}}{\sum_{j=1}^m e^{z_j}}$

$L = \sum_{j=1}^m y_j \log(p_j)$

For multiclass
Softmax regression

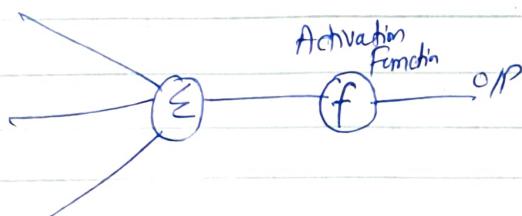
⇒ For Regression

Activation function → Linear function
Loss function → MSE ($y - y_0$)

Linear regression

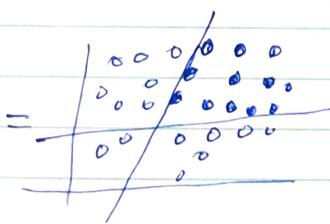
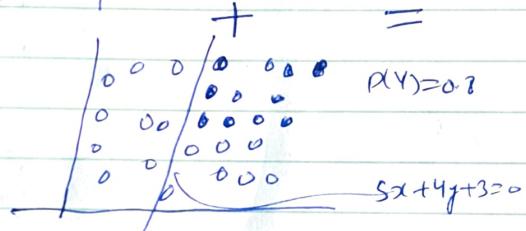
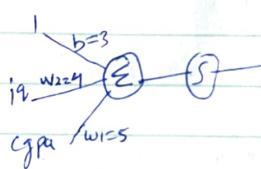
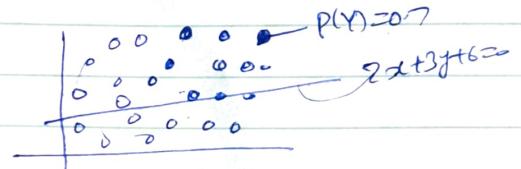
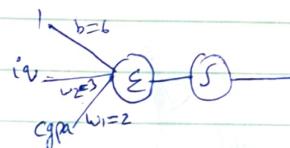
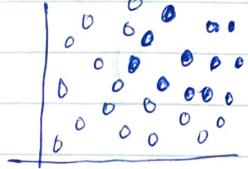
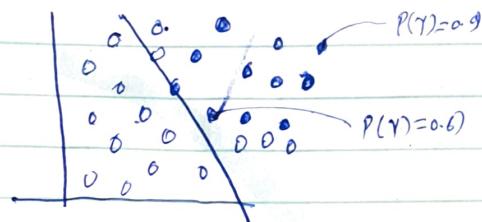
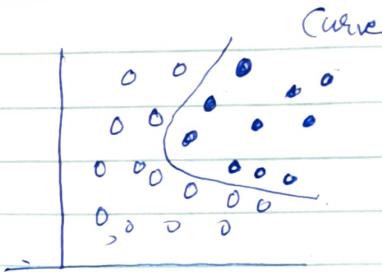
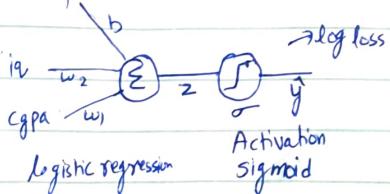
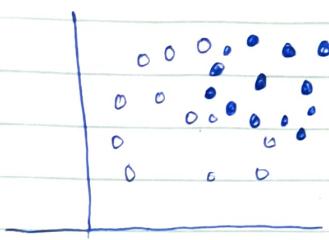
⇒ Perception is a mathematical model which can be used in multiple ways because it is flexible by design

i.e	Loss function	Activation	Output
	Hinge loss	Step	perception → binary classifier
	Log loss (Binary cross entropy)	Sigmoid	Logistic regression → binary classifier
	Categorical cross entropy	Softmax	Softmax regression → multiclass classifier O/P probability
	MSE	linear	Linear regression O/P = number

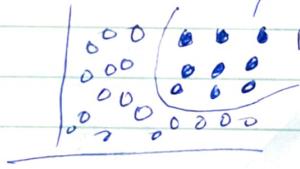


Multi Layer Perceptron (MLP)

Date _____
Page _____

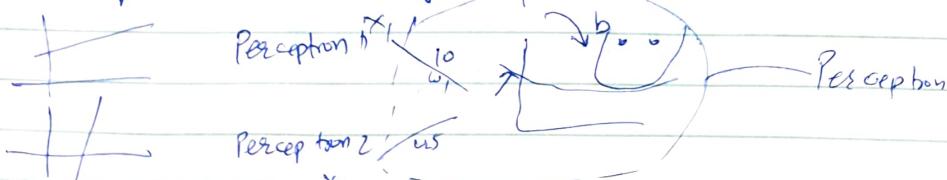


smoothing



$$0.7 + 0.8 = 1.5 \quad \sigma(z) = \frac{1}{1+e^{-z}} = 0.82$$

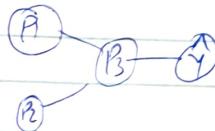
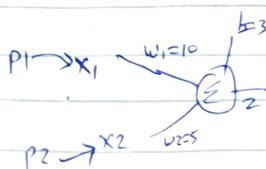
linear Combination of two perception

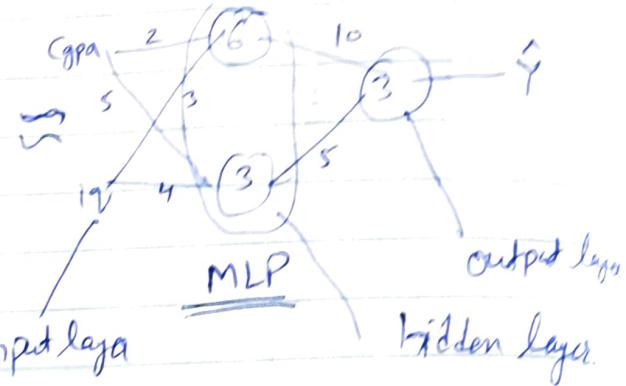
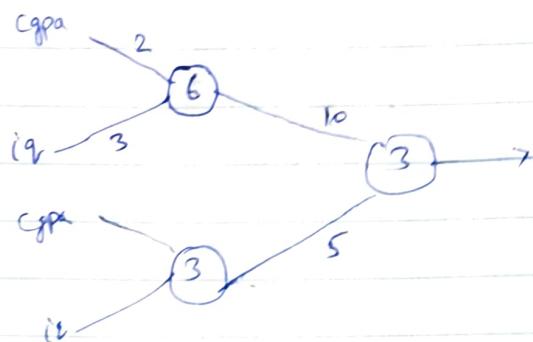


$$\text{weighted } z = 0.7 \times 10 + 0.8 \times 5 + b$$

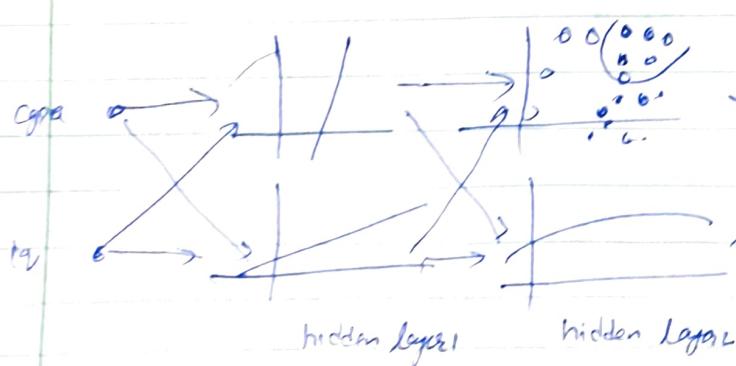
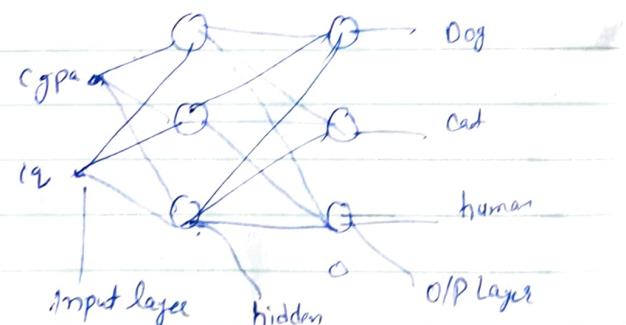
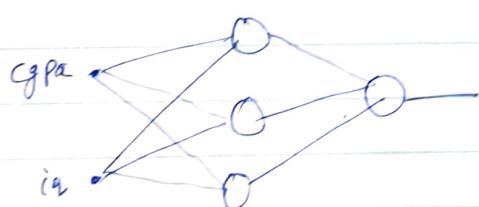
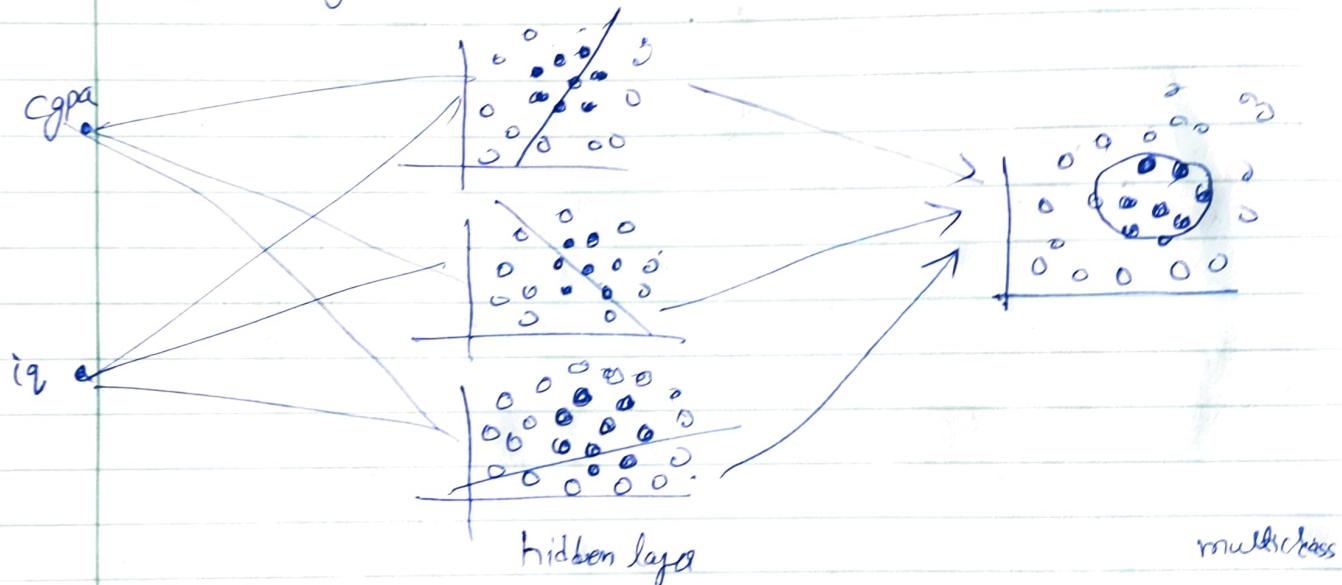
$$\sigma(z) = \frac{1}{1+e^{-z}}$$

weightage of Perception is more





By Using multiple perception we can capture non-linear,

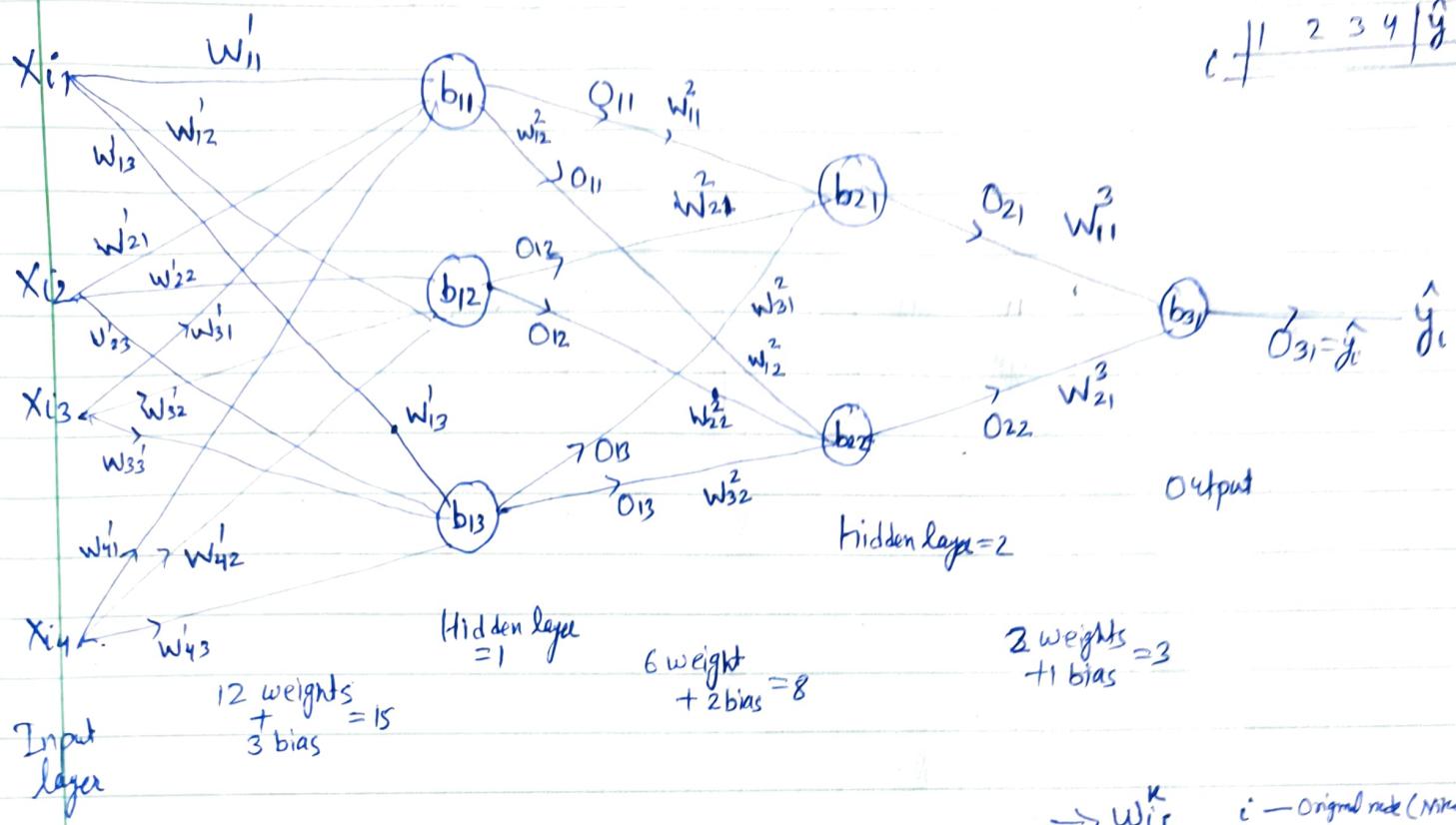


Demo at <https://playground.tensorflow.org/>

Multi-layer Perceptron Notation

Data $\rightarrow \{m \times n\}$
 $m \rightarrow \text{rows}$
 $n = 4$
 each row denoted by X_i

$i \mid 1 \ 2 \ 3 \ 4 \mid \hat{y}$

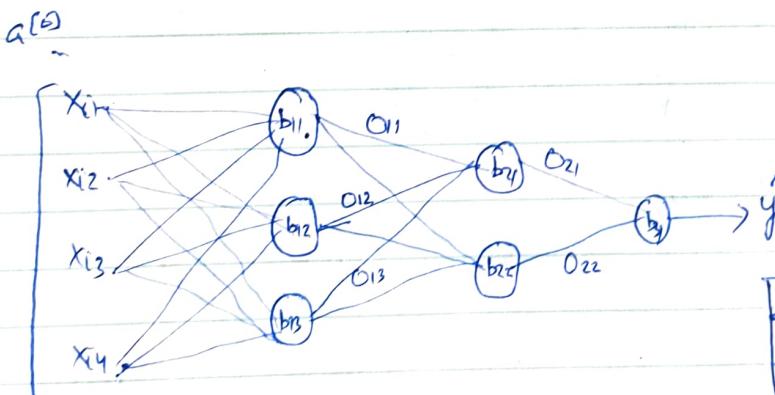


$$\text{Trainable parameters} = 15 + 8 + 3 = 26$$

$\Rightarrow w_{ij}^k$ $i \rightarrow \text{Original node (Middle node)}$
 $j \rightarrow \text{Dest node (Final output)}$

$\Rightarrow b_{ij}^k$ $i \rightarrow \text{layer}, j \rightarrow \text{node}$

$\Rightarrow \text{Output } O_{kj}$



Forward Propagation

How a NN Predicts O/P?

of trainable parameters

CGPA	i ₁	i ₂	i ₃	Predicted
7.2	72	69	81	1
8.1	92	75	76	0

Layer #1

$$\begin{bmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \\ w_{31} & w_{32} & w_{33} \\ w_{41} & w_{42} & w_{43} \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix} = \begin{bmatrix} w_{11}x_1 + w_{21}x_2 + w_{31}x_3 + w_{41}x_4 \\ w_{12}x_1 + w_{22}x_2 + w_{32}x_3 + w_{42}x_4 \\ w_{13}x_1 + w_{23}x_2 + w_{33}x_3 + w_{43}x_4 \end{bmatrix} \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

Prediction $\Rightarrow o(W^T x + b)$

$$= \sigma \left(\begin{bmatrix} w'_{11}x_{i1} + w'_{21}x_{i2} + w'_{31}x_{i3} + w'_{41}x_{i4} + b_{11} \\ w'_{12}x_{i1} + w'_{22}x_{i2} + w'_{32}x_{i3} + w'_{42}x_{i4} + b_{12} \\ w'_{13}x_{i1} + w'_{23}x_{i2} + w'_{33}x_{i3} + w'_{43}x_{i4} + b_{13} \end{bmatrix} \right)$$

$$= \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix} \leftarrow a^{[1]}$$

Layer #2

$$\begin{bmatrix} w_{21}^2 & w_{22}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}^T \begin{bmatrix} O_{11} \\ O_{12} \\ O_{13} \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}_{2 \times 1}$$

$$= \sigma \left(\begin{bmatrix} w_{11}^2 O_{11} + w_{21}^2 O_{12} + w_{31}^2 O_{13} + b_{21} \\ w_{12}^2 O_{11} + w_{22}^2 O_{12} + w_{32}^2 O_{13} + b_{22} \end{bmatrix} \right) = \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix} \leftarrow a^{[2]}$$

Layer #3

$$\begin{bmatrix} w_{11}^3 & w_{21}^3 \\ w_{21}^3 & w_{31}^3 \end{bmatrix}^T \begin{bmatrix} O_{21} \\ O_{22} \end{bmatrix}_{2 \times 1} + b_{31}$$

$$= \sigma \left([w_{11}^3 O_{21} + w_{21}^3 O_{22} + b_{31}] \right) = \hat{y}_i \leftarrow a^{[3]}$$

$$a^{[1]} = \sigma(a^{[0]} w^{[1]} + b^{[1]})$$

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

$$a^{[3]} = \sigma(a^{[2]} w^{[3]} + b^{[3]})$$

$$\left(\left(\sigma \left(\left(\sigma \left(\left(a^{[0]} w^{[1]} + b^{[1]} \right) w^{[2]} + b^{[2]} \right) w^{[3]} + b^{[3]} \right) \right) \right) \right) \hat{y} = a^{[3]}$$

$a^{[1]}$

$a^{[2]}$

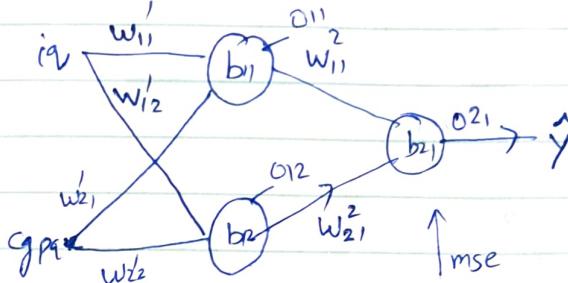
$a^{[3]} = \hat{y}$

Backpropagation

Date _____
Page _____

Backpropagation, short for "backward propagation of errors," is an algorithm for supervised learning of ANN using gradient descent. Given an ANN and an error function, the method calculates the gradient of the error function w.r.t. the neural network's weights.

i ₁	cgpa	lpa
80	8	3
60	9	5
70	5	8
120	7	11



$$O_{21} = w_{11}^2 O_{11} + w_{12}^2 O_{12} + b_{21}$$

9 trainable parameters

(w, b) Random value

w → 1, b → 0

Step 1: You selected a point (row)
↳ student

2. Predict (lpa) → forward propagation [Dot product]

3. Choose a loss function

4. Weight and bias update using

Gradient Descent

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$

$$b_{\text{new}} = b_{\text{old}} - \eta \frac{\partial L}{\partial b_{\text{old}}}$$

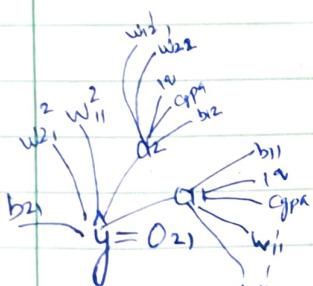
$$\text{eg } \hat{y}_1 = 18 \text{ lpa}$$

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

$$(3 - 18)^2 = 225 \text{ error}$$

$$\hat{y} \leftarrow 18 \downarrow$$

$$\hat{y} \rightarrow 19$$



derivative of Loss w.r.t. weight & bias

g derivative

$$\left| \begin{array}{l} \frac{\partial L}{\partial w_{11}^2}, \frac{\partial L}{\partial w_{12}^2}, \frac{\partial L}{\partial b_{21}}, \frac{\partial L}{\partial w_{11}}, \frac{\partial L}{\partial w_{12}}, \frac{\partial L}{\partial b_{11}}, \frac{\partial L}{\partial w_{21}}, \frac{\partial L}{\partial w_{22}}, \frac{\partial L}{\partial b_{12}} \end{array} \right|$$

$$w_{11\text{new}} = w_{11\text{old}} - \eta \frac{\partial L}{\partial w_{11}}$$

$$w_{21\text{new}} = w_{21\text{old}} - \eta \frac{\partial L}{\partial w_{21}}$$

$$b_{21\text{new}} = b_{21\text{old}} - \eta \frac{\partial L}{\partial b_{21}}$$

$$\min_{w_{11}} \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial g} \times \frac{\partial g}{\partial w_{11}^2} = \frac{\partial L}{\partial g} \times \frac{\partial [o_{11} w_{11}^2 + o_{12} w_{21}^2 + b_{21}]}{\partial w_{11}} = -2(y - \hat{y}) \times o_{11}$$

$$(ii) \frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial g} \times \frac{\partial g}{\partial w_{21}^2} = -2(y - \hat{y}) \times o_{12}$$

$$(iii) \frac{\partial L}{\partial b_{21}} = \frac{\partial L}{\partial g} \times \frac{\partial g}{\partial b_{21}} = -2(y - \hat{y}) \times 1 = -2(y - \hat{y})$$

$$\left| \begin{array}{l} \frac{\partial L}{\partial \hat{y}} = \frac{\partial (y - \hat{y})^2}{\partial \hat{y}} = -2(y - \hat{y}) \end{array} \right|$$

$$\hat{y} = w_{11}^2 o_{11} + w_{21}^2 o_{12} + b_{21}$$

$$\frac{\partial L}{\partial y} = -2(y - \hat{y})$$



$$(iv) \frac{\partial L}{\partial w_{11}} = \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{11}} \right) \cdot \frac{\partial o_{11}}{\partial w_{11}} \\ = -2(y - \hat{y}) \cdot w_{11}^2 \cdot x_{11}$$

$$(vii) \frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{12}} \cdot \frac{\partial o_{12}}{\partial w_{12}} \\ = -2(y - \hat{y}) \cdot w_{21}^2 \cdot x_{12}$$

$$(v) \frac{\partial L}{\partial w_{21}} = \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{11}} \right) \cdot \frac{\partial o_{11}}{\partial w_{21}} \\ = -2(y - \hat{y}) w_{11}^2 \cdot x_{12}$$

$$(viii) \frac{\partial L}{\partial w_{22}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{12}} \cdot \frac{\partial o_{12}}{\partial w_{22}} \\ = -2(y - \hat{y}) w_{21}^2 \cdot x_{12}$$

$$(ix) \frac{\partial L}{\partial b_{11}} = \left(\frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{11}} \right) \cdot \frac{\partial o_{11}}{\partial b_{11}} \\ = -2(y - \hat{y}) \cdot w_{11}^2 \cdot 1$$

$$(ix) \frac{\partial L}{\partial b_{12}} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial o_{12}} \cdot \frac{\partial o_{12}}{\partial b_{12}} \\ = -2(y - \hat{y}) w_{21}^2 \cdot 1$$

$$\therefore \frac{\partial o_{11}}{\partial w_{11}} = \frac{\partial [i_1 w_{11} + c_{pa} w_{21} + b_{11}]}{\partial w_{11}} = i_1 = x_{11}$$

$$\frac{\partial o_{12}}{\partial w_{12}} = \frac{\partial [i_2 w_{12} + c_{pa} w_{22} + b_{12}]}{\partial w_{12}} = i_2 = x_{12}$$

$$\frac{\partial o_{11}}{\partial w_{21}} = \frac{\partial [i_1 w_{11} + c_{pa} w_{21} + b_{11}]}{\partial w_{21}} = x_{12}$$

$$\frac{\partial o_{12}}{\partial w_{22}} = c_{pa} = x_{12}$$

$$\frac{\partial o_{11}}{\partial b_{11}} = 1$$

$$\frac{\partial o_{12}}{\partial b_{12}} = 1$$

Algorithm Steps

0> weights / bias \rightarrow initialization
 $w_{11} = \text{random}$
 $b = 0$

loop over epochs
 until converge
 1> for i in range(4):

1(a) 1 student \rightarrow forward prop \rightarrow predict (lpa)

1(b) Loss calculate (mse)

1(c) Adjust all weights & bias

$$w_{11} = w_{11} - \eta \frac{\partial L}{\partial w_{11}}$$

3 times

epoch = 5

for i in range(epoch):

for j in range(X.shape[0]):
 1(a), 1(b), 1(c)

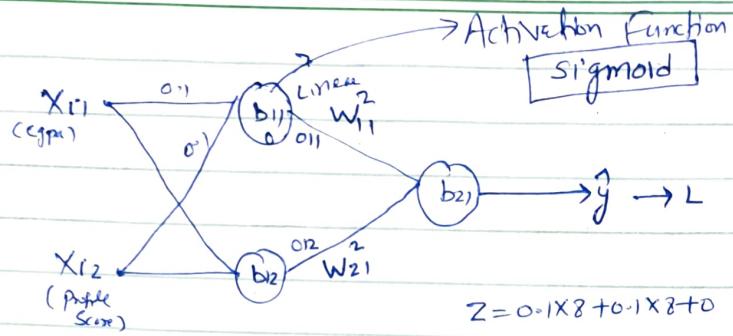
calc avg loss for the epoch

"You can never be overdressed or too well dressed."

Classification Example

Date _____
Page _____

Emp	Profile Score	Placement
8	8	1
7	9	1
6	10	0
5	5	0



$$Z = O_1 \cdot X_1 + O_1 \cdot X_2 + b_1$$

$$\sigma(Z) = O_{11}$$

$$\boxed{\text{Loss} = -Y \log(\hat{y}) - (1-Y) \log(1-\hat{y})} \quad \text{Binary cross entropy}$$

$$y = \sigma(Z)$$

$$(i) \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial w_{11}^2} = \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \times \hat{y} \times \frac{1}{Z} = \frac{-(y-\hat{y})}{\hat{y}(1-\hat{y})} \times \frac{1}{w_{11}^2} \rightarrow \hat{y} \rightarrow Z_{\text{final}}$$

$$= -(y-\hat{y}) O_{11} (w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_2)$$

$$(ii) \frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial w_{21}^2} = -(y-\hat{y}) O_{12}$$

$$(iii) \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial b_2} = -(y-\hat{y}) \otimes$$

$$\therefore \frac{\partial L}{\partial \hat{y}} = \frac{\partial}{\partial \hat{y}} [-Y \log(\hat{y}) - (1-Y) \log(1-\hat{y})] = \frac{-Y}{\hat{y}} + \frac{(1-Y)}{(1-\hat{y})} = \frac{-Y(1-\hat{y}) + \hat{y}(1-Y)}{\hat{y}(1-\hat{y})} = \frac{-Y + Y\hat{y} - \hat{y} + \hat{y}}{\hat{y}(1-\hat{y})}$$

$$\therefore \frac{\partial \hat{y}}{\partial Z} = \frac{\partial (\sigma(Z))}{\partial Z} = \sigma(Z)[1 - \sigma(Z)] = \hat{y}(1-\hat{y})$$

$$Z_F = w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_2$$

$$\frac{\partial Z}{\partial w_{11}^2} = O_{11}$$

$$(iv) \frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial w_{11}^2} = -(y-\hat{y}) \times w_{11}^2 \times O_{11} \times (1-O_{11}) \times X_{i1}$$

$$\hat{y} \rightarrow Z_F \rightarrow O_{11} \rightarrow Z_{\text{prev}}$$

$$(v) \frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial w_{21}^2} = -(y-\hat{y}) w_{21}^2 \times O_{11} \times (1-O_{11}) \times X_{i2}$$

$$\hat{y} = O_{11} = \sigma(Z_{\text{prev}})$$

$$O_{11} = \sigma(Z_{\text{prev}})[1 - \sigma(Z_{\text{prev}})]$$

$$= O_{11}[1 - O_{11}]$$

$$(vi) \frac{\partial L}{\partial b_2} = \frac{\partial L}{\partial \hat{y}} \cdot \frac{\partial \hat{y}}{\partial Z} \cdot \frac{\partial Z}{\partial b_2} = -(y-\hat{y}) w_{11}^2 O_{11} (1-O_{11}) \otimes 1$$

$$Z_{\text{prev}} = w_{11}^2 X_{i1} + w_{21}^2 X_{i2} + b_{11}$$

$$Z_f = w_{11}^2 O_{11} + w_{21}^2 O_{12} + b_{21}$$

$$Z_{\text{prev}} = w_{12}^1 X_{11} + w_{22}^1 X_{12} + b_{12}$$

Date _____
Page _____

$$(vii) \frac{\partial L}{\partial w_{12}} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial Z_f} \cdot \frac{\partial Z_f}{\partial O_{12}} \frac{\partial O_{12}}{\partial Z_{\text{prev}}} \frac{\partial Z_{\text{prev}}}{\partial w_{12}} = -(y - \hat{y}) w_{21}^2 O_{12} (1 - O_{12}) X_{11}$$

$$(viii) \frac{\partial L}{\partial w_{22}} = -(y - \hat{y}) w_{21}^2 O_{12} (1 - O_{12}) X_{12}$$

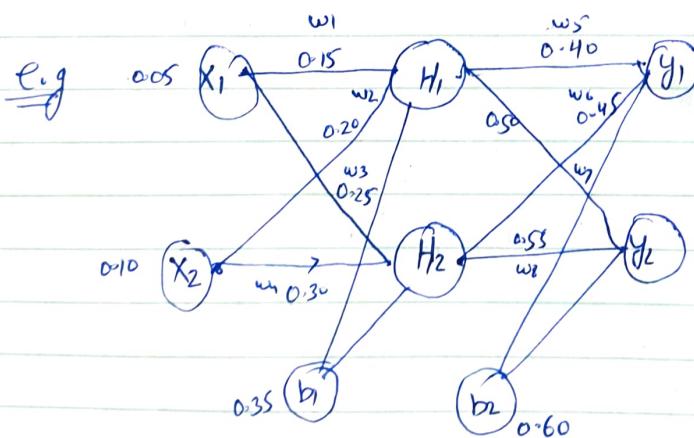
$$(ix) \frac{\partial L}{\partial b_{21}} = -(y - \hat{y}) w_{21}^2 O_{12} (1 - O_{12})$$

$\frac{\partial L}{\partial b}$ = rate of change of L w.r.t b

$$\frac{\partial L}{\partial b} = -ve \quad b \uparrow \quad L \downarrow$$

$$\frac{\partial L}{\partial b} = +ve \quad b \uparrow \quad L \uparrow$$

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial b}$$



Target Values
 $T1 = 0.01$
 $T2 = 0.99$

First calculate the values of H_1 & H_2 by a forward pass

$$H_1 = X_1 \times 0.15 + X_2 \times 0.20 + b_1 \\ = 0.05 \times 0.15 + 0.10 \times 0.20 + 0.35 = 0.3775$$

$$H_1^{\text{final}} = \sigma(H_1) = \frac{1}{1 + e^{-0.3775}} = 0.593269992$$

$$H_2 = 0.05 \times 0.25 + 0.10 \times 0.10 + 0.35 = 0.3925$$

$$H_2^{\text{final}} = \sigma(H_2) = \frac{1}{1 + e^{-0.3925}} = 0.596884378$$

$$y_1 = H_1 \times 0.40 + H_2 \times 0.45 + 0.60 \\ = 0.59326992 \times 0.40 + 0.59688 \times 0.45 + 0.60 = 1.10590597$$

$$\sigma(y_1) = y_1^{\text{final}} = \frac{1}{1 + e^{-1.105905}} = 0.75136507$$

$$y_2 = 0.59326992 \times 0.50 + 0.59688 \times 0.55 + 0.60 = 1.2249214$$

$$\sigma(y_2) = y_2^{\text{final}} = \frac{1}{1 + e^{-1.2249214}} = 0.77292845$$

$$E_{\text{total}} = \frac{1}{2} \sum (\text{target} - \text{output})^2 = \frac{1}{2} [0.01 - 0.751365]^2 + \frac{1}{2} [0.99 - 0.77292845]^2 \\ = 0.29837111$$

Now we will backpropagate this error to update the weights using a backward pass

$$\frac{\partial E}{\partial w_5} = \frac{\partial E}{\partial y_1^{\text{final}}} \times \frac{\partial y_1^{\text{final}}}{\partial y_1} \times \frac{\partial y_1}{\partial w_5} \\ = 0.74136507 \times 0.186815602 \times 0.59326992 \\ = 0.0821670407$$

$$\begin{cases} y_1 = w_5 H_{1, \text{final}} + w_6 H_2^{\text{final}} + b_2 \\ \frac{\partial y_1}{\partial w_5} = H_{1, \text{final}} = 0.59326992 \end{cases}$$

Similarly calc.

$$\frac{\partial E}{\partial w_6}, \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \frac{\partial E}{\partial w_3}, \frac{\partial E}{\partial w_4}$$

$$w_5^{\text{new}} = w_5^{\text{old}} - \eta \frac{\partial E}{\partial w_5} \quad \text{Let } \eta = 0.5 \\ = 0.4 - 0.5 \times 0.0821670 \\ = 0.35891648$$

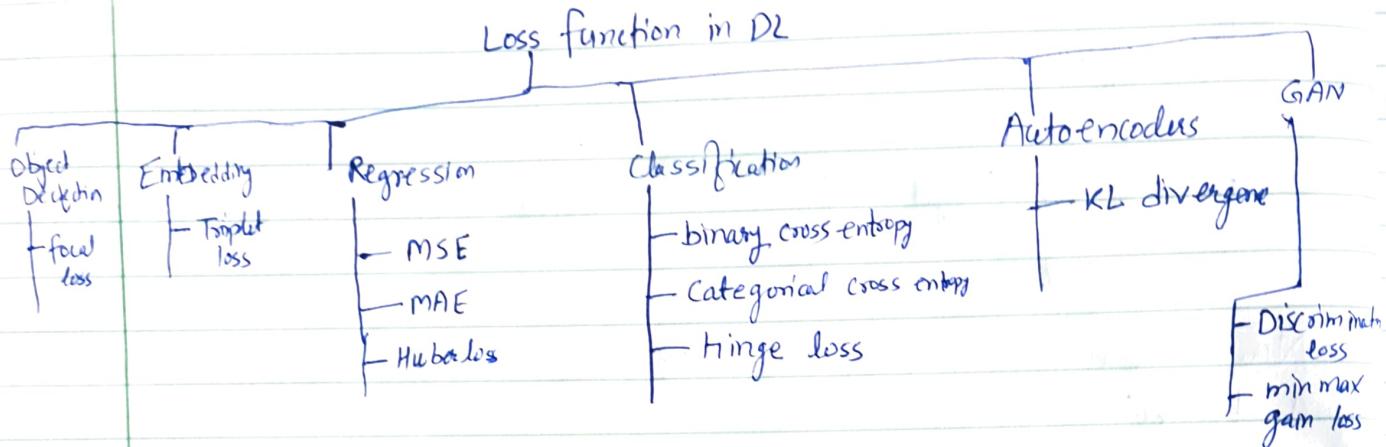
$$\frac{\partial y_1^{\text{final}}}{\partial y_1} = \frac{\partial \sigma(y_1)}{\partial y_1} = \sigma(y_1)E - \sigma(y_1) \\ = y_1^{\text{final}} [1 - y_1^{\text{final}}] \\ = 0.751365 [1 - 0.75136507] \\ = 0.186815602$$

$$\frac{\partial E}{\partial y_1^{\text{final}}} = \frac{1}{2} \sum (y_1 - \hat{y}_1)^2$$

$$= -(y_1 - \hat{y}_1) \\ = -(0.01 - 0.75136507) \\ = 0.74136507$$

LOSS Function in Deep Learning

Loss Function: Loss function is a method of evaluating how well your algorithm is modelling your dataset.



Loss function \rightarrow single training example
 Cost function \rightarrow Whole training example

$$\text{Loss function} = (y_i - \hat{y}_i)^2$$

$$\text{cost function} = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

1. Mean Squared Error (MSE)

squared loss L2 loss

$$= (y_i - \hat{y}_i)^2 = (\text{true} - \text{predicted})^2$$

$$CF = \frac{1}{n} \sum (y_i - \hat{y}_i)^2$$

quadratic

outliers

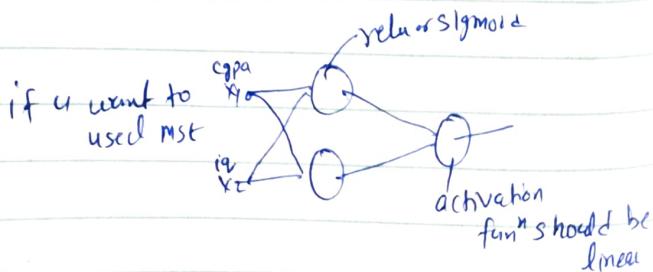


Advantages:

1. Easy to interpret
2. Differentiable (CD)
3. 1 local minima

Disadvantages:

1. Error unit (squared) — different
2. Not Robust to outliers



2 Mean Absolute Error (MAE) \rightarrow L1 loss

$$L = |y_i - \hat{y}_i|$$

$$C = \frac{1}{n} \sum |y_i - \hat{y}_i|$$

Advantage

1. Intuitive and easy
2. Unit is same
3. Robust to outliers

Disadvantage

1. Not differentiable. We cannot apply gradient descent

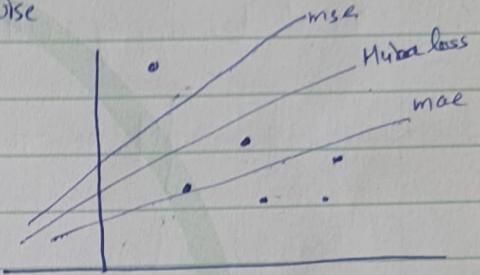
Huber Loss

$$L = \begin{cases} \frac{1}{2} (y - \hat{y})^2 & \text{for } |y - \hat{y}| \leq \delta \\ \delta |y - \hat{y}| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

δ = hyperparameter

It lies b/w MSE & MAE

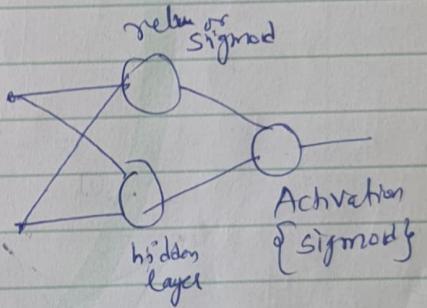
MSE - outlier impact
MAE - normal points

or Log loss

4. Binary Cross Entropy used for classification

$$\text{Loss function} = -y \log \hat{y} - (1-y) \log(1-\hat{y})$$

$$\text{Cost function} = -\frac{1}{n} \left[\sum_{i=1}^n y_i \log \hat{y}_i + (1-y_i) \log(1-\hat{y}_i) \right]$$

Advantage

1. Differentiable

Disadvantage

1. Multiple local minimum
2. Intuitive (not)

5. Categorical Cross Entropy [Used in Softmax Regression]

[Multi-class] {Classification}

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

where k is # classes in the data

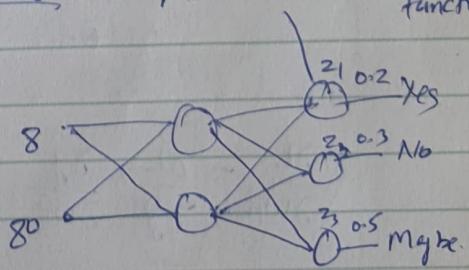
$$L = -y_1 \log \hat{y}_1 - y_2 \log \hat{y}_2 - y_3 \log \hat{y}_3$$

$$\begin{bmatrix} 0.2 & 0.3 & 0.5 \\ 1 & 0 & 0 \end{bmatrix}$$

$$L = -1 \times \log(0.2)$$

$$\begin{bmatrix} 0.3 & 0.6 & 0.1 \\ 0 & 1 & 0 \end{bmatrix} \quad L = -1 \log(0.6)$$

Activation function = Softmax function



CGPA	IQ	Placed
8	80	Yes 1
6	60	No 2
7	70	Maybe 3

ON E mg

6.

Sparse Categorical Cross Entropy

assigning integer value to classes

$$\{0.1 \ 0.4 \ 0.5\}$$

$$L = -1 \times \log(0.1)$$

It will be fast. Because we don't need to apply OHE.

$$L = - \sum_{j=1}^k y_j \log(\hat{y}_j)$$

$$C = -\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^k y_{ij} \log(\hat{y}_{ij})$$

OHE		Sparse CCE	
7	70	Yes	1
8	80	No	2
6	60	Maybe	3

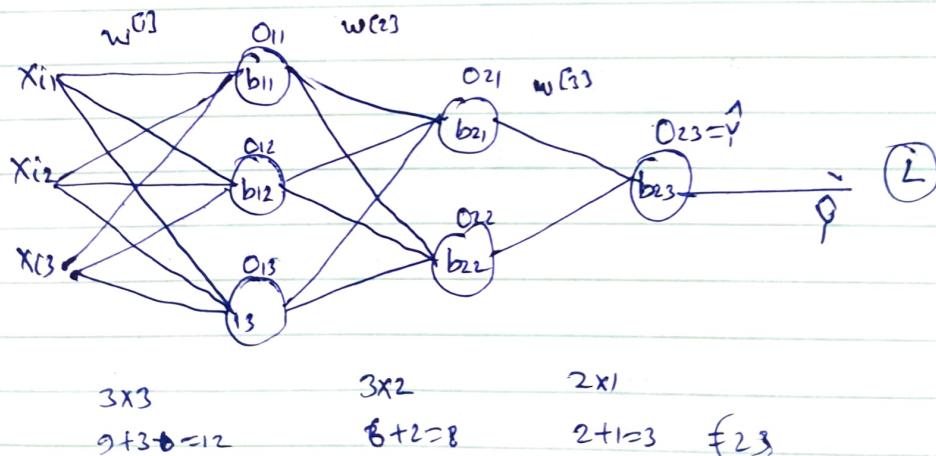
$$\{0.1 \ 0.4 \ 0.5\} \text{ if second class}$$

$$L = -1 \times \log(0.4)$$

3rd

$$\{0.6 \ 0.12 \ 0.2\}$$

$$L = -1 \times \log(0.2)$$

MLP Memoization

$$\frac{\partial L}{\partial w_{11}^3} = \left(\frac{\partial L}{\partial g} \right) \times \frac{\partial g}{\partial w_{11}^3}$$

$$\frac{\partial L}{\partial w_{11}^2} = \left(\frac{\partial L}{\partial g} \right) \times \left(\frac{\partial g}{\partial o_{21}} \right) \times \frac{\partial o_{21}}{\partial w_{11}^2}$$

$$L \rightarrow \hat{y} \rightarrow w_{11}^3$$

$$L \rightarrow \hat{y} \rightarrow o_{21} \rightarrow w_{21}^2$$

$f(x)$ $g(x)$

$$h(f(x), g(x)) = \frac{\partial h}{\partial x} \left[\frac{\partial h}{\partial f(x)} \times \frac{\partial f(x)}{\partial x} + \frac{\partial h}{\partial g(x)} \times \frac{\partial g(x)}{\partial x} \right]$$

$$\frac{\partial L}{\partial w_{ii}} = \frac{\partial L}{\partial j} \left[\frac{\partial j}{\partial o_{21}} \times \frac{\partial o_{21}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{ii}} + \frac{\partial j}{\partial o_{22}} \times \frac{\partial o_{22}}{\partial o_{11}} \times \frac{\partial o_{11}}{\partial w_{ii}} \right]$$

for larger hidden layer it will be more complex.

∴ We will save the intermediate derivative $\frac{\partial L}{\partial j}$ in memory for further processing.

Keras
TF

Backpropagation → Chain diff rule + memoization

Gradient Descent in Neural Network

Batch GD (Vanilla GD) → entire dataset → update

Stochastic GD → For each row → update
help the algo to move out of local minima

Mini batch GD

⇒ which is faster (given same no. of epochs) Batch GD

⇒ which is the faster to converge (given same # epochs) Stochastic

Vectorization :- epoch=10 Batch GD

for i in range(10):
 $y = np.dot(x, w) + b$
dot → smart replacement → loops
↓
Vectorization → faster than loops

Fast
 SGD > MBGD > SGD

Convergence SGD < MBGD < SGD

Vanishing Gradient Problem in ANN

→ As more layers using certain activation function are added to neural networks, the gradients of the loss function approaches zero, making the network hard to train

When there are more layers in the network, the value of the product of derivative decreases until at some point the partial derivative of the loss function approaches a value close to zero, and the partial derivative vanishes. We call this the vanishing gradient Problem.

Exploding Gradient Problem: - If the gradients are large, the multiplication of these gradients will become huge over time. This results in the model being unable to learn & its behaviour becomes unstable. This problem is called the exploding gradient problem.

Solution of Vanishing Gradient Problem

- Change Activation function (use ReLU)
- Reduce model complexity (i.e. reduce hidden layers)
- Proper weight initialization
- Batch Normalization
- Use Residual N/w

Solution of Exploding Gradient

- Use Batch Normalization
- Use less no. of layers
- Carefully initialize weights
- Use gradient clipping

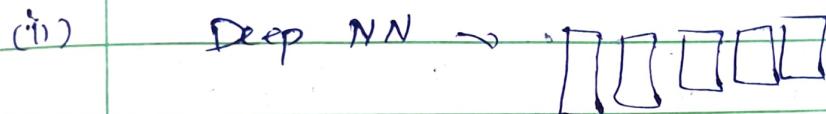
Conclusion:

During the backpropagation in the deep neural networks, the vanishing gradient problem occurs due to the sigmoid & tan activation function & the exploding gradient problem occurs due to large weights.

Vanishing Gradient Problem: In ML, the VGP is encountered when training ANN with Gradient based learning methods & backpropagation. In such methods, during each iteration of training each of the NN weights receives an update proportional to the partial derivatives of the error function with ~~at~~ the current weight. The problem is that in some cases, the gradient will be vanishingly

small, effectively preventing the weight from changing its value
In the worst case, this may completely stop the NN from further training.

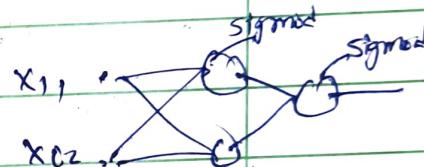
(i) $0.1 \times 0.1 \times 0.1 \times 0.1 = 0.0001 \leftarrow \text{VGP}$



(iii) Sigmoid / tanh \rightarrow Activation function $\boxed{\text{VGP occur}}$

$$w_n = w_0 - \eta \frac{\partial L}{\partial w}$$

\rightarrow Sigmoid derivative



$$\frac{\partial L}{\partial w_{11}} = \frac{\partial L}{\partial y} \times \frac{\partial y}{\partial z_2} \times \frac{\partial z_2}{\partial a_1} \times \frac{\partial a_1}{\partial w_{11}}$$

$\frac{\partial a_1}{\partial w_{11}}$ This will be very small

if $w_0 = 1$
 $w_n = 1 - 0.01 \times 0.0001$
 $w_n = 0.9999$

$$\frac{\partial L}{\partial w_{11}} = 0.00001$$

\rightarrow Sigmoid
 $0 \rightarrow 1$ Input space
 $(0-1) \rightarrow (0-1)$ Output space

Change is not occurs
 $(1-0.9999)$

→ How to recognize the Vanishing Gradient Problem Occur

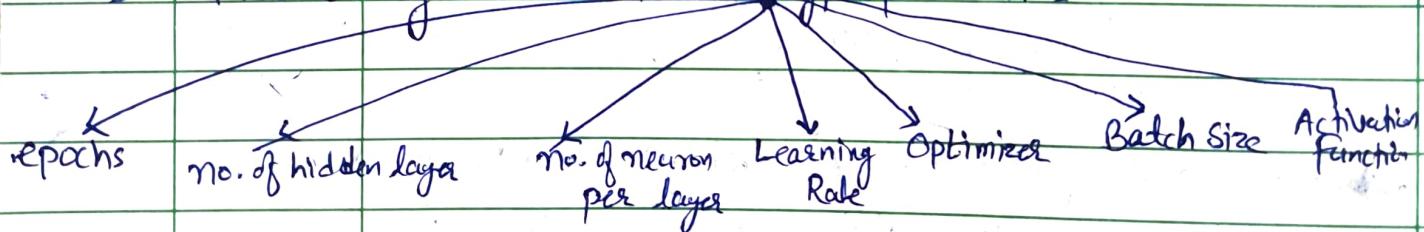
(i) Loss function change not occurring then VGP

(ii) Draw graph b/w weights vs epoch.



HOW TO IMPROVE THE PERFORMANCE OF A NEURAL NETWORK

1. Fine tuning Neural network hyper parameters



(2) By Solving problems — Vanishing/Exploding gradient Problem

