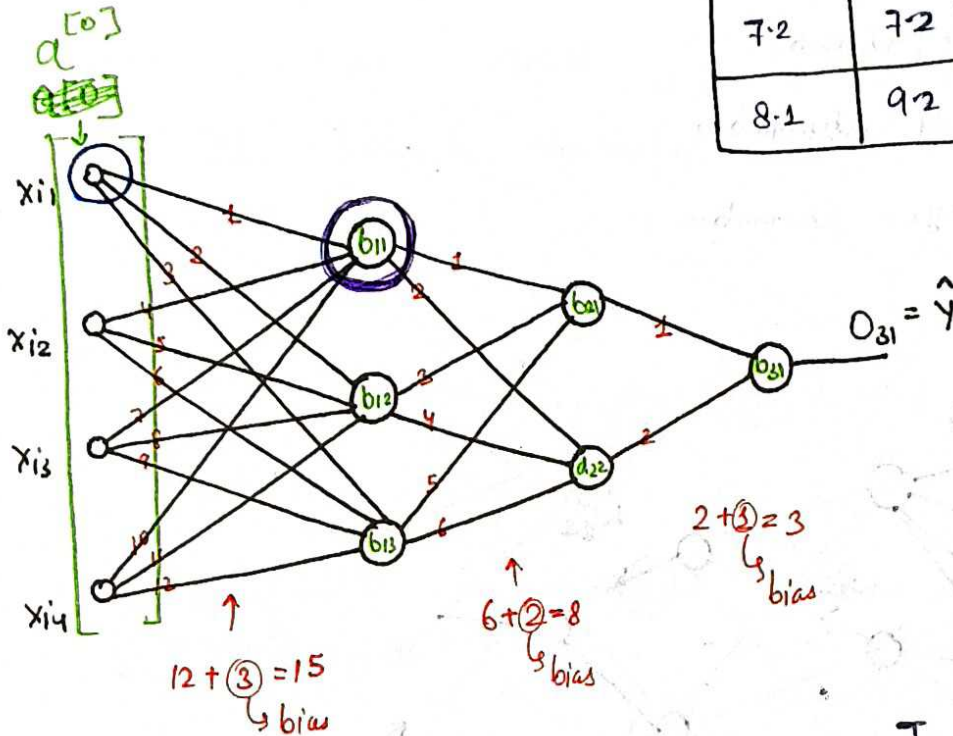


Forward Propagation

gpa	iq	10 th m	12 th m	plaud
7.2	72	69	81	1
8.1	92	75	76	0

of trainable parameters



prediction $\rightarrow \sigma(w^T x + b)$

layer #1

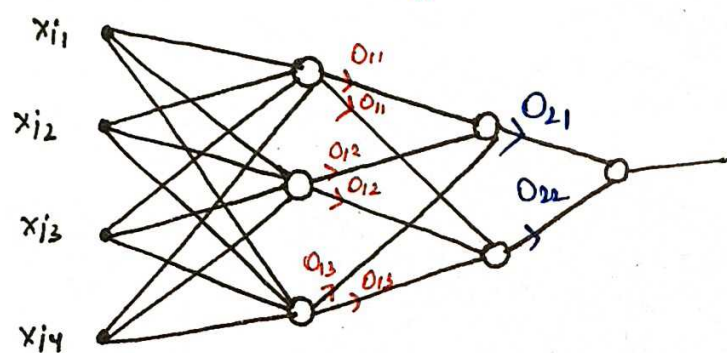
$$\begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix}^T \cdot \begin{bmatrix} x_{11} \\ x_{12} \\ x_{13} \\ x_{14} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

$4 \times 3^T \rightarrow 3 \times 4$ 4×1 3×1

$$= \begin{bmatrix} w_{11}' x_{11} + w_{21}' x_{12} + w_{31}' x_{13} + w_{41}' x_{14} \\ w_{12}' x_{11} + w_{22}' x_{12} + w_{32}' x_{13} + w_{42}' x_{14} \\ w_{13}' x_{11} + w_{23}' x_{12} + w_{33}' x_{13} + w_{43}' x_{14} \end{bmatrix} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}$$

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

$$= \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix} \rightarrow \begin{matrix} a[1] \\ a \end{matrix}$$



Layer #2

$$\begin{bmatrix} w_{11}^2 & w_{12}^2 \\ w_{21}^2 & w_{22}^2 \\ w_{31}^2 & w_{32}^2 \end{bmatrix}^T \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix} + \begin{bmatrix} b_{21} \\ b_{22} \end{bmatrix}$$

$3 \times 2 \xrightarrow{T} 2 \times 3 \quad \quad 3 \times 1 \quad \quad 2 \times 1$

$$= \sigma \begin{bmatrix} w_{11}^2 o_{11} + w_{21}^2 o_{12} + w_{31}^2 o_{13} + b_{21} \\ w_{12}^2 o_{11} + w_{22}^2 o_{12} + w_{32}^2 o_{13} + b_{22} \end{bmatrix} = \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix} \rightarrow \begin{matrix} a[2] \\ a \end{matrix}$$

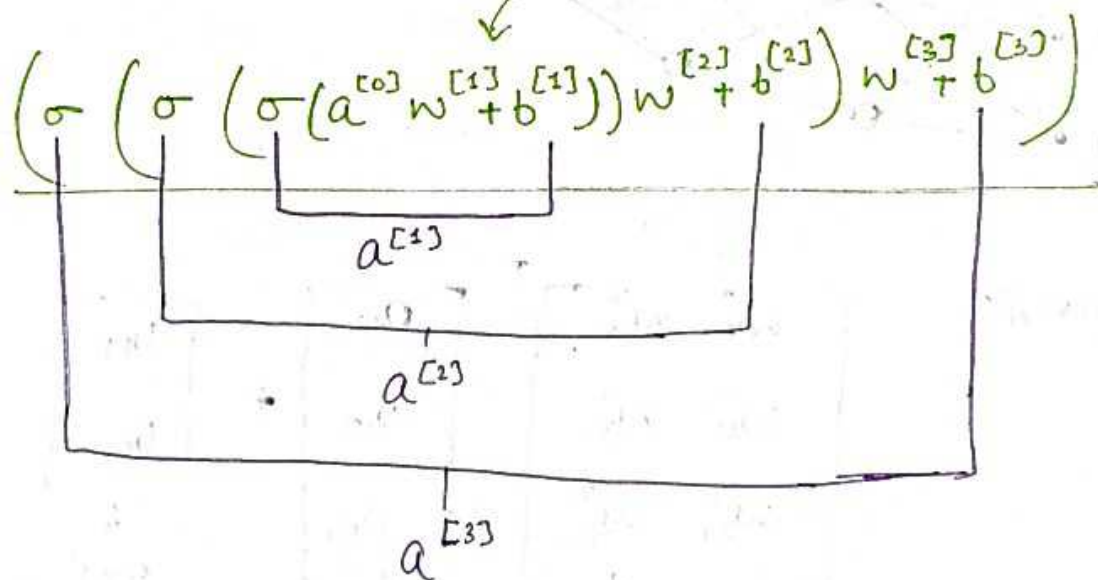
Layer #3

$$\begin{bmatrix} w_{11}^3 \\ w_{21}^3 \end{bmatrix}^T \begin{bmatrix} o_{21} \\ o_{22} \end{bmatrix} + \begin{bmatrix} b_{31} \end{bmatrix}$$

$2 \times 1 \xrightarrow{T} 1 \times 2 \quad \quad 2 \times 1 \quad \quad 1 \times 1$

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

$$\left. \begin{aligned} a^{[1]} &= \sigma(a^{[0]} w^{[1]} + b^{[1]}) \\ a^{[2]} &= \sigma(a^{[1]} w^{[2]} + b^{[2]}) \\ a^{[3]} &= \sigma(a^{[2]} w^{[3]} + b^{[3]}) \end{aligned} \right\} \text{Simple method}$$



What is Loss function?

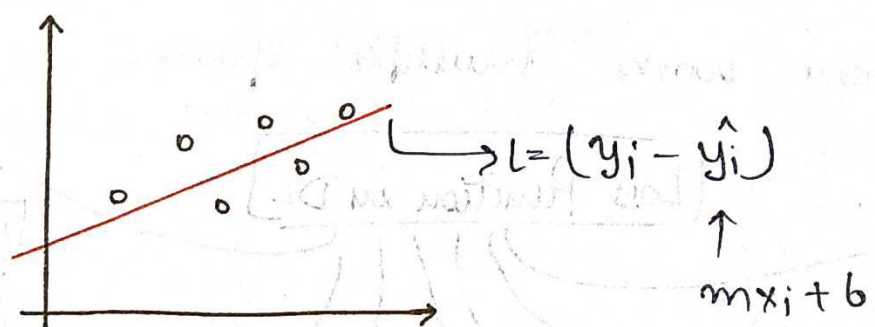
Loss function is a method of evaluating how well your algo is modelling your dataset.

Loss function \rightarrow higher $\xrightarrow{\text{Model}}$ poor
 \searrow low $\xrightarrow{\text{Model}}$ great

Loss function \rightarrow Mathematical function

$L(\text{parameter}) \rightarrow$ parameter change \rightarrow loss function's value change

for eg:



$L(m, b) = \min \rightarrow$ Find value of m, b where loss func is minimum

Why is loss function important?

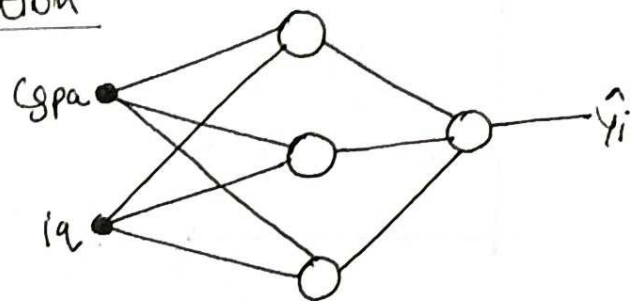
[You can't improve what you can't measure]

\hookrightarrow Peter Drucker

Loss function Vs Cost function

cgpa	iq	y_i packg	\hat{y}_i Prediction
6.3	100	6.3	6.1
7.1	91	4.1	4
8.5	83	3.5	3.7
9.2	102	7.2	7

batch



Loss function \rightarrow single row

$y_i = 6.3$ $\hat{y}_i = 6.1$ $(y_i - \hat{y}_i)^2$

Cost function

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

$\frac{1}{4} [(6.1 - 6.3)^2 + (4.1 - 4)^2 + \dots] = \text{Cost function}$

1. Mean Squared Error (MSE) Squared loss or L2 loss

$\rightarrow (y_i - \hat{y}_i)$

$(\text{true} - \text{predict})^2$

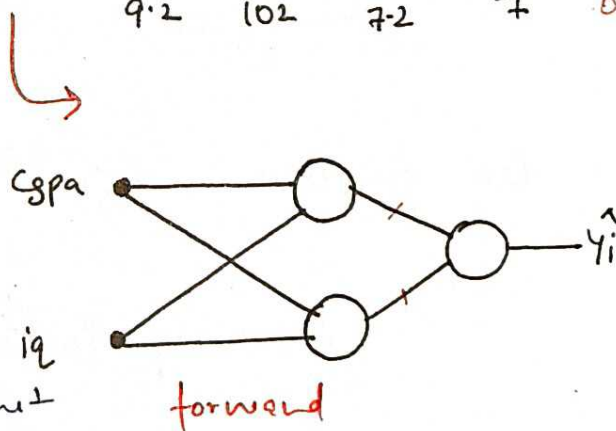
$(6.3 - 6.1)^2 = \dots$

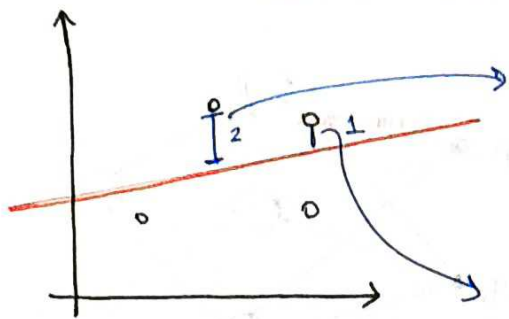
Why we use square?

Cz if error is -ve then we cannot able add so we use square \rightarrow we have +ve but problem?

$(\text{true} - \text{predict}) \rightarrow 1 \text{ unit} \xrightarrow{\text{square}} 1 \text{ unit}$
 $2 \text{ unit} \xrightarrow{\text{square}} 4 \text{ unit}$
 $4 \text{ unit} \xrightarrow{\text{square}} 16 \text{ unit}$

cgpa	iq	y_i packg	\hat{y}_i Pred	$y_i - \hat{y}_i$
6.3	100	6.3	6.1	0.2
7.1	91	4.1	4	0.1
8.5	83	3.5	3.7	-0.2
9.2	102	7.2	7	0.2





error \rightarrow heavy value change in loss function

$2^2 = 4 \rightarrow$ means double the error

$1^2 = 1 \rightarrow$ means double the error 1 but small change in loss function.

\rightarrow When large error \rightarrow double large error

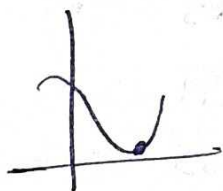
Advantage

1. Easy to interpret
2. This loss function always differentiable (40)
3. 1 local minima

$(w, b) \rightarrow$ one value \rightarrow

minimize

for low loss function \rightarrow global min



Disadvantage

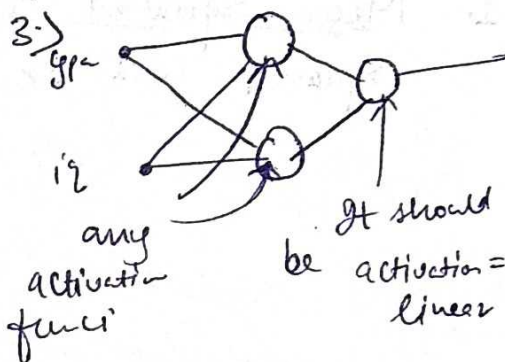
1. Error unit \rightarrow Squared

$lpa \rightarrow$ loss function \rightarrow

$(lpa)^2$

unit \rightarrow diff

2. \rightarrow not Robust to outlier



Cost function of MSB $= \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

{Batch of points} \rightarrow

2. Mean Absolute Error (MAE) \rightarrow L1 Loss

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

$$C = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

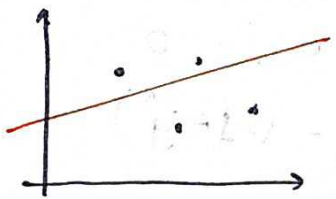
Advantage

1) Intuitive and easy

2) Unit same

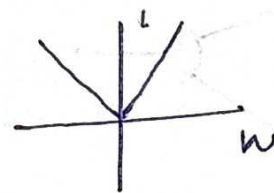
$$\begin{array}{ccc} y_i & \hat{y}_i & = |y_i - \hat{y}_i| \\ \downarrow & & \downarrow \\ \text{Lapa} & \text{LPA} & \text{LPA} \end{array}$$

3) Robust to outliers
 $\odot \rightarrow$ outlier



disadvantage

1) Not differentiable
 \hookrightarrow QD not differentiable
 then we \leftarrow have to use sub QD \rightarrow little bit tricky

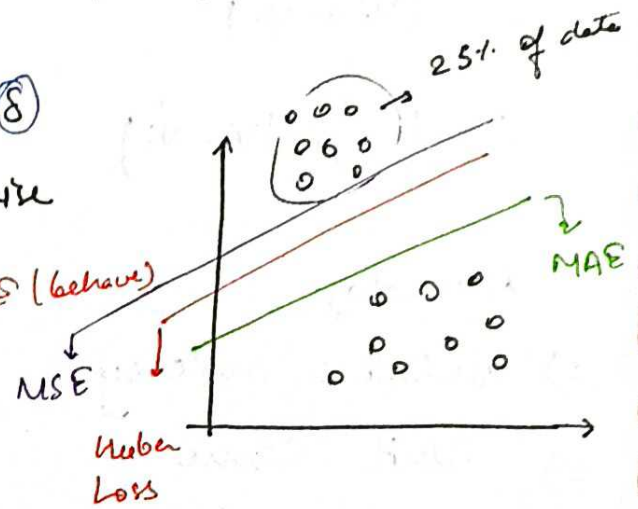


3. Huber loss

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

δ hyperparameter Data Normalized \rightarrow MSE (behaves)
 \rightarrow data outlier \rightarrow MAE (behaves)

This loss func use when data contain many outlier.



4. Binary Cross Entropy

\rightarrow Classification

\rightarrow Two classes

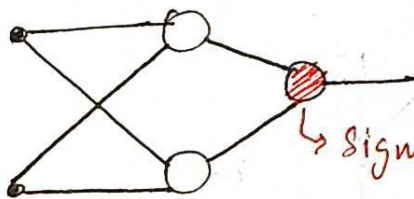
Gpa	iq	Placem ⁺
8	80	1
7	70	0
6	60	0

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

$y \rightarrow$ actual value / target

$\hat{y} \rightarrow$ NN prediction

Condition 3



\rightarrow Sigmoid (Activation)

when we use binary cross entropy loss func then output node activation func always be sigmoid.

Cost function:

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

5. Categorical Cross Entropy [Used in Softmax Regression] (18)

→ Multi-class classification

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

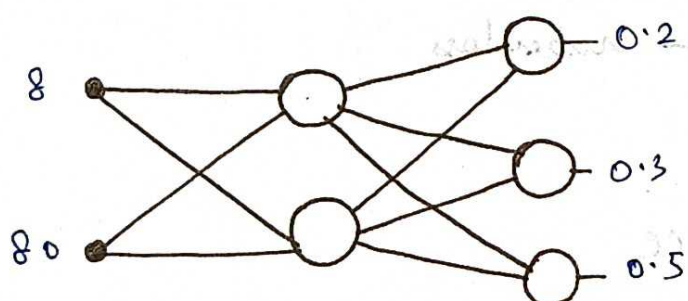
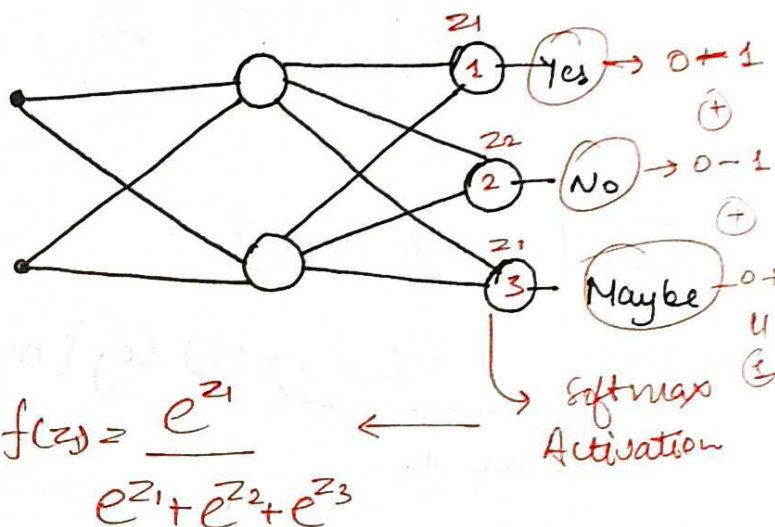
where k is # classes in the data

Maybe $\rightarrow \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$

No $\rightarrow \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$

Gpa	iq	placed
8	80	Yes
6	60	No
7	70	Maybe

Yes	No	Maybe
1	0	0
0	1	0
0	0	1



$$y = [1 \ 0 \ 0]$$

$$\hat{y} = [0.2 \ 0.3 \ 0.5]$$

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

$$L = -1 \log(0.2) - 0 - 0$$

update weight and biased using gradient descent

Sparse categorical cross entropy

All method same but lol diff

$$\hat{y} = [0.1 \ 0.4 \ 0.5] \quad y = 1$$

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

Gpa	iq	placed	OHE
7	70	Yes	(1)
8	80	No	(2)
6	60	Maybe	3

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

Choose only first because original y is 1

* If original y is 2 then we choose $-y_2 \log(\hat{y}_2)$

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

$$\hat{y} = [0.1 \quad 0.2 \quad 0.6]$$

$$y = [2] \quad L = -(1) \log(0.2)$$

we take 1.

Reg \rightarrow mse
 \rightarrow outlier - mae } Huber loss

Classifi \rightarrow binary - bce
 \rightarrow multi - cce
 \rightarrow SLE

*** Why combination of binary cross entropy and sigmoid function work well?

1. The sigmoid function's output fits perfectly into the BCE loss function, which expects a prob value for its calculation.
2. The gradients of the BCE loss function with respect to the model parameters are well-behaved when used with the sigmoid activation, facilitating effective training.

3. The output of the sigmoid function directly represents the probability of the positive class, making it intuitive to interpret the results.