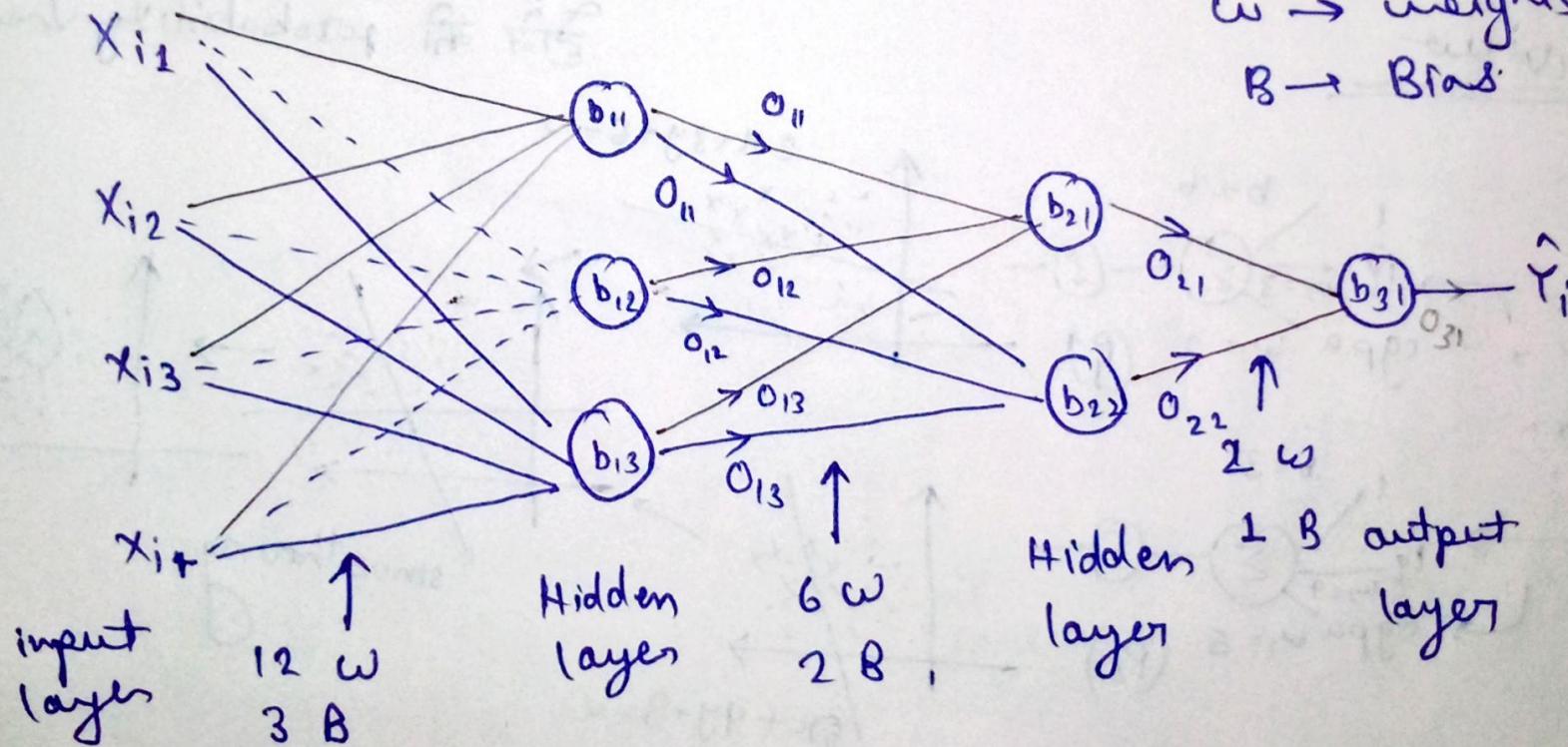


# Multi layer Perceptron (MLP) Notation:

$x_{i1} \rightarrow i = \text{rows}$   $1 = \text{column}$ .

$w \rightarrow \text{weights}$

$B \rightarrow \text{Bias}$



$12+3+6+2+2+1 = 26$  Trainable parameters

### Bias

$b_{ij}$   
i = layer  
j = node

### output

$O_{ij}$   
i = layer  
j = node

### weight

$w_{ijk}$

layer  
PN CN  
PN - previous node  
CN - current node  
i = पिछे layer के किस node से जा रहा है  
j = आने वाले layer के किस node से जा रहा है  
k = किस layer में जा रहा है

### Multilayer Perception (MLP):

#### Sigmoid function

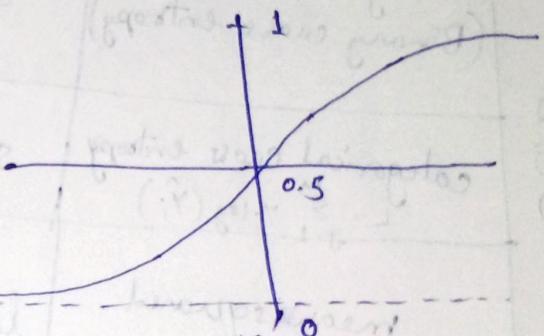
$$\sigma(z) = \frac{1}{1 + e^{-z}} \quad (f)$$

loss function - log loss

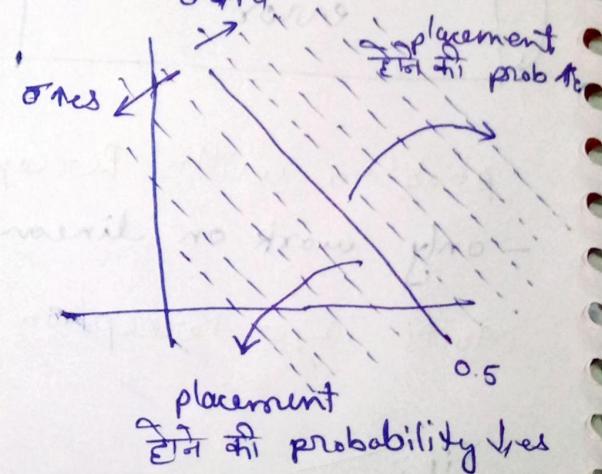
$$z = \sum k x_i X_i$$

If  $z > 0$  then  $\sigma(z) > 0.5$

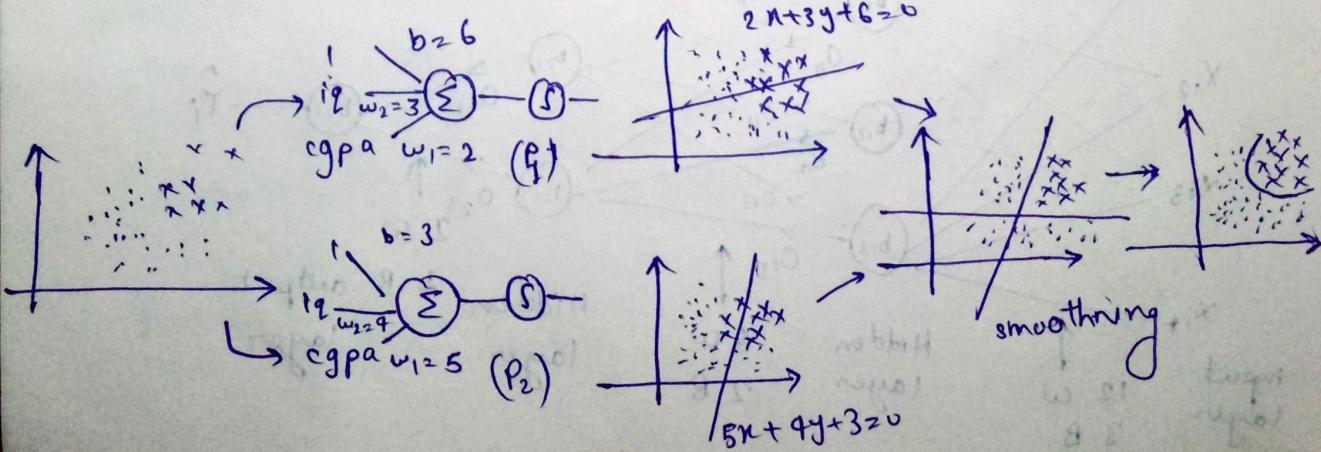
If  $z < 0$  then  $\sigma(z) < 0.5$



If  $\sigma = 0.5$  mean placement होने की probability is 0.5



### Overview

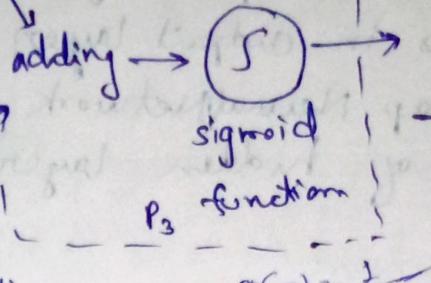


$P_1 \rightarrow$  perception 1  
 $P_2 \rightarrow$  perception 2  
 $P_3 \rightarrow$  perception 3

find  $p(Y)$  for each point  $(P_1)$

find  $p(Y)$  for each point  $(P_2)$

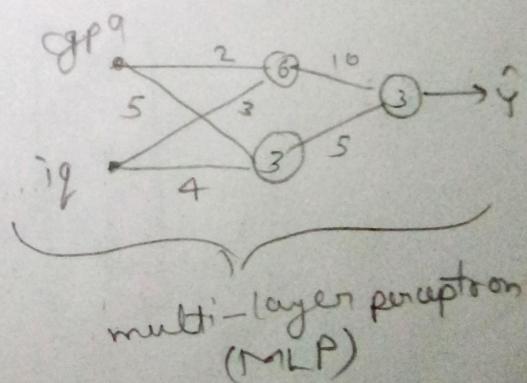
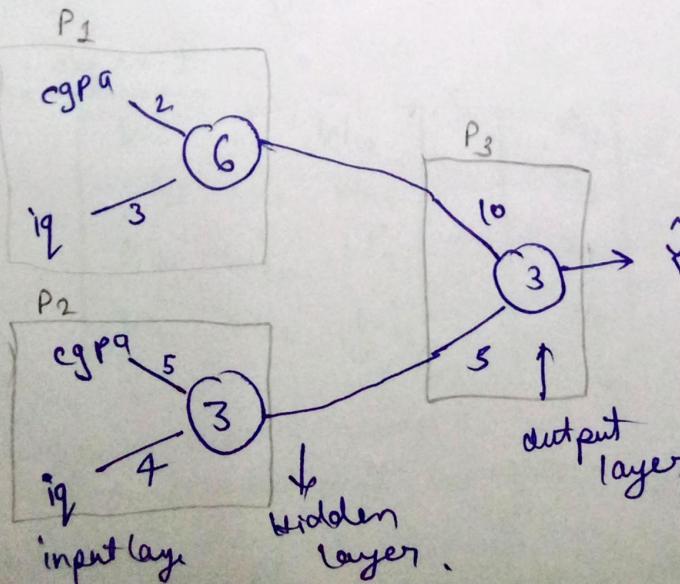
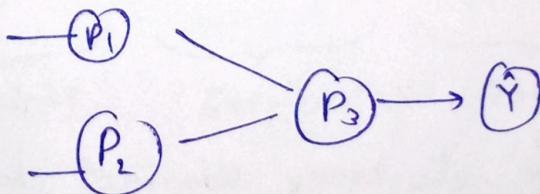
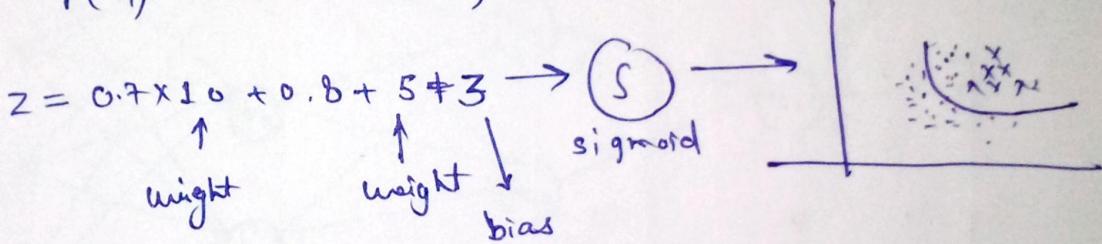
linear combination of two perception.



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

\*Here we can dominate the one perception by other using the weighted of ~~the~~. and also can be add bias

$$p(Y_1) = 0.7 \quad p(Y_2) = 0.8 \quad b = 3$$



How to change the architecture of Neural Network?

- add nodes in hidden layer.
- add nodes in input layer (visualise in 3D)
- add nodes in output layer. (cat, dog classification)
- using Deep Neural Network (increase the number of hidden layers).

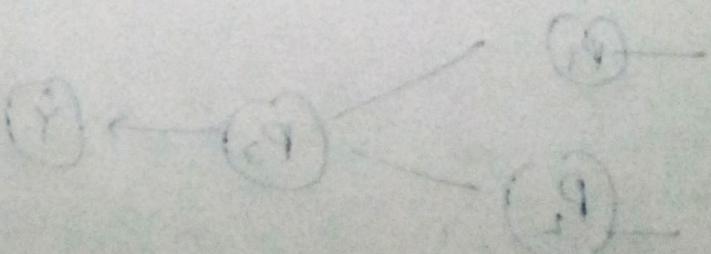
used for multiclass classification  
example predict in images cat, dog, human

\* Neural network also called universal function approximators.

$$\hat{y} = \sigma(\theta^T x + b)$$

$$\hat{y}_1 = \sigma(\theta_1^T x + b_1)$$

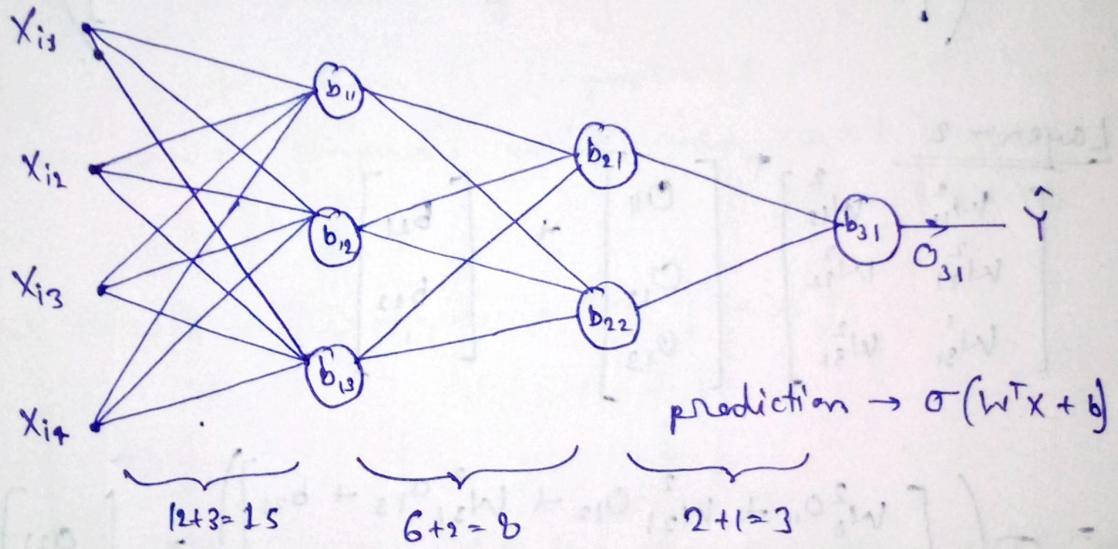
$$\hat{y}_2 = \sigma(\theta_2^T x + b_2)$$



# FORWARD PROPAGATION

Forward Propagation:

c g p a	i q	10 <sup>th</sup>	12 <sup>th</sup>	placed	# for trainable parameters
7.2	72	69	81	1	
8.1	92	75	76	0	



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

# Layer-1

$$\begin{bmatrix}
 w'_{11} & w'_{12} & w'_{13} \\
 w'_{21} & w'_{22} & w'_{23} \\
 w'_{31} & w'_{32} & w'_{33} \\
 w'_{41} & w'_{42} & w'_{43}
 \end{bmatrix}^T
 \begin{bmatrix}
 x_{i1} \\
 x_{i2} \\
 x_{i3} \\
 x_{i4}
 \end{bmatrix}_{4 \times 1} +
 \begin{bmatrix}
 b'_{11} \\
 b'_{12} \\
 b'_{13}
 \end{bmatrix}_{3 \times 1}$$

rows  $\rightarrow$  input weights  
 columns  $\rightarrow$  output weights

$$\begin{bmatrix} w_{11}'x_{i1} + w_{12}'x_{i2} + w_{13}'x_{i3} + w_{14}'x_{i4} \\ w_{12}'x_{i1} + w_{22}'x_{i2} + w_{23}'x_{i3} + w_{24}'x_{i4} \\ w_{13}'x_{i1} + w_{23}'x_{i2} + w_{33}'x_{i3} + w_{43}'x_{i4} \end{bmatrix}_{3 \times 1} + \begin{bmatrix} b_{11} \\ b_{12} \\ b_{13} \end{bmatrix}_{3 \times 1}$$

$$= \sigma \left( \begin{bmatrix} w_{11}'x_{i1} + w_{12}'x_{i2} + w_{13}'x_{i3} + w_{14}'x_{i4} + b_{11} \\ w_{12}'x_{i1} + w_{22}'x_{i2} + w_{23}'x_{i3} + w_{24}'x_{i4} + b_{12} \\ w_{13}'x_{i1} + w_{23}'x_{i2} + w_{33}'x_{i3} + w_{43}'x_{i4} + b_{13} \end{bmatrix} \right) = \begin{bmatrix} o_{11} \\ o_{12} \\ o_{13} \end{bmatrix}$$

# Layer - 2

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

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

# Layer - 3

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

$$\sigma \left( \begin{bmatrix} w_{11}^3 o_{21} + w_{21}^3 o_{22} + b_{31} \end{bmatrix} \right) = \hat{Y}_1 = o_{31}$$

$$\sigma(a^{[0]} W^{[1]} + b^{[1]}) = a^{[1]} \Big|_{a^{[2]} = \sigma(a^{[1]} W^{[2]} + b^{[2]})} \Big|_{a^{[3]} = \sigma(a^{[2]} W^{[3]} + b^{[3]})}$$

# LOSS FUNCTIONS

what is loss function:

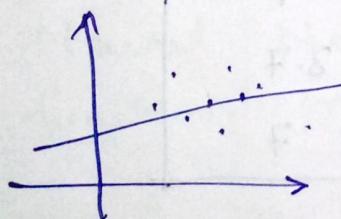
loss function is a method of evaluating how well your algorithm is modelling your dataset.

High  $\rightarrow$  poor } algo performance.  
small  $\rightarrow$  great }

why loss function is important?

"you can't improve what you can't measure"  
Peter Drucker

\* loss function example in ML (MSE)



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

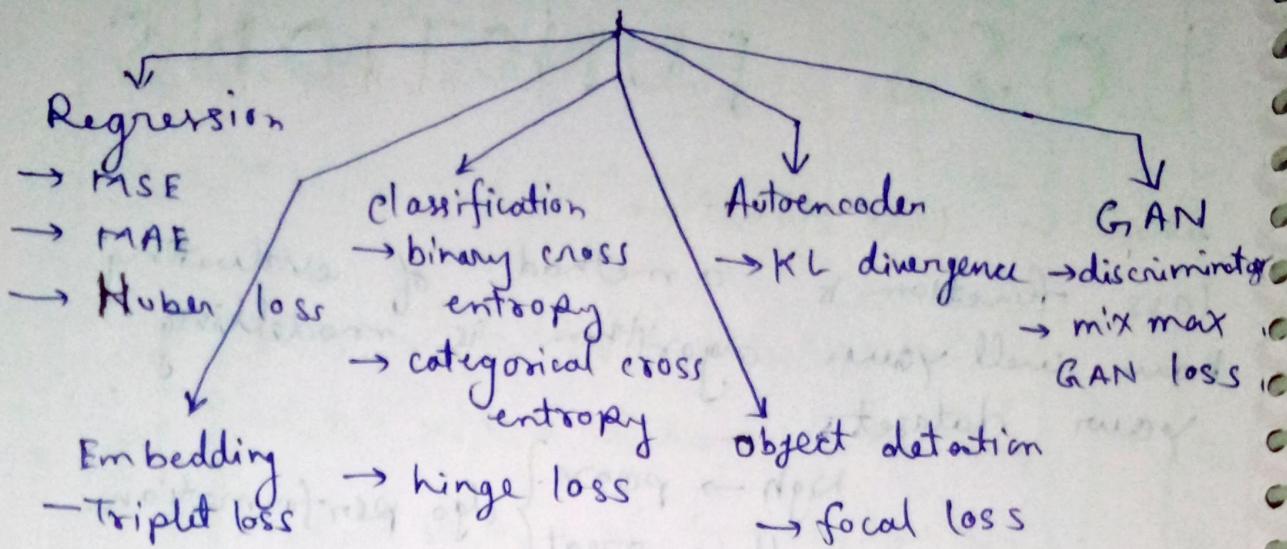
$\downarrow$   
 $m, x; + b$

calculate the combination of  $m$  and  $b$  such that  $L(m, b) = \text{minimum}$ .

for minimizing we use the gradient descent.

"loss function eye of ML algorithm"

## Loss functions in DL



Loss function vs Cost function:

→ loss function single training example  $\rightarrow$  calculate  $\sum \frac{(y_i - \hat{y}_i)^2}{n}$

cgpa	iq	package	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

$$\text{here } y_i = 8.3 \quad \hat{y}_i = 6.1$$

$$L = (y_i - \hat{y}_i)^2 = (8.3 - 6.1)^2 = 0.04$$

→ cost function ~~for~~ training data  $\rightarrow$  calculate  $\sum \frac{(y_i - \hat{y}_i)^2}{n}$

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

$$C = \frac{1}{4} [(6.1 - 6.3)^2 + (9.1 - 4)^2 + (3.5 - 3.7)^2 + (7.2 - 7)^2] =$$

# ① Mean Squared Error (MSE, squared loss, L<sub>2</sub> loss)

$$MSE = (y_i - \hat{y}_i)^2 \quad \text{cost fn} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$MSE = (\text{true value} - \text{predicted value})^2$

$$\text{Step.1} = (6.3 - 6.1)^2 = 0.04$$

for each row to minimize  $\left\{ \begin{array}{l} \text{Step.1} = \text{update the weight and bias} \\ \text{Step.2} = \text{update the weight and bias} \end{array} \right.$

## advantages:

- easy to interpret
- differentiable (used in GD)
- one local minima only.

## disadvantages

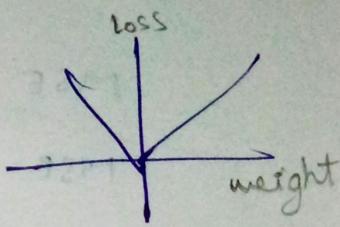
- unit of error is squared
- it is not robust to outliers.

\* if we use this loss function in DL's Neural Network architecture then activation function must be linear. (last one node)

## ②) Mean Absolute Error (MAE or L1 loss)

loss function ( $L$ ) =  $|Y_i - \hat{Y}_i|$

cost function ( $C$ ) =  $\frac{1}{n} \sum |Y_i - \hat{Y}_i|$



### advantages:

→ intuitive and easy to understand

→ unit is same as  $Y$

→ Robust to outlier

### disadvantages

→ not differentiable function

~~if data have no outlier then use MAE~~

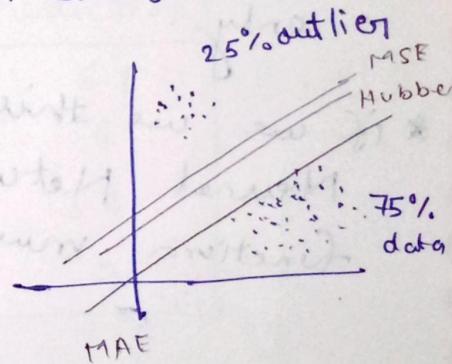
## ③) Huber Loss:

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

if no outlier then Huber = MSE

if outlier then ~~Huber~~ = MAE

Here outlier depends on the ( $\delta$ ) hyperparameter.



## ④) Binary Cross Entropy (log loss):

- used in classification example in logistic regression
- it is used when we have only two class prediction (0 or 1, True or False)

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

$Y$  = actual

$\hat{Y}$  = prediction.

- if we use Binary Cross entropy then in the output layer activation function is only sigmoid, in hidden layer may be change.

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

advantages:

- differentiable
- 

disadvantage:

- multiple local minima
- not a intuitive

### ⑤ Categorical Cross Entropy ~~(used in soft)~~:

- it is used in softmax Regression
- it is used in classification (multi-class classification)

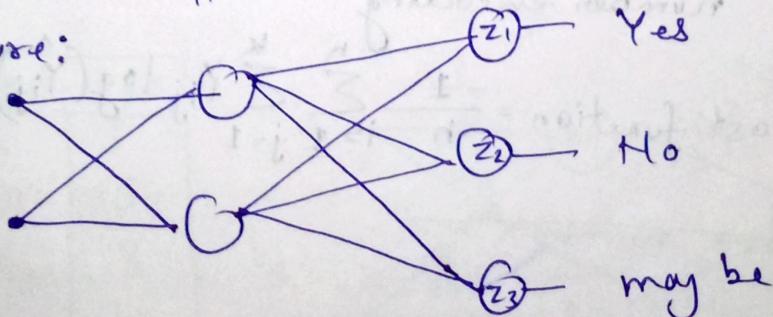
$$L = - \sum_{j=1}^k Y_j \log (\hat{Y}_j)$$

Here  $k$  = number of class.

cgtage	placed	Yes	No	Maybe
8.80	Yes	1	0	0
6.60	No	0	1	0
7.70	maybe	0	0	1

OHE (one hot encode)

Architecture:



↑  
output neurons categories (class)

- in output layer activation function is always softmax.

\* Here the no. of neuron = number of class  
for N=3.  $f(z)_1 = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$

for N=3  $f(z)_2 = \frac{e^{z_2}}{e^{z_1} + e^{z_2} + e^{z_3}}$

for may be  $f(z)_3 = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$

and the loss function for three class.

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

range of  $f(z)_1, f(z)_2$  and  $f(z)_3$  is 0 to 1

$$\text{and } f(z)_1 + f(z)_2 + f(z)_3 = 1$$

⑥ sparse categorical cross entropy

One hot encoding

number encoding

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

- Here the number of

(x,y) as inputs are input tuples

where xi contains activities of input features and yi contains output features