

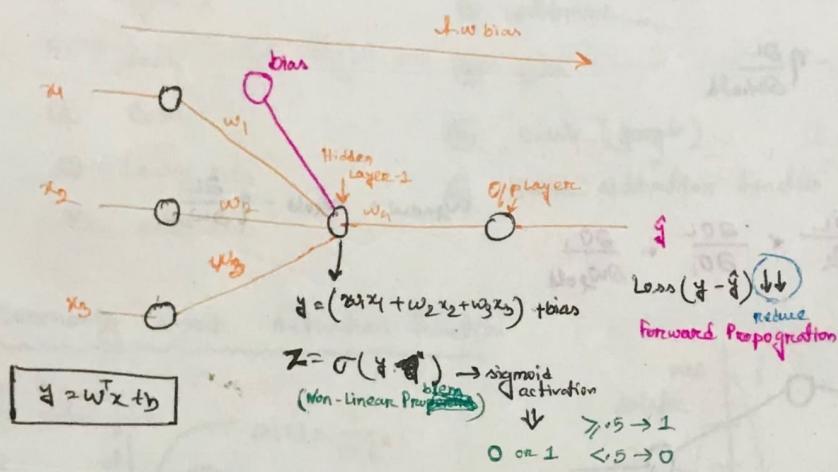
## Tutorial-2

# Forward Propagation, Loss functions, Chain Rule of derivatives

## Day-2 : Deep Learning

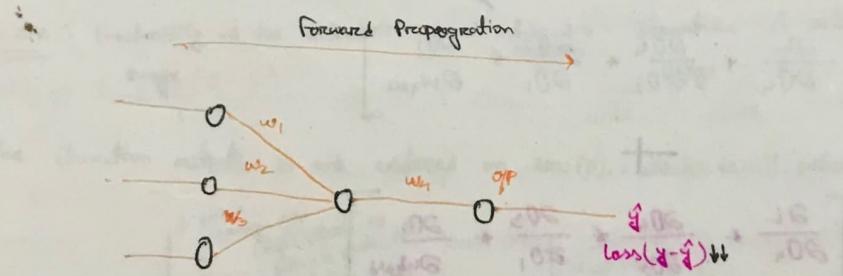
### Agenda

- ① Forward Propagation
- ② Activation function
- ③ Chain Rule of derivative
- ④ Vanishing Gradient Problem
- ⑤ Loss function



### Backward Propagation

- ① Weight update formula
- ② Chain Rule of differentiation



### ③ Backward Propagation

- ④ Weight update formula

$$\frac{\partial L}{\partial w_{old}} =$$

$\eta$  (-ve slope)

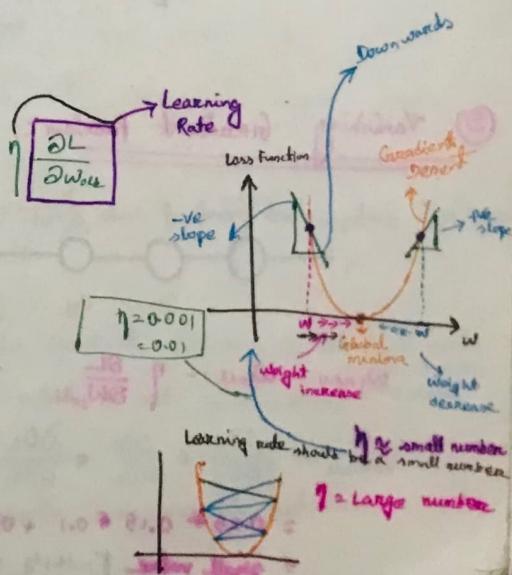
$$w_{new} = w_{old} - \eta (-ve)$$

$$= 2w_{old} + \eta (-ve)$$

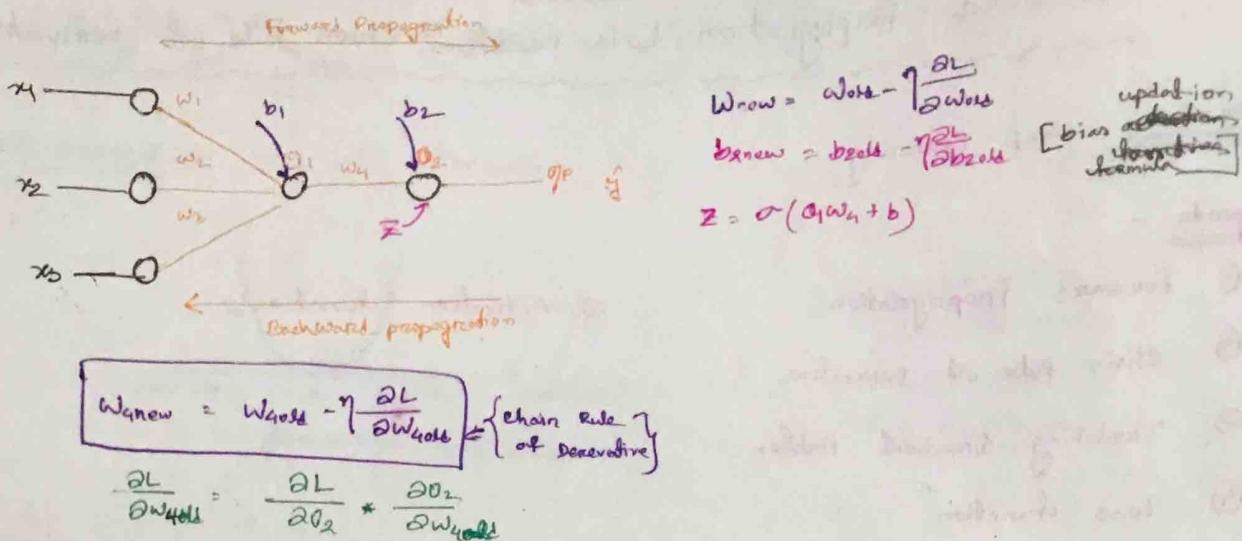
$$w_{new} = w_{old} - \eta (+ve)$$

$$w_{new} >> w_{old}$$

$$w_{new} << w_{old}$$



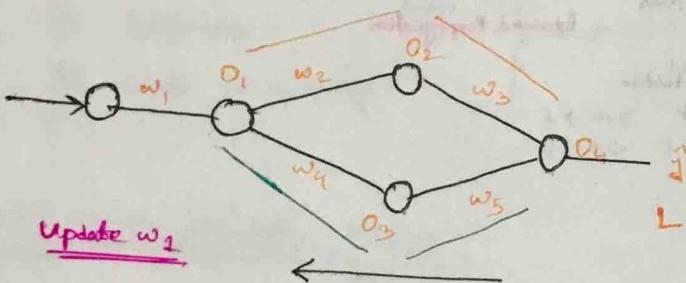
### ④ Chain rule of differentiation:



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

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}}$$

$$\frac{\partial L}{\partial w_{\text{old}}} = \frac{\partial L}{\partial o_2} * \frac{\partial o_2}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}} \leftarrow \frac{\partial o_1}{\partial w_{\text{old}}} \quad w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w_{\text{old}}}$$



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

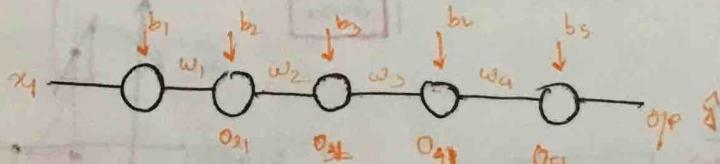
$$\frac{\partial L}{\partial w_{\text{old}}} = \left[ \frac{\partial L}{\partial o_4} * \frac{\partial o_4}{\partial o_2} * \frac{\partial o_2}{\partial o_1} + \frac{\partial L}{\partial o_5} * \frac{\partial o_5}{\partial o_3} * \frac{\partial o_3}{\partial o_1} \right]$$

$$+ \left[ \frac{\partial L}{\partial o_4} * \frac{\partial o_4}{\partial o_3} * \frac{\partial o_3}{\partial o_1} * \frac{\partial o_1}{\partial w_{\text{old}}} \right]$$

chain rule of derivatives

### ⑤ Vanishing Gradient Problem:

Sigmoid value between 0 and 1



$$w_{1,\text{new}} = w_{1,\text{old}} - \eta \frac{\partial L}{\partial w_{1,\text{old}}}$$

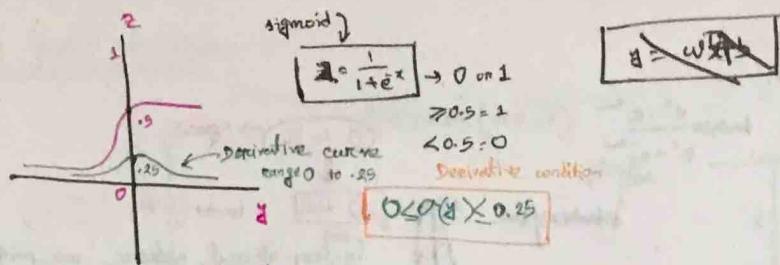
$$\text{Loss} = \frac{1}{2} (\theta - y)^2 \rightarrow \text{MSE (Mean Square Error)}$$

$$o_{51} = \sigma \left[ \frac{\partial o_{51}}{\partial w_{1,\text{old}}} \cdot [ (o_{41} + w_4) + b_3 ] \right] \quad 0 \text{ to } 0.25$$

$$\frac{\partial L}{\partial w_{1,\text{old}}} = \frac{\partial L}{\partial o_{51}} * \frac{\partial o_{51}}{\partial o_{41}} * \frac{\partial o_{41}}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial w_{1,\text{old}}}$$

= 0.25 \* 0.15 \* 0.1 \* 0.05 \* 0.02 [Every hidden neuron derivative loss is decreased]

= small value [multiply with small value results in small]



$$0.5 < \alpha(z) < 0.25$$

$> 0.5 = 1$   
 $< 0.5 = 0$

Derivative condition

$$0.5 < \alpha(z) < 0.25$$

$$W_{\text{new}} = W_{\text{old}} - \eta \quad (\text{small numbers})$$

$$W_{\text{new}} \approx W_{\text{old}}$$

Vanishing gradient  
No change in weights

$$\begin{aligned} W_{\text{new}} &= W_{\text{old}} - \eta \quad (\text{small numbers}) \\ 1000 &= 1000 - .001 \\ 1000 &\approx 999.99 \\ 1000 &\approx 1000 \\ W_{\text{new}} &\approx W_{\text{old}} \end{aligned}$$

to solve use  
Another activation function

Q1 There are many another activation function

① Sigmoid

④ softmax

② tanh

⑤ Elu

③ ReLU

⑥ Swish (google)

④ Leaky-ReLu

⑦ Linear activation function

⑤ pre-ReLu

## Q2 Commonly used Activation function

a. Sigmoid:

Formula

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

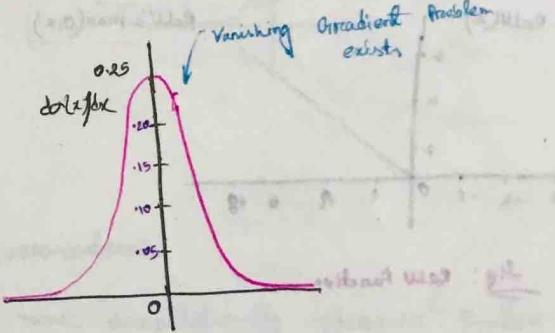
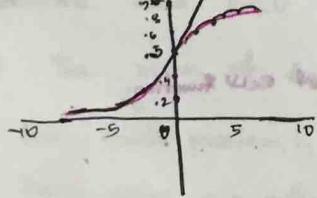
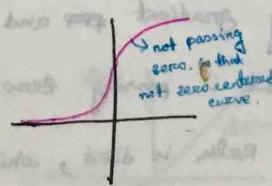
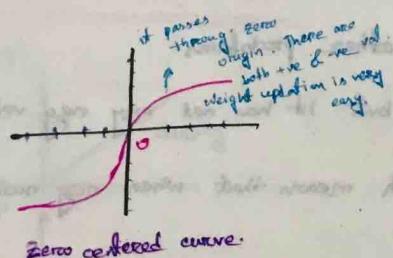


Fig-1-a : Functionality of the activation  
range (0-1)

Fig-1-b : Derivative of activation  
range (0-0.25)

i) The function output is not centered on zero (0), which will reduce the efficiency of weight update.



ii) The sigmoid function performs exponential operation ( $\alpha(z) = \frac{1}{1+e^{-z}}$ ), which is slower for computers. Time increase

### Advantages:

- ① Smooth gradient, preventing 'jumps' in output values.
- ② Output values bound between 0 and 1, normalizing the output of each neuron.
- ③ Clean prediction, that is very close to 1 or 0.

### Disadvantages:

- ① prone to gradient vanishing, so that doesn't create a deep network.
- ② Function output is not zero-centered.
- ③ Power operations are relatively time consuming.  
★ can't create deep neural network.

b. ~~tanh~~: tanh

hyperbolic  
It is a hyperbolic function.

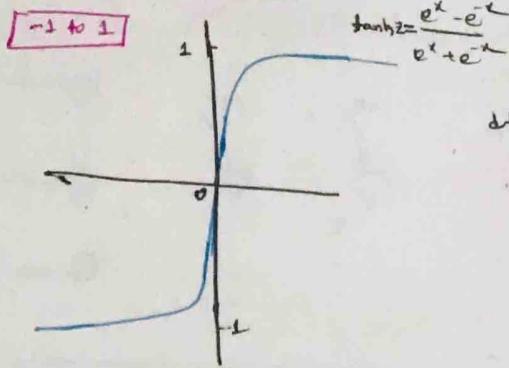


Fig: Functionality of tanh

range  $(-1 \rightarrow 1)$

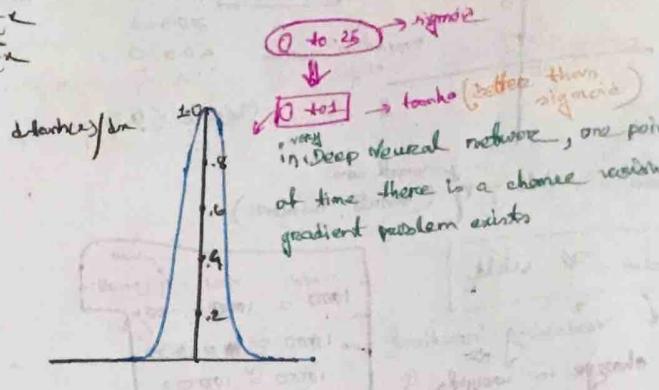


Fig: Derivatives of tanh

range  $(0 \rightarrow 1)$

① zero centric function, hyperbolic function.

② In binary classification, the tanh function is used the hidden layers and sigmoid is used for the output.

c. ReLU function:

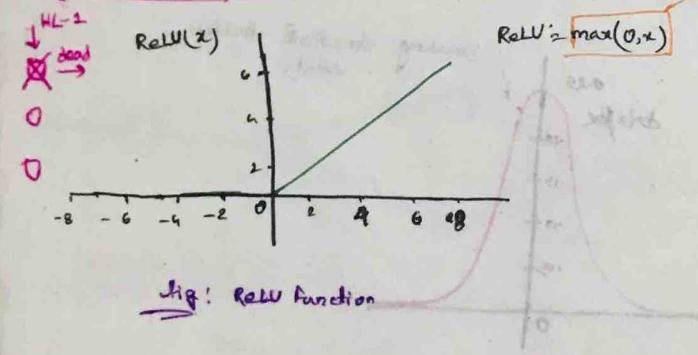


Fig: ReLU function

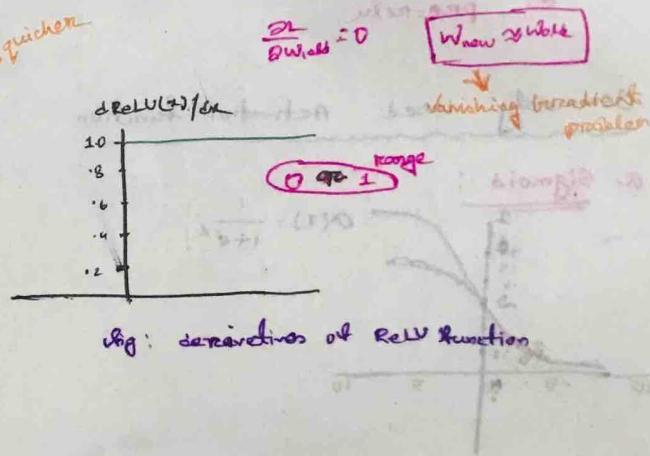
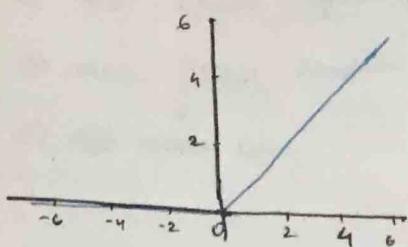


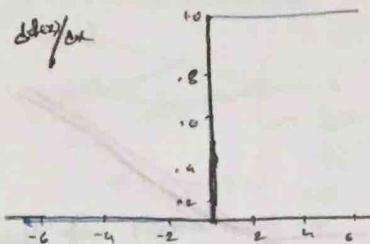
Fig: derivatives of ReLU function

- ① If  $x$  is neg the value becomes zero(0).
- ② When  $\frac{\partial L}{\partial w_{all}} = 0$ , then  $\frac{\partial L}{\partial w_{new}} \approx \frac{\partial L}{\partial w_{old}}$  which is called vanishing gradient problem
- ③ It is better than sigmoid & tanh. It is pretty much quicker.
- ④ It is obviously solving vanishing gradient problem and other problems.
- ⑤ Disadvantage: It is not zero centric. It is passing zero but it has not any neg value.
- ⑥ When input is negative the ReLU is died, which means that when neg number is entered the output val is zero.

### Leaky - ReLU Function:

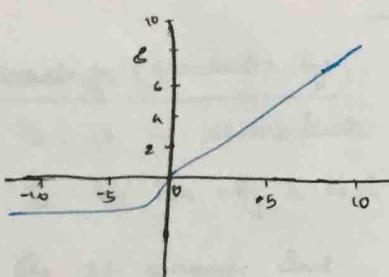


$$f(x) = \max(0.01x, x)$$

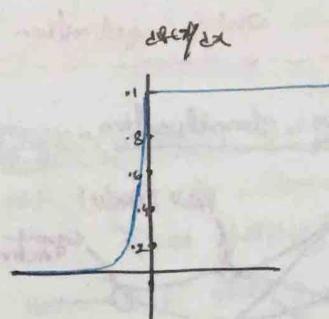


- ① Solves dead neuron problem<sup>(ReLU)</sup>, because it has small neg number
- ② It has been not fully proved.

### ELU (Exponential Linear Units) Function:



$$f(x) = \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{otherwise} \end{cases}$$



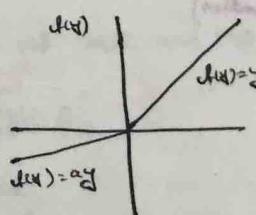
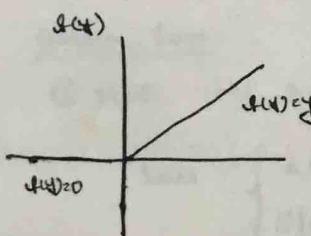
- ① No dead ReLU issue
- ② The mean of the output is closed to 0, zero-centered
- ③ One small problem is that it is slightly more computationally intensive. Similar to Leaky ReLU. Hence using exponential value.

④

### Softmax

⑤

### PReLU (Parametric ReLU) Function:



PReLU is PRELU. For PRELU the coefficient of the negative part is not constant and is adaptively learned.

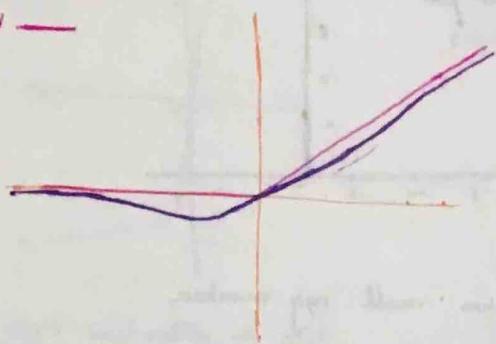
$$f(x_i) = \begin{cases} y_i, & \text{if } y_i > 0 \\ a_i y_i, & \text{if } y_i \leq 0 \end{cases}$$

- ① If  $a_i = 0$ , becomes ReLU
- ② If  $a_i > 0$ , becomes Leaky ReLU
- ③ If  $a_i < 0$ , becomes PRELU. If  $a_i$  is a learnable parameter, it becomes PRELU

## ⑤ Swish

A self ~~gated~~ function. This is brought by google itself.

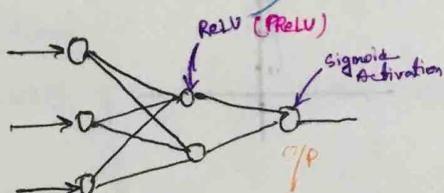
Swish —  
ReLU —



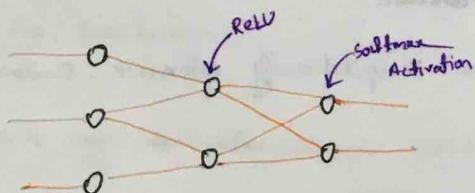
There is a problem in that we cannot find out the derivative of zero.

## ⑥ Technique - which activation function we should use?

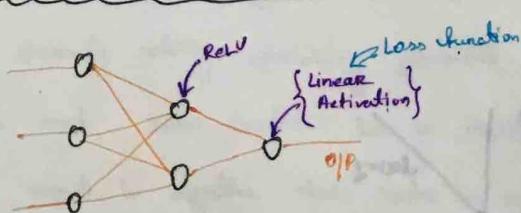
- i) Binary Classification: convergence not happening so that change ReLU to PReLU(FLU)



- ii) Multiclass Classification:



- iii) Regression Classification: it means will be happening continuous value.



## ⑦ Loss Function:

Deep Learning (ANN)

Loss Function

Data:

continuous value Regression

Ex: Degree Salary (C/P)

10 PHD

— —

Regression Problem

Data: Classification

Play

10

9

5

Study

2

3

5

Pass/Fail

Fail

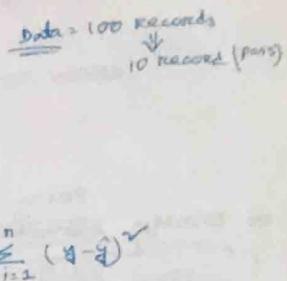
Fail

maybe Pass

## Loss Function in Regression:

- (a) MSE (Mean Square Error)
- (b) MAE (Mean Absolute Error)
- (c) Huber Loss

## Loss & cost function



## MSE (Mean Square Error):

$$\text{Loss function: } \frac{1}{2n}(Y - f)^2$$

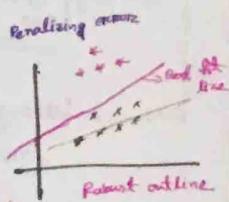
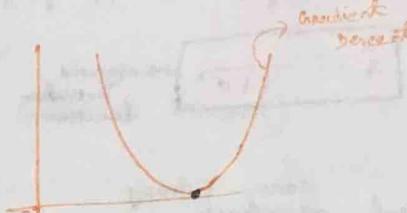
Quadratic Equation

$$(a-b)^2 = a^2 - 2ab + b^2 \quad [\text{quadratic eqn}]$$

$$a^2 + bx + c$$

in cost func provide batch wrt data.

$$\text{cost function: } \frac{1}{2n} \sum_{i=1}^n (Y_i - f_i)^2$$



### Advantages (quadratic eqn):

- (i) it is differentiable
- (ii) it has only 1 local or global minima.
- (iii) it converges fast

### Disadvantages (quadratic equation):

- (i) not robust to outliers  
(line shifting)  
Hence line shifting is huge.

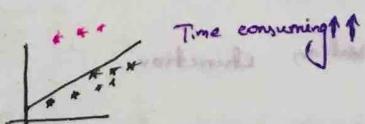
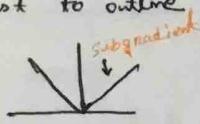
## MAE (Mean Absolute Error):

$$\text{Loss function: } \frac{1}{2}(Y - f)^2$$

$$\text{cost function: } \frac{1}{2} \sum_{i=1}^n |Y_i - f_i|$$

### Advantages:

- (i) Robust to outliers



### Disadvantages:

- (i) Time consuming



## Huber Loss:

- (i) MSE and (ii) MAE (combination of MSE and MAE)

$$\text{Loss} = \begin{cases} \frac{1}{2}(Y - f)^2, & \text{if } |Y - f| \leq \delta \\ \delta |Y - f| - \frac{1}{2}\delta^2, & \text{otherwise} \end{cases}$$

when outliers is present



two functions are selected, the one giving a large loss

and other for small, outliers case.

## Loss function in Classification:

cross entropy → Binary cross entropy → Binary classification

cross entropy → Categorical cross entropy → Multiclass classification

### ④ Binary cross Entropy Classification:

(Binary classification problem)

$$\text{Loss} = -y \log(\hat{y}) - (1-y) \log(1-\hat{y}) \rightarrow \text{logistic regression}$$

$$\text{Loss} = \begin{cases} -\log(1-\hat{y}) & \text{if } y=0 \\ -\log \hat{y} & \text{if } y=1 \end{cases} \rightarrow \text{Binary classification}$$

$$\hat{y} = \frac{1}{1+e^{-x}} \rightarrow \text{sigmoid activation function}$$

### ⑤ Categorical cross entropy:

(Multiclass classification Problem)

$x_1$	$x_2$	$x_3$	opp
1	3	4	Good
2	6	7	Bad
8	9	10	Neutral

number of categories

Good	Bad	Neutral	$y_{it}$
1	0	0	$y_{i1}$
0	1	0	$y_{i2}$
0	0	1	$y_{i3}$

$$\text{Loss}(x_i, y_i) = - \sum_{j=1}^{k_{\text{cat}}} y_{ij} \ln(\hat{y}_{ij})$$

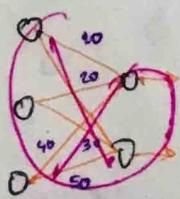
$k_{\text{cat}} = \text{number of categories}$

$$y_{ij} = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ik}]$$

$$y_{ij} = \begin{cases} 1 & \text{if the element is in class} \\ 0 & \text{otherwise} \end{cases}$$

$\hat{y}_{ij} = \text{Softmax Activation function} \leftarrow O/P \text{ Layer}$

$$\alpha(z) = \frac{e^{z_i}}{\sum_{j=1}^k e^{z_j}} \rightarrow \text{softmax activation}$$



$$\begin{aligned} \text{At higher output, probability that class will be coming up} \\ &= \frac{e^{10}}{e^{10} + e^{20} + e^{30} + e^{40} + e^{50}} \\ &= 0.4, 0.5, 0.6 \\ &\quad \{ \} \\ &0/P_1 + 0/P_2 = 1 \rightarrow \text{at higher probability} \end{aligned}$$

## Conclusion:

ReLU, Softmax → Multiclass → categorical cross entropy  
 ReLU, Sigmoid → Binary → Binary cross entropy

## Linear Regression

ReLU, Linear Activation (Loss → MSE, MAE, Huber loss)