

## Day 2

1. forward propagation
2. Chain Rule of Derivative.
3. Vanishing Gradient Problem
4. loss function.

\* Backward Propagation.

weight updating formula.

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

$\eta$  = Learning rate.

$\partial$  = Derivative.

$L$  = Loss

$W$  = Weights.

How to find  $\frac{\partial L}{\partial W_{\text{old}}}$  ← To find this we use

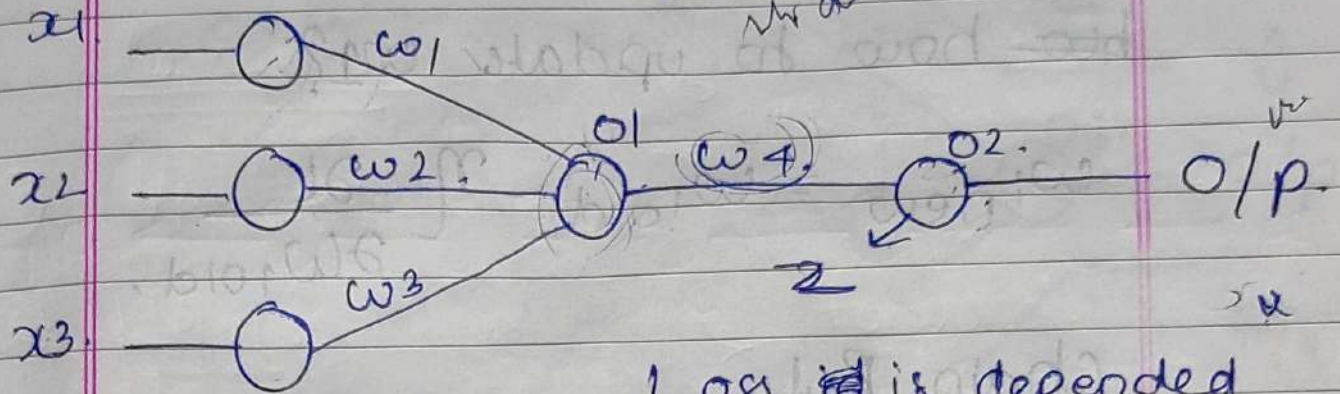
↓  
Chain rule of Derivative



$$z = \sigma(0, w_4 + b)$$

$\sigma = \text{Activation}$

Function.



Loss is depended on  $o_2$ .

$$\frac{\partial L}{\partial w_{4old}} = ?$$

$$w_{4new} = w_{4old} - \eta \frac{\partial L}{\partial w_{4old}}$$

Whenever we want to update  $w_4$ . { chain Rule of Derivative. }

$$\frac{\partial L}{\partial w_{4old}} = \frac{\partial L}{\partial o_2} \times \frac{\partial o_2}{\partial w_{4old}}$$

Bias Updation formula.

$$b_{2w} = b_{2old} - \eta \frac{\partial L}{\partial b_{2old}}$$

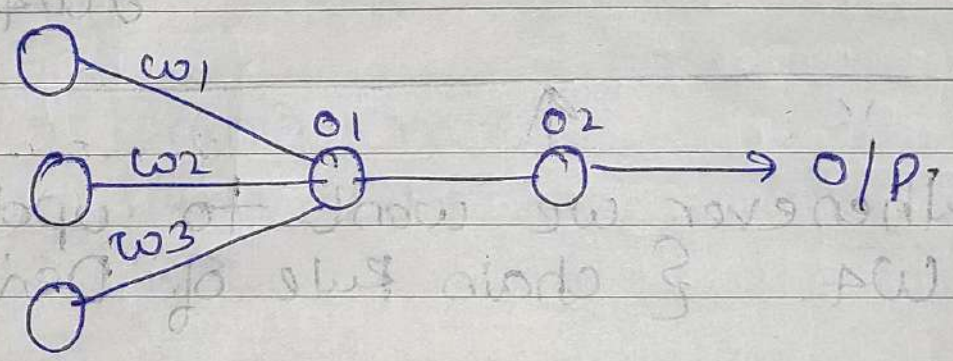


how to update  $w_1$ ?

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

### Chain Rule

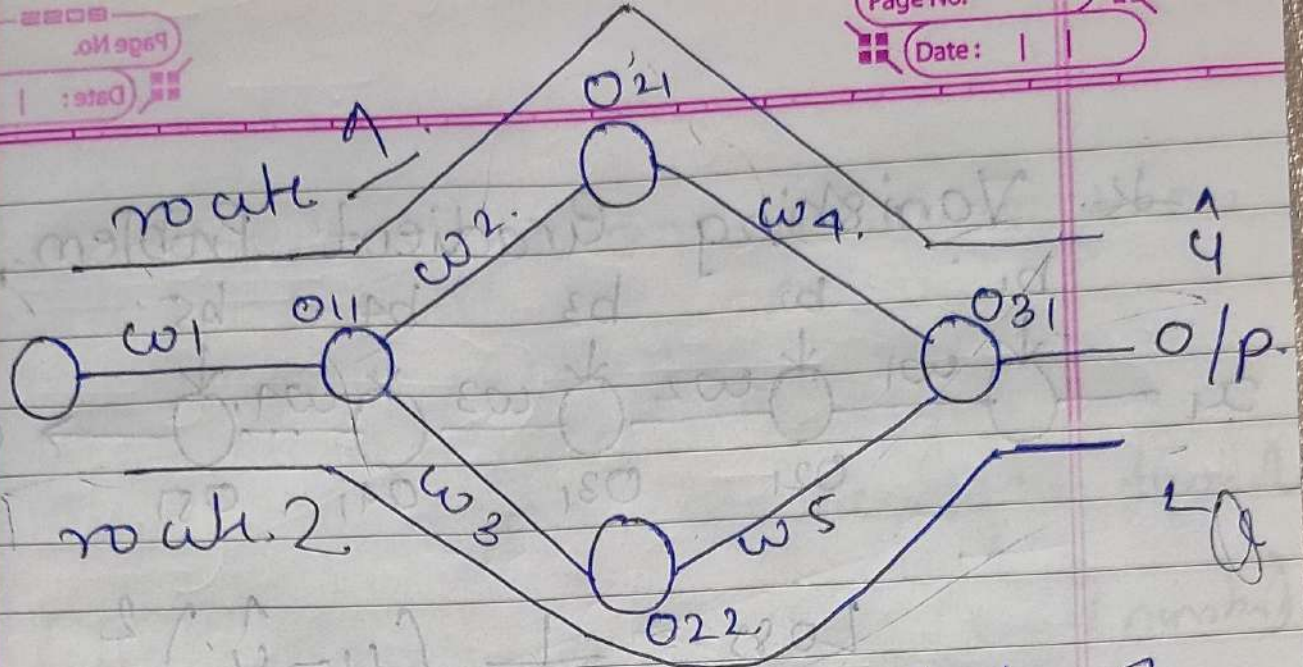
loss is dependable on  $O_2$ , and  
 $O_2$  is dependable on  $O_1$   
 and  $O_1$  is dependable on  
 $w_1$ .



$$\frac{\partial L}{\partial w_{1\text{old}}} = \frac{\partial L}{\partial O_2} \times \frac{\partial O_2}{\partial O_1} \times \frac{\partial O_1}{\partial w_{1\text{old}}}$$

This is the new Chain Rule of differentiation to update  $w_1$ .





how to find  $w_1$  with the help of chain rule.

→ There are 2 routes.

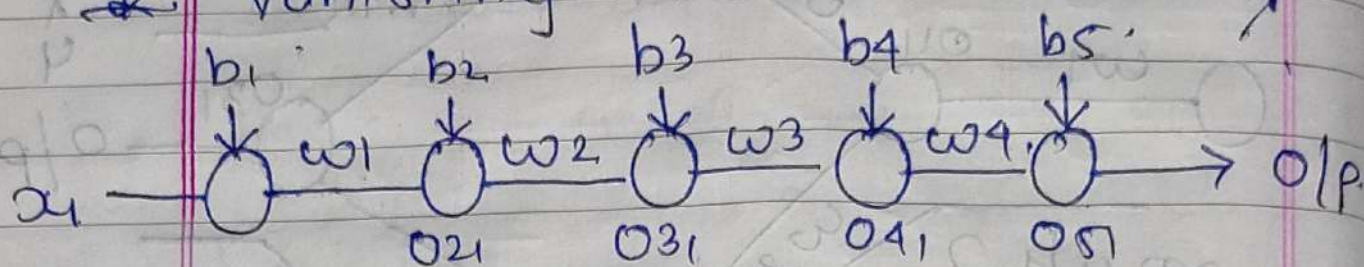
$$w_{1, \text{new}} = w_{1, \text{old}} - \eta \frac{\partial L}{\partial w_{1, \text{old}}}$$

$$\frac{\partial L}{\partial w_{1, \text{old}}} = \left[ \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{21}} * \frac{\partial o_{21}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1, \text{old}}} \right]$$

$$+ \left[ \frac{\partial L}{\partial o_{31}} * \frac{\partial o_{31}}{\partial o_{22}} * \frac{\partial o_{22}}{\partial o_{11}} * \frac{\partial o_{11}}{\partial w_{1, \text{old}}} \right]$$



## \* Vanishing Gradient Problem



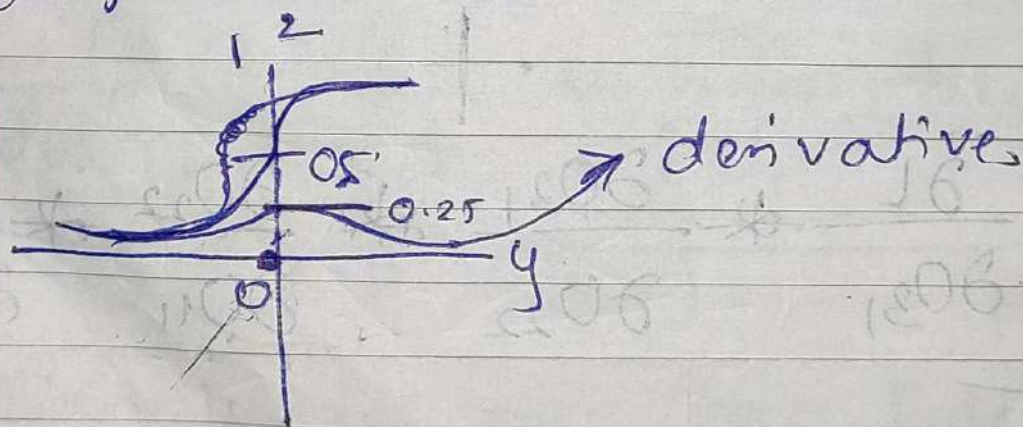
$$\text{Loss} = \frac{1}{2} (y - \hat{y})^2$$

L → MSE

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

$$\frac{\partial L}{\partial w_{\text{new}}} = \frac{\partial L}{\partial o_5} * \frac{\partial o_5}{\partial o_4} * \frac{\partial o_4}{\partial o_3} * \frac{\partial o_3}{\partial o_2} * \frac{\partial o_2}{\partial w_1}$$

## ⊙ Sigmoid Act function



derivative is always range in between 0 to 0.25.



0.25 \* 0.15 \* 0.10 \* 0.05 \* 0.2

↓

= small value

Then we'll update the weight formula

$$w_{\text{new}} = w_{\text{old}} - \eta \frac{\partial L}{\partial w} \quad (\text{small number})$$

→ small number.

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

(Vanishing Gradient Problem)

How to solve this?

→ Use Another Activation function

① Sigmoid

② Tanh

③ ReLU

④ Leaky ReLU

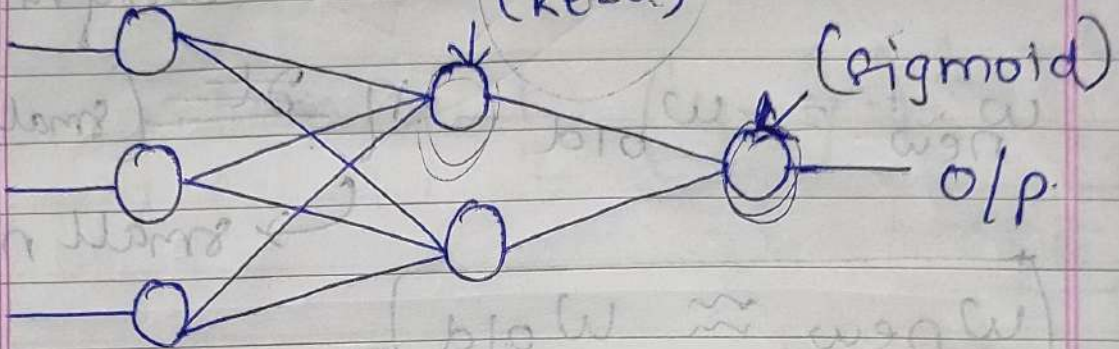
⑤ Pre ReLU

} Activation functions



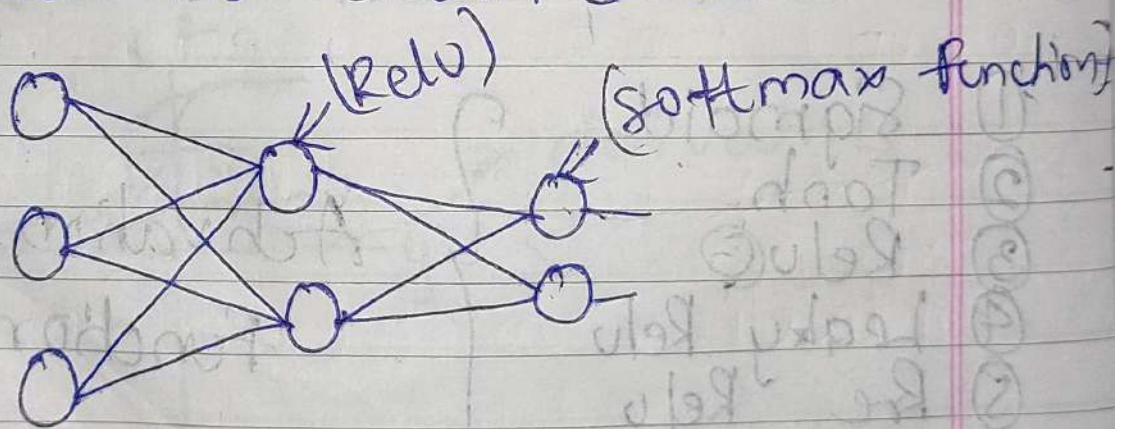
Technique which Activation function we should use.

\* Binary classification.  
(ReLU)



\* Always Remember whenever you are doing a binary classification problem your output should have the sigmoid Activation function.

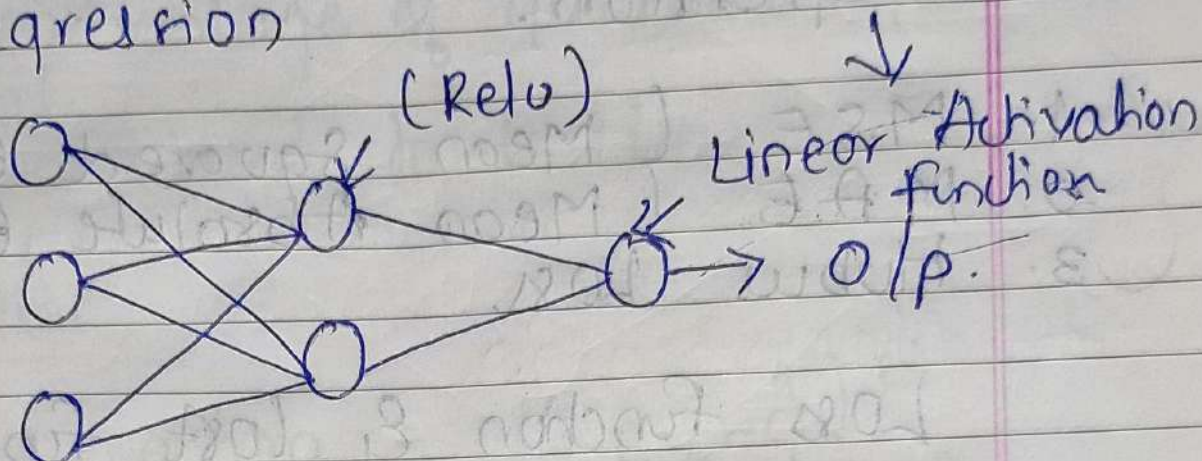
\* Multiclass classification.





# Loss function

## \* Regression



## \* Loss function.

## Deep Learning ANN.

↓  
Regression

↓  
Classification.

Exp.	Deq	Sal.
10	phd	...
—	—	—
—	—	—

Play	Study	P/f
10	2	(F)
2	10	(R)
—	—	(C)

Regression Problem

Classification



## \* Regression.

1. MSE (Mean Squared Error)
2. MAE (Mean Absolute Error)
3. Huber Loss.

Loss function & cost function.

### ① MSE

Loss function =  $\frac{1}{2} (y - \hat{y})^2$

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

→ Also called quadratic Equation.

## \* Advantages of quadratic Equation

- ① Differentiable.
- ② It has only 1 local or Global minima.
- ③ It converges faster.



\* Disadvantage of quadratic Equation.

① ~~Not~~ Robust to Outliers

\* MAE

$$\text{Loss function} = \frac{1}{2} (y - \hat{y})^2$$

$$\text{loss function} = \frac{1}{2} \sum_{i=1}^n (y - \hat{y})$$

\* Advantage.

① Robust to Outlier

\* Huber loss.

hyperparameter

Huber loss is a combination of MSE and MAE

$$\text{Loss fn} = \frac{1}{2} (y - \hat{y})^2 \quad \text{if } |y - \hat{y}| \leq \underline{\quad}$$

↑

Use when outliers are not present.



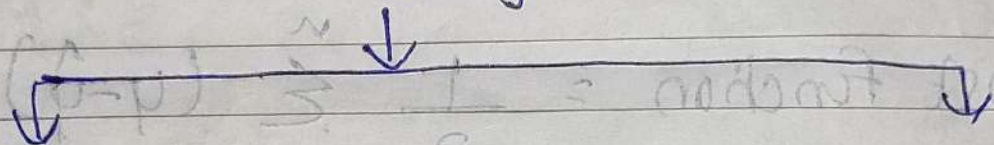
$$\delta |y - \hat{y}| - \frac{1}{2} \delta^2, \text{ otherwise}$$

↑

use when outliers are present

\* Classification.

Cross Entropy



Binary ~~classification~~  
Cross Entropy

Categorical

Cross Entropy



Binary Classification



Multi class Classification

① Binary Cross Entropy.

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



Logistic Regression



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

### ① Categorical Cross Entropy.

$f_1$	$f_2$	$f_3$	O/p.
2	3	4	Good
5	6	7	Bad
8	9	10	Neutral

Convert to One Hot Encoding.

Good	Bad	Neutral
1	0	0
0	1	0
0	0	1

$$L(x_i, y_i) = - \sum_{j=1}^c y_{ij} \cdot \ln(\hat{y}_{ij})$$

$i = \text{rows}$

$j = \text{columns}$

$$y_i = [y_{i1}, y_{i2}, y_{i3}, \dots, y_{ic}]$$

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



Conclusions.

Categorical Cross Entropy.



Relu, Softmax → Multiclass

Relu, Sigmoid → Binary

Linear Regression / Binary Cross Entropy

Relu, Linear Activation



Loss function = MSE, MAE, Huber