

## Module-I Deep Learning

→ What is artificial intelligence?

- It is the ability of a computer to perform tasks commonly associated with intelligent beings.

→ What is machine learning?

- It is the study of algorithms that learn from data.
- P - Performance
- T - Task
- E - Experience
- Supervised Learning
  - Labelled dataset, (<sup>continuous value</sup><sub>categorical value</sub>)
  - Regression, Classification - KNN, SVM, Decision Tree, Random Forest, Artificial Neural Networks.
- UnSupervised Learning:
  - Unlabelled dataset
  - Clustering, Association

- Regression: Trying to fit a line to a set of data points
  - The distance b/w the datapoints and the line is called error.
  - The best fit line will be the one with the less error.

\* In machine learning, we are trying to find the equation of a line

$$y = mx + c$$

- × In the case of regression, we find the best fit line to set of data points
- × In classification, we are trying to find the line that partitions of the data points.

### Regression:

$$y = x + c$$

↓              ↓  
dependant    independant  
variable       variable  
(what  
we're  
predicting)

$$b_1 = \frac{n \sum (x_i y_i) - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

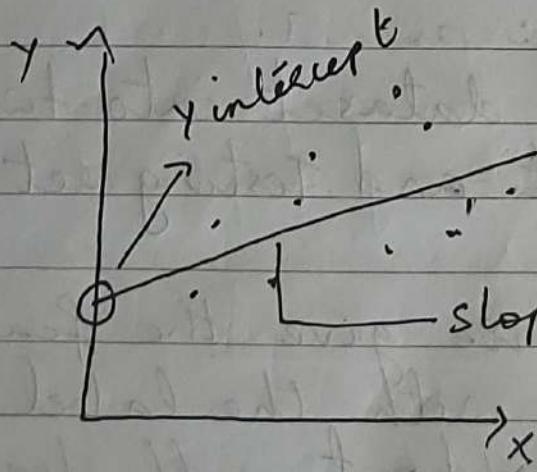
\* If there are three variables (equation of a plane)

$$Y = b_0 + b_1 x_1 + b_2 x_2$$

- \* If its only one variable, it is a threshold.
- \* If there are more than three variables we are trying to find the equation of a hyperplane.

$$y = mx + c$$

↓                    ↓  
slope              Y-intercept



If we keep changing the slope and Y intercept, we can find the best fit line, that has the lowest error.

Training

## Linear Regression:

### Residuals:

Difference b/w actual and predicted value.

- \* The best model is estimated by minimizing the sum of squares of these residuals

$$SS_{res} = \sum_i^n (y_i - f(x_i))^2$$

↓  
cost function used to measure the error of the model.

### Model Training:

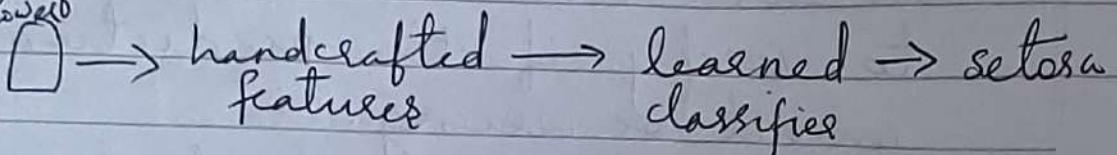
Partition the dataset into two:  
Training set and Testing set.

- \* For training, we give the feature data along with the label to model and the output will be an equation  $y = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n$ .
- \* After training, we will get a learned model.

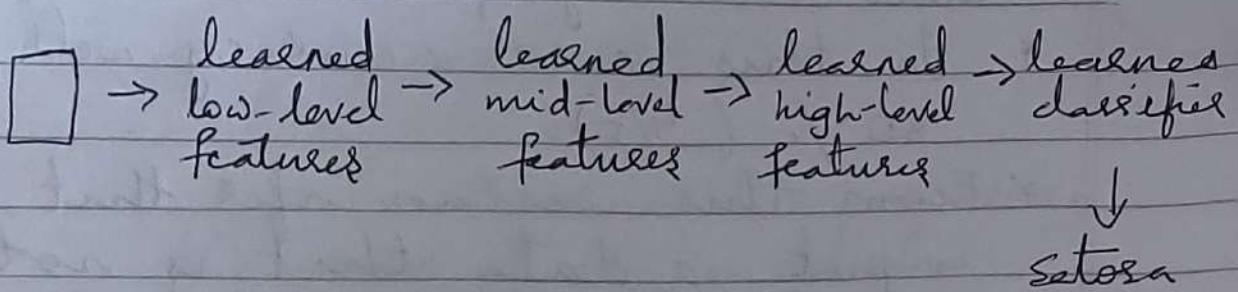
- \* In test data, we will mask the label and the model will predict the category. It may be correct or wrong.
- \* We can cross check with the label to verify if the prediction was correct or wrong.
- \* If 4 out of 5 predictions were right, we use evaluation metrics like accuracy, precision and recall.
- \* In this case accuracy will be  $80\% (4/5)$ .
- \* From this we can infer that, if we input a data that is not in the label set, the model will predict correctly with 80% accuracy.
- \* In machine learning, we make the machines learn the weights associated with the features.

## \* What is deep learning?

Traditional machine learning:  
image of  
a flower



Deep "end-to-end" learning:



\* Deep learning uses artificial neural networks to learn from lots of data without needing explicit programming

Artificial Neural Network:

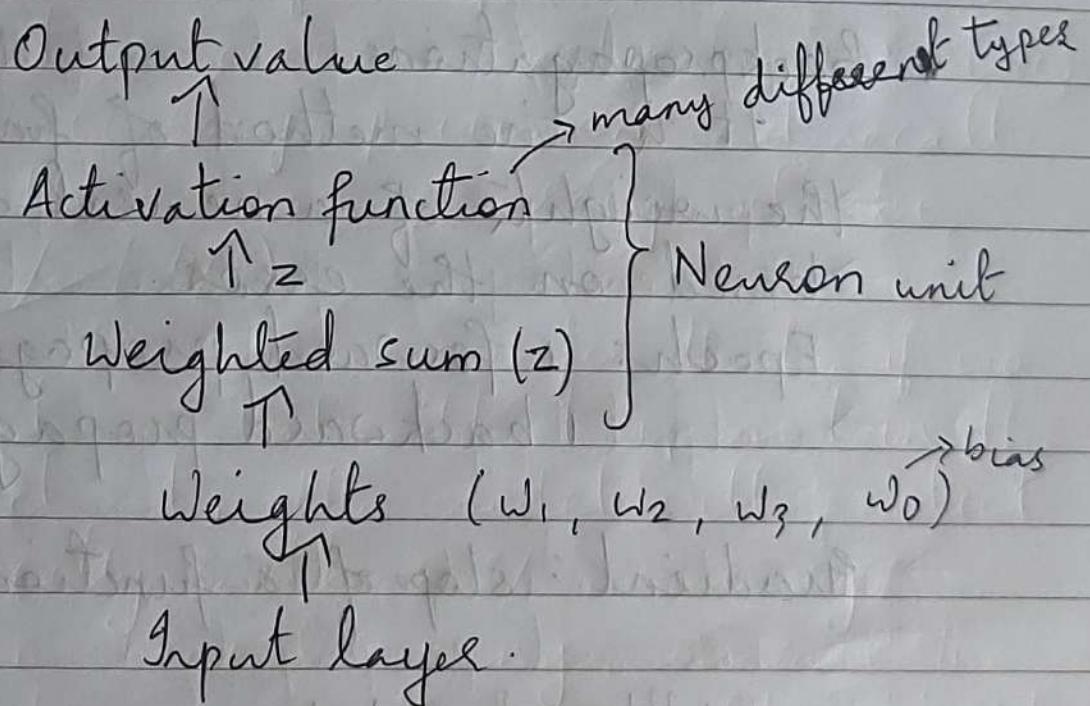
→ ANN are machine learning models that tries to mimic the complex functions of the human brain.

→ Simulate our brain

Why neural networks?

→ Can be used for wide variety of applications

$w_0 \rightarrow$  bias, it is the y intercept. To provide some control for the equation



if  $z \rightarrow$  Regression function (continuous value)

$$z >= 0 \rightarrow 1$$

else

$$z \rightarrow -1$$

step function

(Activation function)

↓  
transforms the weighted sum (z)  
to our desired output.

\* For classification, we use sigmoid and softmax activation function

Sigmoid (0-1):

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

\* Backpropagation:

It is the method of fine tuning the weights of a neural network based on the error.

Epoch: 1 forward propagation and 1 backward propagation.

Gradient: slope of a function.

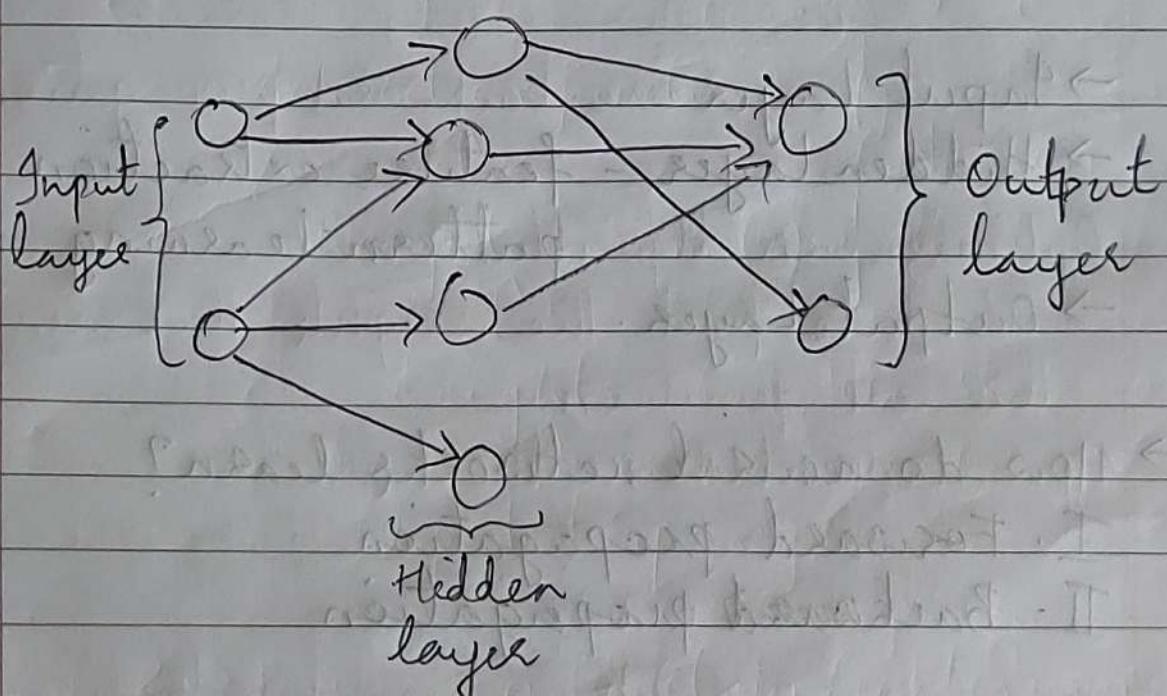
$$w_i' = w_i + \Delta w_i$$

- \* We find the values of  $w$  for which the error is minimal.
- \* We keep changing the weights iteratively till we get the minimum error.
- \* We use optimization algorithm to minimize the error.

- \* Example of optimization algorithm is gradient descent.
- \* Error / cost / loss

Perceptron  $\rightarrow$  a single neuron structure sideways

Multi-layer Perceptron



- \* If there are more than two hidden layers then it is called deep neural network.
- \* Less than two - shallow neural network.

- Why do we need multiple patterns?
  - To get better results.

Hyperparameter : The user chooses:

- ↳ no. of hidden layers
- ↳ no. of neurons to include

Components / Architecture of Neural network -

- Input Layer.
- Hidden layer - feature extraction, pattern learning
- Output Layer.

- How do neural networks learn?

- I. Forward propagation
- II. Backward propagation

Forward Propagation:

- ↳ Each layer accepts input and passes the output to the next layer in the forward direction.
- ↳ Pre-activation: calculating weighted sum

↳ Activation: Passing the sum into an activation function.

How does the computer know there is error in prediction?

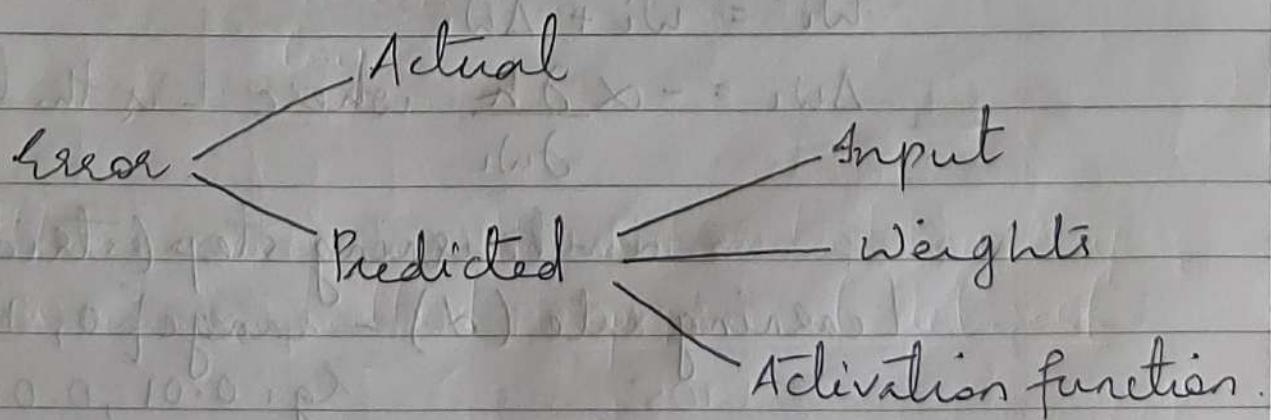
↳ Use cost / loss / error function

Eg: Mean square error.

$$MSE = \frac{1}{n} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$\frac{\text{(actual value} - \text{predicted value})^2}{n}$$

The error happens because of wrong weights learnt.



Now, we update the weights to fine tune the error.

Back propagation:

→ Trying to update weights to minimize the error.

One epoch - one forward propagation

one backward propagation

Optimization Algorithm:

Used to find the weights that minimize the cost function.

Eg: Gradient Descent

$$w_i' = w_i + \Delta w_i$$

$$\Delta w_i = -\alpha \frac{\partial L}{\partial w_i} \text{ where } L \text{ is the loss function}$$

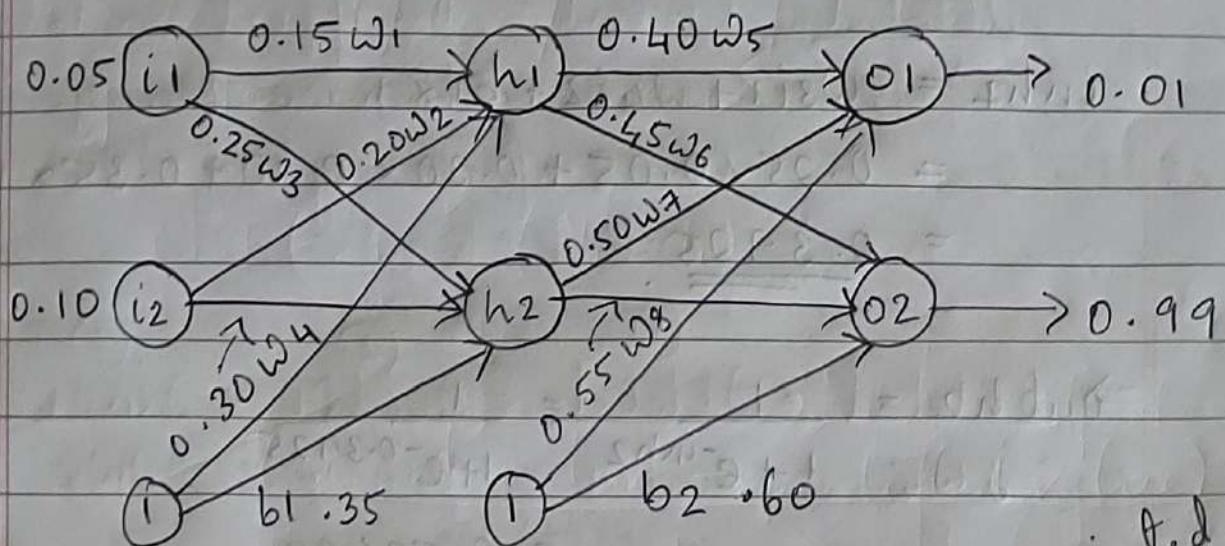
$\alpha$  - learning step (stride)

Learning rate ( $\alpha$ ) - range [0, 1]

Eg: 0.01, 0.001, 0.0001

\* Lesser the  $\alpha$ , the more optimal the solution is.

## How Backpropagation Works?



$h_1$

$$\text{in } h_1 = w_{11}i_1 + w_{21}i_2 + b_1 \rightarrow \begin{array}{l} \text{value associated with} \\ \text{bias} \end{array}$$

$$\text{in } h_1 = w_{11}i_1 + w_{21}i_2 + b_1 \times 1$$

(input to  $h_1$ )

$$= 0.15 \times 0.05 + 0.20 \times 0.10 + 0.35 \times 1$$

$$= \underline{\underline{0.3775}}$$

Assume activation function is sigmoid

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

$$\text{outh}_1 = \frac{1}{1 + e^{-\text{inh}_1}} = \frac{1}{1 + e^{-0.3775}} \\ = \underline{\underline{0.59326}}.$$

h<sub>2</sub>

$$\text{inh}_2 = w_3 i_1 + w_4 i_2 + b_1 x_1 \\ = 0.25 \times 0.05 + 0.30 \times 0.10 + 0.35 \times 1 \\ = \underline{\underline{0.3925}}.$$

$$\text{outh}_2 = \frac{1}{1 + e^{-\text{inh}_2}} = \frac{1}{1 + e^{-0.3925}} \\ = \underline{\underline{0.59688}}.$$

Now, we have finished the hidden layer activation.

O<sub>1</sub>

$$\text{inO}_1 = w_5 \text{outh}_1 + w_6 \text{outh}_2 + b_2 x_1 \\ = 0.40 \times 0.59326 + 0.45 \times 0.59688 \times 0.60 \times 1 \\ = 1.105$$

$$\text{out01} = \frac{1}{1 + e^{-\text{in01}}} = \frac{1}{1 + e^{-1.105}} \\ = 0.7511 \quad (\text{expected is } 0.01)$$

Q2.

$$\text{in02} = w_7^{\text{out}} \text{in1} + w_8^{\text{out}} \text{in2} + b_2 \times 1 \\ = 0.50 \times 0.59326 + 0.55 \times 0.59688 + \\ 0.60 \times 1 \\ = \underline{\underline{1.22484}}$$

$$\text{out02} = \frac{1}{1 + e^{-1.22484}} = 0.7727 \quad (\text{expected is } 0.99)$$

→ Everything till now is forward propagation.

Now compute the error using MSE

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

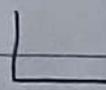
$$\text{Etotal} = \frac{1}{n} [(target_{01} - \text{out01})^2 + (target_{02} - \text{out02})^2] \\ = \frac{1}{2} [(0.01 - 0.7511)^2 + (0.99 - 0.7727)^2] \\ = \underline{\underline{0.2982}} \rightarrow \text{Total Error}$$

→ Now we use backpropagation to adjust the error.

→ First, we update weight  $w_5, w_6, w_7, w_8$  since it's closer.

$\frac{\partial E_{\text{total}}}{\partial w_5}$  (What is the effect on total error when we change  $w_5$ )  
 (Rate of change of error total w.r.t to  $w_5$ )

$$\Rightarrow \frac{\partial E_{\text{total}}}{\partial w_5} = \frac{\partial E_{\text{total}}}{\partial \text{out}_{01}} \times \frac{\partial \text{out}_{01}}{\partial w_5} \times \frac{\partial \text{in}_{01}}{\partial w_5}$$



(Chain Rule)

$$\textcircled{1} \quad \frac{\partial E_{\text{total}}}{\partial \text{out}_{01}} = \frac{\partial}{\partial \text{out}_{01}} \left[ \frac{1}{2} [\text{target}_{01} - \text{out}_{01}]^2 + [\text{target}_{02} - \text{out}_{02}]^2 \right]$$

since the second part doesn't contain  $\text{out}_{01}$

$$\frac{\partial}{\partial \text{out}_{01}} = \frac{1}{2} [(\text{target}_{01} - \text{out}_{01})^2]$$

we can ignore the rest.

$$\frac{\partial}{\partial \text{out}_{01}} = \frac{1}{2} \cdot \frac{1}{2} (\text{target}_{01}^2 - 2 \times \text{target}_{01} \times \text{out}_{01} + \text{out}_{01}^2)$$

net = in

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

$$= \frac{\partial}{\partial \text{out}_{01}} \cdot \left[ \frac{1}{2} (\text{out}_{01}^2 - 2 \times \text{target}_{01} \times \text{out}_{01}) \right]$$

$$= \frac{\partial}{\partial \text{out}_{01}} = \left[ \frac{1}{2} (2\text{out}_{01} - 2\text{target}_{01}) \right]$$

$$= \frac{\partial}{\partial \text{out}_{01}} = \text{out}_{01} - \text{target}_{01}$$

$$= 0.7511 - 0.01$$

$$= \underline{\underline{0.7411}}$$

$$\textcircled{2} \quad \frac{\partial \text{out}_{01}}{\partial \text{in}_{01}} = \frac{\partial}{\partial \text{in}_{01}} \left[ \frac{\text{out}_{01}}{1 + e^{-\text{in}_{01}}} \right] = \text{out}_{01} (1 - \text{out}_{01}) \\ = 0.7511 (1 - 0.7511) \\ = \underline{\underline{0.1869}}$$

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

$\sigma'(x) : \sigma(x) (1 - \sigma(x))$  Derivative

$$\textcircled{3} \quad \frac{\partial \text{in}_{01}}{\partial w_5} = \frac{\partial}{\partial w_5} [\omega_5 \text{out}_1 + \omega_6 \text{out}_2 + b_2 \times 1]$$

$$= \text{out}_1 + 0 + 0$$

$$= \text{out}_1$$

$$= 0.5932$$

$$\frac{\partial E_{\text{total}}}{\partial w_5} = 0.3411 \times 0.1869 \times 0.5932 \\ = \underline{0.0821}$$

## Gradient Descent

$$w_5' = w_5 + \Delta w_5$$

$$\Delta w_5 = -\alpha \frac{\partial L}{\partial w_5} \quad [L \rightarrow \text{loss function}] \\ \text{here it is MSE function}$$

Assume that  $\alpha = 0.5$

$$\Delta w_5 = -0.5 \times \frac{\partial E_{\text{total}}}{\partial w_5} \\ = -0.5 \times 0.0821 \\ = -0.041$$

$$w_5' = w_5 + \alpha \Delta w_5 \\ = 0.40 + -0.041 \\ = \underline{0.359}$$

↳ updated value of  $w_5$

[ $w_6, w_7, w_8$  can be calculated using this procedure]

W1

$$\frac{\partial E_{\text{total}}}{\partial w_1} = \frac{\partial E_{\text{total}}}{\partial \text{out}_01} \times \frac{\partial \text{out}_01}{\partial w_1} \times \frac{\partial \text{in}_01}{\partial w_1} \times \frac{\partial \text{out}_{\text{hi}}}{\partial \text{in}_{\text{hi}}} \\ * \quad \frac{\partial \text{out}_{\text{hi}}}{\partial w_1}$$

① and ② we have already found out.

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_01} = 0.7411 \quad \frac{\partial \text{out}_01}{\partial w_1} = 0.1869$$

$$\begin{aligned} \textcircled{3} \quad \frac{\partial \text{in}_01}{\partial w_1} &= \frac{\partial}{\partial w_1} [w_{\text{south}_1} + w_{\text{south}_2} + b_2 \times 1] \\ &= \text{outh}_1 + 0 + 0 \\ &= \text{outh}_1 \\ &= w_5 + 0 + 0 \\ &= w_5 \\ &= \underline{0.40} \end{aligned}$$

$$\textcircled{4} \quad \frac{\partial \text{outh}_1}{\partial w_1} = \frac{\partial}{\partial w_1} \left[ \frac{1}{1 + e^{-\frac{w_1 + b_1}{0.59326}}} \right] = \text{outh}_1(1 - \text{outh}_1) \\ = 0.59326(1 - 0.59326) \\ = \underline{0.24130}$$

$$\textcircled{5} \quad \frac{\partial h_1}{\partial w_1} = \frac{\partial}{\partial w_1} [w_{11} + w_{21}x_2 + b_1x_1] \\ = \cancel{w_{11}} + 0 + 0 \\ = \underline{\underline{0.05}}$$

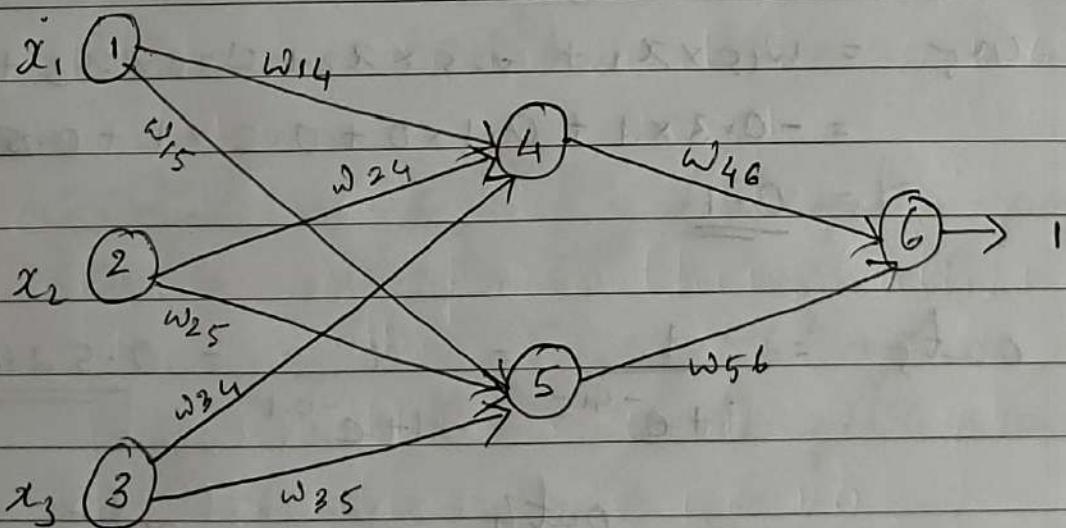
$$\frac{\partial E_{\text{total}}}{\partial w_1} = 0.7411 \times 0.1869 \times 0.40 \times \\ 0.24130 \times 0.05 \\ = \underline{\underline{0.0006685}}$$

$$w_1' = w_1 + \Delta w_1$$

$$\Delta w_1 = -\alpha \frac{\partial L}{\partial w_1} \\ = -\alpha \times \frac{\partial E_{\text{total}}}{\partial w_1} \\ = -0.5 \times 0.0006685 \\ = \underline{\underline{-0.00033425}}$$

$$w_1' = 0.15 + -0.00033425 \\ = \underline{\underline{0.1496}}$$

Q2. For the given multilayer feed-forward neural network, let the learning rate be 0.9. The initial weight and bias values of the network are given in Table along with the first training tuple;  $x = \{0, 1, 0\}$  with a class label of 1. Find the updated weights and biases after 1 epoch.



$x_1$	$x_2$	$x_3$	$w_{14}$	$w_{15}$	$w_{24}$	$w_{25}$
1	0	1	0.2	-0.3	0.4	0.1

$w_{34}$	$w_{35}$	$w_{46}$	$w_{56}$	$\theta_4$	$\theta_5$	$\theta_6$
-0.5	0.2	-0.3	-0.2	-0.4	0.2	0.1

## Forward propagation

$$\begin{aligned}in_4 &= w_{14} \times x_1 + w_{24} \times x_2 + \cancel{\theta_4} w_{34} \times x_3 + \theta_4 \\&= 0.2 \times 1 + 0.4 \times 0 + \cancel{0.2} - 0.5 \times 1 + -0.4 \\&= \underline{-0.7}\end{aligned}$$

$$out_4 = \frac{1}{1+e^{-in_4}} = \frac{1}{1+e^{-(0.7)}} = 0.3318$$

$$\begin{aligned}in_5 &= w_{15} \times x_1 + w_{25} \times x_2 + w_{35} \times x_3 + \theta_5 \\&= -0.3 \times 1 + 0.1 \times 0 + 0.2 \times 1 + \cancel{0.2} 0.2 \\&= \underline{0.1}\end{aligned}$$

$$out_5 = \frac{1}{1+e^{-in_5}} = \frac{1}{1+e^{-0.1}} = 0.5249$$

$$\begin{aligned}in_6 &= w_{46} \times \cancel{out_4} + \cancel{\theta_6} w_{56} \times \cancel{out_5} + \theta_6 \\&= -0.3 \times -0.7 + \cancel{0.2} - 0.2 \times 0.1 + 0.1 \\&= \underline{-0.1045}\end{aligned}$$

$$out_6 = \frac{1}{1+e^{-in_6}} = 0.4738 \cdot (\text{Expected } 1)$$

$n \rightarrow$  no. of output

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

$$\begin{aligned}
 \text{MSE} &= \frac{1}{n} [(target_6 - out_6)^2] \\
 &= \frac{1}{21} [(1 - 0.4738)^2] \\
 &= (1 - 0.4738)^2 = \underline{\underline{0.2768}}.
 \end{aligned}$$

## Backpropagation

$$\frac{\partial E_{\text{total}}}{\partial w_{46}} = \frac{\partial E_{\text{total}}}{\partial out_6} \times \frac{\partial out_6}{\partial w_{46}} \quad \begin{matrix} ① \\ ② \\ ③ \end{matrix} \quad \begin{matrix} \partial E_{\text{total}} \\ \partial out_6 \\ \partial w_{46} \end{matrix} \quad \begin{matrix} \partial out_6 \\ \partial w_{46} \\ \partial w_{46} \end{matrix}$$

$$\begin{aligned}
 ① \frac{\partial E_{\text{total}}}{\partial out_6} &= \frac{\partial}{\partial out_6} \left[ \frac{1}{2} (target_6 - out_6)^2 \right] \\
 &= \cancel{\frac{\partial}{\partial out_6} \left[ target_6^2 + out_6^2 / 2 \times target_6 \times out_6 \right]} \\
 &= \cancel{\frac{\partial}{\partial out_6} [out_6^2 - 2 \times target_6 \times out_6]} \\
 &= \frac{\partial}{\partial out_6} [target_6^2 + out_6^2 - 2target_6 \cdot out_6] \\
 &= \frac{\partial}{\partial out_6} [target_6^2 + out_6^2 - out_6] \\
 &= 2out_6 - 2target_6 \\
 &= 2[0.474 - 1] \\
 &= \underline{\underline{-1.052}}
 \end{aligned}$$

$$\begin{aligned}
 \textcircled{2} \quad \frac{\partial \text{out}_6}{\partial \text{in}_6} &= \frac{\partial}{\partial \text{in}_6} \left( \frac{1}{1+e^{-\text{in}_6}} \right) \\
 &= \text{out}_6 - (1 - \text{out}_6) \\
 &= 0.474 - (1 - 0.474) \\
 &= \underline{-0.052}
 \end{aligned}$$

$$\textcircled{3} \quad \frac{\partial \text{in}_6}{\partial w_{46}} = \frac{\partial}{\partial w_{46}} \left[ \text{out}_4 \cdot w_{46} + \text{out}_5 \cdot w_{56} + b_6 \right]$$

$$= \text{out}_4$$

$$= 0.331$$

$$-1.052 \times \underline{-0.052 \times 0.331} = 0.0181$$

$$w_{46}' = w_{46} + \Delta w_{46}$$

$$\Delta w_{46} = -\alpha \times \delta E_{total}$$

$$= -0.9 \times 0.0181$$

$$= -0.01629$$

$$w_{46}' = -0.3 + -0.01629$$

$$= \underline{-0.31629}$$

## Gradient Descent

→ Optimization algorithms

→ We are trying to find the equation of the best fit line

$$y = \underset{\text{feature}}{\overset{w_1}{\nearrow}} m x_i + c - \underset{w_0}{\swarrow} w_0$$

this  
is what  
we're  
trying to  
predict

→ We use error function to find the error.

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_{pred})^2$$

Replace  $\hat{y}_{pred}$  with  $m x_i + c$

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - (m x_i + c))^2$$

↓  
We have to minimize this cost → gradient descent.

$$\frac{2m^2x_i^2}{2m^2} = \frac{2mx_i^2}{\text{derivative}}$$

$$2m^2 - 1$$

→ Predicted value  $y_{\text{pred}} = m x_i + c$ .

$$\text{MSE (Cost fn L)} = \frac{1}{n} \sum_{i=1}^n (y_i - y_{\text{pred}})^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i - (mx_i + c))^2$$

$$= \frac{1}{n} \sum_{i=1}^n (y_i^2 + (mx_i + c)^2 - 2y_i(mx_i + c))$$

$$L = \frac{1}{n} \sum_{i=1}^n (y_i^2 + m^2x_i^2 + c^2 + 2mx_i c - 2y_i mx_i - 2y_i c)$$

$$\frac{\partial L}{\partial m} = \frac{1}{n} \sum_{i=1}^n (0 + 2mx_i^2 + 0 + 2x_i c - 2y_i x_i - 0)$$

$$\frac{\partial L}{\partial m} = \frac{1}{n} \sum_{i=1}^n 2x_i (\underbrace{mx_i + c - y_i}_{y_{\text{pred}}})$$

$$\checkmark \frac{\partial L}{\partial m} = -2 \sum_{i=1}^n x_i (y_i - y_{\text{pred}})$$

$$\left. \begin{aligned} m' &= m + \Delta m \\ \Delta m &= m - \alpha \frac{\partial L}{\partial m} \end{aligned} \right\}$$

$$\frac{\partial L}{\partial c} = \frac{1}{n} \sum_{i=1}^n (0 + 0 + 2c + 2mx_i - 0 - 2y_i)$$

$$\frac{\partial L}{\partial c} = \frac{1}{n} \sum_{i=1}^n (2c + 2mx_i - 2y_i)$$

$$\frac{\partial L}{\partial c} = \frac{1}{n} \sum_{i=1}^n 2(c + mx_i - y_i)$$

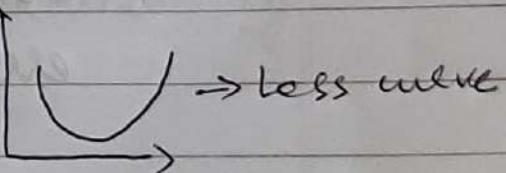
$$\frac{\partial L}{\partial c} = \frac{1}{n} \sum_{i=1}^n 2 \left( \underbrace{(mx_i + c)}_{y_{pred}} - y_i \right)$$

$$\checkmark \frac{\partial L}{\partial c} = -2 \sum_{i=1}^n (y_i - y_{pred})$$

$$\begin{aligned} \checkmark c' &= c + \Delta c \\ &= c - \alpha \frac{\partial L}{\partial c} \end{aligned}$$

on different shapes

- \* if we choose larger  $\alpha$  or learning rate, gradient descent can overshoot.
- \* if we choose lower  $\alpha$ , gradient descent will take small steps to reach local minima and will take a longer time to reach minimum.



- \* Binary - Sigmoid
  - \* Multi class - Softmax
  - \* Regression - linear regression
- } Activation functions

Variants of gradient descent:

- Stochastic - one forward propagation, computing the error and one back propagation
- Mini-batch - using the computed error to back-propagate and update the weights.
- Batch

→ If there are 4 samples, so 4 backpropagation will be there.

→ finding the average error to back-propagate. If there are two samples and combine the errors and backpropagate. Samples will be divided into batches: No. of weight updation will be less.

→ Weight will be updated at the end of the iteration. Only one backpropagation will be there using the combined error.

## Vanishing and Exploding gradient problem

- This happens with deep neural network (more than 2 hidden layers)
- When the updated weight is very minimal, it is negligible.
- When the gradient is very minimal, weight updation is negligible.
- \* → If its a ~~more~~ complex problem, we need to use deep neural network.
- It often happens when we use sigmoid and tanh activation functions.

### Exploding gradient:

- When the weight updation is too large, it is called exploding gradient
- It happens with relu activation function.

### Activations Function

- Transform the weighted sum to required output.
- Purpose: To put some non-linearity to the neural network.

Types of activation fn:  
 → Binary Step function

$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ 1 & \text{for } x \geq 0 \end{cases}$$

- Cannot use for multi-class classification problem

→ Linear Activation function-

$$f(x) = x$$

- Weighted sum is the output.
- Used for regression function.  
 (continuous values)
- Will not be using in hidden layers because the derivative will be 0.

→ Non-Linear Activation Function:

1. Sigmoid activation

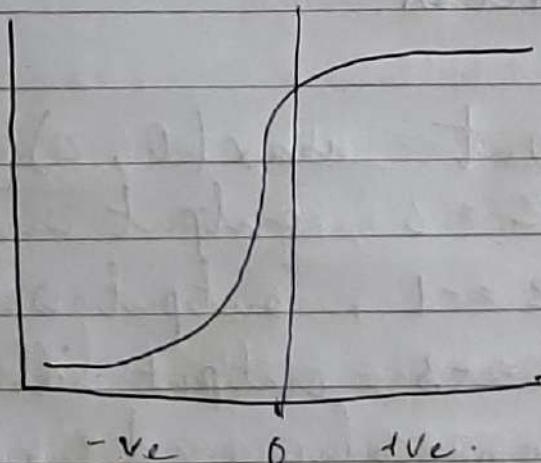
$$= \phi(z) = \frac{1}{1+e^{-z}}$$

- The output will range from 0 to 1.
- Value of  $z$  depends on the value of  $x$ .
- Used for binary classification problems.

$x$	Output
3	sigmoid
1.75	0.95
-2	0.12
0.5	0.62

↓

It is never a probability curve.



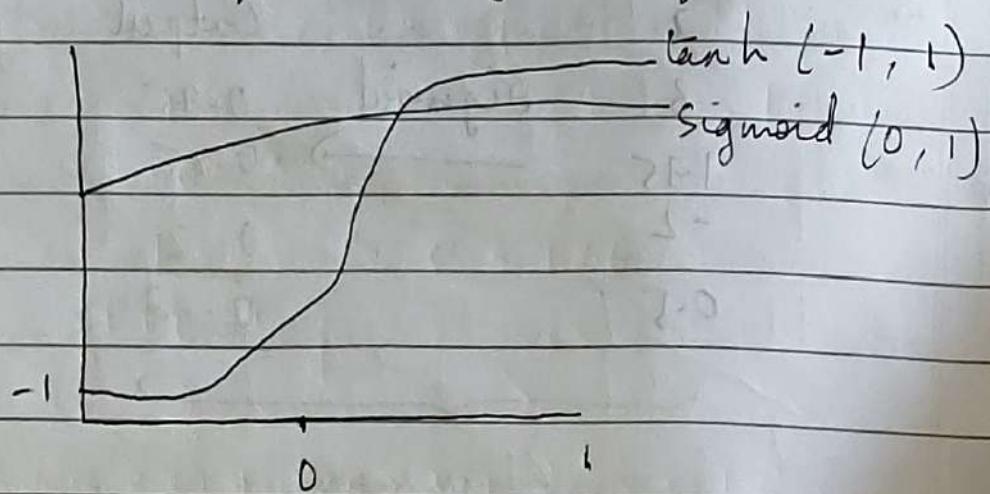
→ If value of  $z$  is -ve, the fn will push it toward 0 and if its +ve, the value is pushed toward 1.

## 2. Tanh Function:

$$= \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

→ Output ranges from -1 to +1.

→ Used for binary classification.



## 3. ReLU function:

→ Output =  $\max(0, z)$

if  $z = 3$ , output will be 3.

$z = -1$ , output will be 0.

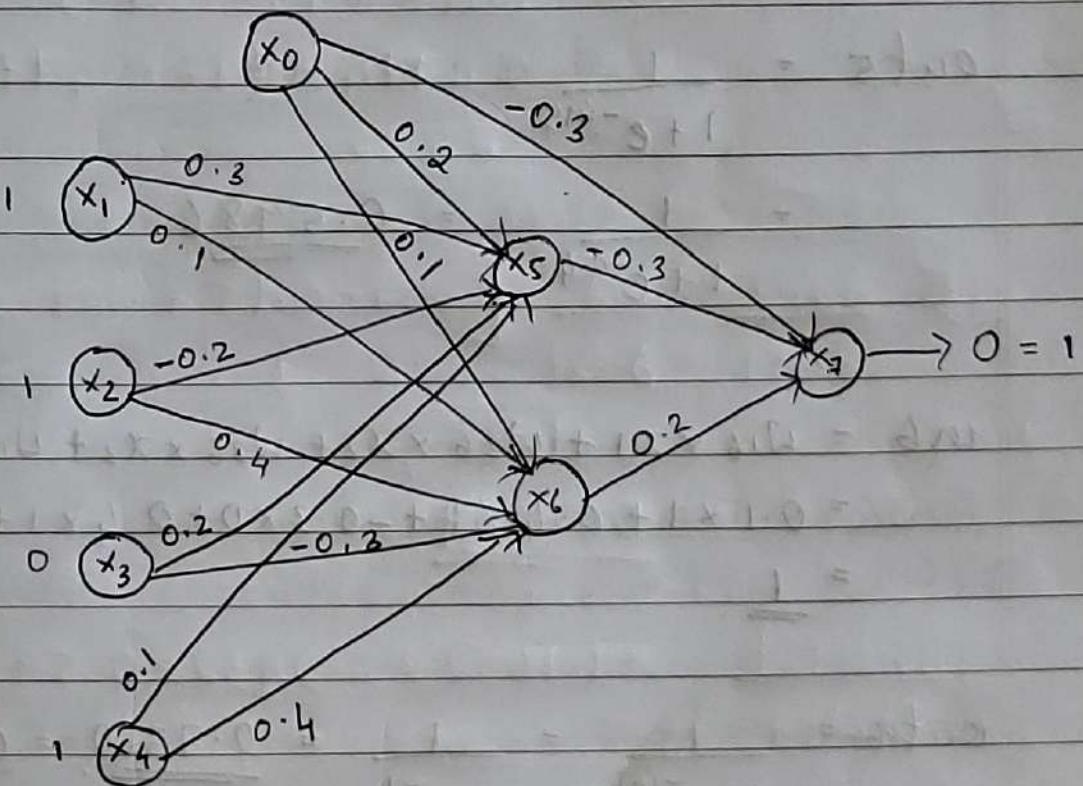
$z = -3$ , output will be 0.

→ Most commonly used hidden layer activation function.

Vismaya Dinesh

Assignment.

- Q. For the ANN shown below, compute the updated weights of  $w_{57}$ ,  $w_{67}$  and  $w_{07}$ . Assume learning rate = 0.8, activation function = sigmoid, output  $O = 1$ .



$x_1$	$x_2$	$x_3$	$x_4$
1	1	0	1

$x_1$	$x_2$	$x_3$	$x_4$
1	1	0	1

## Forward Propagation:

$$\begin{aligned}
 \text{in}_5 &= w_{15} \times x_1 + w_{25} \times x_2 + w_{35} \times x_3 + \\
 &\quad w_{45} \times x_4 + b_5 \\
 &= 0.3 \times 1 + -0.2 \times 1 + 0.2 \times 0 + 0.1 \times 1 + 0.2 \\
 &= \underline{\underline{0.4}}
 \end{aligned}$$

$$\begin{aligned}
 \text{out}_5 &= \frac{1}{1 + e^{-\text{in}_5}} \\
 &= \frac{1}{1 + e^{-0.4}} = \underline{\underline{0.5986}}
 \end{aligned}$$

$$\begin{aligned}
 \text{in}_6 &= w_{16} \times x_1 + w_{26} \times x_2 + w_{36} \times x_3 + w_{46} \times x_4 + b_6 \\
 &= 0.1 \times 1 + 0.4 \times 1 + -0.3 \times 0 + 0.4 \times 1 + 0.1 \\
 &= \underline{\underline{1}}
 \end{aligned}$$

$$\text{out}_6 = \frac{1}{1 + e^{-\text{in}_6}} = \frac{1}{1 + e^{-1}} = \underline{\underline{0.7310}} = \underline{\underline{0.7311}}$$

$$\begin{aligned}
 \text{in}_7 &= w_{57} \times \text{out}_5 + w_{67} \times \text{out}_6 + b_7 \\
 &= -0.3 \times 0.5986 + 0.2 \times 0.7311 + -0.3 \\
 &= -0.17961 + 0.14622 - 0.3 \\
 &= \underline{\underline{-0.3334}}
 \end{aligned}$$

$$\text{out}_7 = \frac{1}{1+e^{-\text{int}}} = \frac{1}{1+e^{-0.3334}} = \underline{0.4174} \cdot (\text{Expected})$$

$$\begin{aligned} \text{MSE} &= \frac{1}{2n} (\text{target}_7 - \text{out}_7)^2 \\ &= \frac{1}{2} (1 - 0.4174)^2 \\ &= \underline{0.3394} \rightarrow E_{\text{total}} \end{aligned}$$

### Back propagation

w<sub>57</sub>:

$$\frac{\partial E_{\text{total}}}{\partial w_{57}} = \frac{\partial E_{\text{total}}}{\partial \text{out}_7} \times \frac{\partial \text{out}_7}{\partial \text{int}} \times \frac{\partial \text{int}}{\partial w_{57}}$$

$$\begin{aligned} \textcircled{1} \quad \frac{\partial E_{\text{total}}}{\partial \text{out}_7} &= \frac{\partial}{\partial \text{out}_7} [(\text{target}_7 - \text{out}_7)^2] \\ &= \frac{\partial}{\partial \text{out}_7} [\text{target}_7^2 + \text{out}_7^2 - 2 \text{target}_7 \text{out}_7] \\ &= (2 \text{out}_7 - 2 \text{target}_7) \\ &= -2 (\text{target}_7 - \text{out}_7) \\ &= -2 (0.4174 - 0.3394) \end{aligned}$$

$$= -2(1 - 0.4174)$$

$$= \underline{-1.1652}$$

$$\textcircled{2} \quad \frac{\partial \text{out}_7}{\partial \text{in}_7} = \frac{\partial}{\partial \text{in}_7} \left[ \frac{1}{1 + e^{-\text{in}_7}} \right]$$

$$= \text{out}_7 * (1 - \text{out}_7)$$

$$= 0.4174 * (1 - 0.4174)$$

$$= \underline{0.2432}$$

$$\textcircled{3} \quad \frac{\partial \text{in}_7}{\partial w_{57}} = \frac{\partial}{\partial w_{57}} [w_{57} \times \text{out}_5 + w_{67} \times \text{out}_6 + w_{07}]$$

$$= \text{out}_5$$

$$= \underline{0.5987}$$

$$\frac{\partial E_{\text{total}}}{\partial w_{57}} = -1.1652 \times 0.2432 \times 0.5987$$

$$= \underline{-0.1697}$$

$$w_{57}' = w_{57} + \Delta w_{57}$$

$$\Delta w_{57} = -\alpha \times \frac{\partial E_{\text{total}}}{\partial w_{57}}$$

$$= -0.8 \times -0.1697$$

$$= 0.13576$$

$$\begin{aligned} w_{57}' &= -0.3 + 0.13576 \\ &= \underline{-0.16424} \end{aligned}$$

w<sub>67</sub>:

$$\frac{\partial E_{\text{total}}}{\partial w_{67}} = \frac{\partial E_{\text{total}}}{\partial \text{out}_7} \times \frac{\partial \text{out}_7}{\partial \text{int}} \times \frac{\partial \text{int}}{\partial w_{67}}$$

$$\frac{\partial E_{\text{total}}}{\partial \text{out}_7} = -1.1652 \quad \frac{\partial \text{out}_7}{\partial \text{int}} = 0.2432.$$

$$\begin{aligned} \frac{\partial \text{int}}{\partial w_{67}} &= \frac{\partial}{\partial w_{67}} [w_{57} \times \text{out}_5 + w_{67} \times \text{out}_6 + w_{07}] \\ &= \text{out}_6 \\ &= \underline{0.7311} \end{aligned}$$

$$\begin{aligned} \frac{\partial E_{\text{total}}}{\partial w_{67}} &= -1.1652 \times 0.2432 \times 0.7311 \\ &= \underline{-0.2072} \end{aligned}$$

$$\begin{aligned} w_{67}' &= w_{67} + \Delta w_{67} \\ &= 0.2 + (-0.8)(-0.2072) \\ &= 0.2 + 0.16576 \\ &= \underline{0.36576} \end{aligned}$$

w<sub>07</sub>

$$\frac{\partial E_{\text{total}}}{\partial w_{07}} = \frac{\partial E_{\text{total}}}{\partial \text{out}_7} \times \frac{\partial \text{out}_7}{\partial w_{07}} \times \frac{\partial \text{in}_7}{\partial w_{07}}$$

$$\frac{\partial \text{in}_7}{\partial w_{07}} = \frac{\partial}{\partial w_{07}} [w_{57} \times \text{out}_5 + w_{67} \times \text{out}_6 + w_{07}] \\ = \underline{\underline{1}}$$

$$\frac{\partial E_{\text{total}}}{\partial w_{07}} = -1.1652 \times 0.2432 \times 1 \\ = -\underline{\underline{0.2834}}$$

$$w_{07}' = w_{07} + \Delta w_{07} \\ = -0.3 + (-0.8)(-0.2834) \\ = -0.3 + 0.22672 \\ = -\underline{\underline{0.07328}}$$

	old value	updated value
w <sub>57</sub>	-0.3	-0.16424
w <sub>67</sub>	0.2	0.36576
w <sub>07</sub>	-0.3	-0.07328

#### 4. Leaky RELU:

- Attempt to solve the dying RELU problem.
- Helps increase the range of RELU.
- A small portion of negativeness is bypassed.
- leaky factor  $a = 0.01$
- When  $a$  is not 0.01, it is called Randomized RELU.

#### Hidden layers activation fns:

- RELU (Rectified Linear Activation)
- Logistic (Sigmoid)
- Tanh (Hyperbolic tangent)

#### Output layer activation fns:

- Linear
- Sigmoid
- Softmax

- For multi-class classification problems, we use Sigmoid Softmax.
- For multi-label classification problem, we use Sigmoid.

$$\text{Softmax fn } \pi_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}} \text{ for } j=1 \dots k$$

→ ~~After~~ The output can be  
interpreted as a probability  
distribution.

→ Whichever class having the  
higher probability will be  
the output or predicted class.

→ Sum of the probabilities should be 1.

### Loss / Cost Functions:

→ Choice of cost function depends on  
classification or regression

→ For regression, MSE (Mean Squared  
Error) is used.

→ Mean Absolute Error (MAE), also  
used for regressions.

↳ This is used when there are  
outliers in the data.

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

↳ MAE treats all errors equally.

→ For binary classification problem, Binary cross entropy (Log loss) is used

$$\hookrightarrow \frac{1}{n} \sum_{i=1}^n [y_i \log(y_{\text{pred}}) + (1 - y_i) \log(1 - y_{\text{pred}})]$$

→ For multi-class classification, Categorical cross entropy. (one-hot encoding is used)

$$\hookrightarrow - \sum_{i=1}^n y_i \log(\pi_i y_{\text{pred}})$$

→ Another loss function is Sparse categorical cross entropy, used for multi-class classification. But instead of one-hot encoding, the labels are represented as integers itself.  
↳ Used when there is a huge number of classes.

\* One-hot encoding is used to eradicate any differences in between the class labels.

1 - 1 0 0    3 - 0 0 1

2 - 0 1 0

→ Loss function  
→ Optimizes  
→ Evaluation metrics } needed for defining a model.

→ Kullback - Leibler Divergence ( $\ell_2$  Divergence).

↳ Used for unsupervised learning.  
↳ finds the difference b/w two probabilities.

→ Hinge Loss.

↳ Used for SVM required tasks.  
↳ Where margin maximization is important.

Optimizers.

→ Used to adjust weights.

→ Eg: Gradient descent.

→ Stochastic Gradient descent (SGD)  
↳ Used when data is plentiful and simple.

↳ Slow convergence but leads to good generalization. It should fit different or various data sets.

→ SGD with Momentum.

↳ Used when the optimization is slow or get stuck in local minima.

\* We are trying to find global minima with lowest error.

→ RMSprop (Root Mean Square)

↳ Great for mini batch training and noise.

→ Adam (Adaptive Moment Estimation)

↳ General-purpose and widely used.

→ Adagrad

↳ Useful for sparse data

→ Adadelta - for preventing slowdowns

→ Nadam - when you need fast convergence.

## Module - II

### Deep Learning:

→ Uses many ANNs to learn from data.

→ Why DL?

↳ Can learn complex data

↳ Large dataset training

→ Data-driven learning

↳ Automatically learns from data. Manual assistance is very low.

### Convolutional Neural Network (CNN)

→ Working with images and manually extracting features from them is very tedious.

→ That is why automatic extraction of features was introduced.

→ CNN model works in two steps:

↳ Feature extraction

↳ Prediction

→ Fully connected Layer / Dense : Every neuron is connected to every other neuron in the next layer.

Vennote

Building blocks of CNN:

- Input Layer
- Convolution Layer
- Flattening Layer
- Pooling Layer
- Fully Connected Layer

Input Layer: values

→ Gray scale - bits will range from 0 - 255. If 8 bits are stored.

00000000 - 0

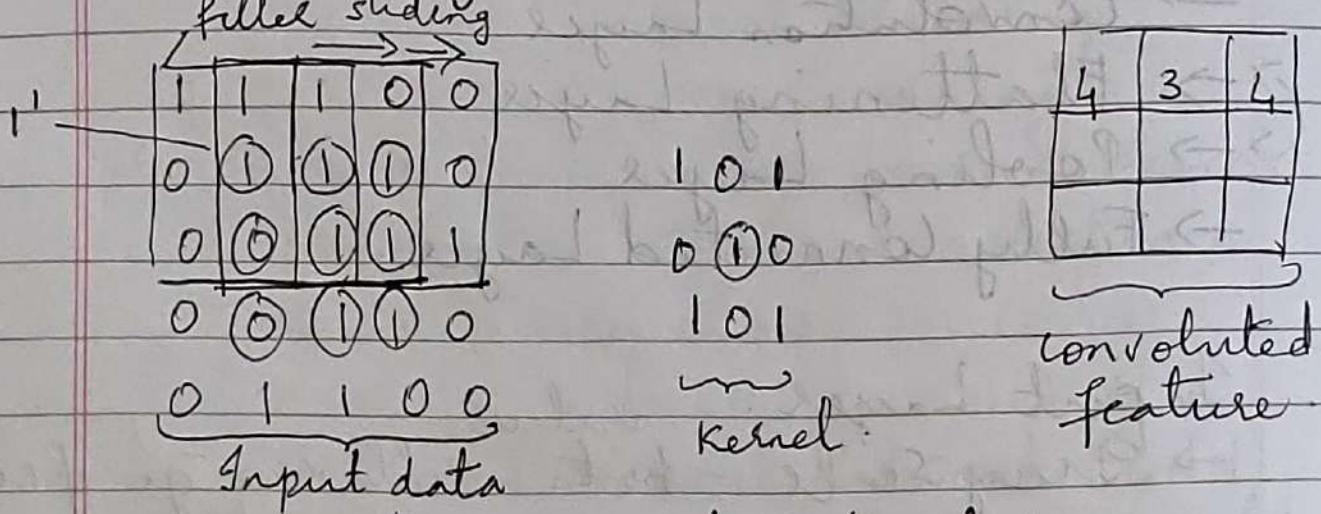
11111111 -  $2^8$  - 255.

→ Black and white binary pixel - 0/1

→ If it's a color image, a full red image will have pixels indicating the shade of red from 0 to 255.

→ Convolutional Layers:

- Applies filters over the image (kernel or masks) to extract features.
- Output of this will be edges, corner or texture.
- Sliding a kernel over the input.



→ element-wise multiplication

$$\begin{aligned}
 l_1' &= 1 \times 1 + 1 \times 0 + 1 \times 1 + 0 \times 0 + 1 \times 1 + 1 \times 0 + \\
 &\quad 0 \times 1 + 0 \times 0 + 1 \times 1 \\
 &= \underline{\underline{4}}
 \end{aligned}$$

$$l_2' = \underline{\underline{3}}$$

$$l_3' = \underline{\underline{4}}$$

Why should we move this kernel?

→ Based on type values in the mask, we will extract different types features.

→ Sobel vertical and sobel horizontal are edge detection masks or kernel.

Padding:

→ Entire image in the boundary will have zero - zero-padding.

→ This is done to make the response map and the actual image ~~is~~ the same size.  
→ convolved or convulated feature.

Stride:

→ How quickly window/filter/kernel slides

→ Stride 2 means, window moves by 2 pixels at a time

\* Filter bank: A set of filters are applied on the image to extract features.

→ Activation Function

↳ ReLU : Rectified linear unit.

It sets all negative values to zero.

$$R(z) = \max(0, z)$$

Introduces non-linearity.

and output is a feature map.

→ Pooling Layer:

↳ Trying to aggregate features

↳ Downsamples the feature

maps to reduce dimensionality.

↳ Here also a mask is applied

↳ max pooling and average pooling

↳ A blank mask is applied.

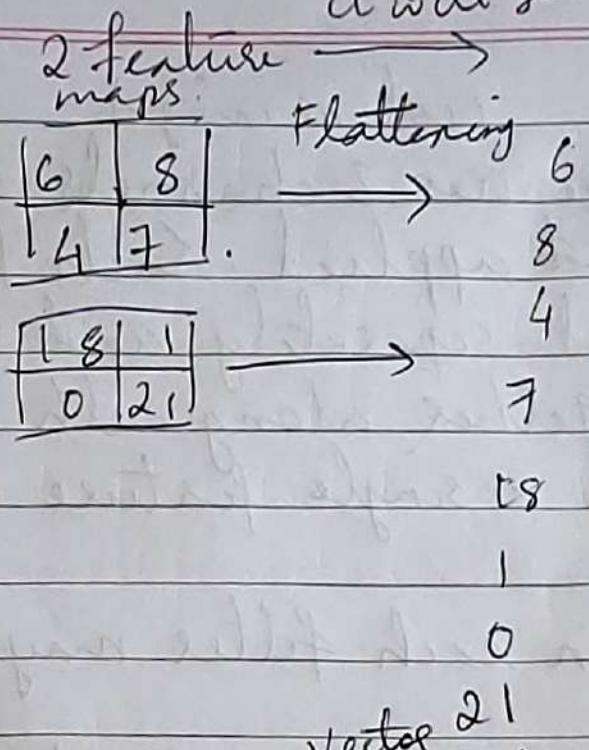
→ Flattening Layer:

↳ Now we have many feature maps after pooling.

↳ The 2D array is converted to a 1D array vector

↳ All the 2D feature maps are converted.

it will get appended.



- ↳ The feature that we get will become the "input to ANN".
- ↳ It will pass through the fully connected layers.

- \* Feature extraction is automated in CNN.
- \* For text input, it will be 1D array. So the mask & size will also be one dimensional.
- Convolution in colour image:
- \* Multi channel image will have multiple a filter bank (no. of filters according to the no. of channels) and the outcome will be a single feature map.

(color image)

- \* If an image has 3 channels, 3 kernels will be applied. Each kernel is calculated separately and is added together along with a bias to produce a single feature map.
- \* The values in each filter may not be same.
- \* Algorithm will decide the values in the filter. This is similar to the weights in ANN. The algorithm will learn on its own if there is any change.

What all parameters determine the convolution operation?

- Size of the ip image
- stride
- padding.

To find the size of convolution image:

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

P - no. of zero-padding

S - stride

F - filter size

W - size of the i/p image

So if there are n-number  
of channels,

then  $W_{out} \times n$  is used.

Pooling layer.

$$W_{out} = \frac{W - F}{S} + 1$$

W - response of the convolution

F - Pooling filter

Conv 1

(Q) Input size =  $28 \times 28 \times 1$ .

Four hyperparameters

No. of kernels,  $k = 16$

size of each kernel = 5 (F)

stride,  $s = 1$

padding,  $P = 2$ .

$$W_2 = W - F + 2P + 1$$

S

$$= 28 - 5 + 2 \times 2 + 1$$

1

$$= \underline{\underline{28}}$$

$$H_2 = 28$$

Output will be  $28 \times 28 \times 16$ . } input to the  
 pooling layer  
 because 2 filters  
 are applied.

Pool 1

$$\text{Pooling: } \frac{w-f+1}{s} = \frac{28-2+1}{2} = \underline{\underline{14}}$$

$F=2$  because its given that  
a  $2 \times 2$  pooling filter.

So output of pooling is  $14 \times \underline{\underline{14}} \times 16$ .

This will be the input to the next convolutional layer: Conv 2.

$$k = 32$$

$$\text{size}(k) = 5$$

$$S = 1$$

$$P = 2$$

$$\frac{w-f+2p+1}{1} = \frac{14-5+2(2)+1}{1} = \underline{\underline{14}}$$

$$\text{Output} = 14 \times 14 \times 32$$

↓  
32 filters  
are applied

Pooling 2.  $F = 2, S = 2$ .

$$\frac{W - F + 1}{S} = \frac{14 - 2 + 1}{2} = \frac{7}{2}$$

Output =  $7 \times 7 \times 32$ . [  $32$   $7 \times 7$  images will be the input to the next layer, flattening layer].

→ The flattening layer will have  $7 \times 7 \times 32$  elements because it is a one-dimensional array.

→ This will be the input to the fully connected layer.

Computation of number of parameters  
least:

For a convolution layer:

→ Will depend on the filter size.

Parameters =  $(F_w \times F_h \times \text{no. of channels} + \text{bias}) * K$ .

Example, if it's a colour image

Parameters =  $(3 \times 3 \times 3 + 1) \times 5 = \underline{140}$ :

A  $3 \times 3$  color image with 3 channels.  
and 5  $3 \times 3$  filters are applied.

Parameters  $\rightarrow$  weights

classmate

Date \_\_\_\_\_

Page \_\_\_\_\_

For a pooling layer:

$\rightarrow$  No parameters are learned because we are applying blank filters.

For a fully connected layer:

Parameters = (input units  $\times$  output units)  
+ output units

every neuron  
is associated  
with a bias.

128 input 64 output

$$\begin{aligned} \text{Parameters} &= (128 \times 64) + 64 \\ &= \underline{\underline{8256}}. \end{aligned}$$

Q. Compute the output shape and #parameters of each layer.

Layers:

Image ( $32 \times 32 \times 1$ )

Convolution

Kernel size	No. of kernels	Stride
-------------	----------------	--------

$5 \times 5$  6 1

Avg. pooling

$2 \times 2$  6 2

Convolution

$5 \times 5$  16 1

Avg. pooling

$2 \times 2$  16 2

Convolution

$5 \times 5$  120 1

FC

(Neurons = 84)

FC

(Neurons = 10)

\* If padding is not given, it will be 0.

Conv 1

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

$$= \frac{32 - 5}{4} + 1 = \underline{\underline{28}}$$

Output of conv1 =  $28 \times \underline{28} \times 6$ .

### Pool1

$$W_{out} = \frac{W - F + 1}{S}$$

$$= \frac{28 - 2}{2} + 1 = \underline{\underline{14}}$$

Output of pool1 =  $14 \times \underline{14} \times 6$ .

### Conv2

$$W_{out} = \frac{W - F + 2P}{S} + 1$$

$$= \frac{14 - 5}{1} + 1 = \underline{\underline{10}}$$

Output of Conv2 =  $10 \times 10 \times 16$ .

### Pool2

$$W_{out} = \frac{W - F}{S} + 1$$

$$= \frac{10 - 2}{2} + 1 = \underline{\underline{5}}$$

Output of Pool2 =  $5 \times \underline{5} \times 16$

Conv3

$$\text{Want} = \frac{W - F + 2P}{S} + 1$$

$$= \frac{5 - 5 + 1}{1} = 1$$

$$\text{Output} = 1 \times \underline{1} \times 120$$

Parameters learnt for conv1.

$$= (5 \times 5 \times 3 + 1) * 6$$

$$= \underline{156}$$

Parameters learnt for conv2:

$\overset{6}{\rightarrow}$  no. of channel of the previous layer

$$= (5 \times 5 \times 6 + 1) \times 16$$

$$= \underline{2416}$$

Parameters learnt for conv3:

$$= (5 \times 5 \times 16 + 1) \times 120$$

$$= \underline{8120} \cdot \underline{120} = \underline{48120}$$

Parameters learnt for FC-1

$$= (120 \times 84) + 84$$

$$= \underline{\underline{10164}}.$$

Parameters learnt for FC-2.

$$= (84 \times 10) + 10$$

$$= \underline{\underline{850}}.$$

$$\begin{aligned} \text{Total no. of parameters} &= 156 + 2416 \\ &\quad + 48120 + 10164 + 850 \\ &= \underline{\underline{61706}}. \end{aligned}$$

\* For exam, write the entire table with output feature map and no. of parameters.

$$* \text{ Padding} = (F - 1)/2$$

↳ to find the padding size.

F → filter size

## Q. Layers

Image ( $32 \times 32 \times 3$ ) Kernel size No. of kernel stride Padding

Conv 1  $11 \times 11$  96 4 Same

Max Pooling  $2 \times 2$  ? 2

Conv 2  $5 \times 5$  256 1 Same

Max Pooling  $2 \times 2$  ? 2

Conv 3  $3 \times 3$  384 1 Same

FC  
(Neurons = 84)

FC  
(Neurons = 10)

Conv 1

$$= W - F + 2P + 1 \quad P = (F-1)/2$$

$$= (11-1)/2 = 5$$

$$= 32 - 11 + 2(5) + 1$$

4

$$= 8 \cdot 75 = \underline{\underline{8}}$$

$$\text{Output} = 8 \times \underline{\underline{8}} \times 96.$$

Max Pooling 1

$$= \frac{W - F}{S} + 1$$

$$= \frac{8 - 2}{2} + 2 = \underline{\underline{4}}$$

$$\text{Output} = \underline{\underline{1}} \times \underline{\underline{4}} \times 96$$

Conv 2

$$= \frac{W - F + 2P}{S} + 1$$

$$P = (F - 1) / 2$$

$$= (5 - 1) / 2 = \underline{\underline{2}}$$

$$= \underline{\underline{4}} - 5 + 2(2) + 1$$

$$= \underline{\underline{4}}$$

$$\text{Output} = \underline{\underline{4}} \times \underline{\underline{4}} \times 256$$

Max Pooling 2

$$= \frac{W - F}{S} + 1 = \frac{4 - 2}{2} + 1 = \underline{\underline{2}}$$

$$\text{Output} = \underline{\underline{2 \times 2}} \times \underline{\underline{256}}$$

Conv3-

$$\begin{aligned}
 & W - F + 2P + 1 \\
 & = \frac{s^2 - 3 + 2(1) + 1}{1} \\
 & = \frac{4^2 - 3 + 2(1) + 1}{1} \\
 & = \frac{16 - 3 + 2 + 1}{1} \\
 & = 14
 \end{aligned}$$

$$\text{Output} = \frac{2 \times 2 \times}{4 \times 4 \times 384}$$

Parameters learnt in conv1:

$$= (11 \times 11 \times 3 + 1) * 96.$$

$$= \underline{\underline{38234944}}.$$

Parameters learnt in conv2:

$$= (5 \times 5 \times 96 + 1) * 256$$

$$= \underline{\underline{614656}}.$$

Parameters learnt in conv3:

$$= (3 \times 3 \times 256 + 1) * 384$$

$$= \underline{\underline{885120}}$$

Parameters learnt in FC-1

$$= (384 \times 84) + 84$$

$$= (2 \times 2 \times 384 * 84) + 84 \quad (\text{output of conv 3}) \\ = \underline{\underline{129108}}$$

Parameters learnt in FC-2

$$= (84 \times 10) + 10$$

$$= \underline{\underline{850}}$$

Total parameters learnt

$$= 34944 + 614656 + 885120 + 129108 \\ + 850$$

$$= \underline{\underline{1664678}}$$

Sparse connectivity in CNN:

→ In convolution layer, if we're using a  $3 \times 3$  filter, then we are trying to learn 9 weights only.

→ This is called sparse connectivity in CNN

→ In ANN, it is fully connected layer and depends on the input.