



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

#### UNIT-II RECOMMENDED SYSTEM

RECOMMENDED SYSTEM:- Recommender system is an information-filtering technique that provides user with recommendation for which they might be interested in.

► Recommender system acts as a solution for your day-to-day choices.

For eg- which websites will you find interesting.  
- which degree and university are best  
- for your future.

- which is the best investment for supporting the education of your children.

- which digital camera should you buy.

- which is the best holiday for you and your family

- which trending course will be helpful for your future.

- "RS are important class of ml algos that offers relevant suggestions to users."

Features - Recommender system recommends items that are most popular among all users.

- They divide users into multiple segments based on their preferences and recommends items to them based on the segment they belongs to.

### Purpose of Recommender System -

#### 1). Predictive Perspective:-

- They predicts to what degree a user likes an item.
- Most popular scenario in research.

#### 2). Retrieval Perspective -

- Reduces search costs.
- provides correct proposals.

#### 3). Interaction Perspective -

- Educate users about the complete product domain.

#### ④ Recommendation Perspective -

- identifies items from a long tail.

Recommender system reduces the information overload by estimating relevance.



# GURUKUL KANGRI VISHWAVISHTA COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO.

collaborative Filtering Recommendation:- Recommendation

what's popular among your peers.

- it uses the rating data to calculate the similarity between items.
- it does not consider the semantic relationship between different items.

collaborative method work with interaction matrix that can also be called as rating matrix. Here the task of ml algm is to learn a function that predicts the utility of item to each user.

	DCG	DC4	3Dlat
John	5	1	3
Tom	?	?	?
Nia	4	?	5

Here set of users  $U$  and items  $I$  to be recommended. To the user  $u \in U$ , learn a func based on the past data that predicts the utility of each item  $i \in I$  to each user  $u \in U$ .

- It uses User Behaviour to predict the recommended items.

content-based recommendation:- display products that are similar to the product you have liked before. i.e.

↳ it works with the data provided by user. either explicitly through rating or implicitly (by clicking on a link).

↳ It generates ~~checkbox~~ - user profile.

# LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Class/Section: III Yr. Section- A

Date: .....

Name of Faculty: Praveen Kumar Yadav

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*CONTENT BASED RECOMMENDED SYSTEM.*

## IMPORTANT & RELEVANT QUESTIONS:

## FEED BACK QUESTIONS (AFTER 20 MINUTES):

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

## REFERENCES: Text/Ref. Book with Page No. and relevant Internet Websites:

DETAILED LECTURE NOTES

Self Revision: 14:

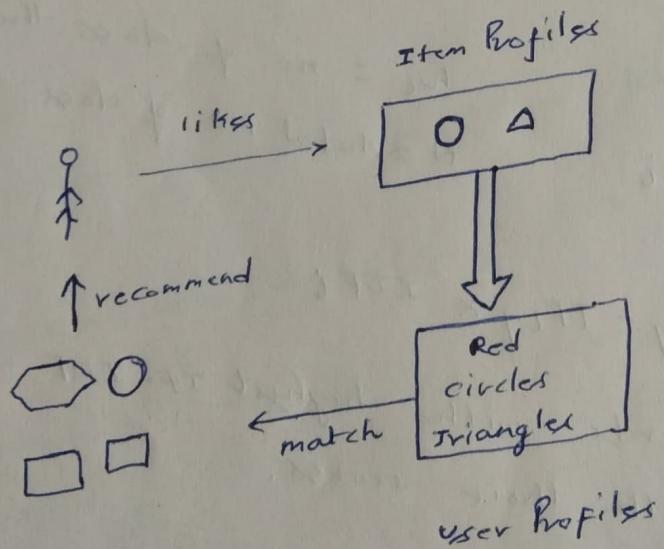
Content-Based Recommendations System:-

Main Idea - To Recommend items to customer & similar to previous items rated highly by u.

e.g. Movies -  
same Actor, director, genre

website, blogs, news  
- Articles with "similar" contents.

People  
- Recommends people with many common friends.



(Item Profiles) - Here for each item, create an item profile

- Profile is a set of features.

- Movies - author, title, actor, director.

- Images, videos - meta data and tags.

- People - set of friends.

- convenient to think of the item profile as a vector.

- one entry per feature

- vector might be boolean or real valued.

- vector might be boolean or real valued.

Text-Features - Profile - set of "important" words in item (document).

TF (Term Freq vs Inverse document Freq) :-  
if  $f_{ij}$  = freq. of an feature i in a document J.

$$TF_{ij} = \frac{f_{ij}}{\max_k f_{kj}} \quad (\text{feature } i \text{ appearing in other documents})$$

$$IDF_i = \log \frac{N}{n_i}$$

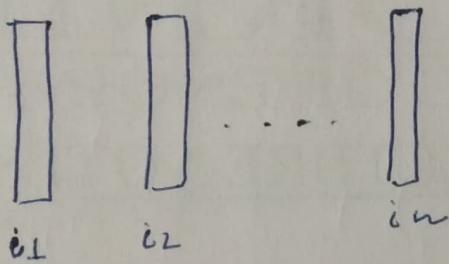
$n_i$  = no of docs that mention term i  
 $N$  = total no of docs

$$\text{TF-IDF score } w_{ij} = TF_{ij} * IDF_i$$

Doc profile: set of words with highest TF-IDF scores, together with their scores.

## User Profiles -

Here User has rated items with profiles  $i_1, \dots, i_n$ .



simple (weighted) average of rated item profile.

variant - Normalize weights using average rating of a user.

## Boolean Utility Matrix

e.g. 1. Boolean utility matrix  
eg items are movies, only feature is "actor".  
- item profile - vector with 0 or 1 for each vector.

5 movies - 2 movie - actor A  
3 " " - " B

- User Profile = mean of item profiles

$$\text{Feature 'A' weight} = 2/5 = 0.4$$

$$\text{Feature 'B' " } = 3/5 = 0.6$$

# DETAILED LECTURE NOTES

PAGE NO.

e.g.-2. star Rating -

1-5 star ratings

higher rating

- Actor A's movies rated 3 and 5

1, 2 and 4.

lower rating  
c-ve "

- Actor B " " by subtracting

useful step - Normalize rating by user's mean rating (3)

- Actor A's normalized rating = 0, +2

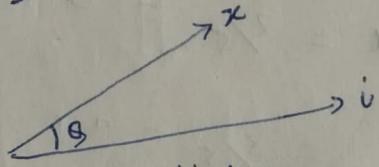
$$\text{Profile weight} = (0+2)/2 = 1$$

- Actor B's normalized rating = -2, -1 +1

$$\text{Profile weight} = -2/3$$

Making Predictions - User Profile  $\mathbf{x}$ , item profile  $\mathbf{i}$

$$\text{Estimate } v(\mathbf{x}, \mathbf{i}) = \cos(\theta) = \mathbf{x} \cdot \mathbf{i} / (\|\mathbf{x}\| \|\mathbf{i}\|)$$



Technically, the cosine distance is actually the angle  $\theta$  and the cosine similarity is the angle  $180-\theta$ .  
we use  $\cos(\theta)$  as our cosine similarity measure.

if  $\theta$  is less which means  $\mathbf{x}$  and  $\mathbf{i}$  are more similar in nature.

cons: Content-based Approach

- Finding the appropriate features is hard.  
eg. images, movies, music
- Overspecialization
  - Never recommends items outside user's content profile.
  - People might have multiple interests.
  - Unable to exploit quality judgements of other users.
  - Cold-start problem for new users
    - How to build a user profile.

Pros:

- ① No need for data on other users
- Able to recommend to users with unique tastes.
- Able to recommend new and unpopular items.
- No first-rater problem.
- Explanations for recommended items.

**Campus: PCE**      **Course: BTECH in CSE**      **Class/Section: III Yr. Section- A**  
**Name of Faculty: Praveen Kumar Yadav**      **Name of Subject: Machine Learning**  
**Date (Prep.): .....**      **Date (Del.): .....**      **Unit No.: .....V**      **Lect. No: .....**  
**Date: .....**      **Code: 60**

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

## COLLABORATIVE FILTERING

### **IMPORTANT & RELEVANT QUESTIONS:**

### **FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

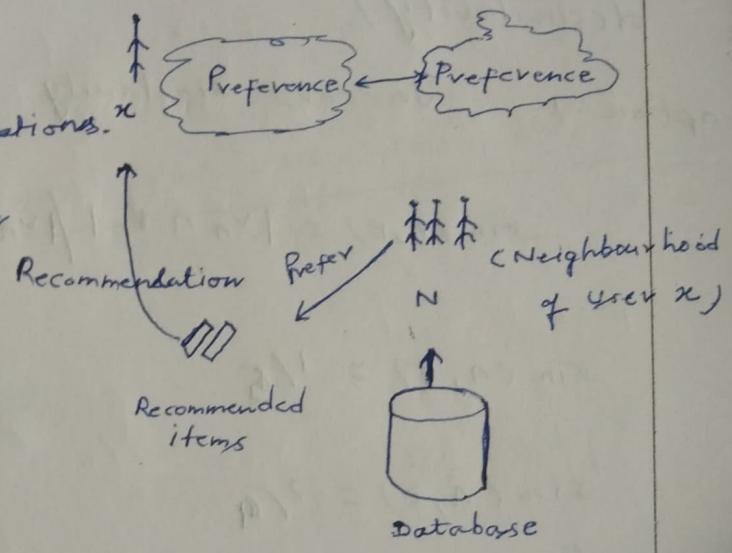
**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:

DETAILED LECTURE NOTESCOLLABORATIVE FILTERING:-

Also known as User-User collaborative filtering.

- Consider a user  $x$ , to whom we want to recommendations.

- Find a set of  $N$  of other users whose likes or dislikes (ratings) are similar to the likes and dislikes of user  $x$ .



- Estimate  $x$ 's ratings based on ratings of users in  $N$ .

Here we have to find the notion of similarities to find the neighbourhood of user  $x$ .

eg -

	HPL	HP2	HP3	TW	swL	sw2	sw3
A	4			5	L		
B		5					
C							
D			3				3

rating vector  $v_A$

rating vector  $v_B$

rating vector  $v_D$

consider user  $x$  and  $y$  with rating vectors  $v_x$  and  $v_y$ .

we need to find a similarity metric  $\text{sim}(u, y)$  capture intuition that  $\text{sim}(A, B) > \text{sim}(A, C)$ . (Correctly classified the user based on similarity and dissimilarity).

option-1. JACCARD similarity.

$$\text{sim}(A, B) = |v_A \cap v_B| / |v_A \cup v_B|$$

$$\text{sim}(A, B) = 1/5$$

$$\text{sim}(A, C) = 2/4$$

- Here  $\text{sim}(A, B) < \text{sim}(A, C)$

Prob- it ignores rating values.

option-2. cosine similarity

$$\text{sim}(A, B) = \cos(v_A, v_B)$$

Here all unknown values are initialized with zero.

$$\text{sim}(A, B) = 0.38$$

$$\text{sim}(A, C) = 0.32$$

- $\text{sim}(A, B) > \text{sim}(A, C)$ , but not by much.

Prob- Treats missing rating as negative. (Bad assumption)

option-3. Centered cosine

Normalize ratings by subtracting row mean.

	HPL	HP2	HP3	TW	$\sum w_1$	$\sum w_2$	$\sum w_3$	Avg Row Mean
A	4			5	1			$10/3$
B	5	5	4					$14/3$
C				2	4	5		$11/3$
D		3					3	$6/2$



	HPL	HP2	HP3	TW	$\sum w_1$	$\sum w_2$	$\sum w_3$	Sum = 0
A	$4 - 10/3 = 2/3$				$5/3$	$-7/3$		
B	$5/3$	$5/3$	$-2/3$		$-5/3$	$4/3$	$4/3$	
C								0
D		0						

Here 0 is average rating.

$$\text{Now } \text{sim}(A, B) = \cos(\vec{v}_A, \vec{v}_B) = 0.09$$

$$\text{sim}(A, C) = -0.56$$

$$\therefore \text{sim}(A, B) > \text{sim}(A, C)$$

captures intuition better

treated as average.

- Missing rating

Handles "tough raters" and "easy raters".

- Also known as Pearson correlation coefficient.



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO.

#### RATING PREDICTIONS:-

Let  $v_k$  be the vector of user  $x$ 's ratings.

Let  $N$  be the set of  $k$  users most similar to  $x$  who have also rated item  $i$ .

prediction for user  $x$  and item  $i$ .

option -1.  $v_{xi} = \frac{1}{k} \sum_{y \in N} v_{yi}$

option -2 - weighted Average  $v_{xi} = \sum_{y \in N} s_{xy} v_{yi} / \sum_{y \in N} s_{xy}$

where  $s_{xy} = \text{sim}(x, y)$ .

That approach is also known as User-User collaborative filtering.

## Item-item Collaborative Filtering:-

Another view - item-item

- For item  $i$ , find other similar items.

- Estimate rating for item  $i$  based on ratings for similar items.

can use same similarity metrics and predictions functions as in user-user model.

$$v_{ui} = \frac{\sum_{j \in N(i; u)} s_{ij} \cdot r_{uj}}{\sum_{j \in N(i; u)} s_{ij}}$$

↑ Neighbourhood  
size = 2

$s_{ij}$  ... similarity of item  $i$  and  $j$   
 $r_{uj}$  - rating of user  $u$  on item  $j$   
 $N(i; u)$  - set of items rated by  $u$  similar to  $i$ .

e.g. item-item CF ( $N=2$ ) find two similar movies.

	1	2	3	4	5	6	7	8	9	10	11	12	Sim( $i, m$ )
1	1		3		?	5		5		2	1	3	-0.15
2				5	4			4					
3	2	4		1	?		3		4	3	5	0.41	
4		2	4	1	5			4		2			-0.10
5			4	3	?	4	2			2	5		-0.31
6	1	3	3	2	?	3	1	2		4			0.59

□ - Unknown rating

Rating between 1 to 5



# POORNIMA COLLEGE OF ENGINEERING DETAILED LECTURE NOTES

PAGE NO.

goal - estimate the rating of Movie 1 by user 5

Step 1 - Find other movie that are similar to movie 1.

Here we use Pearson correlation as similarity.

Subtract mean rating  $m_i$  from each movie i.

$$m_1 = (1+3+5+5+4)/5 = 3.6$$

$$\text{row}_1 = [-2.6, 0, -0.6, 0, 0, 1.4, 0, 0, 1.4, 0, 0.4, 0]$$

Step 2. Compute cosine similarities between rows.

$$s_{13} = 0.91 \quad s_{16} = 0.59$$

Step 3 - Predict by taking weighted Average

$$v_{15} = (0.41 * 2 + 0.59 * 3) / (0.41 + 0.59) = 2.6$$

## LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Class/Section: III Yr. Section- A

Date: .....

Name of Faculty: Praveen Kumar Yadav

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Prep.): .....

Date (Del.): ..... Unit No.: 5 Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

Neural Network.

**IMPORTANT & RELEVANT QUESTIONS:****FEED BACK QUESTIONS (AFTER 20 MINUTES):**

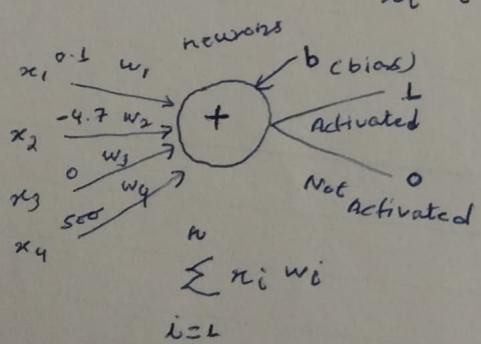
**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:

In 1950 or earlier the people came up with the idea that we can probably trying to replicate/reproduce the functionality of brain i.e. what is happening of a human brain).

The neuron receives the input, then neurons somehow figure out what is important / what's not, so it can activate that signal or ignore that signal. You have chains of all those neurons connected to each other and somehow we are processing all the input data that are coming in neurons and we can do something on that processed information. So same process is being

replicated in the machine through the algorithm.



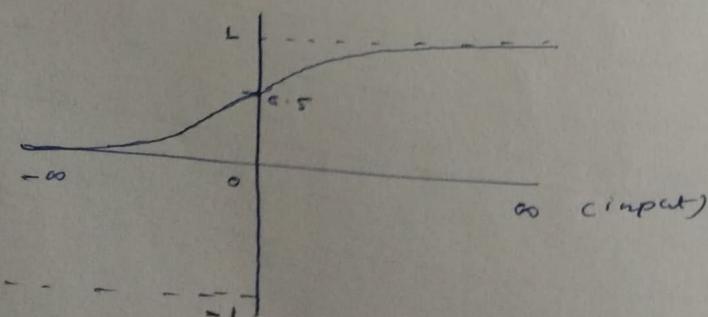
$n_i w_i =$  if  $w_i = 0$  then o/p is 0  
else ignore the signal

if  $w_i >= 0$   
 $w_i = 500$  then o/p - amplify that signal

$w_i < 0$   
then - deamplify the signal.  
 $w_i = 0.0001$

( $-\infty$  to  $+\infty$ )  $\xrightarrow{\text{Transformation done by sigmoid function}}$  ( $0, 1$ )  
input range

These are rules that what we defined.





# POORNIMA

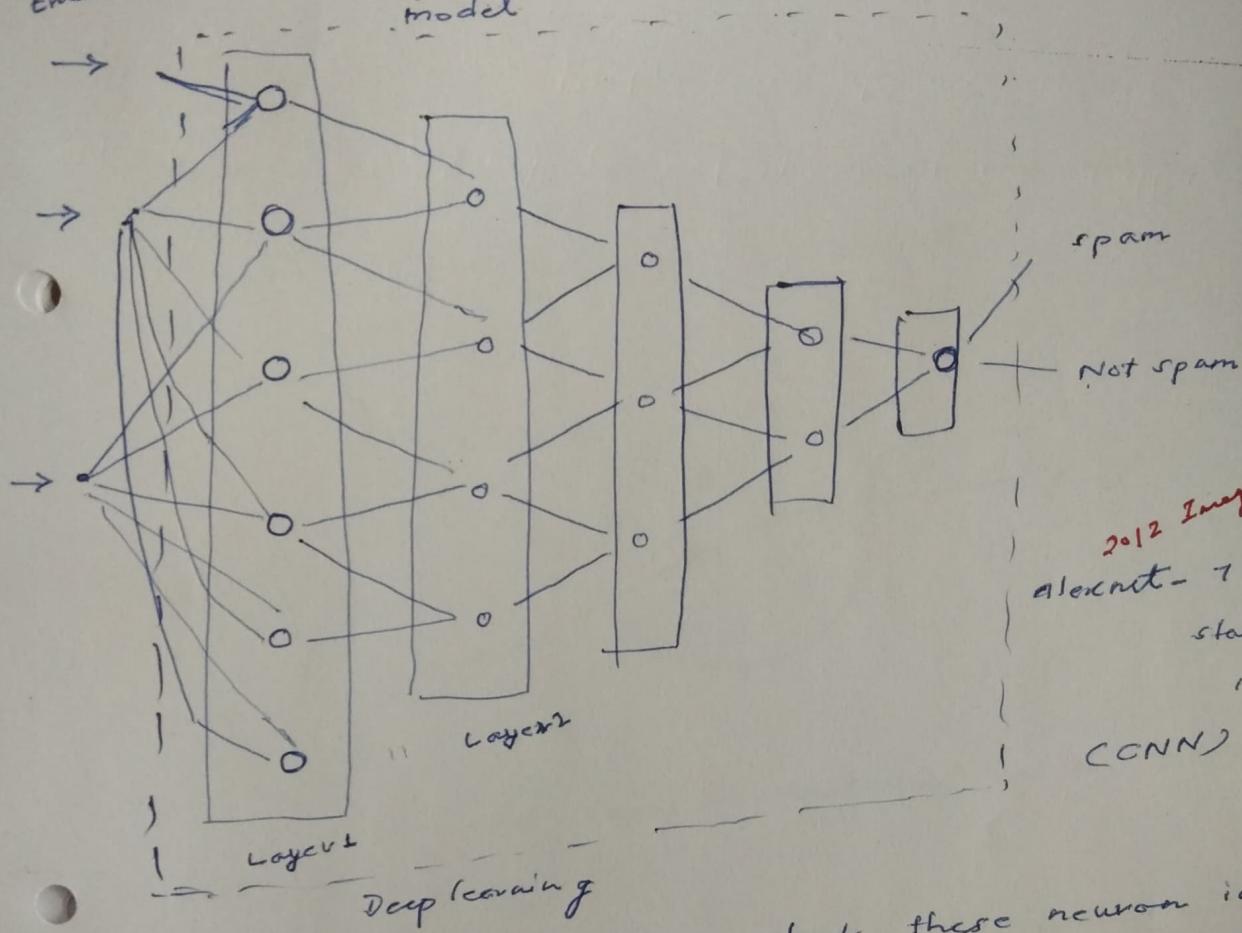
## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

Email text

model

PAGE NO.



so In Deep learning you can stack these neuron in one layer and processing simultaneously and on top of that we have another layers of neuron of partially or fully connected layer and you can have multiple layers like that and last output layer produces output of all processed signals.

That kind of architecture is known as deep learning neural network.

2012 ImageNet Competition  
alexnet - 7 or 8 layer  
standford University  
1.4 Million Images  
(CNN) 95%



# POORNIMA

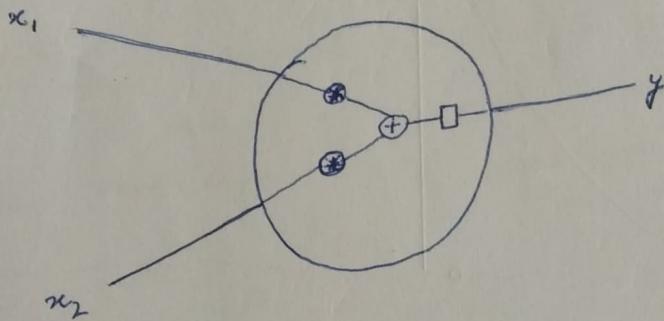
## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

#### Artificial Neural Network :-

**Neuron:** A neuron takes input, does some math with them and produces one output.



$$x_1 = w_1 * x_1$$

↳ - each input is multiplied by a weight

$$x_2 = w_2 * x_2$$

Next, all the weighted inputs are added together with a bias  $b$ -

i.e

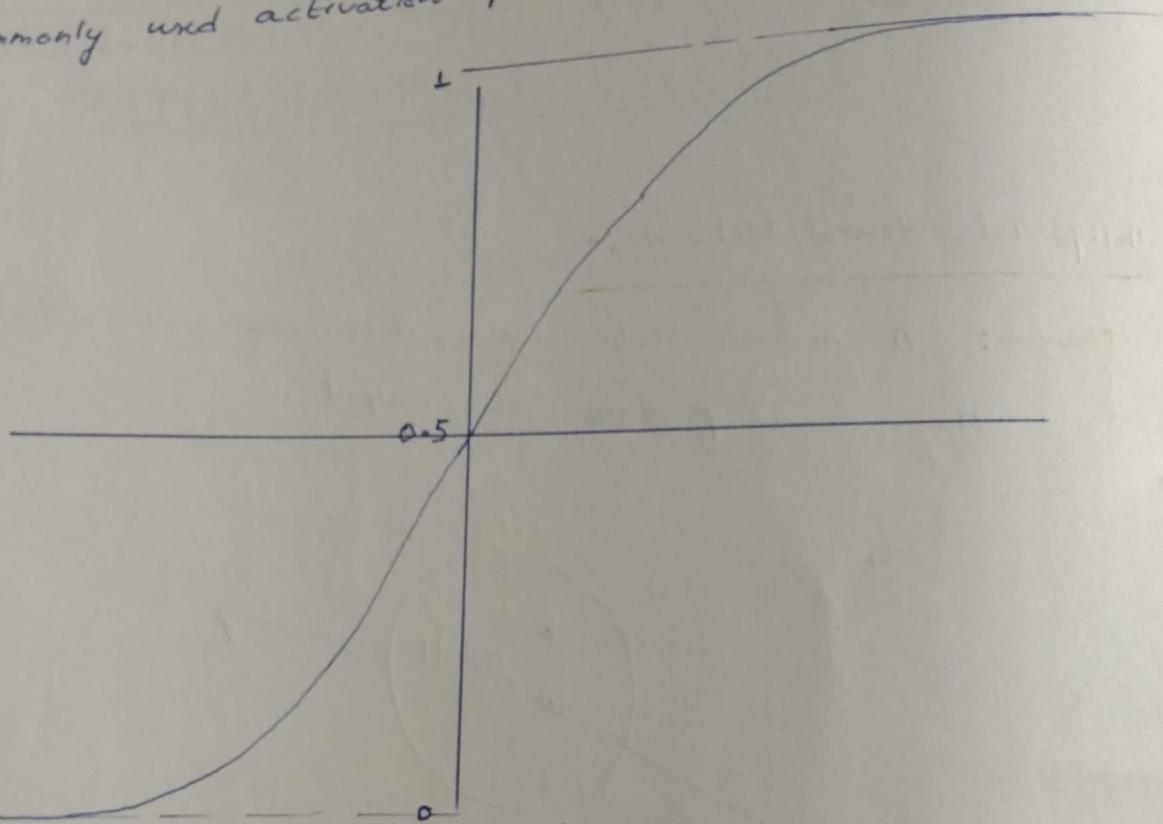
$$(x_1 * w_1) + (x_2 * w_2) + b$$

↳ the sum is passed through an activation

Finally, the function-

$$y = f(x_1 * w_1 + x_2 * w_2 + b)$$

The activation function is used to turn an unbounded input into an output that has a nice, predictable form. A commonly used activation function is the sigmoid function.



The sigmoid function only outputs numbers in the range  $(0, L)$ . You can think of it as compressing  $(-\infty \text{ to } +\infty)$  to  $(0, L)$  - big negative numbers become  $\approx 0$  and big positive numbers become  $\approx L$ .

$$\text{For eg. } w = [w_1, w_2] \quad (w \cdot x + b) = (w_1 * x_1 + w_2 * x_2) + b$$

$$b = 4 \quad = 0 * 2 + 1 * 3 + 4$$

$$x = [x_1, x_2] \quad = 7$$

$$y = f(w \cdot x + b) \\ = F(7) = 0.9999$$



# POORNIMA

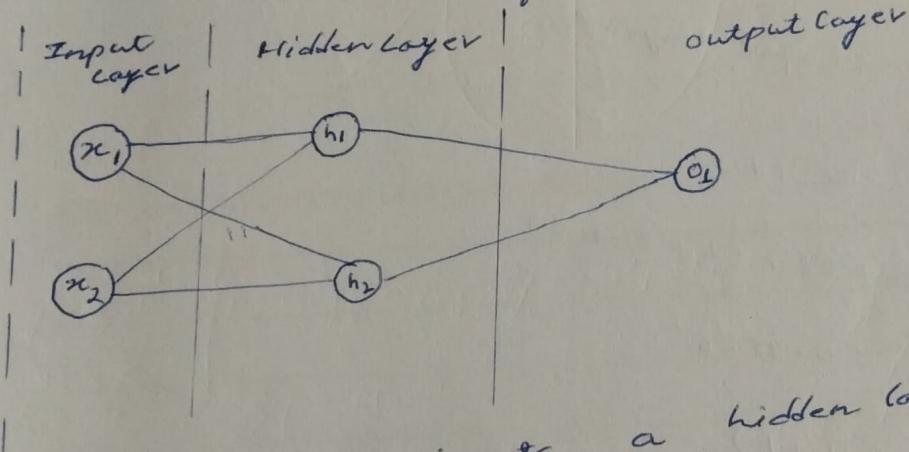
## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO.

Here the neurons outputs 0.999 given the inputs  $x = [2, 3]$ .  
This process of passing inputs forward to get an output is known as Feed-Forward Neural Network.

Neural Network:- A neural network is nothing more than a bunch of neurons connected together.



This network has 2 inputs, a hidden layer with 2 neurons ( $h_1$  and  $h_2$ ) and output layer with 1 neuron ( $o_1$ ). Here input for  $o_1$  are the outputs from  $h_1$  and  $h_2$  - that's what makes this a network.  
A hidden layer is any layer between the input (first) layer and output (last layer).

Feed-Forward Neural Network :- A neural network can have any number of layers with any number of neurons in those layers. So in the feed-forward neural network we are feed the input(s) forward through the neuron's in the network to get the output(s) at the end. This is called as **Feed-Forward Neural Networks.**

eg -  $w_1, w_2$   
 $w = [0, 1]$  for all neurons

bias  $b = 0$  and same sigmoid function.

$h_1, h_2, o_1$  denote the outputs of neurons they

represent.

$x_1, x_2$   
 $x = [2, 3]$  (Input Layer)

$$\begin{aligned} h_1 &= h_2 = f(w \cdot x + b) \\ &= f((0 \cdot 2) + (1 \cdot 3) + 0) \\ &\approx f(3) = 0.9526 \end{aligned}$$

$$\begin{aligned} o_1 &= f(w \cdot [h_1, h_2] + b) \\ &= f((0 \cdot h_1) + (1 \cdot h_2) + 0) \\ &\approx f(0.9526) \end{aligned}$$

$$= 0.7216$$

so the output of a neural network for input  $x = [2, 3]$  is  $0.7216$ .

# LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Class/Section: III Yr. Section- A

Date: .....

Name of Faculty: Praveen Kumar Yadav

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Perception*

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

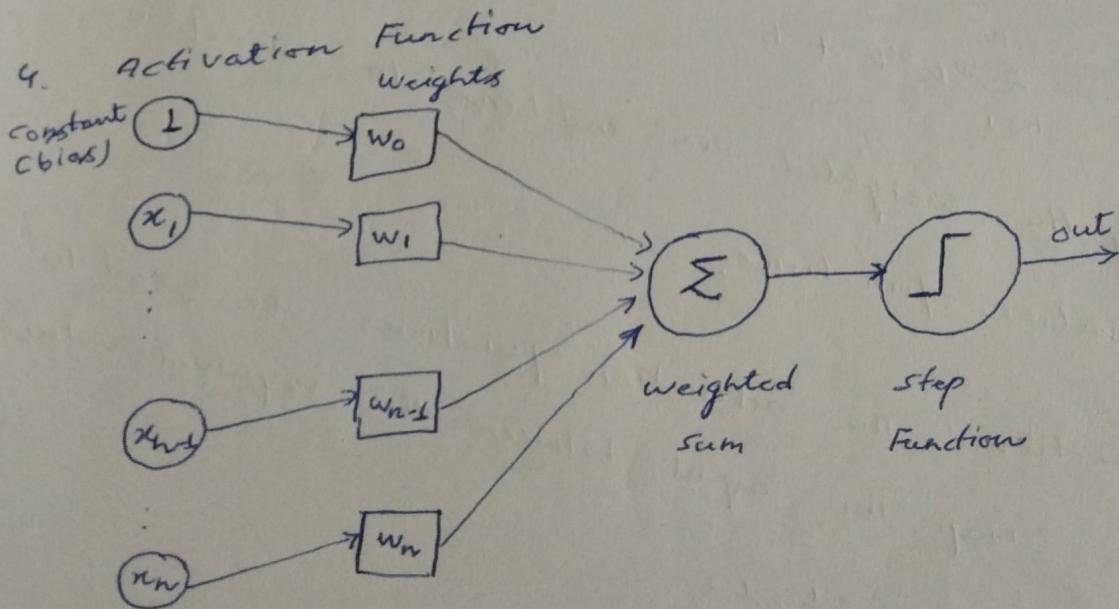
PAGE NO. ....

**Perception:-** A perception is a single layer network and a multilayer perceptron is called as **Neural networks.**

Perception is a linear classifier (binary). Also it is used in supervised learning.

The Perception consists of 4 parts.

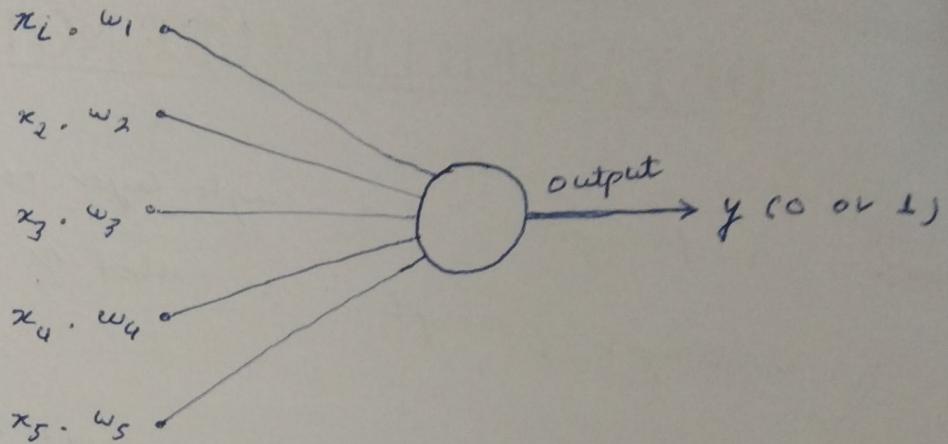
1. Input values or one input layer
2. weights and bias.
3. Net sum



(Perception)

## How Perceptron works:-

- All input  $x$  are multiplied with their weights  $w$ . Let's call it  $K$ .



- Add all the multiplied value and call them weighted sum.
- ↓ end term.  
or  $w_K x_K$   
 $\sum_{K=1}^n K$  (formula for  $K^{th}$  term)
- Sigma →  
for summation       $K=1$

- Add bias ( $b$ ) to the weighted sum

$$\sum_{K=1}^n w_K x_K + b$$

- apply the weighted sum and bias to correct activation functions.

**Activation Function.** The activation functions are used to map the input between the required values like  $(0, 1)$  or  $(-1, 1)$ .



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

**weights :-** weight shows the strength of the particular node.

**bias :-** bias is a constant which helps the model in a way that it can best fit for the given data.

bias unit appended to the start/end of the input and each hidden layer and it is not affected by the values in previous layer.

It allow you to move the line up and down to fit the prediction with the data better, without to the line always goes through origin (0,0) and you may get a poorer fit.

## Coding a Neuron -

```
import numpy as np
```

```
def sigmoid(x):
```

# our activation function:  $f(x) = \frac{1}{1 + e^{-x}}$

```
return 1 / (1 + np.exp(-x))
```

```
class Neuron:
```

```
def __init__(self, weights, bias):
```

```
    self.weights = weights
```

```
    self.bias = bias
```

```
def feedforward(self, inputs):
```

```
    total = np.dot(self.weights, inputs) + self.bias
```

```
    return sigmoid(total)
```

```
weights = np.array([0, 1]) #  $w_1 = 0, w_2 = 1$ 
```

```
bias = 4 #  $b = 4$ 
```

```
n = Neuron(weights, bias)
```

```
x = np.array([2, 3])
```

#  $x_1 = 2, x_2 = 3$

```
print(n.feedforward(x))
```

## Multilayer Neural Network -

Coding a Neural Network:-

```
import numpy as np
```

```
class OurNeuralNetwork:
```

```
    def __init__(self):
```

```
        weights = np.array([0, 1])
```

```
        bias = 0
```

```
        self.w1 = Neuron(weights, bias)
```

```
        self.w2 = Neuron(weights, bias)
```

```
        self.o1 = Neuron(weights, bias)
```

```
    def feedforward(self, x)
```

```
        out_w1 = self.w1.feedforward(x)
```

```
        out_w2 = self.w2.feedforward(x)
```

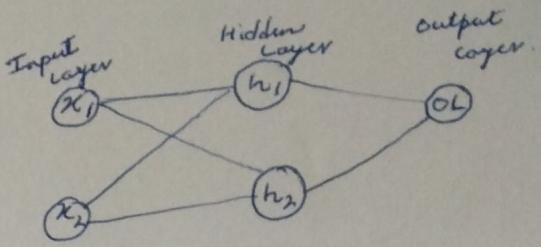
```
        out_o1 = self.o1.feedforward(np.array([out_w1,  
                                              out_w2]))
```

```
    return out_o1
```

```
network = OurNeuralNetwork()
```

```
x = np.array([2, 3])
```

```
print(network.feedforward(x))
```



# LECTURE NOTES

Campus: PCE

Course: BTECH in CSE

Class/Section: III Yr. Section- A

Date: .....

Name of Faculty: Praveen Kumar Yadav

Name of Subject: Machine Learning

Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... <sup>5<sup>th</sup></sup> Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

Multilayer Network

Multilayer Perceptron.

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:

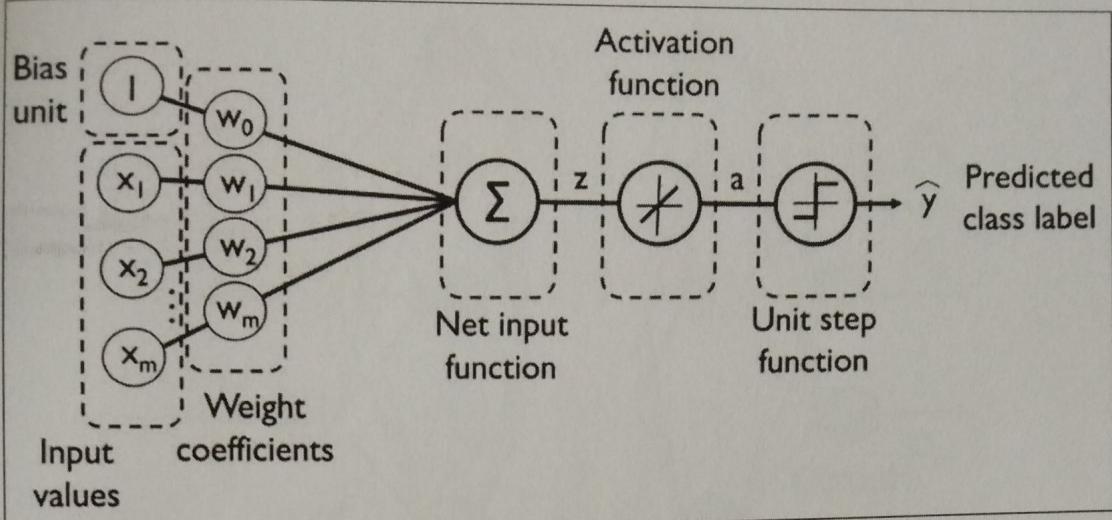
Multilayer Neural Network:-

A fully connected multi-layer

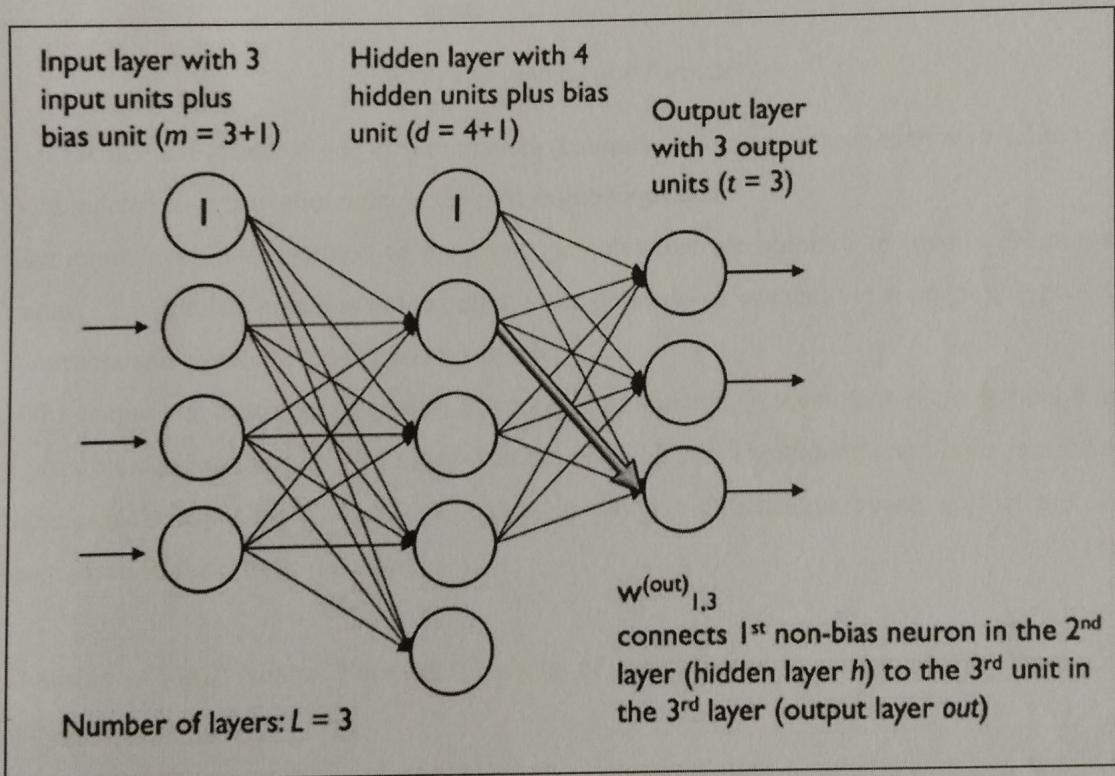
neural network is called a Multilayer

Percptron.

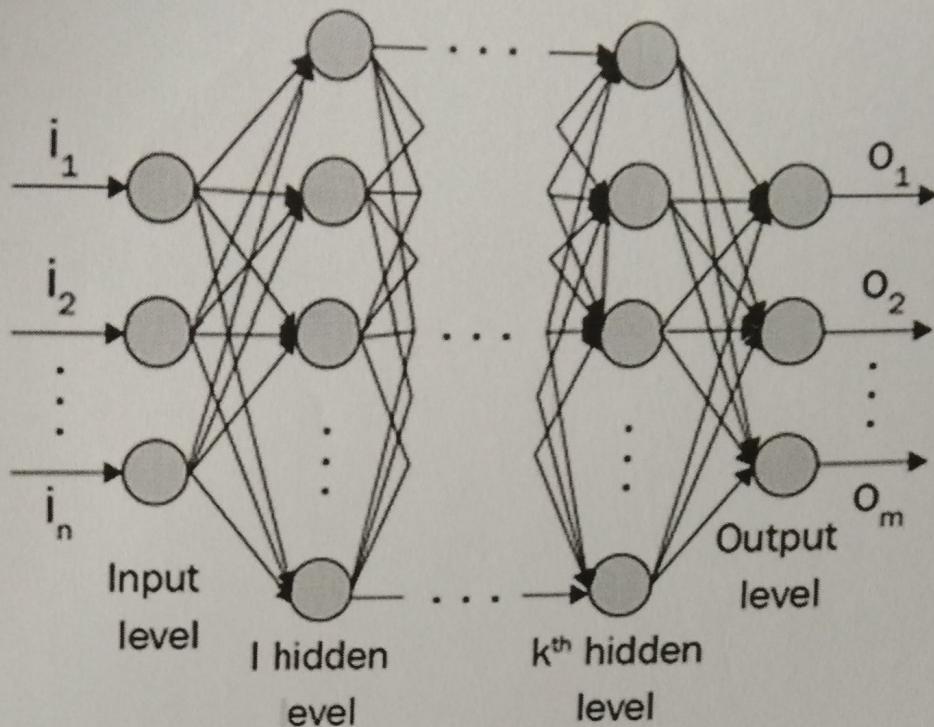
- MLP is a typical example of a feed forward artificial neural network.
- MLP is basically used to solve non-linear separable problems.



Single Layer Neural Network



A Multilayer Perceptron



Feed-Forward Neural Network

A layer is a collection of one or more nodes (computation units), where each node in a layer is connected to every other node in the next immediate layer.

The input layer is constituted of the input variables that are required to predict the output values. The number of nodes in the output layer depends on whether we are trying to predict a continuous variable or a categorical variable.

If the output is a continuous variable, the output has one unit. If the output is categorical with  $n$  possible classes, there will be  $n$  nodes in the output layer. The hidden level/layer is used to transform the input layer values into values in a higher-dimensional space, so that we can learn more features from the input.

**Training of MLP Neural Network:** Training of Neural network can be done by means of two algorithm-

1. **Feed Forward Propagation Algorithm:** In forward-propagation, we apply a set of weights to the input data, pass it through the hidden layer, perform the nonlinear activation on the hidden layer output, and then connect the hidden layer to the output layer by multiplying the hidden layer node values with another set of weights. For the first forward-propagation, the values of the weights are initialized randomly.

2. **Back Propagation Algorithm:** In back-propagation, we try to decrease the error by measuring the margin of error of output and then adjust weight accordingly. Neural networks repeat both forward- and back-propagation to predict an output until the weights are calibrated.



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE Course: BTECH in CSE Class/Section: III Yr. Section- A Date: .....

Name of Faculty: Praveen Kumar Yadav Name of Subject: Machine Learning Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

Deep Neural networks -

MLP Training- chain-rule Implementation.

#### IMPORTANT & RELEVANT QUESTIONS:

#### FEED BACK QUESTIONS (AFTER 20 MINUTES):

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



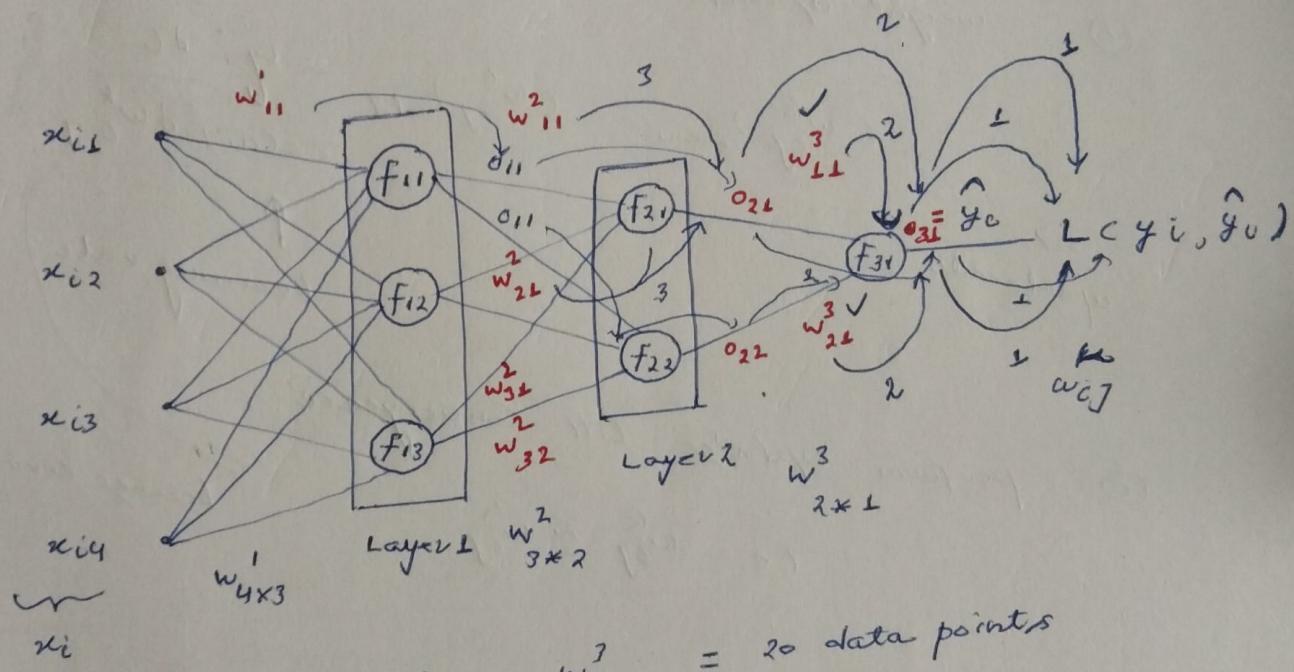
# Poornima COLLEGE OF ENGINEERING

## DETAILED LECTURE NOTES

PAGE NO. ....

Training on MCP-chain-rule:-

$$D = \{x_i, y_i\} \quad x_i \in \mathbb{R}^4 \quad \left\{ \begin{array}{l} \text{reg. prob.} \rightarrow \text{square loss} \\ y_i \in \mathbb{R} \end{array} \right.$$



determine  $w'_{4 \times 3}, w^2_{3 \times 2}, w^3_{2 \times 1}$  = 20 data points

step 1 - define the cost function -

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg. (regularization term)} \quad \text{on weights}$$

SQUARE LOSS

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

$$L = \sum_{i=1}^n L_i + \text{reg.}$$

$$\sum_{i,j,k} |w_{ijk}|^2$$

Step 2 - minimize  $L$  ↗ square loss  
↗ regularization pm.  
 $w^1, w^2, w^3$  i.e.

$$\min_{w_{i,j}^k} L$$

Step 3 - use SGD or GD

(a) initialization of  $w_{i,j}^k$  randomly.

(b) weight update rule

$$(w_{i,j}^k)_{\text{new}} = (w_{i,j}^k)_{\text{old}} - \gamma_i \frac{\partial L}{\partial w_{i,j}^k}$$

$\uparrow$   
learning rate

(c) perform updates till convergence.

$$\text{i.e. } w_{i,j}^k \approx w_{i,j}^k \quad \text{becomes zero}$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

$w_{11}^3$

$$\frac{\partial L}{\partial w_{11}^3} = \frac{\partial L}{\partial O_{3L}} \cdot \frac{\partial O_{3L}}{\partial w_{11}^3} \quad (\text{chain rule of differentiation})$$

$$\frac{\partial L}{\partial w_{21}^3} = \frac{\partial L}{\partial O_{3L}} \cdot \frac{\partial O_{3L}}{\partial w_{21}^3}$$

$w_{11}^2$

$$\frac{\partial L}{\partial w_{11}^2} = \frac{\partial L}{\partial O_{3L}} \cdot \frac{\partial O_{3L}}{\partial O_{2L}} \cdot \frac{\partial O_{2L}}{\partial w_{11}^2}$$

$$\frac{\partial L}{\partial w_{21}^2} = \frac{\partial L}{\partial O_{3L}} \cdot \frac{\partial O_{3L}}{\partial O_{2L}} \cdot \frac{\partial O_{2L}}{\partial w_{21}^2}$$

$$\frac{\partial L}{\partial w_{3L}^2} = \frac{\partial L}{\partial O_{3L}} \cdot \frac{\partial O_{3L}}{\partial O_{2L}} \cdot \frac{\partial O_{2L}}{\partial w_{3L}^2}$$

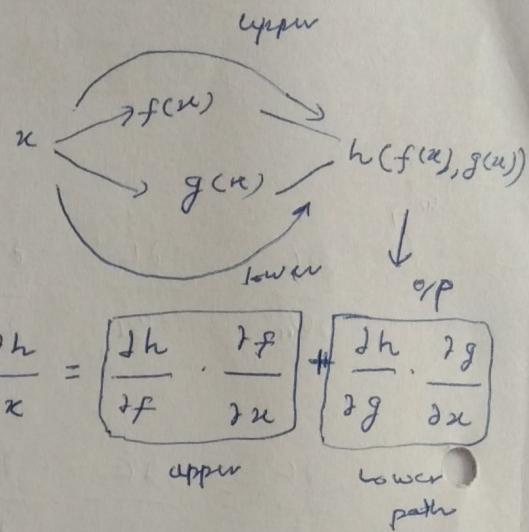
$w'_{11}$ 

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial o_{3L}} \cdot \frac{\partial o_{3L}}{\partial o_{2L}} \cdot \frac{\partial o_{2L}}{\partial o_{11}}$$

$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial o_{3L}} \cdot \frac{\partial o_{3L}}{\partial w'_{11}}$$

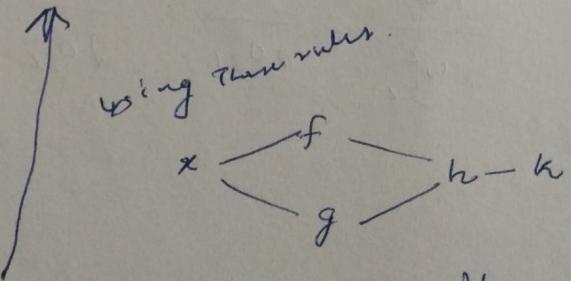
$$\frac{\partial o_{3L}}{\partial w'_{11}} = \frac{\partial g_1}{\partial o_{2L}} \cdot \frac{\partial o_{2L}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11}}$$

$$+ \frac{\partial o_{3L}}{\partial o_{2L}} \cdot \frac{\partial o_{2L}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11}}$$



$$\frac{\partial L}{\partial w'_{11}} = \frac{\partial L}{\partial o_{3L}} \cdot \left( \frac{\partial o_{3L}}{\partial o_{2L}} \cdot \frac{\partial o_{2L}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11}} + \right.$$

$$\left. \frac{\partial o_{3L}}{\partial o_{2L}} \cdot \frac{\partial o_{2L}}{\partial o_{11}} \cdot \frac{\partial o_{11}}{\partial w'_{11}} \right)$$



$$\frac{\partial K}{\partial x} = \frac{\partial K}{\partial h} \times \left\{ \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x} \right\}$$

$$\frac{\partial K}{\partial x} = \frac{\partial K}{\partial h} \cdot \frac{\partial h}{\partial f} \cdot \frac{\partial f}{\partial x} + \frac{\partial K}{\partial h} \cdot \frac{\partial h}{\partial g} \cdot \frac{\partial g}{\partial x}$$



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE Course: BTECH in CSE Class/Section: III Yr. Section- A Date: .....

Name of Faculty: Praveen Kumar Yadav Name of Subject: Machine Learning Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Back propagation algorithm for Perceptron*

#### IMPORTANT & RELEVANT QUESTIONS:

#### FEED BACK QUESTIONS (AFTER 20 MINUTES):

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Backpropagation Algorithm in Neural network :-

Notations :-

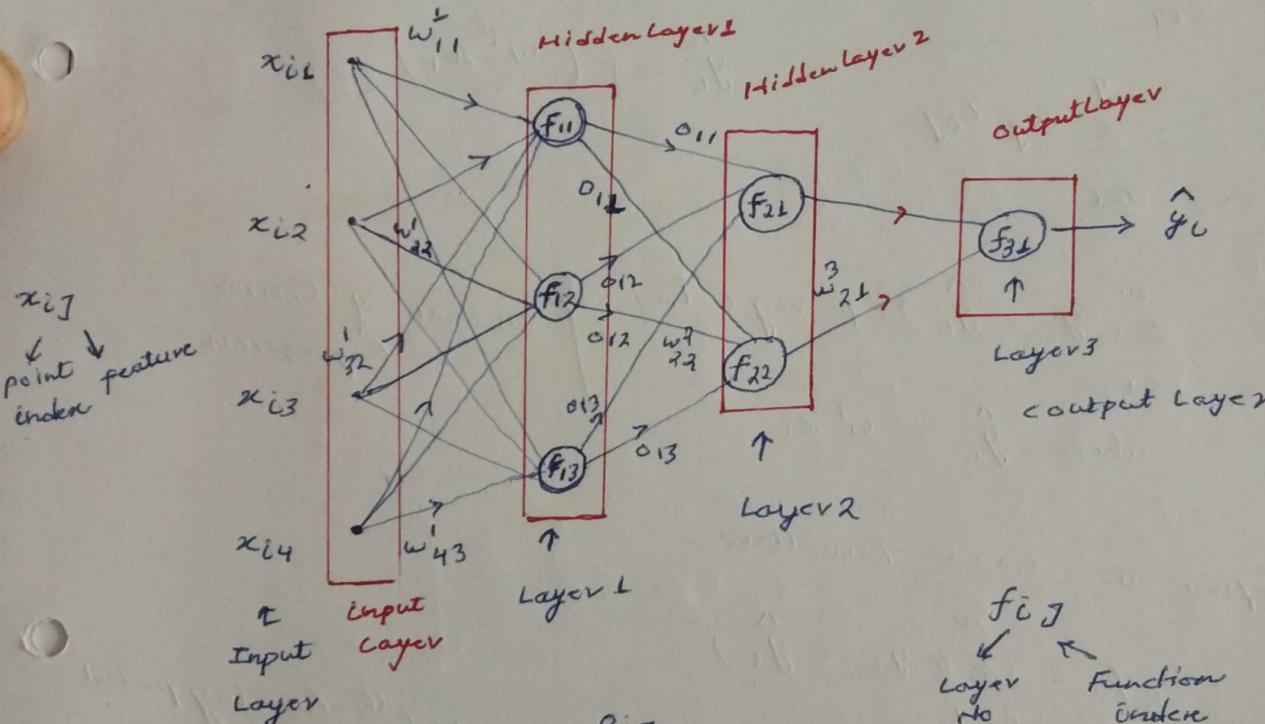
$$W^1_{4 \times 3}$$

$$W^2_{3 \times 2}$$

$$W^3_{2 \times 1}$$

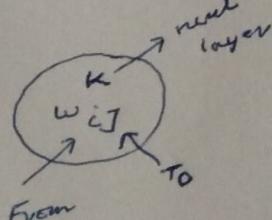
$$D = \{x_i, y_i\}$$

$$x_i \in \mathbb{R}^4; y_i \in \mathbb{R}$$



$o_{ij}$   
i<sup>th</sup> layer      j<sup>th</sup> neuron  
under

$f_{i,j}$   
layer no  
Function  
order  
(neuron index)



Fully Connected Multilayer Perceptron.

$$\text{input } \times 3 = 12$$

4

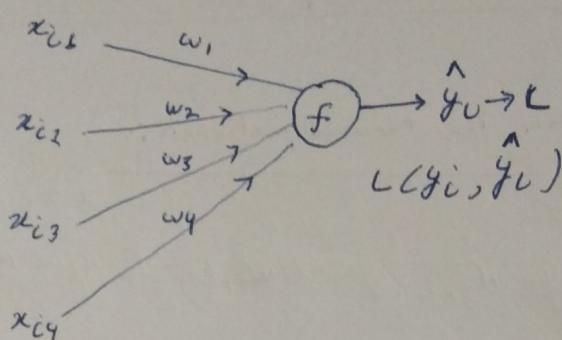
$$W^L =$$

$$\begin{bmatrix} w_{11}^1 & w_{12}^1 & w_{13}^1 \\ w_{21}^1 & w_{22}^1 & w_{23}^1 \\ w_{31}^1 & w_{32}^1 & w_{33}^1 \\ w_{41}^1 & w_{42}^1 & w_{43}^1 \end{bmatrix}_{4 \times 3}$$

## Training a single Neuron model:-

↳ To find the best edge-weights using  $D_{Tr}$  (Training Data)  
 Perceptron & LR  $\rightarrow$  single Neuron model for classification

$$D_{Tr} = \{x_i, y_i\}_{i=1}^n$$



LR. Regression  $x_i \in \mathbb{R}^{d_{CR}}$

$$\hat{y}_i = \sum_{j=1}^d w_j x_{ij} \quad y_i \in \mathbb{R}$$

$$\hat{y}_i = w^T x_i$$

$$\hat{y}_i = f(w^T x_i)$$

$$= \sum_{j=1}^d w_j x_{ij}$$

$$\min_{w_i} \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \text{reg}(P_w)$$

where  $\hat{y}_i = w^T x_i$

*In case of Linear Regression*

Step 1 - Define a loss Function -

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$L_i = (y_i - \hat{y}_i)^2 \quad \text{square loss of } i^{\text{th}} \text{ training point}$$

$$L = \sum_{i=1}^n L_i$$

Step 2 optimization problem -

$$\min_{w_i} \sum_{i=1}^n (y_i - w^T x_i)^2$$

$$\min_{w_i} \sum_{i=1}^n (y_i - f(w^T x_i))^2 + \text{regularization}$$



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

$$w^* = \underset{w}{\operatorname{argmin}} \sum_{i=1}^n (y_i - f(w^T x_i))^2 + \text{reg.}$$

3) solve the optimization problem.

a) initialization of  $w^{(0)} \rightarrow$  random initialization

b) Compute the Derivative of  $L$  w.r.t to  $w$  i.e.

$$\nabla_w L = \begin{bmatrix} \frac{\partial L}{\partial w_1} \\ \frac{\partial L}{\partial w_2} \\ \vdots \\ \frac{\partial L}{\partial w_4} \end{bmatrix} \quad w \in \mathbb{R}^4 \quad x_i \in \mathbb{R}^4$$

① c)  $w_{\text{new}} = w_{\text{old}} - n \nabla_w L$

for  $\text{iter} = 1 \text{ to } k$ .

$$(w_i)_{\text{new}} = (w_i)_{\text{old}} - n \left[ \frac{\partial L}{\partial w_i} \right]_{(w_i)_{\text{old}}}$$

weight update rule.

$$\nabla_{\omega} L = \left[ \frac{\partial L}{\partial \omega_1}, \frac{\partial L}{\partial \omega_2}, \frac{\partial L}{\partial \omega_3}, \frac{\partial L}{\partial \omega_4} \right]^T$$

$$\frac{\partial L}{\partial \omega_1} = \frac{\partial L}{\partial f} \cdot \frac{\partial f}{\partial \omega_1} \rightarrow \text{chain rule of differentiation}$$

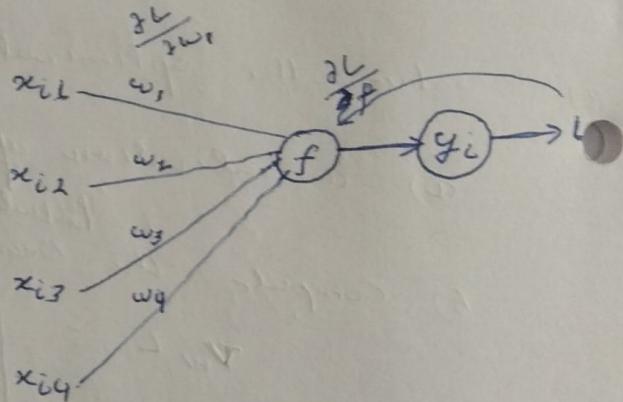
partial  
chain rule

$$\text{cf. } \frac{\partial L}{\partial x} = \frac{\partial L}{\partial f_1} \frac{\partial f_1}{\partial x} + \frac{\partial L}{\partial f_2} \frac{\partial f_2}{\partial x}$$

$$L = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$= \sum_{i=1}^n (y_i - f(\omega^T x_i))^2$$

$$\frac{\partial L}{\partial f} = - \sum_{i=1}^n 2(y_i - f(\omega^T x_i))$$



$$\frac{\partial f}{\partial \omega_1} = x_i$$

$$\boxed{\frac{\partial L}{\partial \omega_1} = \sum_{i=1}^n -2x_i(y_i - \hat{y}_i)}$$

$$\frac{\partial L}{\partial \omega_2} \quad \dots \quad \frac{\partial L}{\partial \omega_4}$$



# POORNIMA

## FOUNDATION

### LECTURE NOTES

Campus: PCE      Course: BTECH in CSE      Class/Section: III Yr. Section- A      Date: .....

Name of Faculty: Praveen Kumar Yadav      Name of Subject: Machine Learning      Code: 6CS4-02

Date (Prep.): ..... Date (Del.): ..... Unit No.: ..... Lect. No: .....

**OBJECTIVE:** To be written before taking the lecture (Pl. write in bullet points the main topics/concepts etc., which will be taught in this lecture)

*Back propagation Algorithm*

**IMPORTANT & RELEVANT QUESTIONS:**

**FEED BACK QUESTIONS (AFTER 20 MINUTES):**

**OUTCOME OF THE DELIVERED LECTURE:** To be written after taking the lecture (Pl. write in bullet points about students' feedback on this lecture, level of understanding of this lecture by students etc.)

**REFERENCES:** Text/Ref. Book with Page No. and relevant Internet Websites:



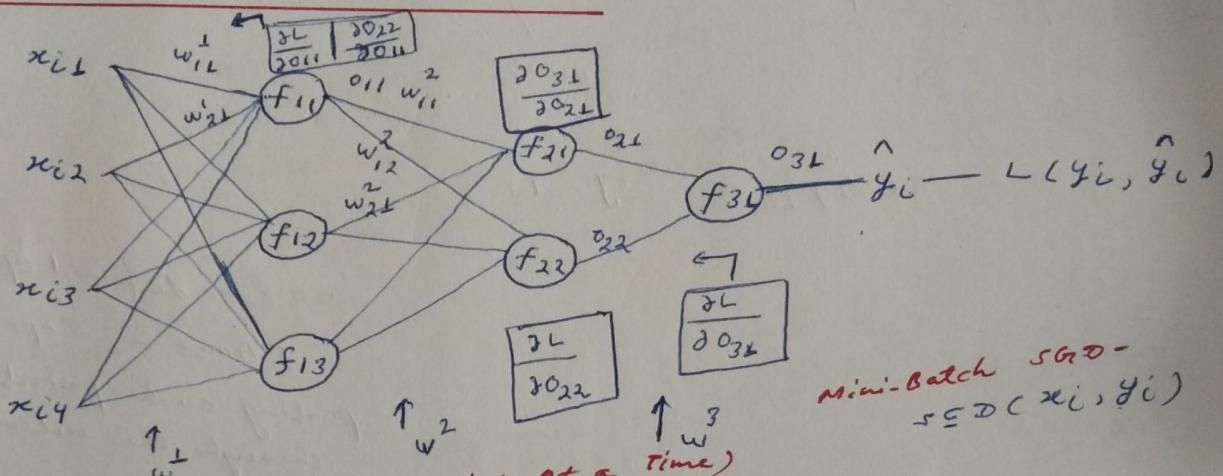
# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

PAGE NO. ....

Backpropagation Algorithm in MLP:-



mini-Batch SGD -  
 $\Sigma \in \mathcal{D}(x_i, y_i)$

SGD algorithm: (only one point at a time)

$$\mathcal{D} = \{x_i, y_i\}$$

step 1 :- Initialize  $w_{ij}^k$ 's randomly.

step 2 :- for each  $x_i$  in  $\mathcal{D}$

(a) - pass  $x_i$  forward through the network  
 (Forward propagation)

and calculate  $L_{SGD}$

(b) - compute  $L(\hat{y}_i, y_i)$

(c) Compute all derivatives using chain rule of  
 derivatives, and Memorization  $(\frac{\partial L}{\partial w_{ij}})$

(d) update weight from end of the network  
 to the start. (Backward propagation)

For example -

$$(\omega_{11}^3)_{\text{new}} = (\omega_{11}^3)_{\text{old}} - \eta \frac{\partial L}{\partial \omega_{11}^3}$$
$$\downarrow$$
$$\frac{\partial L}{\partial o_3 L} \cdot \frac{\partial o_3}{\partial \omega_{11}^3}$$

Step 3 - repeat step 2 till convergence to solution.  
↓  
means.

$$(\omega_{ij}^k)_{\text{new}} \approx (\omega_{ij}^k)_{\text{old}}$$

$D = \{x_i, y_i\}$  - 1 epoch

5 times - 5 epoch

passing each points  
through the network.

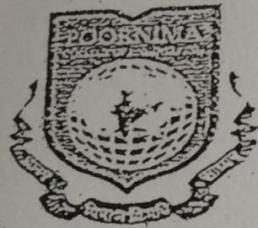
Backprop -

- ① init forward prob
- ② 1 epoch / compute loss
- / compute derivatives
- / update weight.

③ repeat step 2 till convergence

so basically a backpropagation is a basically multiple epoch training methodologies where we leverage the concept of chain rule and memorization to update the weights of NN.

v. very important points.  
Back propagation is only work when the activation function is differentiable otherwise training of NN is not possible.



# POORNIMA

## COLLEGE OF ENGINEERING

### DETAILED LECTURE NOTES

(SGD Approach)

PAGE NO. ....

- ② Instead of sending one point we can also send the batch of points in the network (Mini-batch SGD).  
and if we send all the data points through the network then it is called as Gradient Descent Approach).

- ③ All of the datapoints in RAM and computing " " using GD is extremely time consuming.  
(in case of Gradient descent Algorithm).