

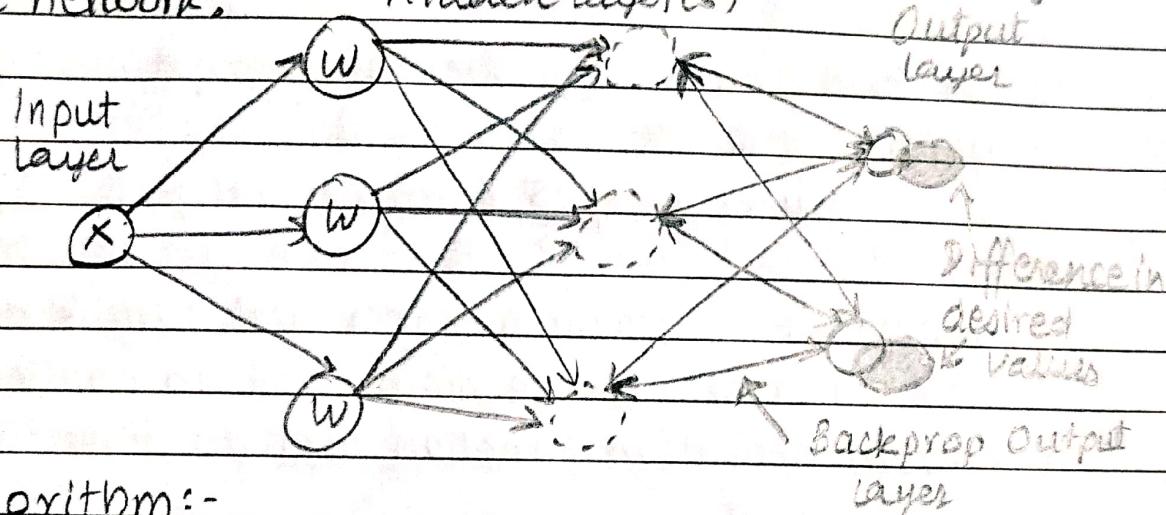
\* UNIT 3 already studied in ML \*

We will only cover concepts not studied

## UNIT 4 : Artificial Neural Network

\* Back Propagation algorithm :-

- Method of fine tuning weights of a neural network based on errors obtained in previous epoch.
- Proper tuning helps in increasing model generalization.
- Helps calculate loss function with respect to all weights in the network.



\* Algorithm:-

- 1) Inputs  $X$  arrive through preconnected path.
- 2) Input is modelled using real weights (randomly selected)
- 3) Calculate output for every neuron from the input layer to the hidden layer, to the output layer.
- 4) Calculate error in the outputs:-

$$\text{Error} = \text{Actual Output} - \text{Desired Output}$$

- 5) Travel back from output layer to the hidden layer and adjust weights such that error is decreased.

\* Advantages:-

Fast and simple, flexible & standard method

\* Disadvantages:-

Sensitive to noise, Performance dependent on input data.

- \* Generalized Delta learning Rule:
    - Most common method for training backpropagation network.
    - Iterative gradient-descent method that minimizes the least-mean-squares (LMS) error in the output.
    - Consider  $i_n$ ,  $o_n$  and  $w$  to be input, output and weights.
- If we use an activation function such as sigmoid a single layer network output is given by
- $$out_i = \text{Sigmoid} \left( \sum_i i_n \cdot w_{ii} \right)$$

Due to derivative properties of sigmoid, weight update equation is

$$\Delta w_{ki} = \eta \sum_p (\text{targ}_p - out_i) \cdot f' \left( \sum_i i_n \cdot w_{ii} \right) \cdot i_k$$

where  $\eta$  is learning rate and  $\text{targ}$  is desired output. This equation can be simplified to contain only neuron activation function and no derivatives.

$$\Delta w_{ki} = \eta \sum_p (\text{targ}_p - out_i) \cdot out_i \cdot (1 - out_i) \cdot i_k$$

This is known as generalized Delta learning rule.

- \* Limitations of MLP:
  - MLP is Multi-layer perceptron, a class of feed-forward network. It has  $\leq 3$  layers. Every neuron uses a non-linear activation  $f^n$ . Uses backpropagation for training.
  - Limitations:
    - 1) Computations are difficult & time consuming.
    - 2) Not known to what extent independent variable is affected by dependent variable.

- 3) Proper function of model depends on quality.
- 4) It includes too many parameters being fully con-

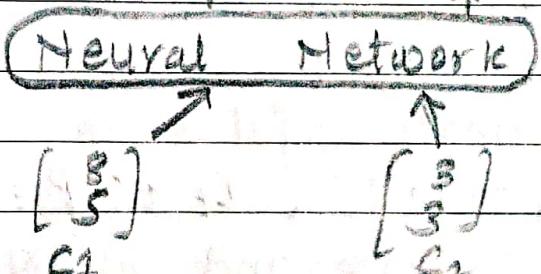
## UNIT V: Convolution Neural Network

### \* Recurrent Neural Network:-

- To work with phrases and sentences we need relation between word & multi-word expressions.
- For example "consider" == "take into account".
- To solve this problem we use recurrent neural network.
- It uses binary tree to identify phrases or sentences.
- We recursively merge pairs of representation of segments to get representation of bigger segments.
- The tree is created based on 2 rules:-

  - Semantic representation if two nodes are merged.
  - Score of how matching two words are.

$\text{Score} = \frac{1}{2} \cdot 3$        $\{S\}$  = parent

### \* Recurrent Neural Network:-

- Done in ML

### \* Long-Short Term Memory:-

- An advanced RNN, sequential network that allows

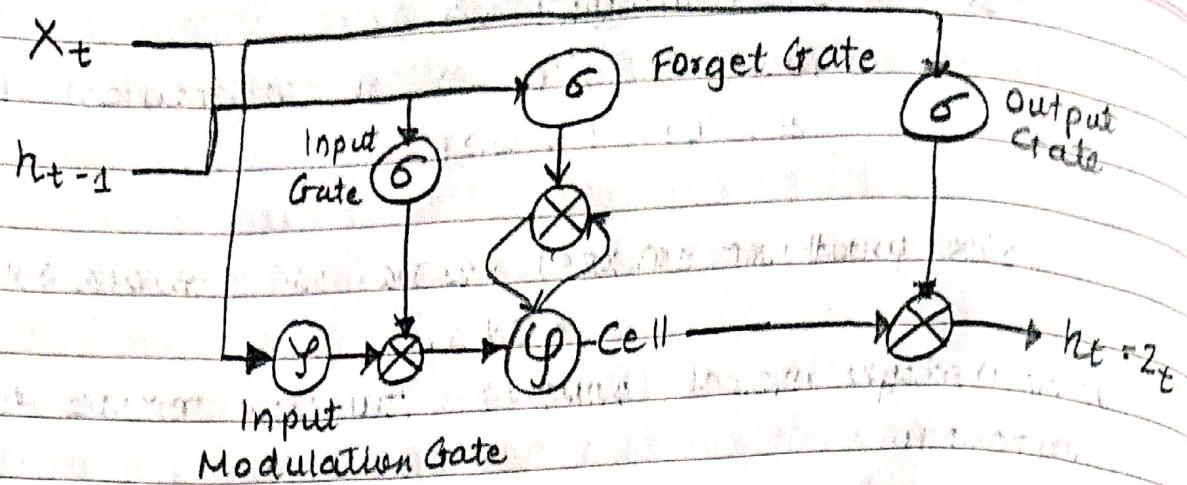
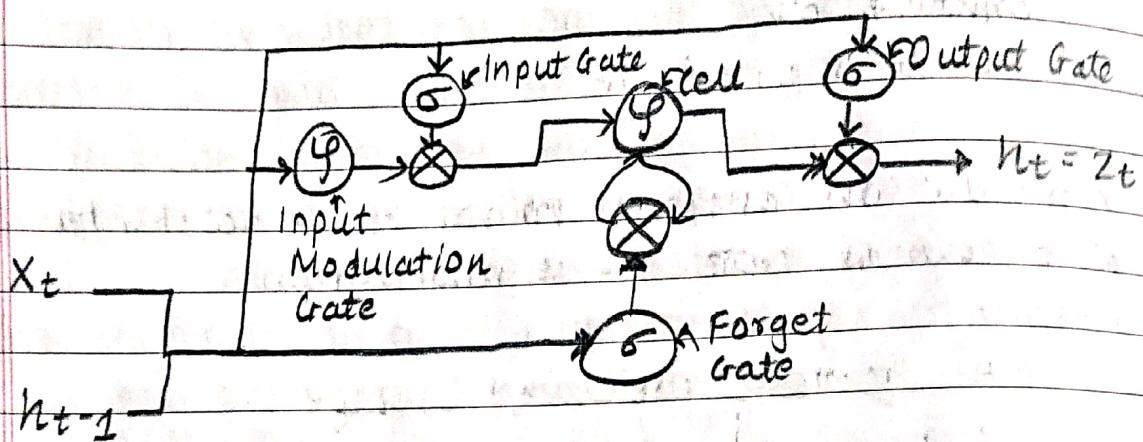


Fig: LSTM

Let's redraw to simplify



### 1. Forget Gate:

- Info that's no longer useful is removed with this gate.
- $x_t$ : Input at a time "t",  $h_{t-1}$ : Previous cell output
- These are fed to gate and multiplied with weights.
- Result is passed through a binary activation function.
- If output is 1 info is retained, if 0 it is forgotten.

### 2. Input Gate:

- Info is regulated and filtered out values that need to be remembered.
- Modulated input and normal input are multiplied together to get useful information.

### 3. Output Gate:

- Extracts useful info from current cell state
- output is sent as an input to next cell.

## \* Gradient Descent Optimization :-

- There are various methods that make gradient descent more optimized. Let us see :-

### 1) Momentum Method:

- To make gradient descent faster, we consider exponentially weighted average of gradients.
- Using average makes converging towards minima faster.
- Pseudocode:-

$$v = 0$$

for each iteration  $i$ :

compute  $dw$

$$v = \beta v + (1 - \beta) dw$$

$$w = w - \alpha v$$

- $v$  and  $dw$  are analogous to velocity & acceleration respectively
- $\alpha$  is learning rate,  $\beta$  is momentum

### 2) RMSprop:

- Apply exponentially weighted average to second moment of gradients ( $dw$ ).

- Pseudocode:-

$$s = 0$$

for each iteration  $i$ :

compute  $dw$

$$s = \beta s + (1 - \beta) dw^2$$

$$w = w - \alpha [dw / (\sqrt{s} + \epsilon)]$$

### 3) Adam Optimization:

- Adam Optimization = RMSprop + Momentum + Bias correction

- Pseudocode:

$$v = 0$$

$$s = 0$$

for each iteration  $i$ :

compute  $dw$

$$v = \beta_1 v + (1 - \beta_1) dw$$

$$s = \beta_2 s + (1 - \beta_2) dw^2$$

$$v = v / (1 - \beta_1^t)$$

$$s = s / (1 - \beta_2^t)$$

$$w = w - \alpha [v / (\sqrt{s} + \epsilon)]$$

## UNIT VI: Applications Perspective

### \* Text Preprocessing :-

#### \* Tokenization :-

- Process of dividing text into meaningful pieces.

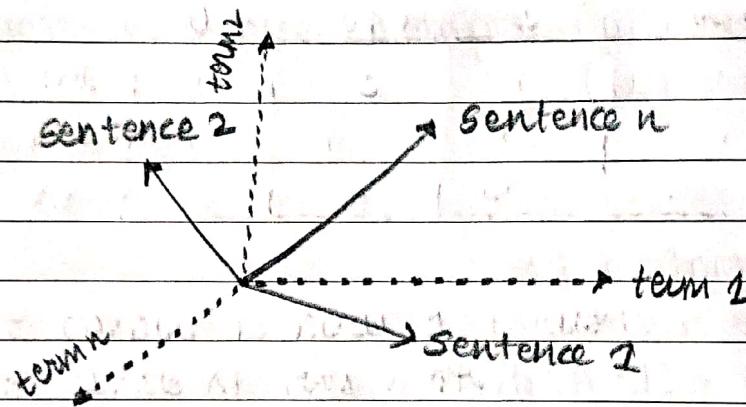
- Pieces are called tokens.

- Most common way of this is whitespace/unigram tokenization in which entire text is split into words by splitting them at whitespaces.

- Another type is regular expression tokenization which will use regular expression pattern to get tokens.

#### \* Document representation :-

- In vector space model each term/word is an axis/dimension. The text/document is represented as a vector in multi-dimensional space.



- Some ways :-
- 1-Hot encoding :-
- We represent words with one hot vectors i.e we associate each unique word with an index in the vector.

Ex :

The	sky	the	great
1	0	1	0
0	1	0	0
0	0	0	0
0	0	0	1

- n-grams :- n-words sequence, we take  $n$  words from sentence and consider it as a token.
- n-gram models estimate the following word probability given previous words.
- Bag of words:
- Used in sentiment analysis uses term frequency (TF).
- There are many more ...

#### \* Feature Selection / Extraction:

- We convert the text in form of features.
- Common techniques to perform feature extraction are:-
- 1) Bag of Words:
  - keep frequency count of all unique words and consider it as feature.
  - Steps involved:
    - i) Identify unique words from document.
    - ii) Find frequencies of unique words from a single sentence.
- 2) Term Frequency - Inverse Document Frequency (TF-IDF):
  - TF : Frequency of each word in document
  - IDF : Assign lower weight to words that appear more frequently, basically depicts rarity of word in document.

$$TF = \frac{\text{Term } i \text{ frequency in document } j}{\text{Total words in document } j}$$

$$IDF = \log_2 \left[ \frac{\text{Total documents}}{\text{documents with term } i} \right]$$

- #### \* Topic Modeling algorithm :- Latent Dirichlet Allocation.
- Latent Dirichlet Allocation (LDA) helps extract topics from a given corpus.
  - LDA classifies text into a document and words per topic, these are modelled based on dirichlet distributions and processes.

- Assumptions made by LDA:
  - i) Documents are a mixture of topics.
  - ii) Topics are a mixture of tokens.
- Working of LDA:  
Consider a corpus [collection of documents] has 3 documents:

Document 1: Are you watching closely?

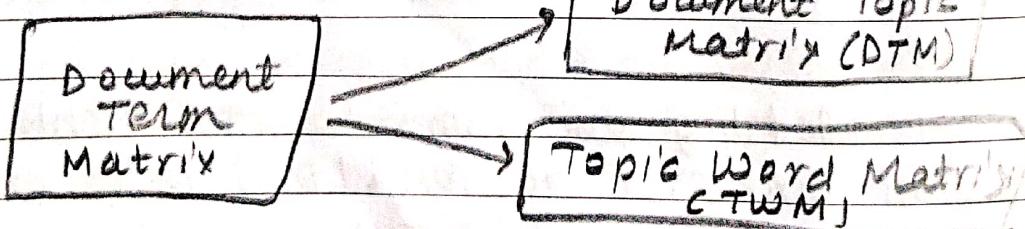
Document 2: You either die a hero or live long enough to see yourself become the villain.

Document 3: Wake up to reality. Nothing ever goes as planned in this accursed world.

$D_1, D_2, D_3$  are examples of 3 documents and consider words are represented as  $W_s$  and for simplicity we consider there are 5 unique words

	$w_1$	$w_2$	$w_3$	$w_4$	$w_5$
$D_1$	0	0	1	1	1
$D_2$	0	1	1	1	1
$D_3$	1	0	0	0	0

The corpus is now a document-word matrix  
LDA converts this matrix into : Document-topic matrix and topic-word matrix



DTM: Possible topics

TWM: Words that the topic will contain

## \* Text Similarity Measures:-

- Used to find similar text.
- Google, Quora, Stack Overflow use it to find similar questions.
- Let us see some similarity measures:-

### • Jaccard Similarity:

- Ratio of common words to total words

$$\text{Jaccard Similarity} = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

- It does not give importance to duplication of words.

### • Cosine Similarity:-

$$\text{Similarity}(A, B) = \frac{A \cdot B}{|A| \times |B|} = \frac{\sum (A_i \times B_i)}{\sqrt{\sum A_i^2} \times \sqrt{\sum B_i^2}}$$

- Suitable when words are repeated and are important.
- It is the ratio of dot product of two vectors of words to their product of magnitude.