

# NATURAL LANGUAGE PROCESSING

→ TEXT DATA

NLP can be used in both machine learning and deep learning.

## Why NLP?

Till now we have just seen numerical data which can be easily solved by ML and DL algorithm. but when it comes to textual data it becomes very hard to solve that's why NLP was introduced.

Human → Communication → Machine

→ Here we need NLP

## Roadmap:



## Tokenization

Converting sentences into words

Eg: Hello, How are you? → Sentence

Hello How are you → Word

## Stopwords

The words which help to form sentences

Eg: tea with toast  
↳ with Stopword.

## Stemming

Process of reducing words to their base word stem.

Eg: Creation  
Creativity  
Creating  
→ Create  
↳ Stem word

## Cons of stemming

↳ Removing meaning of the word

going  
goes  
gone  
→ go  
↳ Stem word

## Pros of stemming

↳ Faster

Stemming usecase → Spam classification  
Review classification

## Lemmatization

process of reducing the words to their base form with meaning

Eg:



### Cons of lemmatization

↳ Slower than stemming

### Pros of lemmatization

↳ Meaningful word



## Basic Terminology used in NLP

1. CORPUS → paragraph
2. Documents → sentence
3. Vocabulary → unique words
4. Words → all words

## Text Preprocessing

Steps:



## Convert word into Vectors

### 1. One hot encoding

[ A man eat food  
  cat eat food   → Corpus ,  
  people watch movie ]   vocabulary - 8

Document 1 → A mat eat food

[ [ 1 0 0 0  
      0 1 0 0  
      0 0 1 0  
      0 0 0 1 ] ]

↳ Sparse matrix

### Advantage      Disadvantage

- Simple to implement - Sparse matrix

↳ majority of zero present

- Out of vocabulary → extra word not handle

- Size is not fixed

- Semantic meaning between  
word is not captured.

Document 2 → Cat eat food → Vocabulary size

[ 1 0 0    is decreasing ↓  
      0 1 0  
      0 0 1   ↓  
                  cannot train model  
                  ↓  
                  Input size not same

## 2 Bag of words (Bow)

→ Removed  
Stopwords

Document 1 → He is @ good boy → good boy

Document 2 → She is @ good girl → good girl

Document 3 → Boys and girls are good → Boy girl good

Vocabulary      Frequency

Good	3
Boy	2
girl	2

feature<sub>1</sub>      feature<sub>2</sub>      feature<sub>3</sub>  
Good              Boy              Girl

document 1	1	1	0	→ Binary Bow
document 2	1	0	1	
document 3	1	1	1	

### Advantage

Simple and Intuitive

### Disadvantage

- Sparsity will be there
- Out of vocabulary
- Ordering of word get change
- Semantic meaning between word not captured

↳ Cosine Similarity

In order to capture the semantic information

↳ we use ngrams

N-GRAM → Bigram, Trigram, - Ngram

	feature <sub>1</sub>	feature <sub>2</sub>	feature <sub>3</sub>	feature <sub>4</sub>	feature <sub>5</sub>	
	Good	Boy	Girl	good boy	good girl	will check combo of two words are present or not in document 1, 2, 3 if present than we will fill value as 1 otherwise 0.
document 1	1	1	0	1	0	
document 2	1	0	1	0	1	
document 3	1	1	1	0	0	

Example :

Anirudh eats food → Bigram?

Anirudh eats , eats foods } → 2 possible bigram

I am not feeling well → Trigram?

I am not , am not feeling , not feeling well } → 3 Possible Trigram

ngram(1,3)

Anirudh is not feeling well      unigram → Bigram → Trigram

f<sub>1</sub> f<sub>2</sub> f<sub>3</sub> f<sub>4</sub> f<sub>5</sub>      Bigram      Trigram

Anirudh is not feeling well      — — —

TF-IDF

Bow is not able capture exact meaning or semantic meaning so to solve this we use concept of TF-IDF

Term frequency - inverse document frequency (TF-IDF)

Example

Senti - Good boy

Sent2 - good girl

Sents - boy girl good

Good present in all

boy and girl present in only few sentences

we will give some weights  
to boy and girl in that  
sentences only where the  
word is present

Term frequency =  $\frac{\text{no of repetition of words in sentences}}{\text{no of words in sentence}}$

Inverse document frequency =  $\log_e \left( \frac{\text{no of sentences containing the word}}{\text{no of sentences}} \right)$

$$TF * IDF = TFIDF$$

Sent<sub>1</sub> - good boy

Sent<sub>2</sub> - good girl

Sent<sub>3</sub> boy girl good

## Term frequency

## Inverse Document frequency

	Sent <sub>1</sub>	Sent <sub>2</sub>	Sent <sub>3</sub>	Words	IDF
good	1/2	1/2	1/3	good	$\log_e(3/3) = 0$
boy	1/2	0	1/3	boy	$\log_e(3/2)$
girl	0	1/2	1/3	girl	$\log_e(3/2)$

$f_1$        $f_2$        $f_3$   
good      boy      girl

Sent <sub>1</sub>	$\frac{1}{2} * 0$	$\frac{1}{2} * \log_e(3/2)$	$0 * \log_e(3/2)$
Sent <sub>2</sub>	$\frac{1}{2} * 0$	$0 * \log_e(3/2)$	$\frac{1}{2} * \log_e(3/2)$
Sent <sub>3</sub>	$\frac{1}{3} * 0$	$\frac{1}{3} * \log_e(3/2)$	$\frac{1}{3} * \log_e(3/2)$

Over here some amount of semantic meaning is captured.

Common in all sentences

## Advantage

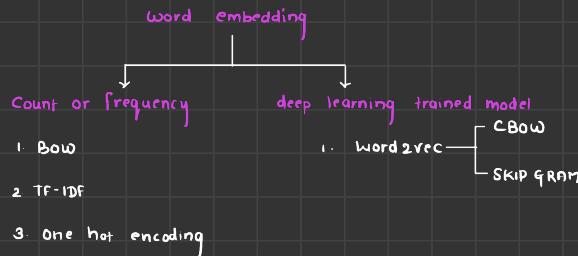
- Intuitive
- word importance is getting captured

## Disadvantage

- sparsity
- out of vocabulary

## Word Embedding

technique which converts the words into vectors

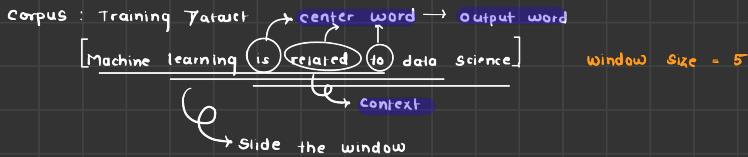


## Word2vec

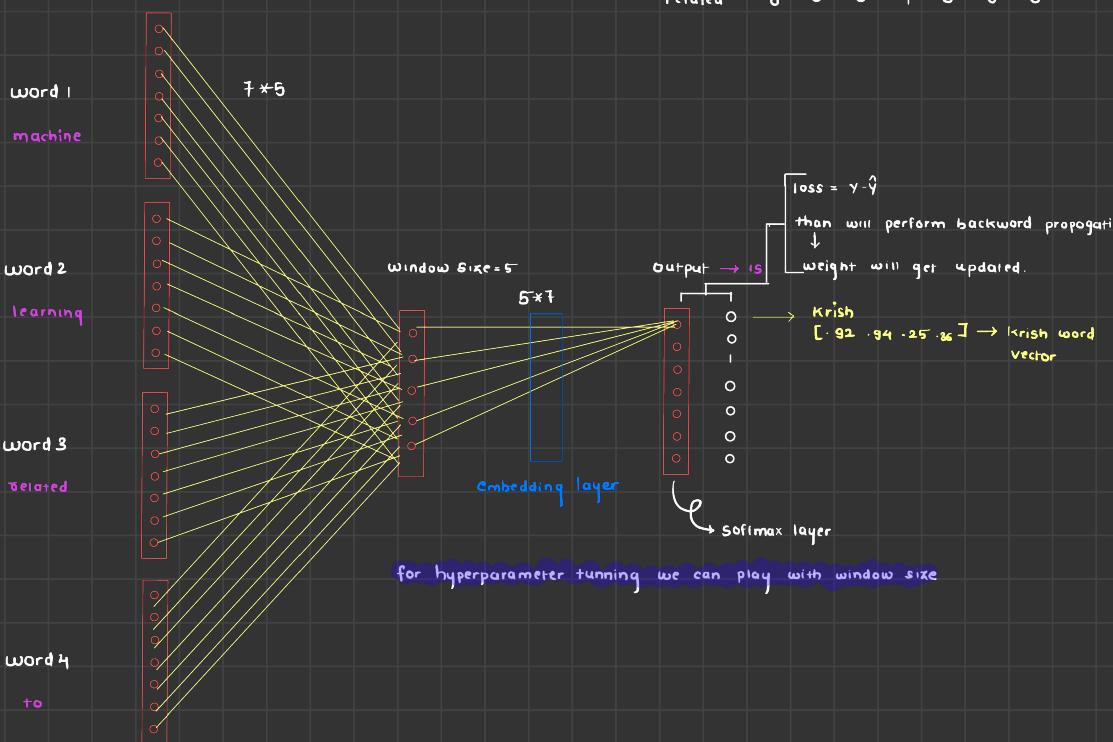
## feature representation

- In word2vec we create vector for each word with limited dimension.
  - Sparsity is reduced
  - Semantic meaning is captured

## CBOW - Continuous Bag of words



Independent features	output	Bag of word
Machine, learning, related, to	is	Machine 1 0 0 0 0 0 0 learning 0 1 0 0 0 0 0 is 0 0 1 0 0 0 0 related 0 0 0 1 0 0 0
learning, is, to, data	Related	
is, related, data, science	TO	



## Skipgram

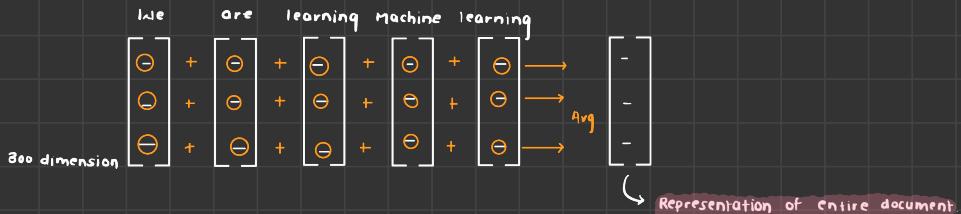
Reverse the CBOW network will get skipgram

Input	Output
IS	Machine, learning, related, to
Related	learning, is, to, data
To	is, related, data, science

## Avgwordvec

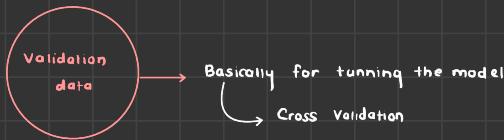
→ used to represent the vector of entire sentence or document not single word.

Example: We are learning Machine learning



## Interview Question

A. Train vs Test vs Validation



B. Why Random forest instead of decision tree?

• Reduce overfitting and bias

• Generalized

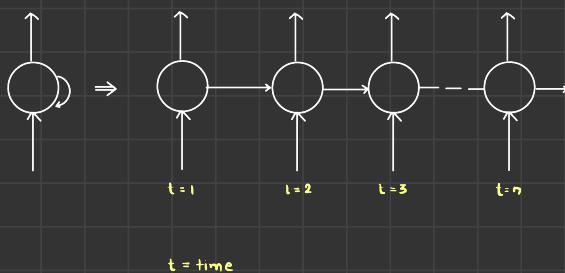
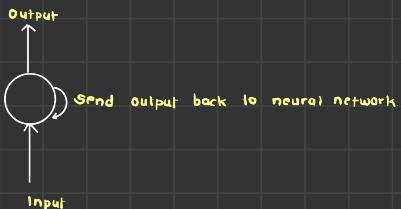
• Higher Accuracy

# RECURRENT NEURAL NETWORK

Use Case:



Basic Architecture of RNN

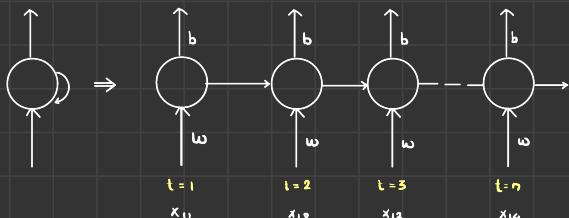


Example

Sentiment Analysis      Output  
The food is good      positive

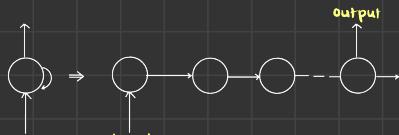
At time  $t_1$  will pass  $x_{11}$  and weight  $w$  will initialize.

$x_{11} \rightarrow \text{Word2vec} \rightarrow \text{Vector}$   
 $\hookrightarrow \text{dimension} = d = 3$



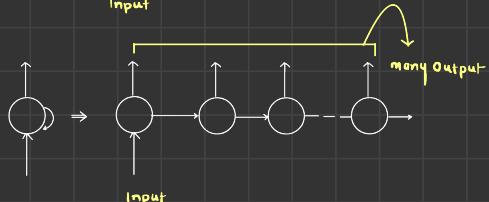
## Types of RNN :

### 1. One to One Rnn



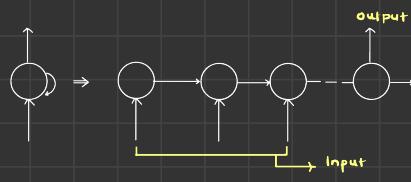
### 2. One to many Rnn

Eq: Text generation



### 3. Many to one Rnn

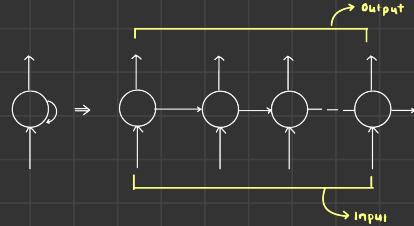
Eq: Sentiment analysis



### 4. Many to many Rnn

Eq: language translation

chatbot



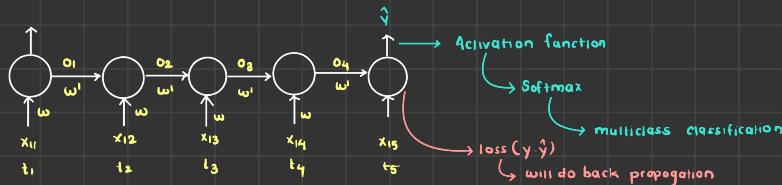
## Forward Propagation in Rnn :

Eq: Sentiment Analysis

The food is very good

positive

$x_{11} \quad x_{12} \quad x_{13} \quad x_{14} \quad x_{15}$

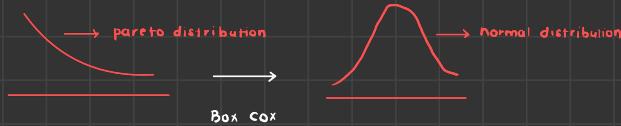


$$o_1 = f(x_{11} * \omega) + b$$

$$o_2 = f(x_{12} * \omega + o_1 * \omega') + b$$

$$o_3 = f(x_{13} * \omega + o_2 * \omega' + o_1 * \omega'') + b$$

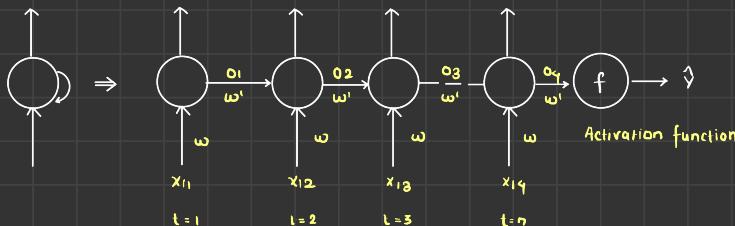
### Interview Question



Transformation can convert pareto to normal



### Back Propagation



$$o_1 = f(x_1 * \omega_1)$$

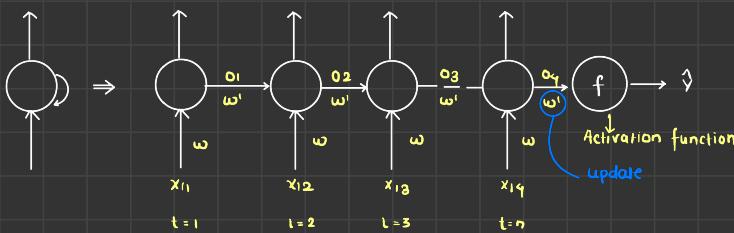
$$o_4 = f[(x_{14} \omega_1) + (o_3 \omega_1)]$$

$$o_3 = f[(x_{13} \omega_1) + (o_2 \omega_1)]$$

$$o_2 = f[(x_{12} \omega_1) + (o_1 \omega_1)]$$

In back propagation we update the weights

first we will calculate the loss than we will update weights.



weight updation formula

$$w^i_{\text{new}} = w^i_{\text{old}} - \eta \frac{\partial L}{\partial w^i}$$

$$\frac{\partial L}{\partial w^i} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w}$$

→ Chain Rule

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

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial o_4} * \frac{\partial o_4}{\partial w}$$

→ Chain Rule

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

$$\frac{\partial L}{\partial w} = \frac{\partial L}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial o_4} * \frac{\partial o_4}{\partial w^i}$$

→ Chain Rule

In back propagation we may face problem of vanishing gradient so to solve this we use LSTM RNN

LSTM stand for long short term memory

In LSTM RNN we have four main section →

- 1 memory cell
- 2 forget cell
- 3 input cell