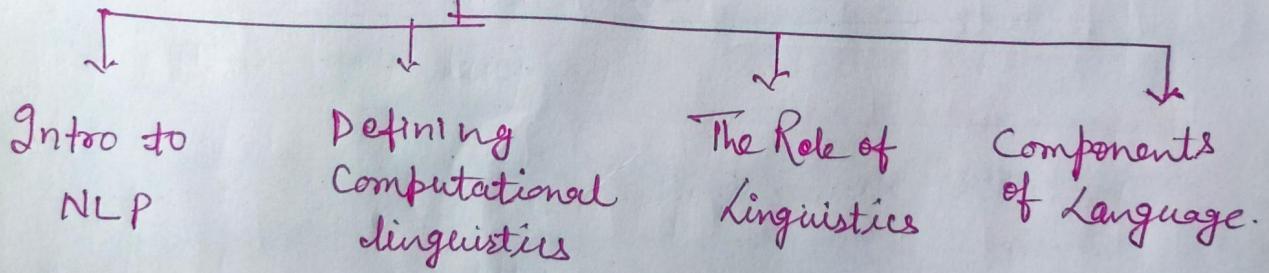


# \* NLP (Natural Language Processing)

- NLP stands at the intersection of linguistics, computer science, and artificial intelligence.
- It enables machines to understand, interpret and respond to human language.
- NLP's applications range from chatbots and language translation to sentiment-analysis and content generation.
- It is a computational techniques for processing and analyzing human language.

## the Art and Science of NLP



- Defining Computational Linguistics →
  - It combines linguistic theories with computational methods.
  - It involves developing algorithms for automating language analysis, processing, and generation.
  - This field bridges the gap between language and machines, enabling meaningful interactions.

# Role of Linguistics →

## ① Semantic Understanding →

Linguistic insights enable NLP models to comprehend word meanings, nuances, and contextual subtleties, enhancing accurate interpretation of text.

## → Syntax Analysis →

→ Linguistic theories guide NLP systems in structurally analyzing sentences, identifying grammatical components, and understanding relationships between words.

## → Named Entity Recognition (NER) →

Linguistic patterns assist NER algorithms in identifying and categorizing entities like names, dates and locations, aiding in information extraction.

## \* Components of Natural Language →

→ Human language comprises various components, including syntax (sentence structure), semantics (meaning), morphology (word forms), and phonetics (speech sounds).

→ NLP algorithms must decipher each of these components to comprehend and generate language accurately.

## → Benefits of Computational Linguistics →

### Efficient Language Processing:

Computational linguistics streamlines the analysis, understanding, and generation of human language, enhancing the efficiency of various tasks.

### Improved Communication: -

It bridges language barriers, enabling seamless communication between individuals speaking different languages.

### Enhanced Data Insights: →

Computational linguistics extracts valuable insights from textual data, supporting informed

### Innovation and Automation: →

By automating language-related tasks, it paves the way for innovative applications such as chatbots, sentiment analysis, and content generation.

### Applications of NLP: →

- (1) Sentiment Analysis
- (2) Language Translation
- (3) Named Entity Recognition
- (4) Deep learning for NLP.

## History of NLP ⇒

### Various NLP Tasks ⇒

Language Modeling → This is the task of predicting what the next word in a sentence will be based on the history of previous words.

→ It is useful for building solutions for a wide variety of problems, such as speech recognition, OCR, handwriting recognition, machine translation.

### Text Classification →

This is the task of classifying the text into a known set of categories based on its content.

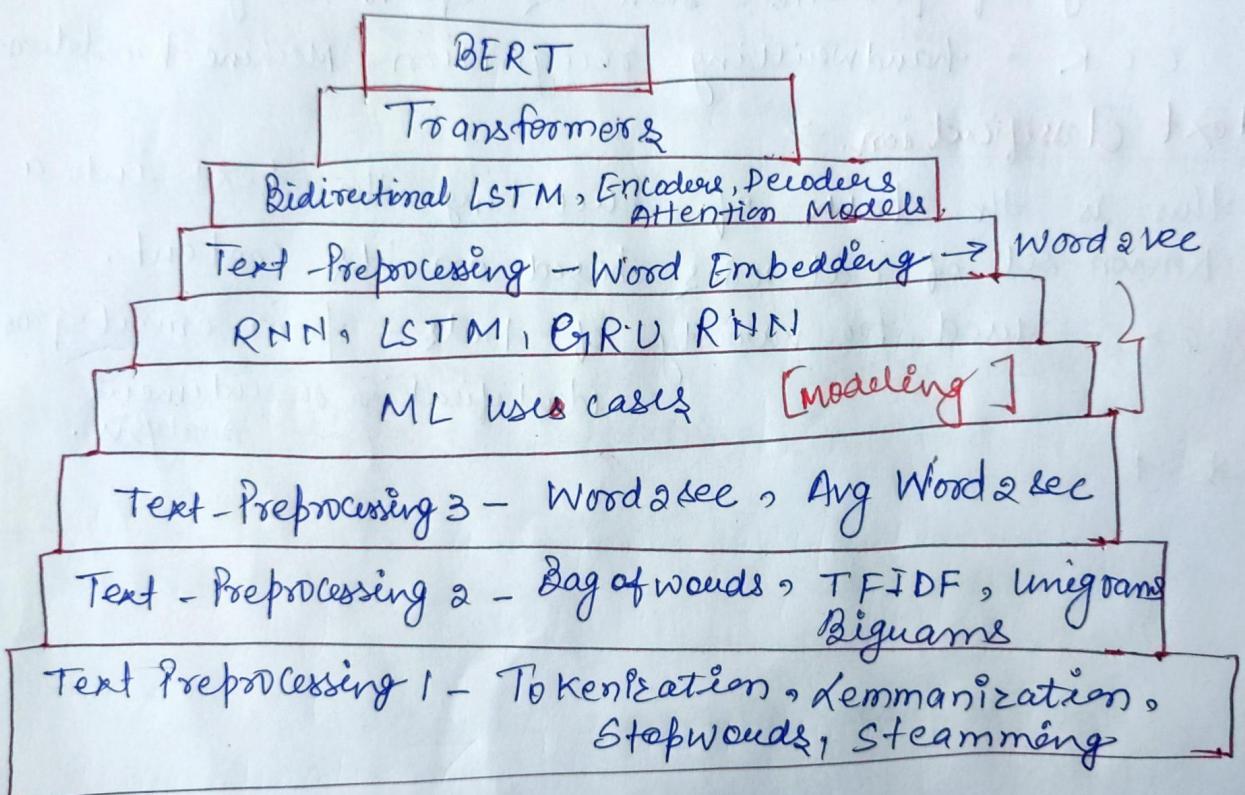
Ex - Used for variety of tools - from email spam identification or sentiment analysis.

Why NLP? → NLP used on text based dataset.

→

Examples of NLP →

- Chat-bots
- Email Filters
- Sentiment Analysis
- Machine Translation



Libraries used in NLP → NLTK, spaCy, TextBlob

① Tokenization ⇒ (Breaking text into words or phrases)

- It is converting sentence into words.
- It is the process of tokenizing or splitting a string, text into a list of tokens.

Stemming → It is the process of reducing a word to its word stem that affixes to suffixes and prefixes or the Roots.

Ex:-

going }  
goes } → go (meaningful word)  
gone }

finally }  
final } ⇒ final (meaning is gone)  
finalized }

Advantage → It is really fast.

Disadvantage → It is removing the meaning.

→ Stemming is a NLP technique that is used to reduce words to their base form as root.

\* Lemmatization → (Reducing words to ~~to~~ their base form)  
It is the process of grouping together the different inflected forms of a word so they can be analyzed as a single item.

history  
historical > history

→ Lemmatization does morphological analysis of the words.

Ex :- better → good  
rocks → rock  
Corpora → Corpus

Applications →

- Used in compact indexing.
- Used in comprehensive retrieval systems like search engines.

Advantage → Meaningful words

Disadvantage → It is slow.

Use cases in stemming → (1) Spam classification  
(2) Review classification

Use cases of lemmatization → (1) Text summarization  
(2) Chat bot  
(3) Language translation.

Text Preprocessing (1) →

- (1) Tokenization
- (2) Stop words
- (3) Stemming
- (4) Lemmatization

Stop words Removal

Removing common words with no significant meanings like "the", "and", "a".

Parts-of-Speech Tagging

Identifying the syntactical properties of words.

## Text preprocessing Step 2 : → [Feature engineering]

Convert words into vectors.

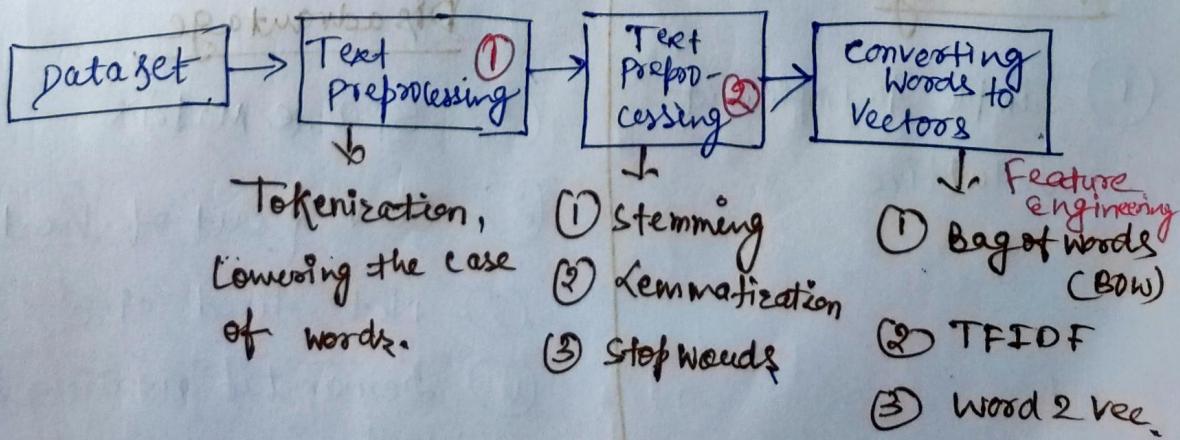
words → vectors.

Techniques used in step 2 :

- ① Bag of words
- ② TFIDF
- ③ Word2Vec.
- ④ One hot Encoding
- ⑤ Inverse document Frequency
- ⑥ Ngrams.

\* Basic Terminologies used in NLP →

- ① CORPUS → Paragraph (Data points)
  - ② Documents → It is like sentences
  - ③ Vocabulary → Unique words in dictionary
  - ④ Words.
- ① One hot encoding → output should be in binary (0/1) format.



## \* (1) ~~bag of words~~

## \* ① One hot encoding →

Ex :- [ A man eat food  
① Cat eat food } Coopers  
People Watch KRISH YT ]

Paragraph ( no. of words ) = 9

A man eat food cat people watch KRISH YT

\* Represent as one hot encode,  $\rightarrow$

D<sub>1</sub> [ [ 10000000 ], [ 01000000 ] ]  
 [ 00100000 ], [ 00010000 ] ]  
east      food

In this encode there are lots of issues -

- ① Sparsity
  - ② When Vocab size decreases then can't train the model.

## Advantage

- ① Simple to Implement
  - ② Intuitive

## Disadvantage

- (1) Sparse Matrix
  - (2) OOV {out of Vocabulary}
  - (3) Not fixed size
  - (4) Semantic meaning between word is not capture.

## (2) Bag of words →

- $D_1 \rightarrow$  He is a good boy → good boy  
 $D_2 \rightarrow$  She is a good girl → good girl  
 $D_3 \rightarrow$  Boys and girls are good → Boys girls good

Remove these words

Vocabulary → Frequency

good → 3 times

boy → 2 times

girl → 2 times

Doc1  
Doc2  
Doc3

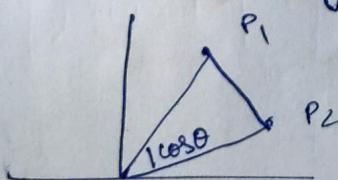
	$f_1$	$f_2$	$f_3$	
→ good	1	1	1	0/p
→ boy	1	1	0	→ not present
→ girl	1	0	1	1
$D_3 \rightarrow$	1	1	1	→ All present

→ In BOW, we can make binary BOW.

### Advantages

① It is simple and Intuitive

② Cosine Similarity



$$\cos \theta = 0.53$$

### Disadvantages

① sparsity

② OOV

③ Ordering of the words

$$1 - 0.53 = \text{Cos-Similarity}$$

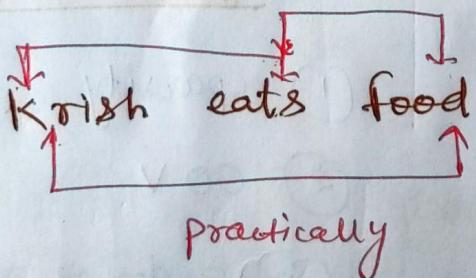
\* Ngrams → It helps to find semantic meanings.

	$f_1$ good	$f_2$ boy	$f_3$ girl
Sentence 1	1	1	0
Sent 2	1	0	1
Sent 3	1	1	1

shown with words

Bigrams → Apart from using single features we will be using common features also.

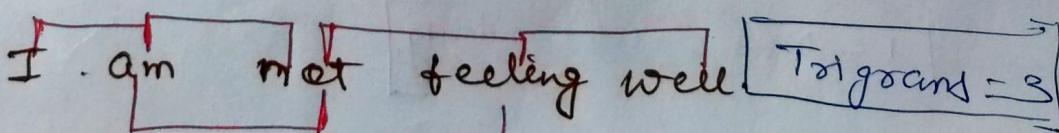
	$f_4$ Good boy	$f_5$ Good girl
S <sub>1</sub>	1	0
S <sub>2</sub>	0	1
S <sub>3</sub>	0	0



Krish eats

Eats food

Trigrams →



I am not, am not feeling, not feeling well.

~~f<sub>1</sub>~~ f<sub>2</sub> f<sub>3</sub> f<sub>4</sub> f<sub>5</sub>  
Krish is not feeling well

### NLTK. sent\_tokenize() →

Return a sentence-tokenized copy of "text" using NLTK's recommended sentence tokenizer.

T.e. re.sub(pattern, repl, string, count=0, flags=0) ⇒

Return the string obtained by replacing the leftmost non-overlapping occurrences of the pattern in string.

1.Q If we remove the words like a, the, is, an  
we use ?

Ans → Stop word

2.Q Bag of words in text preprocessing is a -  
Ans → Feature extraction Technique

3.Q. For the sentence "cat eat food", how many bigrams can be created?

Ans → 2

cat eat food

4.Q Collections of documents is called ,

Ans → Corpus

5.Q tf-idf is used in  
→ Text processing & Page ranking in search engine

① Bagwords → Text → Vectors

Sentence 1 → He is a good boy. → good boy

Sentence 2 → She is a good girl → good girl

Sentence 3 → Boy and girl are good → Boy girl  
good

②

Stop words → used to remove unnecessary words (. is, a, he, she, are)

Sent 1 good boy

Sent 2 good girl

Sent 3 boy girl good

③ Frequency (Vocabulary) →

Frequency

good 3

boy 2

girl 2

	Good	boy	girl
	$f_1$	$f_2$	$f_3$
Sent 1	1	1	0
Sent 2	1	0	1
Sent 3	1	1	1

### Disadvantage

- ① Out of vocabulary
- ② Computation is high.
- ③ Semantic meaning is missing.

→ So T F F DF used to remove semantic meaning missing.

# \* TF - IDF (Term Frequency-Inverse Document frequency)

Sent 1 : good boy

Sent 2 : good girl

Sent 3 : boy girl good

} which ever words are rarely present in the sentences we should give more weighted:

Term - Frequency  $\rightarrow$  (TF)

$\frac{\text{No. of repetitions of words in sentence}}{\text{No. of words in sentence}}$

Inverse Document frequency (IDF)

$\log \left( \frac{\text{No. of sentences}}{\text{No. of sentences containing the word}} \right)$

Term frequency (TF)  $\rightarrow$

	Sent 1	Sent 2	Sent 3
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$

TF x IDF

\* Inverse Document Frequency (IDF)  $\rightarrow$

$$\text{good } \log_e \left( \frac{3}{2} \right) = 0$$

$$\text{boy } \log_e \left( \frac{3}{2} \right)$$

$$\text{girl } \log_e \left( \frac{3}{2} \right)$$

## TF \* IDF

	$f_1$	$f_2$	$f_3$
	good	boy	girl
Sent 1	$\frac{1}{2} \times 0 = 0$	$\frac{1}{2} \times \log(3/2)$	$\frac{1}{2} \times \log(3/2) = 0$
Sent 2	$\frac{1}{2} \times 0 = 0$	$0 \times \log(3/2) = 0$	$\frac{1}{2} \times \log(3/2)$
Sent 3	$\frac{1}{2} \times 0 = 0$	$\frac{1}{3} \log(3/2)$	$\frac{1}{3} \log(3/2)$

here, some amount of semantic information is present.

→ Words that are frequently present and words that are only present in a specific sentence.

### Advantages

① It is Intuitive

② Word importance is getting captured

→ Overcome sparse matrix.  
→ Sparsity problem removed

### Disadvantages

① Sparsity

② Out of Vocab

③ Not capture semantic meaning.

## \* Word Embeddings $\Rightarrow$ [Words $\rightarrow$ Vectors]

→ It is a technique which converts words into vectors.

### Word Embedding

Count or Frequency

Deep learning train  
Models.

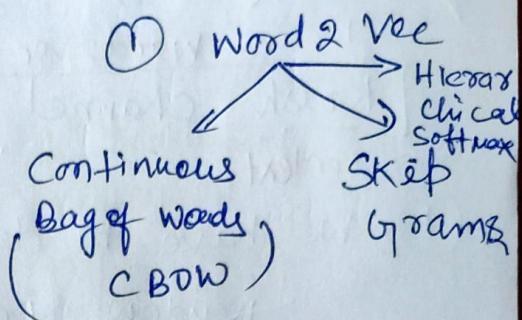
① Count or Frequency →

② Deep Learning Trained  
model

① Bag of words (BOW)

② TF-IDF

③ One Hot Encoding



② Word2Vec → [Feature Representation]

→ Every word is used to create vectors.

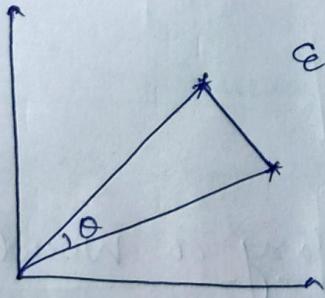
① → That vectors will be in limited dimensions.

② → Sparsity is reduced in Word2Vec.

③ → Semantic meaning of vectors are maintained.

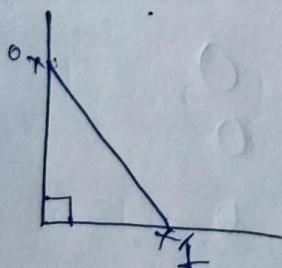
\* Cosine Similarity →

Euclidean Distance -



$$\text{Distance} = \sqrt{1 - \cos(\theta)}$$

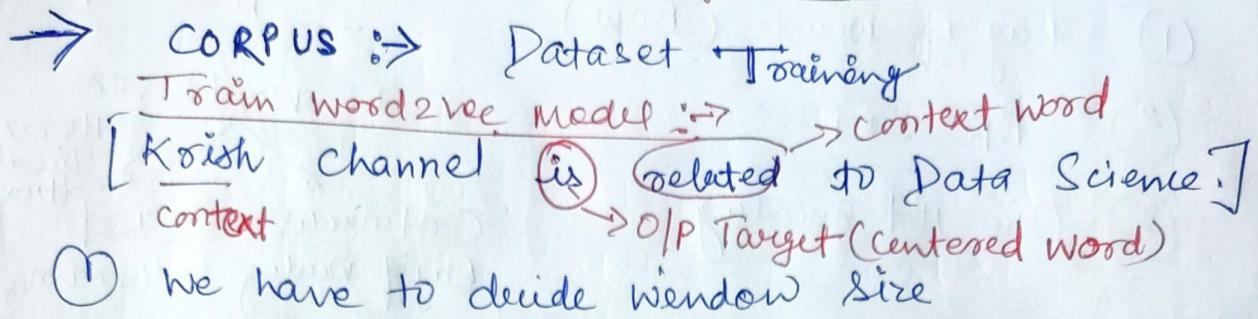
$$\begin{aligned} \cos(\theta) &= \frac{\text{adjacent}}{\text{hypotenuse}} = \frac{1}{\sqrt{2}} \\ \cos(\theta) &= 0.7071 \end{aligned}$$



$$\text{Cosine-Similarity} = \cos(90^\circ) = 0$$

$$\text{Distance} = \sqrt{1 - 0} = 1$$

## \* CBOW (Continuous Bag of words) $\Rightarrow$



### Training Data

#### Independent Feature

- $\rightarrow$  Krish, channel, Related
- $\rightarrow$  To
- $\rightarrow$  Channel, Is, To, Data
- $\rightarrow$  Is, Related, Data, Science

O/P  
 $\downarrow$   
Is

Related  
 $\downarrow$   
To

### \* Bag of words (BOW) $\Rightarrow$

#### Vector representation of vocabulary or word $\rightarrow$

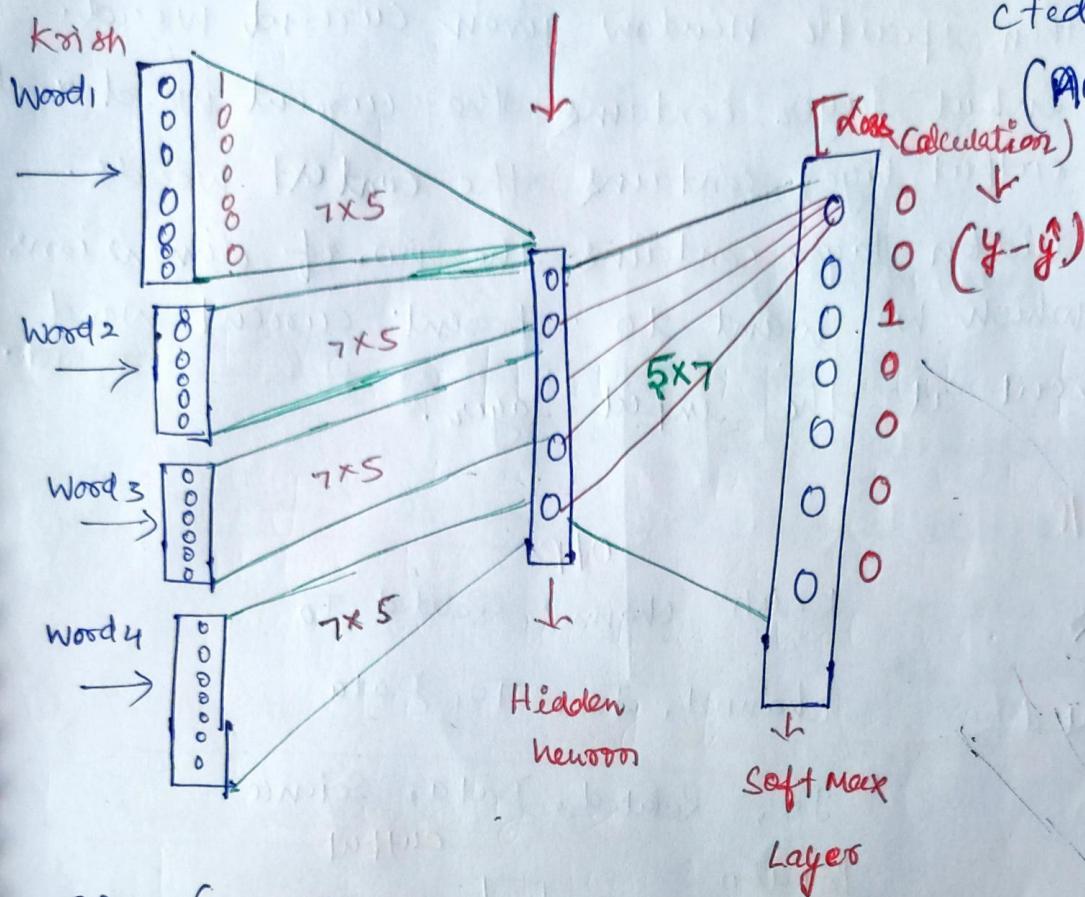
Krish  $\rightarrow$  1 0 0 0 0 0 0

channel  $\rightarrow$  0 1 0 0 0 0 0

is  $\rightarrow$  0 0 1 0 0 0 0

related  $\rightarrow$  0 0 0 1 0 0 0

→ Data representation, independent I/Ps and  
Word size = 5  
ops give fully connected layers.  
(ANN)



CBOW (Continuous Bag of Words) : CBOW model

predicts the current word given context words within a specific window.

• the input layer contains the context words and the output layer contains the current word.

→ the hidden layer contains the no. of dimensions in which we want to represent the current word present at the output layer.

→ CBOW → [ Predicts the target word based on its surrounding context words ]

② Skip Gram  $\Rightarrow$  Predicts the context words based on the target word

- Skip Gram predicts the surrounding context words.
- Within specific window given current word.
- the input layer contains the current word and the output layer contains the context words.
- the hidden layer contains the no. of dimensions in which we want to represent current word present at the input layer.

I/P

Is

Related

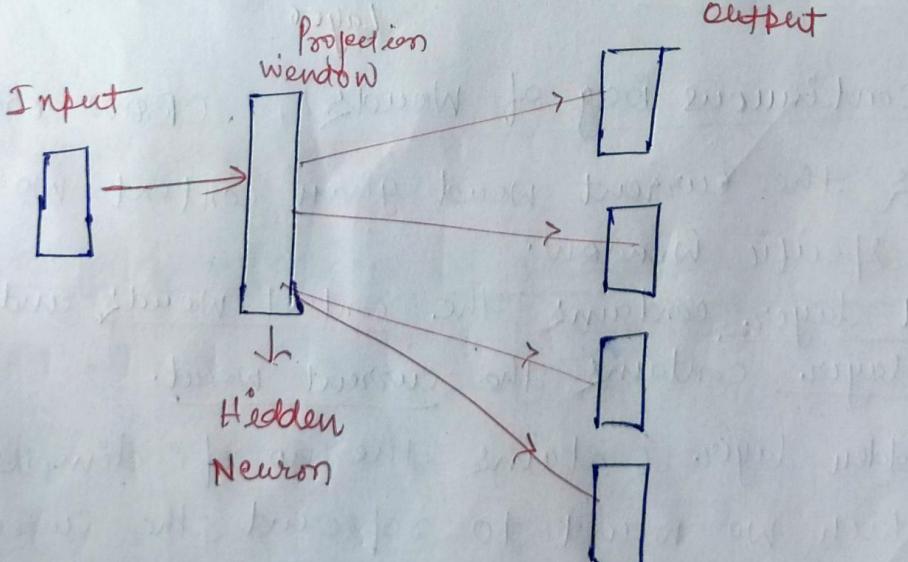
To

O/P

Krish, channel, Related, To

channel, Is, To, Data

output



③ Hierarchical Softmax  $\rightarrow$

- Transforms the problem of predicting one word into a set of binary decisions.

- Word2Vec Applications in NLP → [Enhancing Semantic Understanding]
- (1) Semantic search → Matching synonyms, antonyms, and related concepts.
  - (2) Question answering → Providing answers to questions by understanding the meaning behind them.
  - (3) Chatbots → Teaching machines to understand human language better.
  - (4) Sentiment Analysis → Classifying text as positive, negatives or neutral based on the meaning of the words.

### \* Amplifying NLP with Word Embeddings:

#### Real-world use cases:

- (1) Amazon Recognition → Uses Word2Vec and other algorithms to recognize objects and faces in images.
- (2) Google Translate → Uses Word2Vec to improve machine translation by understanding the meaning of words.
- (3) Spotify → Uses Word2Vec to recommend songs based on the meaning behind the lyrics.
- (4) Google News → Uses Word2Vec to group related articles based on the meaning of the words.

## Word2Vec's Limitations & challenges

### Limitations

- ① Ignorance of concepts and semantic structure
- ② Difficulty handling rare words or complex phrases
- ③ Lack of accounting for different word senses

### Challenges

- ① Optimizing algo performance
- . Reducing the dimensionality of vectors-
- . Dealing with out-of-vocab words and misspellings.