

NLP

[Natural Language Processing]

Text Mining :-

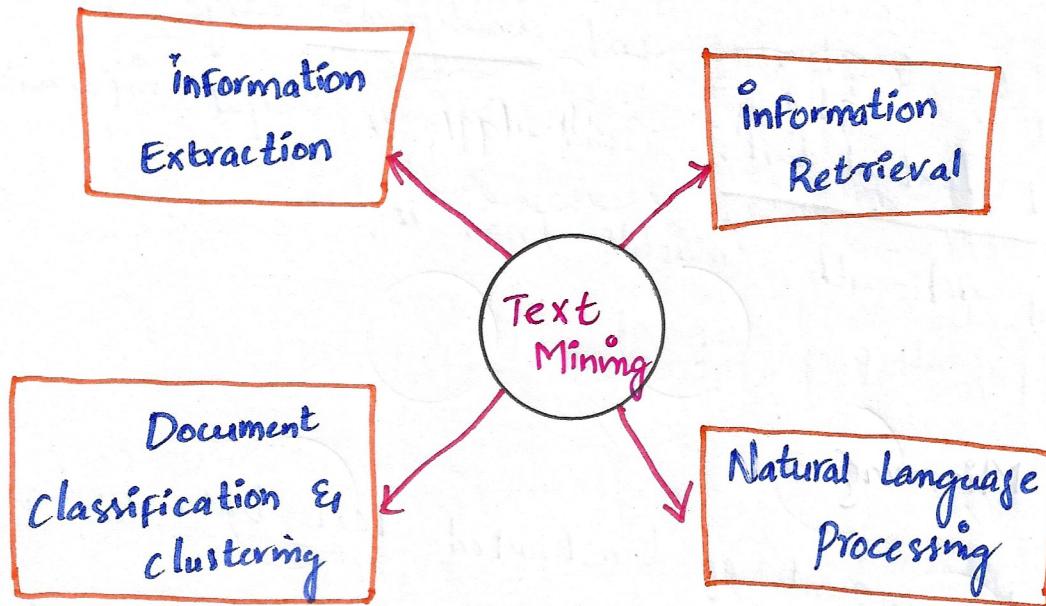
The Process of Exploring and analyzing Large Collection of "Unstructured Textual Data" and Deriving Useful information, patterns, and Actionable insights from it.

Need of Text Mining

- Unprecedented amount of Unstructured text Data being Generated Every Day
- Text mining transforms Unstructured Data into Structured Data to further be used for Analysis
- Identify business insights to feed business processes and reduce risks.
- Make informed Decisions, automate processes, Market research using Sentiment analysis, etc.



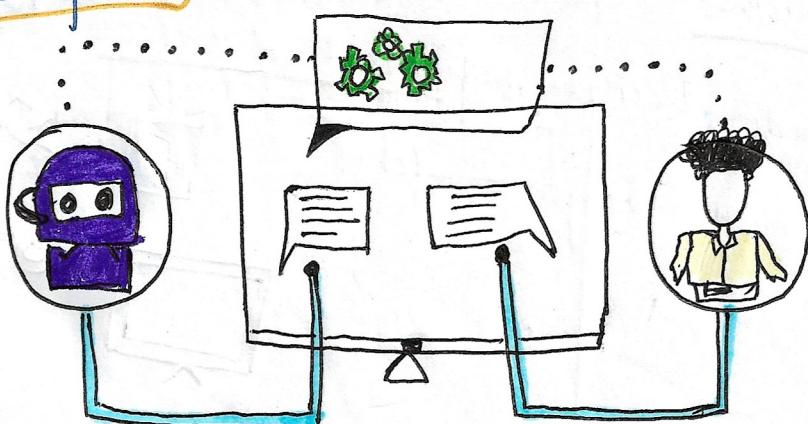
Areas of Text Mining



* NLP [Natural Language Processing]

it is the field of study that deals with the understanding, interpreting and manipulating human language using computers.

Using Computers.



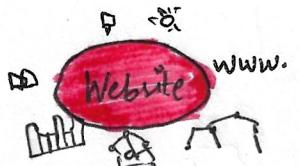
NLP in Text Mining

- Since most of the significant information is written down in Natural language such as English, French, German etc... and it is not conveniently tagged.
- So, after identification and Extraction of the Content needed for text analysis, we can use different NLP techniques to Extract meaningful information from it.
- NLP helps computers to communicate with humans in their own language and perform other language related tasks.
- NLP makes it feasible for computers to read text, interpret etc. it measures Sentiment and determine which party are important.

Applications of NLP

- * Sentimental analysis
- * Chat Bot
- * Speech Recognition
- * Machine translation
- * Spelling checking

- * Information Extraction



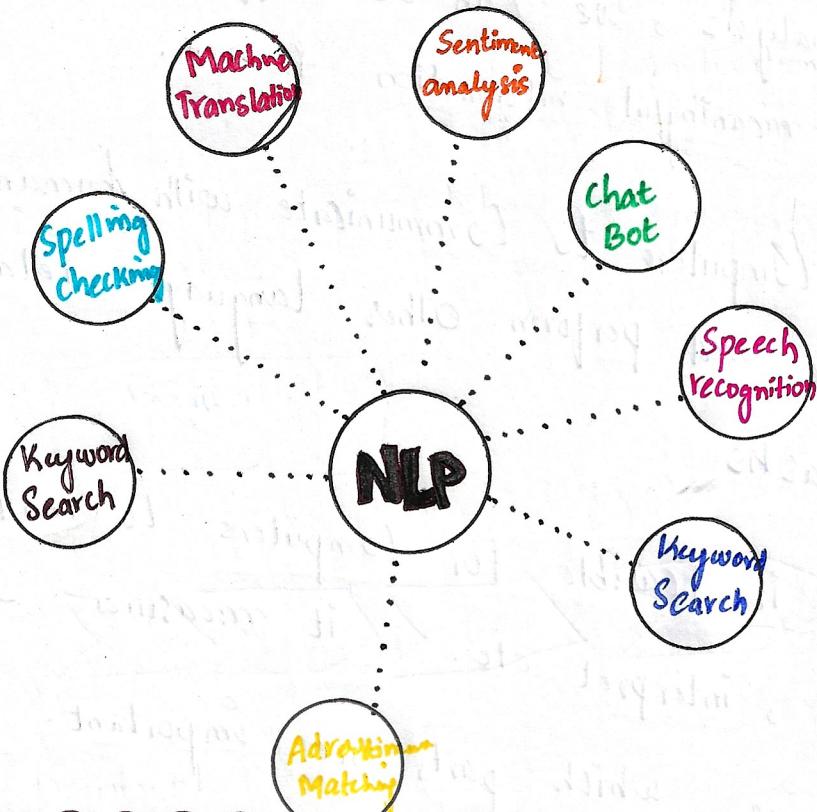
- * Keyword Search



- * Advertisement Matching



Applications of NLP



Natural language Tool Kit

- Python [NLTK, corpora]



- spaCy

- Hugging Face

- BERT

- GPT3

NLTK Corpora

A Corpus is a huge collection of written texts. A Compi
lation of Corpora is called a Corpus. It is a body of
written (or) Spoken Texts used for Linguistic analysis and
the development of NLP tools.

* Corpus :- Multiple Sentences in a Document / File is
called Corpus.

* Corpora :- Multiple "Corpus" is called as Corpora
Ex:- Combination of different G-mails, Ms-word Documents

CODE

```
# import nltk
```

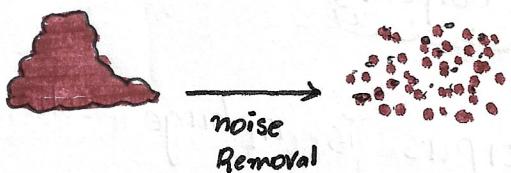
```
# nltk.download()
```

```
# Paragraph = "He is a good boy. She is good girl. boy &  
# girl are good."
```

Steps

- * Read Raw Text
- * tokenize
- * clean text
 - stop words
 - Stemming
 - Lemmatization
- * Vectorize
- * Build Machine Learning model

Text Cleaning



Different methods in Text Cleaning

- * Filter Out Punctuation
- * Split into words [Tokenization]
- * Filter out Stop words [is, the, are]
- * Stemming Words] Base word
- * Lemmatization
- * Normalizing Case

Noise Removal ←

NLT Task



Text Processing

Text Parsing & EDA

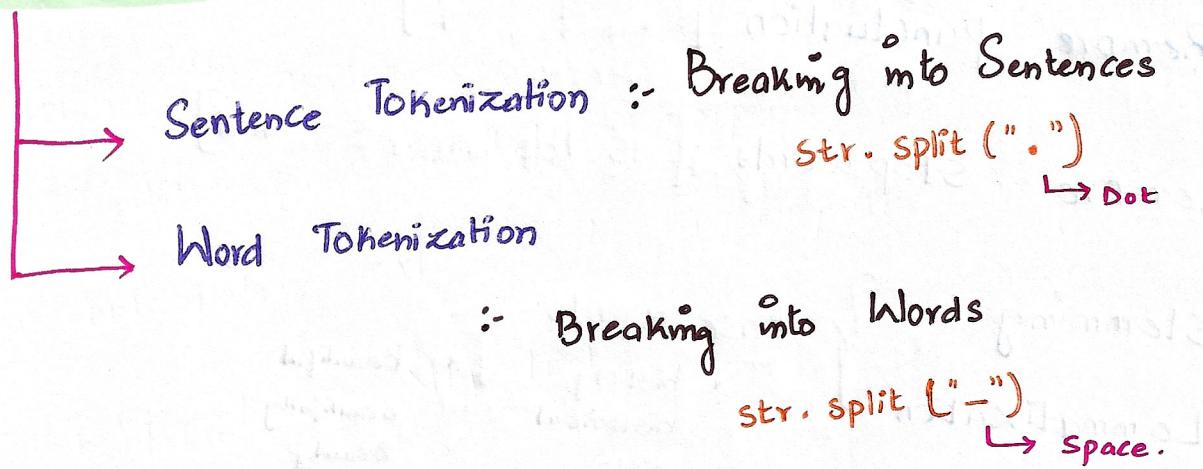
Text Representation & Feature Engineering

Modelling and/or pattern mining

Evaluation & Deployment

- Text Document
- Text processing
- Text Parsing & EDA
- Text Representation & Feature Engineering
- Modelling & Pattern mining
- Evaluation & Deployment

→ **Tokenization** : Breaking into Parts



CODE

Sentence tokenization

Sentences = nltk. sent_tokenize (Paragraph)

Sentences

Out ["he is a good boy."
"she is a good girl."
"boy & girl are good."]

Word tokenization

words = nltk. word_tokenize (paragraph)

words

Out ["he", "she", "boy",
"is", "is", "is",
"a", "a", "a",
"good", "good", "good",
"boy", "girl", "girl",
".", ".", "."]

Text Cleaning

- Remove Punctuation [. , ! " " &]
- Remove Stop words [Is, the, are, so, and]
- Stemming } Base words
- Lemmatization * history
 historical
 histo
- * Beautiful
 Beautifully
 Beauty

⇒ Vectorization [Words → Numeric]
{ Converting words to vectors }

* Text cleaning }
+
Vectorization

Text preProcessing

Remove Punctuation [. , ! " " &]

• Regular Expression [Regex] # Website : regex101.com

import re

Extraction of words, sentences, numbers, punctuations

* \d → matches a digit [0-9]

* \d{3} → matches Exactly Three times { 777 000 999 }

* \d{3}

* (\d{3}) - \d{3} - \d{4} | \d{10}

* We have phone number 8886669990. Using regex how we extract data.

Pattern : \d \d \d \d \d \d \d \d \d
8 8 8 6 6 6 9 9 9 0

(or)

\d{10}

8886669990

* Suppose we have USA number (888)-666-9990

Pattern : \(\d{3}\) - \d{3} - \d{4}
(888) - 666 - 9990

∴ () are special character. So, we use \ before () & End

code

import re

text = "Jack is very Faraway, we can contact him by 7989900099. And he had another USA number (799)-899-0877."
pattern = "\d{10} | \(\d{3}\) - \d{3} - \d{4}"

re.findall(Pattern, text)

Out: ['7989900099', "(799)-899-0877"]

If we want to Extract headline of a paragraph

text = "Note 1 - Overview"

The project is all about

Note 2 - Memory

memory is all about

capturing within

Pattern = "Note Id - [^\n]+)" → including one (or) more
Starting word digit (newline)
(Starting of string)

re.findall(pattern, text)

[out]: ["Overview", "Memory"]

∴ If we don't mention () in pattern.

pattern = "Note Id - [^\n]+"

[out]: ["Note 1 - Overview", "Note 2 - Memory"]

If we want Exact any Series of code in text

FY2022 Q1

Pattern : FY Id{4} Q Id → digits
Starting letters Up to 4 digits * Q' means in question

But, any Non-linear Context like Q5, Q6... From text we don't those Q5, Q6, to Extract

Text = "The year consist FY2021 Q1 was \$21 billions.
FY2022 Q5 \$ 2 billion, FY2023 Q2 \$

Pattern = " FY $\backslash d\{4\}$ Q[1-4]"
6billions" \rightarrow digit
 \rightarrow up to 1 to 4

re.findall ["pattern", text]

[Out] [FY2021 Q1, FY2023 Q2] \therefore Because, we have taken range [1-4], so, it didn't get FY2022 Q5

If we want Ignore Capital and Small letters.

Ex = fy 2021 Q3, FY 2022 Q5, fy 2026 Q6
Pattern = FY $\backslash d\{4\}$ Q[1-4] \rightarrow it ignore case in the text.

re.findall [Pattern, Ex, flags = re.IGNORECASE]

[Out] : fy 2021 Q3, FY2022 Q5, fy2026 Q6

IF we want Only year and Q values.

Ex = fy 2021 Q3, Fy 2022 Q5

Pattern = Fy ($\backslash d\{4\}$ Q\1d)

Here we use Brackets to Exact information only inside Brackets.

re.findall [Pattern, Ex, flags = re.IGNORECASE]

[Out] : [2021 Q3, 2022 Q5]

Extract only \$ 4.85 billions, \$ 3 billion from text

Text = "Jack has \$ 4.85 billions of cash in
Dreams, but in hand he has \$ 3 billion"

Pattern = "\\$[0-9].+in hand"

Starting dollar is a special character so, we mention slash before it

range dot also special character so, mention / before it

include all after the dot

re.findall (Pattern, text)

out ["\$ 4.85", "\$ 3"]

pattern = "\\$ ([0-9].+) " Here we mentioned Bracket, so, it Extract with in the Brackets

out : [4.85, 3]

Extract Email's from the text Using regex

chat 1 = regarding flight ticket Contact abcP@xyz.com

chat 2 = Sorry I lost Contacting you my mail address Sad_29@xyz.in

chat 3 = Your said my mail of mine is dog@que.com

Pattern = "[a-zA-Z]*@[a-zA-Z]*\\.[a-zA-Z0-9A-Z]*"

matches = re.findall(Pattern, chat1)

matches

[out]: abcP@xyz.com

abc_29@mad.in

Badboy@gmail.com

[Same for chat2, chat3]

Email Extraction
: Regex

Extracting Order Issue Details from the Database.

Extracting Order delivery # 412823

Issue 1 = Problem with Order delivery # 412823

Issue 2 = MyOrder 283282 charged 80\$ for product

Issue 3 = maintenance charges details Order 200823

pattern = Order [a-zA-Z0-9A-Z] [^d]* (\d*)

matchv = re.findall(pattern, issue1, flags=re.IGNORECASE)

[out]: Order delivery # 412823

Same format for issue 2, issue 3.

Text cleaning

remove Punctuation.

import re → regex

corpus = []

Len(3 list) → range(3) ⇒ [0, 1, 2, 3]
i = 0

for i in range (len(sentences)):

rp = re. sub ("[^a-zA-Z]", " ", sentences[i])
Subtract ↗
corpus.append (rp) ↗ remove space
remove everything
other than
a-z, A-Z

print (corpus)

Out: ('He is a good boy', 'She is a good girl', 'boy & girl are good')
are removed by "RegEx"

Stopwords :- Words which doesn't have any meaning / use in Sentence.

They are just used for framing Sentences.

Ex:- For Example

∴ Example is Enough
but, we are using "For" ↗ Stopwords

Examples of Stopwords

From nltk.corpus import Stopwords

stopwords.words("English")

Out: i you him it's which is each own will
me you're his it'sself who are few same just
my you'll himself they whom was more so don
myself you'd she them this were most than don't
we yourself she's their that why other now
our your selves her their that'll how some too very
ours he herself themselves these all such if
ourselves it what those any no nor only s m could
ourselves he it both can t couldn't

179 elements -----

remove stopwords

corpus = []

for i in range [len(Sentences))]:

 rp = re.sub("[^a-zA-Z]", " ", Sentences[i])

 rp = rp.lower()

 rp = rp.split()

 rp = [word for word in rp if not word
 in set(stopwords.words("English"))]

 rp = " ".join(rp)

corpus.append(rp)

print(corpus)

Out: [good boy, good girl, boy girl good]

Stemming.

→ Base Word : In order to process Data, we are going work on "WORDS".
that's the Reason we bring it on **Base word**.

Lemmatization

Application :- No human interaction. We use stemming
Ex:- Classification
Gmail

* Stemming

: bring To Base word , which has **no meaning**
(or) not a dictionary word

* Lemmatization

: brings to Same **Base word**, But, it has **meaning**

Application :- human interaction, we use Lemmatization.

Ex:- ChatBOT

Alexa

Siri

Google Assistant

Stemming

code

from nltk.stem import PorterStemmer

PS = PorterStemmer()

ps.stem ("history")

Out: "histori"

Lemmatization

```
from nltk.stem import WordNetLemmatizer  
# wnl = WordNetLemmatizer()  
# wnl.lemmatize("history")
```

out: "history"

Backend.
it opens "Dictionary".

only one should be used
while writing a code
(or)
Building a model.

Stemming

- Converting words to their word stem or root.
- Chopping off the end of a word, to leave only the base.
- beautiful and beautifully words to beauti.

Lemmatization

- reducing words into their Lemma (or) dictionary
- Remove inflectional endings and return to dictionary form of a word
- beautiful and beautifully words to beautiful.

Code

- We are working on Sentiment analysis and Text classification. So, we use "Stemming" for this model.
- Lemmatization is used for the chatbot, Alexa, which has interactions.

Text Cleaning [stopwords, stemming, lemmatization]

```
# Corpus = []
```

```
# for i in range(len(sentences)):
```

```
    rp = re.sub("[^a-zA-Z]", " ", sentences[i])
```

```
    rp = lower()
```

```
    rp = split()
```

```
    rp = [ps.stem(word) for word in rp  
          if not word in set(stopwords.
```

```
words("english"))
```

```
    rp = " ".join(rp)
```

```
    corpus.append(rp)
```

```
# Corpus
```

```
Out: ["good boy", "good girl", "boy girl good"]
```

Vectorization

Convert tokens into to a format that a machine

learning understands
"Converting the Text to Numeric Representation"

⇒ Methods of Vectorization

1. Count Vectorization [Bag of words]
2. N-grams Count Vectorization
3. Term frequency - inverse document frequency
[TF-IDF]

* Bag of Words [Bow] / Count vectorization :-

Counting the number of "words" in the Sentence
is known as "Bow" / count vectorizer.

Ex:- Sentence 1 = good boy

Sentence 2 = good girl

Sentence 3 = boy girl good

	boy	girl	good
Sentence 1	1	0	1
Sent 2	0	1	1
Sent 3	1	1	1

This is how Bow [Bag of words] works. and also

Count vectorizer works.

* it just counts the words in the sentence
and separate them.

Ex:2

Sentence 1 = Data good Future Good

Sentence 2 = Data best Present

Sentence 3 = Data poor Past

	best	Data	future	good	Present	Poor	Past
Sentence 1	0	2	1	1	0	0	0
Sentence 2	1	1	0	0	1	0	0
Sentence 3	0	1	0	0	0	1	1

code # bow [bag of words]

from sklearn.feature_extraction.text import
CountVectorizer

CV = CountVectorizer(binary=True, ngram_range=(3,3))

bow = CV.fit_transform(corpus).toarray()

bow

[out]: array([[1, 0, 1],
[0, 1, 1],
[1, 1, 1]]).

	Boy	girl	good
Sent 1	1	0	1
Sent 2	0	1	1
Sent 3	1	1	1

* Problem with [Bag of words] : it is going to give Only count/frequency of all words.

Same weightage

Count/Frequency

we can't identify which word is important in the sentence

Advantage

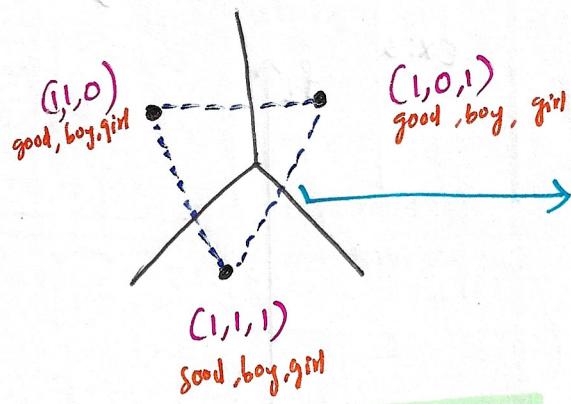
* Simple & Intuitive

DisAdvantage

- * Sparsity * Semantic meaning lost
- * OOV [Out of vocabulary]
- * Ordering of words are completely changed.

If we want to plot in 3 [three] dimension

Ex:-



Distance Between them
can be calculated by.

- * Cosine Similarity
- * Euclidean distance.

Tells How similar sent 1, sent 2
sent 3 are similar to each other.

* So, what is Cosine Similarity

& Cosine Distance

?
∴ it is Widely Used

in "Recommendation Systems".

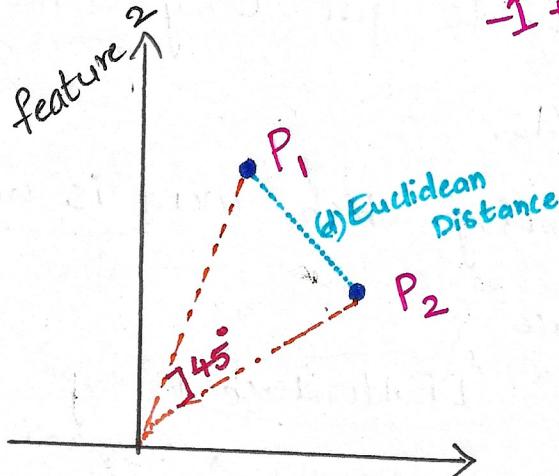
We have Two points : P_1 & P_2

Similarity = $\frac{\text{Distance}}{\text{Distance increases} \uparrow}$
 Decreases \downarrow \Rightarrow Distance increases \uparrow
 Increases \uparrow \Rightarrow Distance Decreases \downarrow

Formula.
 $1 - \text{Cosine-Similarity} = \text{Cosine Distance}$

∴ According to
Geometric Term.

Ex:-



it ranges b/w
-1 to +1

$$\therefore \cos - \sin = \cos \theta$$

Angle b/w P_1 & P_2

$$\therefore \cos 45^\circ = 0.707$$

$\Rightarrow 0.70 \Rightarrow 70\%$. Similar.

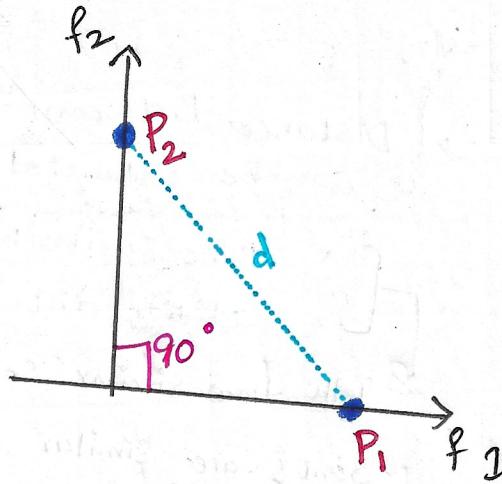
Based on Distance &
angle b/w them

But, In Cos-sine Similarity to find the Similarity

We, have to find "angle" between the Data

points.

Ex:2



$$\therefore \cos 90^\circ = 0$$

These Data points are not similar

To Each Other.

Ex:3

Distance b/w
two records

$$\text{Ex: } A = [5, 1, 3, 2] \\ B = [0, 1, 2, 3]$$

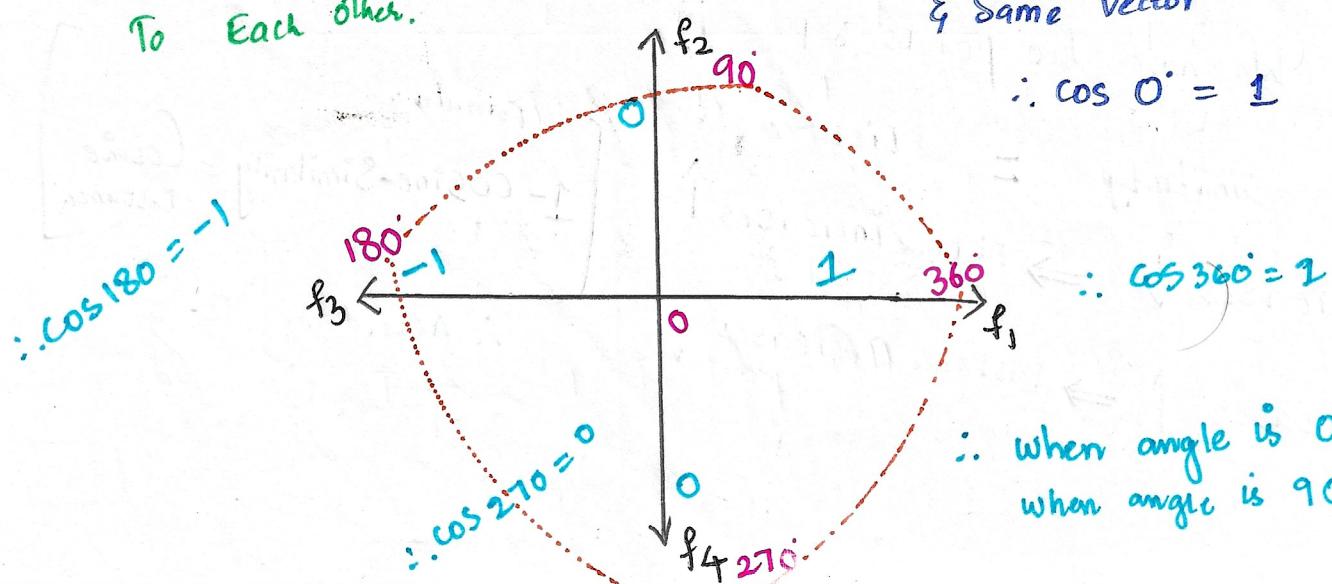
$$= \frac{(5 \times 0) + (1 \times 1) + (3 \times 2) + (2 \times 3)}{\sqrt{5^2 + 1^2 + 3^2 + 2^2} \times \sqrt{0^2 + 1^2 + 2^2 + 3^2}}$$

$$= \frac{5 + 1 + 6 + 6}{\sqrt{5^2 + 1^2 + 3^2 + 2^2} \times \sqrt{0^2 + 1^2 + 2^2 + 3^2}} \\ = 0.7703$$

P_1 & P_2 are on the same plane

& Same vector

$$\therefore \cos 0^\circ = 1$$



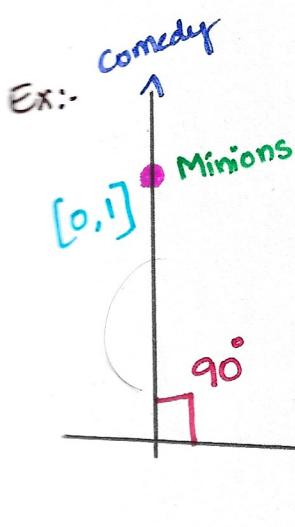
∴ So, we want to calculate ∴ Cosine Distance

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

$$\rightarrow \cos 45^\circ = 1 - \frac{1}{\sqrt{2}} = 0.3$$

$$\rightarrow \cos 90^\circ = 1 - 0 = 1$$

$$\rightarrow \cos 0^\circ = 1 - 1 = 0$$



$$\therefore \cos 90^\circ = 0$$

$$\begin{aligned} \text{Distance} &= 1 - \cos 90^\circ \\ &= 1 - 0 \end{aligned}$$

$$\text{Distance} = 1$$



$$\therefore \cos 0^\circ = 1$$

$$\begin{aligned} \text{Distance} &= 1 - \cos 0^\circ \\ &= 1 - 1 \\ &= 0 \end{aligned}$$

$$\therefore \text{Distance} = 0$$

* In Bag of words [Bow]/Count vectorize. We can't find

The "semantic meaning" (or) ? information.

↳ Make words
Meaningful.

So, we Use Method Called as

"N-grams".

→ Combination of "n" words/features.

N-grams

- * Unigrams = it uses only single word (or) single features
- * Bigrams = Combination of two words/features
- * trigrams = Combination of three words/features

Ex:- Bi-grams.

			Feature F ₁	F ₂	F ₃
D ₁	Good	boy	1	1	0
D ₂	Good	girl	1	0	1
D ₃	Boy	girl	1	1	1

Sentence
1
Sent 2
Sent 3

F₁ F₂ F₃ F₄ F₅
Good boy Girl Good boy Good girl

	Sent 1	1	1	0	1	0
Sent 2	1	0	1	0	1	
Sent 3	1	1	1	0	0	

Bi-Grams

Ex:- Jack Eats Food = Bi-grams

∴ Jack Eats
Eats Food

Trigrams

I am not Feeling well = Trigrams

∴ I am not
am not Feeling
not Feeling well

Ex:- The food is good] → opposite
The food is not good

Food good is not The

1	1	1	0	1
1	1	1	1	1

But, by the Bag of words, it gives Every individual option
has Equal importance.

So, we use method called "TD-IDF"

TF & IDF

TF-IDF: [Term Frequency & Inverse Document Frequency]

∴ Which are words are rarely present in the Sentences
We should try to give more weightage

EX:- { Cat Eat Food } These common words captured
Dog Eat Food by the
PIG Eat Food IDF: inverse Document
frequency.

These are rare / Unique.
So, we calculate

TF : Term Frequency

when we combine both TF & IDF
we can get the most rare words in The
Sentences

Formula :

$$\therefore \text{Term Frequency} = \frac{\text{No. of repetition of words in Sentence}}{\text{No. of words in Sentence}}$$

*

$$\text{Inverse Document Frequency} = \log_e \left[\frac{\text{No. of Sentences}}{\text{No. of Sentences Containing Word}} \right]$$



TF-IDF

\Rightarrow Term Frequency :- For Every Sentences

	Sentence 1	Sent2	Sent 3
Unique words	$\frac{\text{No. of repetition of words in Sentences}}{\text{No. of words in Sentences}}$ $\text{Formula} = \frac{\text{No. of words in Sentences}}{\text{No. of words in Sentences}}$		
good	$\frac{1}{2} \left[\frac{1-\text{good}}{2-\text{good boy}} \right]$	$\frac{1}{2} \left[\frac{1-\text{good}}{2-\text{good girl}} \right]$	$\frac{1}{3} \left[\frac{1-\text{good}}{3-\text{good girl boy}} \right]$
boy	$\frac{1}{2} \left[\frac{1-\text{boy}}{2-\text{good boy}} \right]$	$\frac{0}{2} = 0$	$\frac{1}{3}$
girl	$\frac{0}{2} = 0$	$\frac{1}{2} \left[\frac{1-\text{good}}{2-\text{good girl}} \right]$	$\frac{1}{3}$

Sent 1 : good boy
 Sent 2 : good girl
 Sent 3 : boy girl good

3 sentences

\Rightarrow Inverse Document Frequency : For Every words

Words	IDF
Good	$\log_e \left[\frac{3}{3} \right] = 0$
boy	$\log_e [3/2]$
girl	$\log_e [3/2]$

$\log_e \left[\frac{\text{No. of Sentences}}{\text{No. of Sentences Containing word}} \right]$

TF : Term Frequency \times IDF : Inverse Document Frequency

	Sent1	Sent2	Sent3		Words	IDF
good	$\frac{1}{2}$	$\frac{1}{2}$	$\frac{1}{3}$	*	good	$\log_e [3/3] = 0$
boy	$\frac{1}{2}$	0	$\frac{1}{3}$		boy	$\log_e [3/2]$
girl	0	$\frac{1}{2}$	$\frac{1}{3}$		girl	$\log_e [3/2]$

TF-IDF [Calculation]

	Feature 1	Feature 2	Feature 3
	good	boy	girl
Sent 1	$\frac{1}{2} \times 0 = 0$	$\frac{1}{2} \times \log_e [3/2] = 0.2027$	$0 \times \log_e [3/2] = 0$
Sent 2	$\frac{1}{2} \times 0 = 0$	$0 \times \log_e [3/2] = 0$	$\frac{1}{2} \times \log_e [3/2] = 0.2027$
Sent 3	$\frac{1}{3} \times 0 = 0$	$\frac{1}{3} \times \log_e [3/2] = 0.1351$	$\frac{1}{3} \times \log_e [3/2] = 0.1351$

importances of word Captured Using TF-IDF

"ngram" "code"

1. ngram-range : tuple (\min_n, \max_n)
2. max-features : If not None, build a vocabulary that Consider, top max-features across corpus Ordered by term frequency

Ex [1,1] means "unigram"

[1,2] \Rightarrow unigram & Bigram
[2,2] \Rightarrow

```
from sklearn.feature_extraction.text import TfidfVectorizer
```

```
# cv = TfidfVectorizer(ngram_range=[3,3],  
#                      max_features=10)  
  
# x = cv.fit_transform(corpus)
```

Another Example

```
# corpus[0]
```

Out: "narendra damodardas modi Gujarathi n End

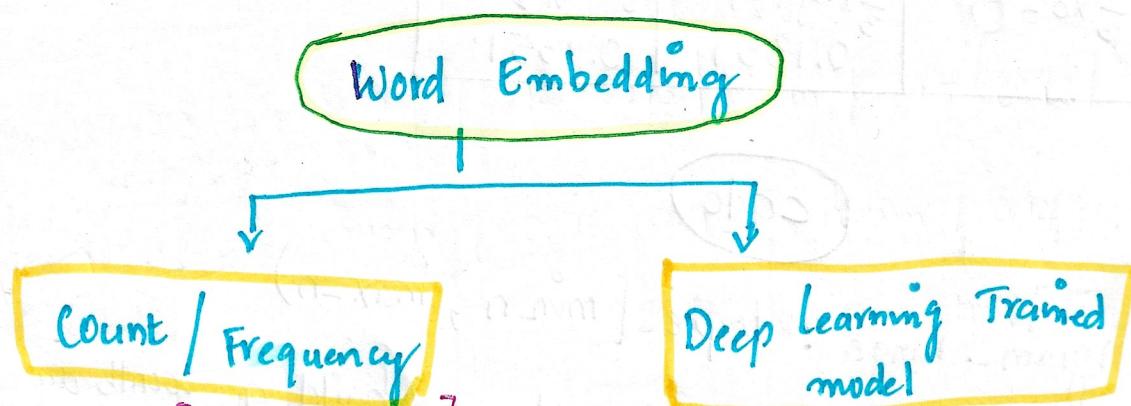
```
# x[0].toarray()
```

Out: array([[0., 0., 0., 0., 0.70710678,
0., 0.70710678, 0., 0.]])

* Word Embedding

Word Embedding

it is a Techniques which converts words into vectors



1. Bow [Bag of words]

2. TF-IDF

3. OHE [One hot Encoding]

1. word 2 vec

* CBOW

* skip grams

2. Avg word 2 vec

Word2Vec

Feature representation

f_1	f_2	f_3	f_4	f_5	f_6	...
Boy	Girl	King	Queen	apple	Mango	

- * Problem with :- **Bow**, **TFIDF**
 - [Semantic Meaning]
 - [Missing]
 - [Synonyms]
 - Similar words [Ex :- good honest]
- * They are creating **Sparse matrix**

$[0, 1]$ (or) $[0.21, 0]$
BOW TDIDF

- * Because, These create "Huge dimensions" and They have more vectors. it leads to "overfitting".

Word2Vec :

- * Limited Dimensions are created
- * Sparsity is Reduced
- * Semantic remained Constant

∴ it is divided into two types

- 1. CBOW [Continuous Bag of words]
- 2. SkipGram

Word2Vec

Feature Representations

Created by model

$f_1, f_2, f_3, f_4, f_5, f_6, \dots$

	Boy	Girl	King	Queen	apple	Mango.
--	-----	------	------	-------	-------	--------

relation b/w

Gender

-1 1 0.92 +0.93 0 -0.02 0.01

0.01 0.02 0.95 0.96 -0.02 0.01

age 0.03 0.02 0.7 0.6 0.95 0.96

food 0.01 0.02 0.95 0.96 -0.02 0.01

;

;

;

300 dimensions vectors

[features]

These dimensions are

Created by "word2vec"

algorithm.

"Google trained with 3 billion words"

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

;

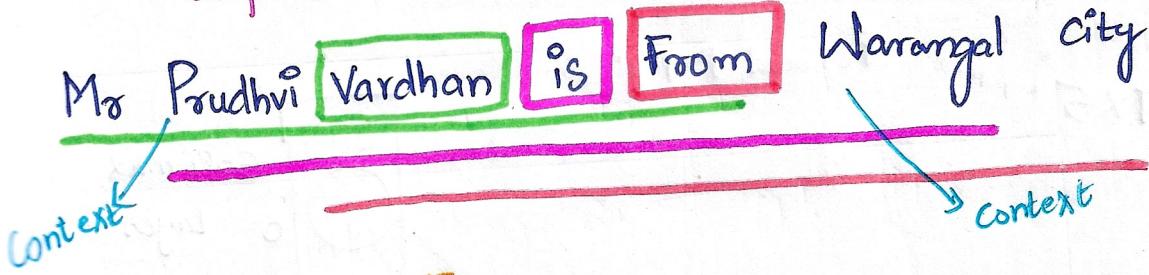
;

;

;</

* CBOW : {continuous Bag of words}

Corpus :-



* Window Size = 5

↳ two words before & two words after the input word.
in addition to input word itself.

Training Data

Independent Feature

Mr, Prudhvi, is, From

Prudhvi, Vardhan, From, Warangal

Vardhan, is, warangal city

O/P

Vardhan

is

From

∴ Here Semantic meaning being captured.

By [Text → vectors]

Given To
Fully Connected Layer
"ANN"

[Artificial Neural Network]

∴ Bow [Bag of words]

• Mr - 1 0 0 0 0 0 0

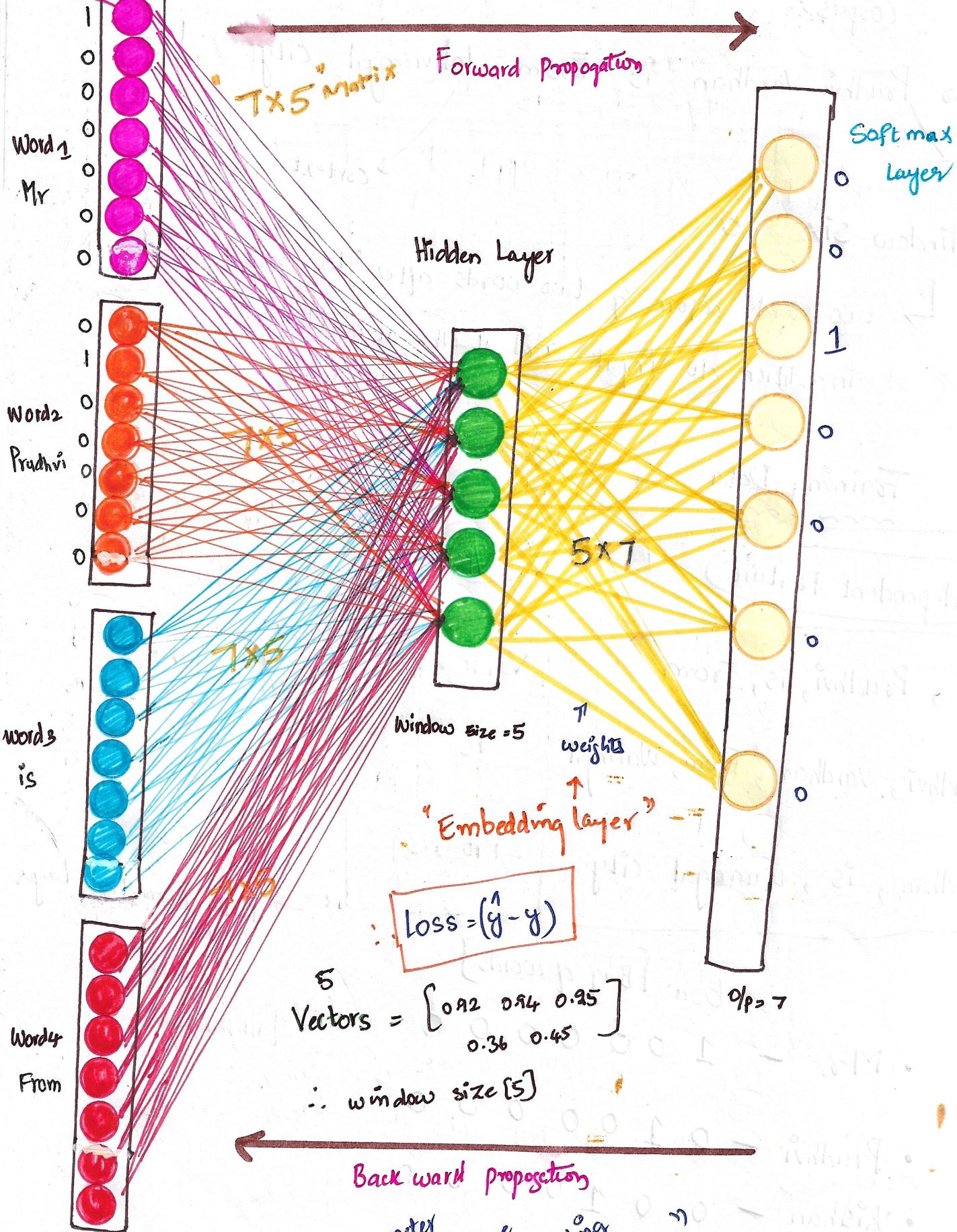
• Prudhvi - 0 1 0 0 0 0 0

• Vardhan - 0 0 1 0 0 0 0

• is - 0 0 0 1 0 0 0

Here input is $s=4$, window size = 5, o/p = 7 → vectors

"CBow" [Genism]



\therefore Hyperparameter Tuning is Done By Extending "window size"

Skip Gram

In The CBOW model, the Distributed representations of Context [or Surrounding words] are combined to Predict the Word in middle.

* But..., In SkipGram model, the distributed representation of input word is used to "Predict the Context".

Ex:- it Exact Opposite of "CBOW"

corpus : Mr Prudhvi Vardham is From Warangal city
Content ↴ Window size = 5

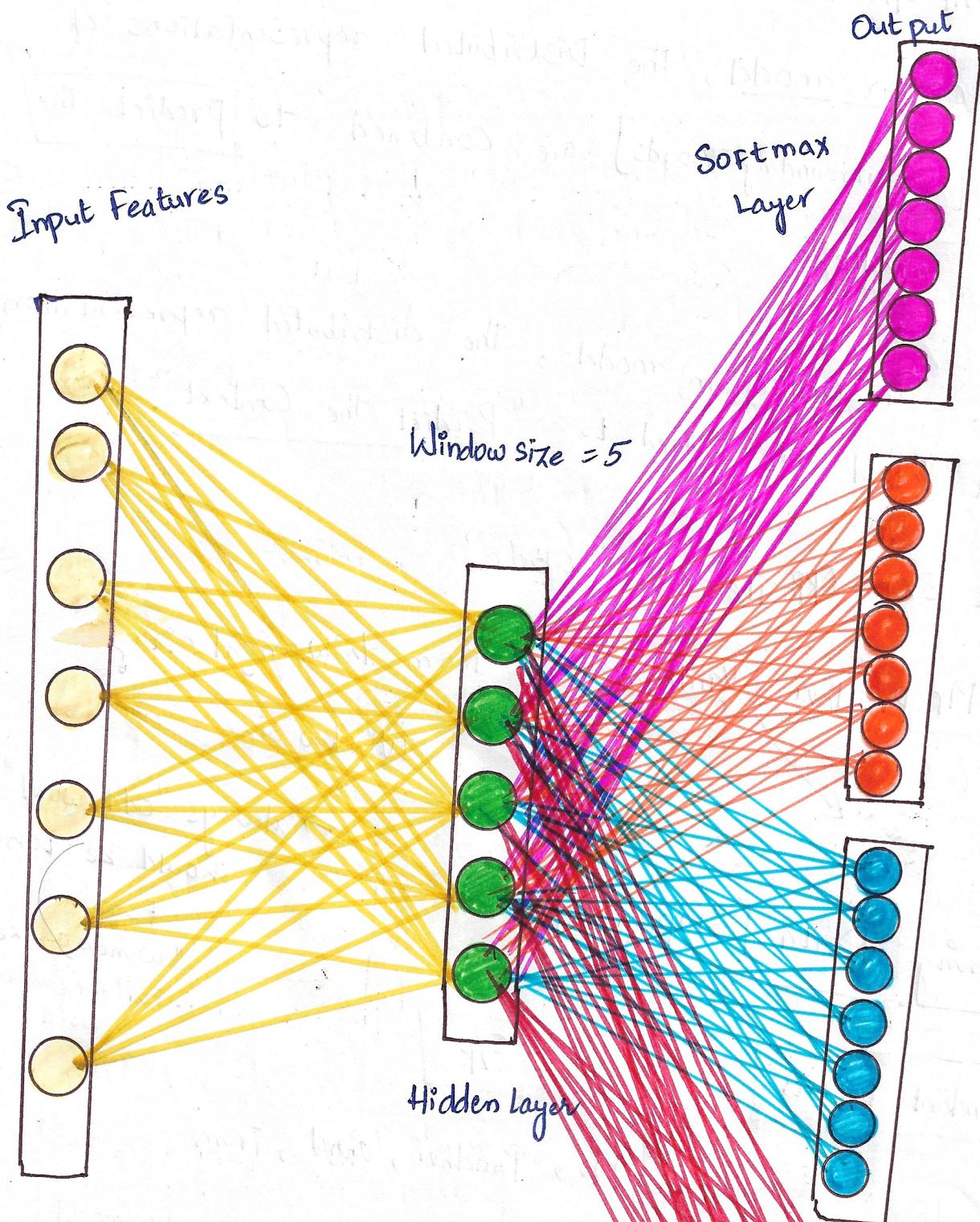
↳ always "odd word" should be taken

Training Data [skip Gram]

Independent Features	O/P
Vardham	Mr, Prudhvi, is , From
is	Prudhvi, Vardham, From , Warangal
From	Vardham, is , Warangal , city

\therefore BIGGER DATA = SKIPGRAM , Short/small Data = CBOW

Here Input features = 7 , Window Size = 5 , O/p = 4



$$\text{Loss} = (y - \hat{y})^2$$

{ If Google is creating
300 dimensions \rightarrow output
Hidden layers has 300 Neurons]
 \therefore window size \rightarrow 4 }

Code

* CBOW / Skipgram

Based on
Transformer
Architecture

Based on
LSTM

- Word2Vec
- Glove
- fastText

→ BERT

→ GPT

→ ELMo

→ Gensim

topic modelling, similarity retrieval
with large corpora.
in NLP, IR [information retrieval]

! pip install gensim

import gensim

from gensim.models import Word2Vec, KeyedVectors

downloader as api

Google
Pre-trained model

has trained
with
billion of words
with 300 dimensions

import gensim.

WV = api.load("word2vec-google-news-300")

WV["King"]

Vec =

Vec

out: array([1.2159
...])

2.97851	8.609	1.396
...

[300 vectors].

Similarity b/w words.

Wv. similarity ("Great", "good")

[out]: 0.729

Wv. similarity ("jesus", "christ")

[out]: 0.637

Semantic meaning can be done :- Option by
"most-similar".

Wv. most-similar ("good")

[out]: [("great": 0.729
"bad": 0.719
"terrific": 0.683
"decent": 0.688
"nice": 0.683
"Excellent": 0.644)]

Wv. most-similar ("jesus christ")

[out]: [("jesus christ": 0.639
"christ": 0.637
"cuzz-u": 0.634
"jesse": 0.632
"lmfao": 0.621)]

Vec = wv["King"] - wv["man"] + wv["woman"]

Vec

[out]: 300 vectors ..
0.236 0.834 0.638 0.843 ..

wv.most_similar([vec])

[out]: [(“King”: 0.647

“queen”: 0.535

“woman”:

King

wv.most_similar(positive = [“King”, “woman”],
negative = [“man”], topn=5) ↳ top 5 numbers

[out]: [(“queen”: 0.7118
“monarch”: 0.6182
“princess”: 0.590
“crown-Prince”: 0.549
“Prince”: 0.537)]

Project :- Spam Classification

```
import Pandas as pd
```

```
# df = pd.read_csv("Spam_Collection", sep="\t", names=[  
# df.head(5)  
    "Label", "messages"])
```

out:

Label	messages
ham	Go until point , crazy, Available Only
ham	OK jar, joking
spam	Free Entry, Coupons Available
ham	You don't know brother
ham	Want Prizes Will Set Friend

Text cleaning :

```
import re
```

```
import nltk
```

```
From nltk.stem import PorterStemmer
```

```
From nltk.corpus import stopwords
```

```
# ps = PorterStemmer()
```

```
# Corpus = []
```

```
for i in range(len(df)):
```

```
    tp = re.sub("[^a-zA-Z0-9]", " ",
```

```
df["messages"][i])
```

→ Separate
By "TAB"

```
rp = rp.lower()
```

```
rp = rp.split()
```

```
rp = [ps.stem(word) for word in rp if not  
      word in set(stopwords.words("english"))]
```

```
rp = " ".join(rp)
```

```
corpus.append(rp)
```

Corpus

```
out: [go journey point - - - -  
-----]  
-----]
```

Vectorization

```
from sklearn.feature_extraction.text import CountVectorizer
```

```
# cv = CountVectorizer()
```

```
# X = cv.fit(corpus).toarray()
```

```
# X.shape()
```

```
out: (5572, 7248)
```

```
# y = pd.get_dummies(df["label"], drop_first = "True")
```

```
# y.shape()
```

```
out: (5572, 2)
```

Here, "ham" column

got deleted,
So, we have only
"spam" column

Train Test Split

```
from sklearn.model_selection import train_test_split  
# x-train, x-test, y-train, y-test = train_test_split (x, y,  
test_size = 0.3, random_state = 29)
```

x-train. shape

Out : $(3900, 1248)$

n-train

out: array ([[0, 0, 0 ..., 0, 0, 0],
[0, 0, 0 ..., 0, 0, 0],
[0, 0, 0 ..., 0, 0, 0].

α -test

Out: array([[0, 0, 0 ..., 0, 0, 0],
[0, 0, 0 ..., 0, 0, 0],
[0, 0, 0, ..., 0, 0, 0],

y-train. Shape

Out: (3900, 1)

$y\text{-train}$

$y\text{-test}$

[out]:

spam

4480

0

3887

0

1653

0

2607

1

2731

1

3900 rows x 1 column

[out]:

spam

534

0

2279

1

:

4159

0

1504

1

2946

1

1672 rows x 1 column

Modelling

from sklearn.naive_bayes import MultinomialNB

model = MultinomialNB()

model.fit(x-train, y-train)

[out]: MultinomialNB()

Prediction

yPred-train = model.predict(x-train)

ypred-test = model.predict(x-test)

Evaluation

```
from sklearn.metrics import accuracy_score
```

```
# print("Accuracy:", accuracy_score(y-train, ypred-train))
```

```
# print("test Accuracy:", accuracy_score(y-test, ypred-test))
```

[out]: Train Accuracy : 0.9928

Test Accuracy : 0.9856

✓