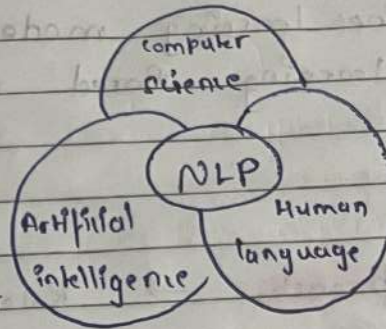


Natural language processing.

• Real world Application.

- Contextual Advertisement
- Email client - spam filtering
- Social media - removing adult content



↳ Common NLP tasks

- (1) Text / Document classification
- (2) Sentiment Analysis
- (3) Information Retrieval
- (4) Part of speech tracking.
- (5) Language Detection and Machine Translation
- (6) Conversational Agents.
- (7) Knowledge graph & QA System.
- (8) Text Summarization
- (9) Topic modelling
- (10) Text generation.
- (11) spell checking and Grammar correction
- (12) Text Parsing
- (13) Speech to Text.

Rules are made
By Human itself

PAGE NO.:

DATE: / /

• Approaches to NLP:

- Heuristic Approaches.
- Machine learning models
- Deep learning Based models.

• ML approach.

→ Rules are made by model itself
after seeing the data.

→ 171 flow

Algorithms used.

- Naive Bayes
- Logistic regression
- SVM
- LDA
- Hidden Markov model.

• Deep learning model approach. :-

Advantage → Retain sequential information.

Architecture used.

- RNN
- LSTM
- GRU / CNN

(Gated Recurrent Unit) → Text generation.

- Transformers
- Auto encoders.

Challenges in NLP

- It is difficult to make machine understand all these

(1) Ambiguity

- I saw the boy on the beach with my binoculars
- I have never tasted a cake quite like that one before

(2) Contextual words.

- I ran to the store because we ran out of milk.

(3) Colloquialisms and slang.

- Piece of cake, - pulling your leg.

(4) Synonyms.

(5) Spelling error

(6) Creativity

(7) Irony, sarcasm

• NLP pipeline.

(1) Data Acquisition

(2) Text Preparation

Text cleanup.

Basic Preprocessing

Advanced Preprocessing

(3) Feature Engineering.

(4) Modelling.

model Building

model evaluation

(5) Deployment.

Deployment

monitoring

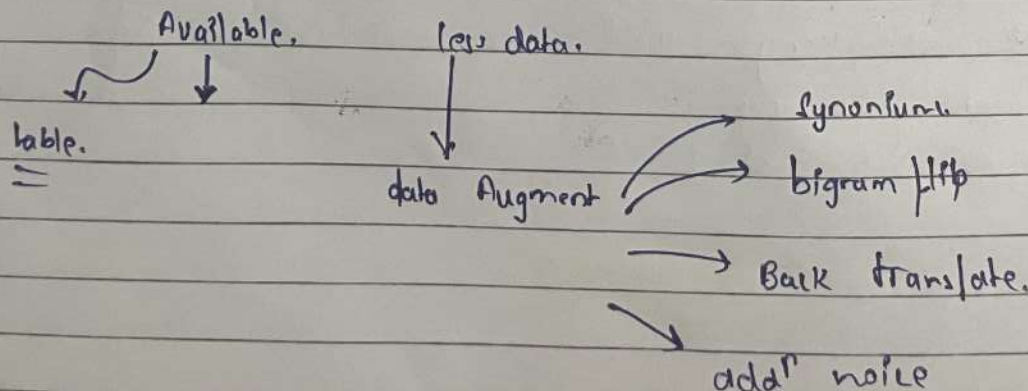
model updates.

Points.

→ this pipeline is not universal

→ Deep learning pipelines are more difficult

* Data Acquisition :- Data,



Other.

→ public dataset

→ web scraping. → beautiful soup

→ API → rapid API → request → json.

* Text preparation

Cleaning

- + html tags cleaning
- + emoji (unicode normalization)
- + spell checking

Basic preprocessing

Basic

optional

- + tokenization
- sentence
- coord.

- + stop word removal
- + stemming
- + Removing digits
- + lower casing
- + language detection

* Advance preprocessing

- + POS tagging → (Parts of speech tagging)
- + Parsing. → understand syntax structure. (eg. chatbot)
- + Coreference Resolution eg. (Lapin directed ^{both some person} of his ^{composed music of most} films.)

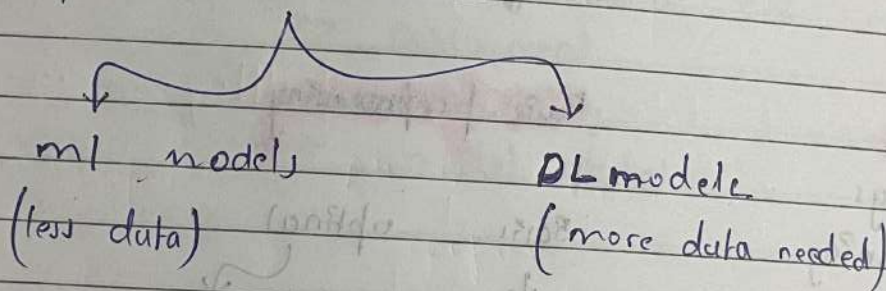
* Feature Engineering :

→ Task is to convert words into numbers.
 Various techniques can be.

- ① Bag of words
- ② Tf-idf
- ③ OHE (one Hot encoding)
- ④ word 2 vec

} depends on Application

* modelling.



* evaluation

→ Intrinsic

extrinsic

(Rule Based)

→ perplexity.

* Deployment

Example :-

PAGE NO. :

DATE : / /

→ Input

chaplin wrote , directed and composed the music for most of his film

→ Tokenization with lemmatization

(stemming)

chaplin wrote , directed and composed the music for most of his film

→ POS Tagging.

NNP VBD VBD CC VBN DT NN IN PRP IN PRP S

chaplin wrote , directed and composed the music for most of his film

→ Coreference Resolution

mention

coref

mention

chaplin wrote , directed and composed the music for most of his

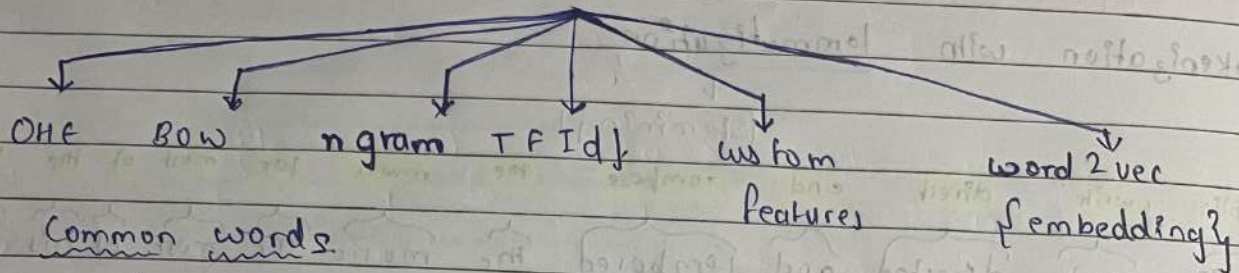
Text Representation :-

PAGE NO.:

DATE: / /

{ Hello , how are you ? }

⇒ Topic is to convert "text meaning into numbers."



* Common words

→ (1) Corpus → Combination of all the words in the Dataset. { words can repeat }

→ (2) Vocabulary → All the unique words of corpus is vocabulary.

→ (3) Document → Each "review" in Dataset is Document.

→ (4) word. → any word in Dataset

(1) One hot Encoding :- Not used These days

PAGE NO.:

DATE: / /

D_1 - people watch campus x

(sparse array)

D_2 - campus x watch campus x

(no fixed size)

D_3 - people write comment

(no capturing of semantic meaning)

D_4 - campus x write comment

Corpus

people watch campus x campus x watch
campus x ~~watch~~ ~~campus x~~ people
write comment campus x write comment

Vocabulary

(V=5)

people watch campus x write comment

now.

every sentence or document D_1, D_2, D_3, D_4 will be converted "V" dimensional array.

	people	watch	campus x	write	comment
people →	1	0	0	0	0
watch →	0	1	0	0	0
campus x →	0	0	1	0	0
write →	0	0	0	1	0
comment →	0	0	0	0	1

$$D_1 = \begin{bmatrix} [1, 0, 0, 0, 0], \\ [0, 1, 0, 0, 0], \\ [0, 0, 1, 0, 0] \end{bmatrix}$$

$$D_2 = \begin{bmatrix} [0, 0, 1, 0, 0], \\ [0, 1, 0, 0, 0], \\ [0, 0, 1, 0, 0] \end{bmatrix}$$

Bag of words :

text

output

$D_1 =$ people watch campus 1
 $D_2 =$ campus watch campus 1
 $D_3 =$ people write comment 0
 $D_4 =$ campus write comment 0

 $V = 5$

	People	watch	campus	write	comment
D_1	1	1	1	0	0
D_2	0	1	1	0	0
D_3	1	0	0	1	1
D_4	0	0	1	1	1

How many time the word is coming in document

Better than "ONE"

Advantages.

Disadvantages.

(1) Sparsity

(2) OOV

(3) Do not consider ordering.

N-Grams.

D_1	people	watch	campus	1
D_1	campus	watch	campus	1
D_2	people	write	comment	1
D_4	campus	write	comment	0
				0

Bi-Gram \rightarrow

Vocabulary of two words.

Tri-Gram \rightarrow

Vocabulary of 3 words.

Vocabulary of Bi-gram \Rightarrow

① people watch, watch campus
 ② campus watch, people write
 ③ write comment, campus write

	①	②	③	④	⑤	⑥
D_1	0	1	0	0	0	0
D_2	0	1	1	0	0	0
D_3	0	0	0	1	1	0
D_4	0	0	0	0	1	1

Bi-Gram is better than Bag of words.

eg.

 D_1 This movie is very good D_2 This movie is not good.Vocabulary \Rightarrow This movie is very good not

① BOW	$\rightarrow D_1$	1	1	1	1	0
	$\rightarrow D_2$	1	1	1	0	1

These two vectors are not far apart.
in 6D space.

60 space.

B5 gram.

vocabulary \Rightarrow This movie / movie is / is very / very good / is not / not good D_1 D_2

vectors are far apart than Bow
 \rightarrow able to capture more semantics of the sentence.

 \rightarrow Dis advantage \rightarrow slows down the algo \rightarrow dimension increases \rightarrow out of vocab**Tf-Idf** :-

Concept is when a word is rare in the whole corpus but is coming more times in a document then it will allot more importance to that word to that document

Tf \rightarrow Term frequencyIdf \rightarrow Inverse dot frequency

$$TF (\text{term frequency}) = \frac{(\text{Number of occurrence of term "t" in document "d"})}{(\text{Total number of term in document d})}$$

$$IDF (\text{inverse doc frequency}) (t) = \log_e \left(\frac{\text{Total number of document in the corpus}}{\text{number of documents with term t in them}} \right)$$

D_1	people watch compux	1		
D_2	compux watch compux	1		
D_3	people write comment	0	people.	$\log(4/2)$
D_4	compux write comment	0	watch	$\log(4/2)$
			compux	$\log(4/3)$
			write	$\log(4/2)$
			comment	$\log(4/2)$

if the "term" is coming in every ~~com~~ Document then
(IDF = 0)

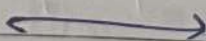
so we add 1 in it

so that they should not be ignored entirely.

$$TF (\text{people}, D_1) = 1/3$$

$$TF (\text{watch}, D_1) = 1/3$$

$$\text{Value of vector of every word in Document} = (TF \times IDF)$$



	people	watch	laptop	write	comment
D_1	$\frac{1}{3} \times 0.3$	$\frac{1}{3} \times 0.3$	$\frac{1}{3} \times 0.125$	0	0
P_2	0	$\frac{1}{3} \times 0.3$	$\frac{2}{3} \times 0.125$	0	0
D_3	$\frac{1}{3} \times 1.3$	0	0	$\frac{1}{3} \times 0.3$	$\frac{1}{3} \times 0.3$
D_4	0	0	$\frac{1}{3} \times 0.125$	$\frac{1}{3} \times 0.3$	$\frac{1}{3} \times 0.3$

* Advantages

→ good in information retrieval

* Disadvantages

- sparsity
- too simple meaning (cannot take out)
- too low dimension

* Word Embeddings :-

→ words are represented in the form of vectors that encodes the meaning of the word such that the words that are close in the vector space are expected to be similar.

* Word & Vec :- (Deep learning Based embedding Technique)

→ (can capture symmetric meaning very properly.)

→ eg. can capture the closeness of the word "joy" and "happy" very properly accurately.

→ low dimensional → (300)

→ makes computation easy.

→ Dense vector

→ very less non-zero values, which reduces overfitting.

→ Demo of word & Vec

• Intuition :- word & Vec

→ It creates features to them and apply probabilistic values.

features ↓	King	Queen	Man	woman	monkey
gender	1	0	1	0	1
wealth	1	1	0.3	0.3	0
power	1	0.7	0.2	0.2	0
Speak.	1	1	1	1	0
weight	0.8	0.4	0.6	0.5	0.3

features of King = $[1 \ 1 \ 1 \ 1 \ 0.8]$

lets do

(king - man + women

$$1 - 1 + 0 = 0$$

$$1 - 0.3 + 0.3 = 1$$

$$1 - 0.2 + 0.2 = 1$$

$$0.8 - 0.1 + 0.5 = 0.7$$

$$1 - 1 + 1 = 1$$

① → Note → The model donot create feature that the example shows above

→ It create features like f_1, f_2, f_3, \dots

② → Who word2vec decide, what the feature will be??

⇒ The underlying assumption of word2vec is that two words sharing similar context also share a similar meaning and consequently a similar vector representation from the model.

word	man	woman	king	queen
man	1	0	0	1
woman	0	1	1	0
king	0	0	1	1
queen	0	0	1	1

$$[0.8 \ 1 \ 1 \ 1 \ 1] = \text{word2vec of king}$$

word2vec
(CBOW) (skip gram)

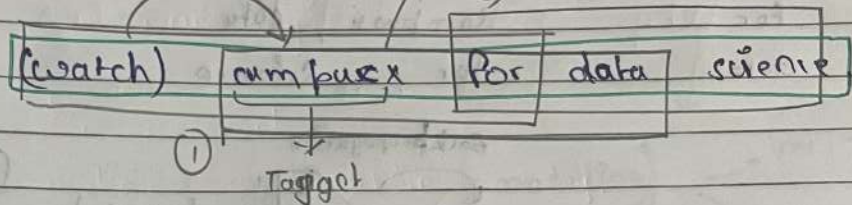
PAGE NO.:

DATE: / /

(1) CBOW :- (Continuous Bag of words)

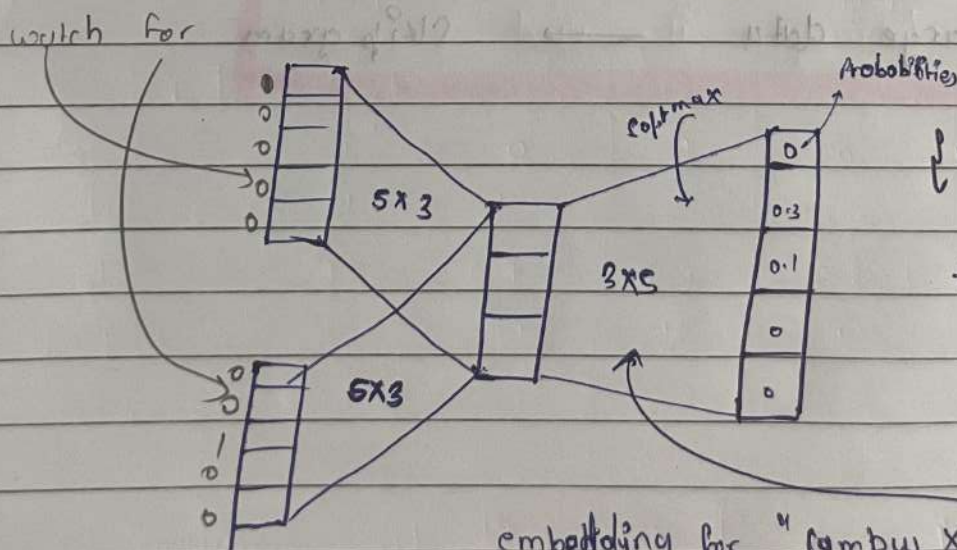
→ It creates a dummy problem to create context of the words.

eg. :-



X	Y	vector for each word.
① watch, for	campus X	watch → [1 0 0 0 0]
campus, data	for	campus X → [0 1 0 0 0]
for, science	data	for → [0 0 1 0 0]
		data → [0 0 0 1 0]
		science → [0 0 0 0 1]

now creating a neural network from scratch.



Values will be compared with output vector then Backpropagation will happen and weights and biases will be rearranged.

embedding for "campus X" will be weights of this layer.

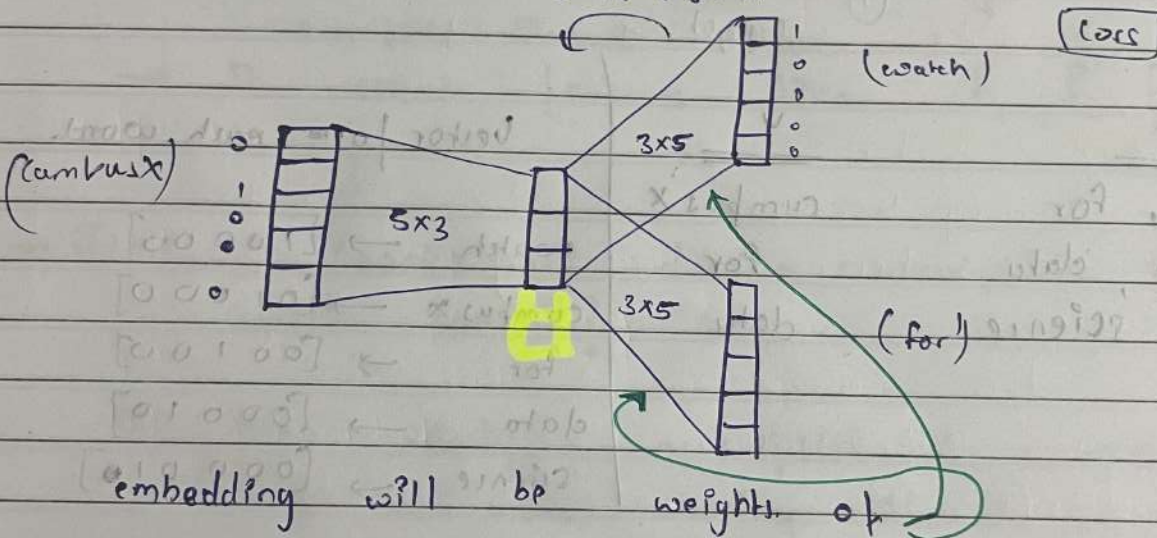
* Skip - Gram :-

window = 3

watch complex for data science

x y
 complex watch for
 for complex data
 data for science

Backpropagation

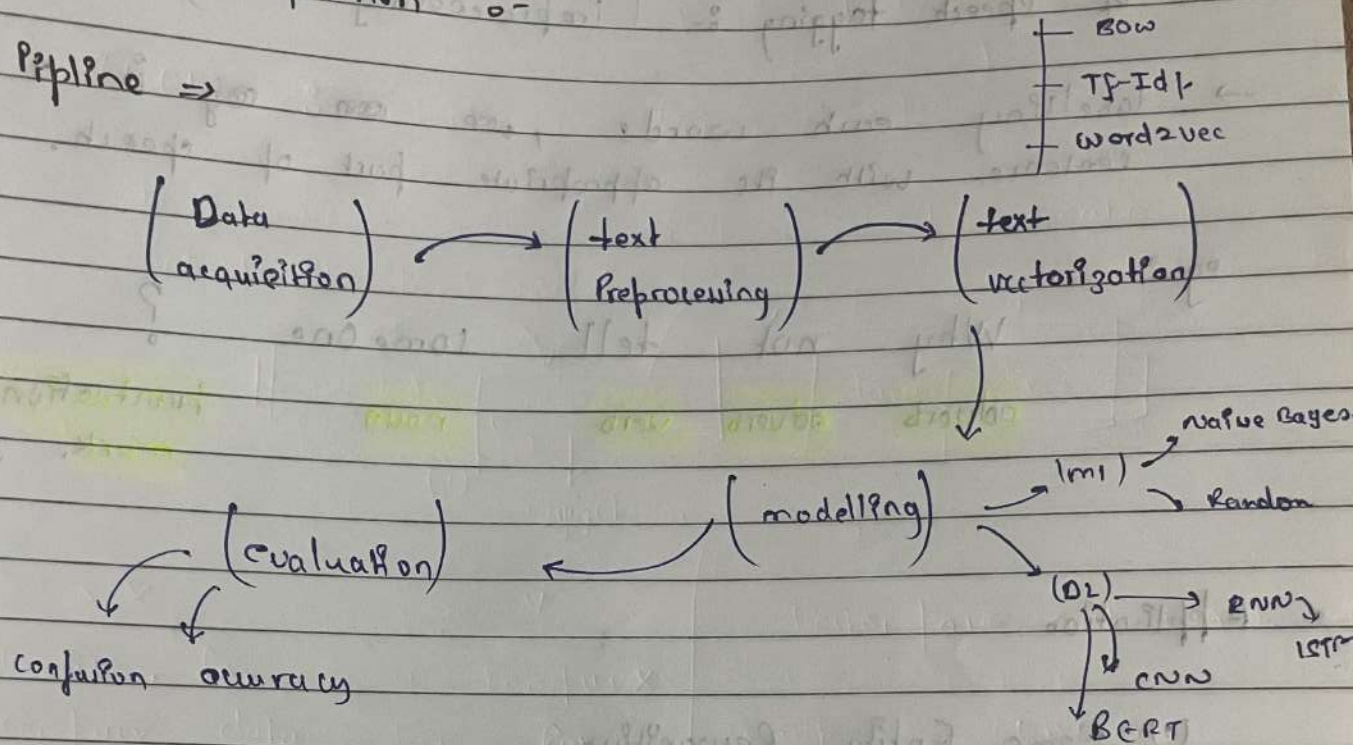


Note :-

small data \longrightarrow CBOWlarge data \longrightarrow Skip gram

Text classification :-

Pipeline \Rightarrow



* Parts of speech tagging :- (Preprocessing task)

→ labelling each words, one can tag in a sentence with its appropriate part of speech.

eg:-

Why not tell some one ?
 adverb adverb verb noun punctuation mark.

• Application

- (1) Name Entity Recognition
- (2) Question Answering System
- (3) word sense disambiguation.
- (4) chat bot

→ POS works on "Hidden Markov model"

→ Spacy library is used