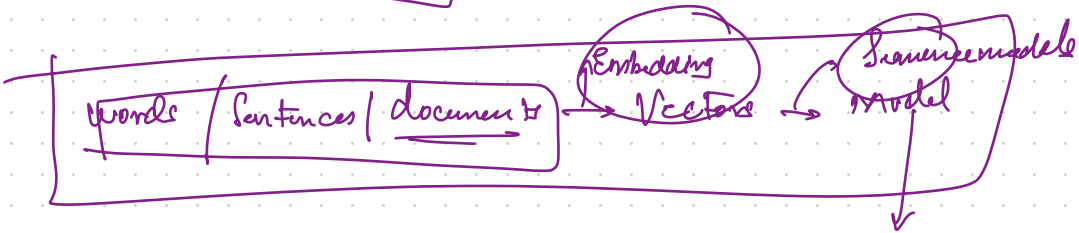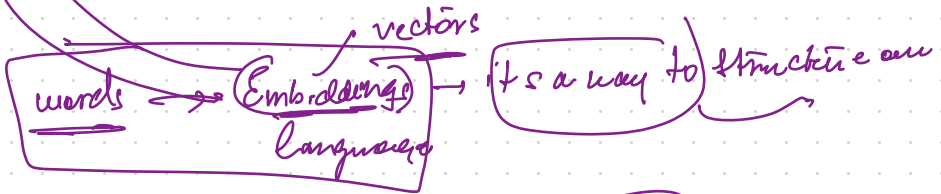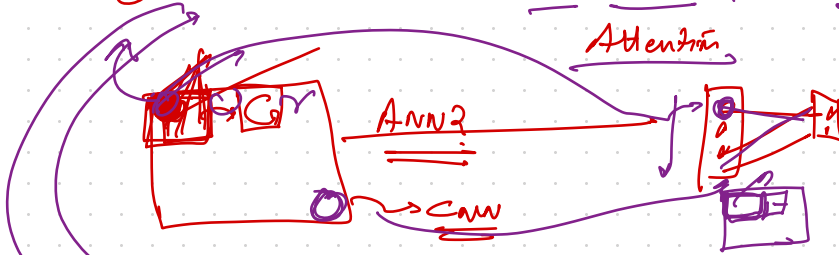# Natural Language Processing

1) NLP problems use case.

- How do we structure words/sentence

↓

Feature engineering

○ Sequence models → RNN, LSTM, GRU, Transformer

Attention

ANN?

GRU

CNN

words → Embeddings → vectors → it's a way to structure our language

| words | Sentences | documents | → Embedding Vectors → Sequence models Model

↓

Large Language models

RLHF → ChatGPT

Img → Text

Multimodal

Question/Prompt → Encoder → □ → Decoder → Answers

Natural language — is a Sequence



This (i) (@) great batch

Great this is batch

Natural Language Processing

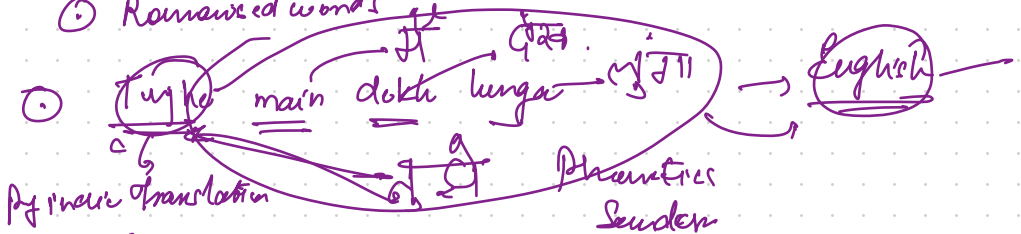Understanding information from text

- Unstructured →
- Content / Sentiments / Slangs / Sarcasm

- Indize — — — — — — —
- Romanized words
- Tuj ke main dekh lunga → मैं → देख → लूंगा → English
  Phonetic translation
  Phonetics Sender

- Structured

Use Cases:
Translation
gmail auto type → auto fill
Search engines

Sentiment analysis

Chatbots
Text classification, Sum
NER, information

Gmail Spam → Spam
→ HAM

$$\bigcirc = \large{f}(x) \quad \text{Language model,}$$

sequential

Language model

Sentences are very cnclear → Social media / Human Conversations

Preprocessing techniques:

Parsues

Special Characters @ http....

Document clasification → sports / Entertainment

Removal of unwanted Characters or patterns of characters from our Source sentence.

Ticket Classifier
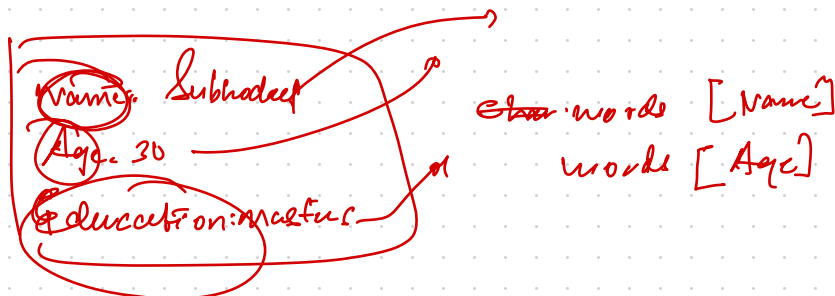
Sentiment classification
☺ → + ve
→ - ve

Ticket
☺ → IT OS
→ IT Apps

Regex → a process of defining certain patterns using certain rules.

logic

" I am very hungry, I will order pizzas and give ✓ to my friend."

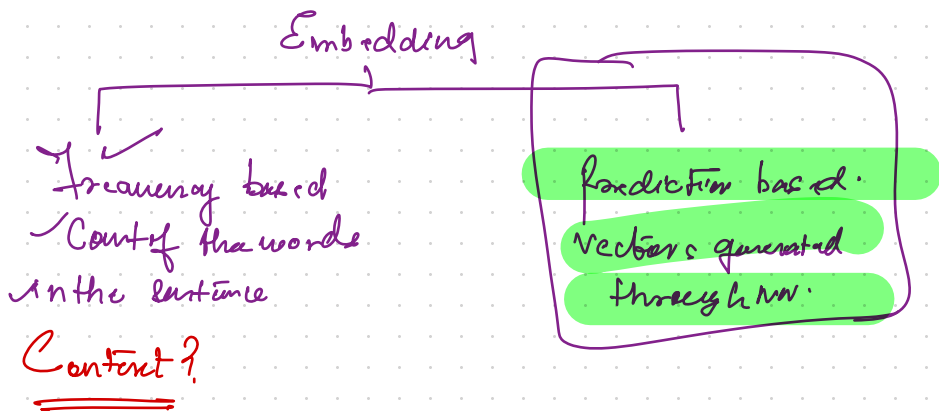Regex → Extract some defined patterns
from sentences

[w]
[s]

Name: Subhadeep
Age: 30
Education: Master

Star words [Name]
words [Age]

Stopwords (an), a, the
because and Theen

Translation , Chatbot → Tomo tur 25°C
→ I am an — — — — —
What's the temp tomorrow

मैं घर जाता हूँ , Tom I am going home

Steming / Lemmatization

Embedding
words → vectors → model

# Embedding

Frequency based
Count of the words
in the sentence

Prediction based.
Vectors generated
through nn.

## Content?

If am going home, home is where heat is

If : 1 , am : 1, going : 1 home 2 , is : 2
                           where 1, heat : 1

Run , Running , Ran & Run

Go , gone , going

Go

helpful

Sentences
Vector
form

1

10:33 → 10:40pm

How to represent sentences/documents in a vectorised form using frequency embeddings.

Text | Category → nFertainment

(1) ------- ✓
(2) ------- ✓    → 5,000 words  ML
100 ✓

→ document

Union of all the words → Vocabulary
Union of all the documents → Corpus
one word → Token

## Bag of words model

### Document classification

→ features/words

flow

(Document) → Outcome

vem
doc2
doc1
BOW

w2
w1
w3

(•) Cristiano Ronaldo got transferred Cristiano is
→ Sports                                   great

(•) Puca hepa performs in Cannes.

| Cristiano | Ronaldo | got | transferred | Dua | ripa | pulpme in | Come |
|-----------|---------|-----|-------------|-----|------|-----------|------|
| ① 1 | 1 | 1 | 1 | 0 | 0 1 | 0 0 | 0 |
| ② 0 | 0 | 0 | 0 | 1 | 1 | 1 1 | 1 |



Ronaldo

doc

Cristiano

got

**Problem:** ① too many dimensions

② does not maintain content

↓ order not maintained
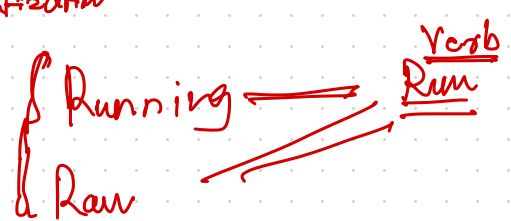
By definition given and
all words are independent
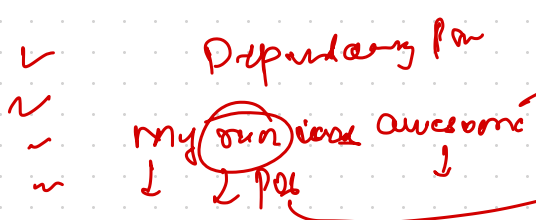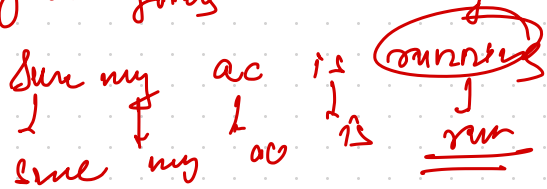


doc2  w
w1
doc1
doc3

Stemming/ lemmatization



Go
going
went
→ GO

Run
Running
Ran
→ Ran

Stemming & Lemmatization Explanation          Porter Stemmer

Lemmatization

Running —— Verb
                Run
Ran

                                    Speedily

                                    Speedy

Going —— GO        GO         Caressed
Gone                           Cary  → Cari

I   am  GO verb  to  visit my  house  and  will
I  am  going  to  visit my  home  and  will make → make

Sure my  ac  is  running
sure  my  ac  is  run

            Dependency Par
            root
My running ioss awesome
        POS

"I  am  great  at  public  speaking"

Stemming &     Lemmatization
(word)
        Sentence  → words
                word tokenizer

Sentence tokeniser

documents → Sentence

• split ("")

"I | am | going | to | my | home" . split (":")

[ 'I', 'am', 'going', 'tu', 'ing', 'homa' ]

nltk → word. tokeniser ( " " )

tweet → words → (count of the word in +ve tweets, count of the word in negative tweet)



−ve counts

+ve counts

Count of +ve words    count of −ve words , 1

tweet →

( )

+ve 50      0      100

60

+ve      −ve

The   Game

tweet :

tweet ( 1 |v , 1 | 0 )

freqs

$$\begin{cases} (game, 0) : 20 \\ (game, 1) : 40 \\ (awesome, 0) : 50 \\ (awesome, 1) : 90 \end{cases}$$

tweet: game is awesome

game awesome

pos = freqs [(game, 1)] = 40 + freqs [(awesome, 1)] →, 90

neg = freqs [(game, 0)] = 20 + freqs [(awesome, 0)] → 50

pos = 130

neg = 70

tweet = $(1, 130, 70)$