

# Natural Language Processing

## History

'Can machines think?'  
— Alan Turing.

## Turing Game

- 3 Participants
- 2 People & 1 computer.
- one is an interrogator & other → <sup>compt.</sup> person.
- computer → will try to make believe that it's a human response.
- person → will try to make believe that it's a machine.

## ELIZA

- early NLP system capable of carrying a limited conversation

ELIZA: what's your name?

RAM: My Name is Ram.

ELIZA: Hi I'm Eliza. what do you want to talk about?

Mid 1950's - 1960's

Birth of NLP & Linguistics

- Main Aim ↓
- Mostly hand coded rules
  - No major advancements

*Machine Translation*

---

1960 - Mid 70's → Dark Era of NLP.

because of no proper development in machine translation in previous era → (1950-1960's)

- people started to believe it as hype & machine / comp. to be incapable of language translation.

↙ 1970's - 1980's Slow Revival of NLP.

- some research activities revived.
- emphasis was still on linguistics.
- work on toy problems

① 1980's & 1990's Statistical Evolution

- by this time computational power increases substantially.
- data-driven / statistical approach

↙ 2000's

② Computational Power ↑↑.

↙ 2010's

Evolution of Deep Learn.  
Neural N/W.

word2vec

Glove

FastText

BERT

2017 +

Transformers  
self atten  
GPT . . .

## NLP.

### Natural Language Processing .

This is a **NLP** class.  
↑  
which class

**NLP**

NLP is a field of Computer Science, Artificial Intelligence & linguistics concerned with the interaction b/w computers & human languages.

specially, the process of a computer extracting meaningful information from Natural language I/P and/or producing Natural Language O/P.

- wikipedia.

## why NLP ?

### ① Answering Questions

→ what is the next train from the city after 3:00 PM Train?

→ I'm a Masters D.S student, which classes do I have today?

→ who is Alan Turing?

## ② Information Extraction

→ we decided to meet for our class at 10:00 AM tomorrow in Learning class.

To do: clan Meeting

Time : 10:00 AM

Venue : Learnbey Clm.

## ③ Translation

मेरा नाम राम है  $\Rightarrow$  My Name is Ram  
Hindi  $\longleftrightarrow$  English

## ④ Text Summarization

→ extract keywords

# Notebook LM.

- extract
- create an abstract of the entire article / Book etc.

## ⑤ Sentiment Analysis

Processing reviews of Product/Services  
& extract sentiment of people and it.

## Level Of Analysis in NLP.

## 1.) Morphological Analysis -

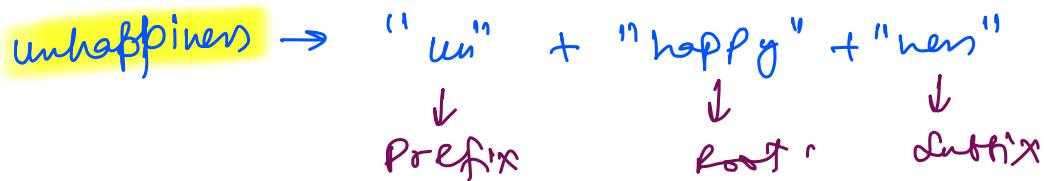
→ breaking down the words into smallest meaningful units  
units → **morphemes**

→ It analyzes structure of word - Prefix, Suffix, root.

eg →

Cats → "Cat" + "s" → Plural

running → "run" + "ing" → suffix



## ② Syntactic analysis

- This examines the grammatical structure of sentences.
- how words combine to form a phrase or a sentence.
- AIC to grammar.

(9)

Eg. → ① The cat sat on the mat.

↓      ↓      ↓  
subject   verb   prepositional  
phrase.

② determining Parts of speech.

POS Tagging

③ Sentence Structure.

## ③ Semantic Analysis

- focuses on meaning of words, phrases & sentences
- going beyond just grammar to understand what the text actually means.

Eg. →

① River bank } different meanings.  
Money bank }

② Resolve pronoun references.

Amit

wrote the exam. He scored 90%.

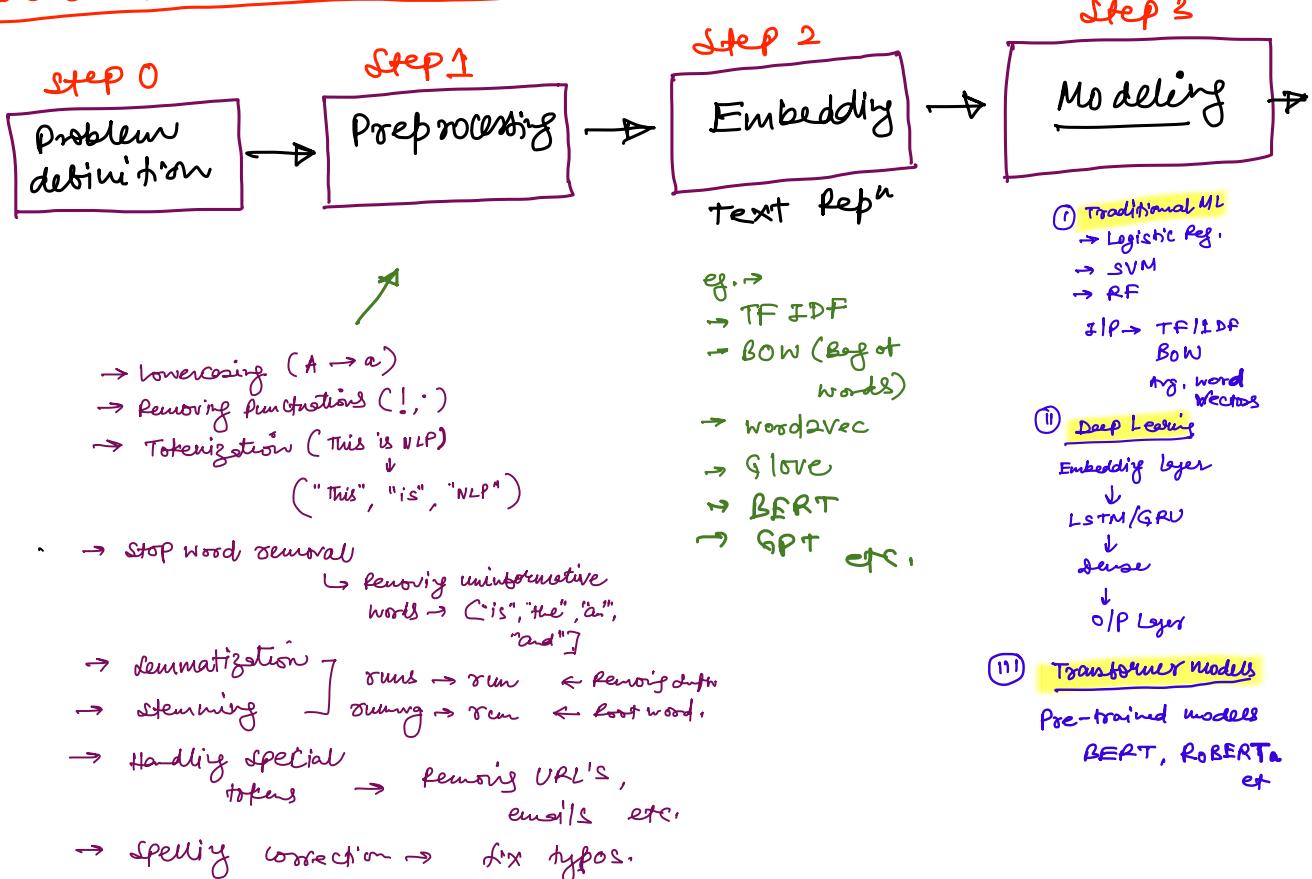
### ③ word sense disambiguation and understanding context.

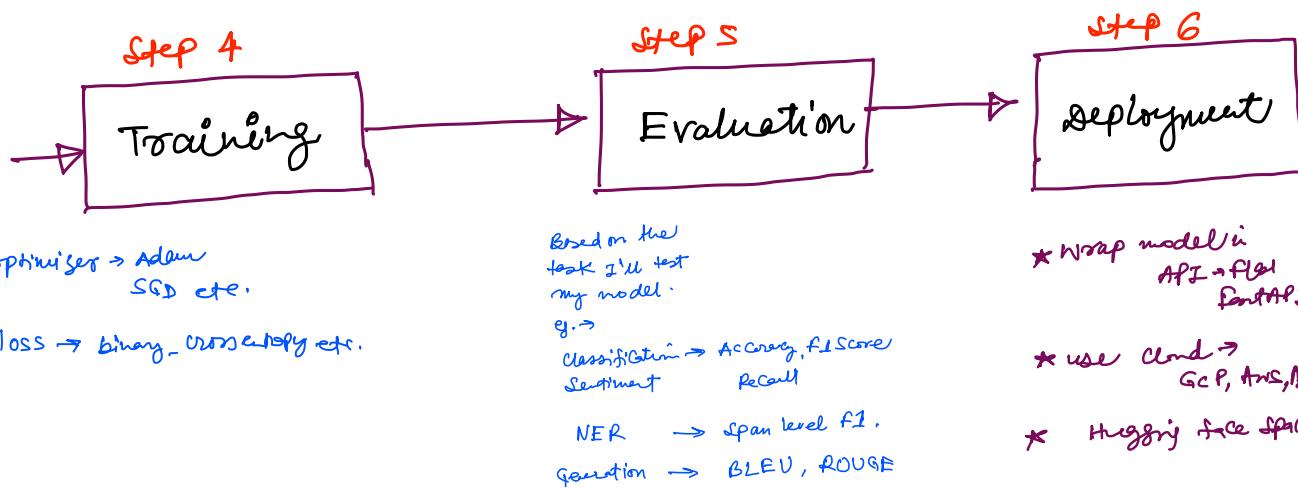
#### ④ discourse analysis

- Examines how meaning is constructed across larger units of texts.
- looks at the relationship betw sentences & paragraphs.
  - ↓  
to understand overall coherence & context.

eg. → how sentences connect in a conversation or a docnt.

#### General flow of NLP.





## # Preprocessing Techniques

### 1. Lowercasing

- each character in your text is converted to lower case.
- "Apple" → "apple"
- "ORANGE" → "orange"
- .lower() → uses unicode mapping.

$$\begin{matrix} A & B & C & D \\ \downarrow & \downarrow & \downarrow & \downarrow \\ a & b & c & d \end{matrix}$$

This step normalizes text because most NLP models treat "APPLE" and "apple" as different tokens unless explicitly trained otherwise.

### 2. Tokenization

- Breaking text into smaller units → tokens.
- tokens can be words / subwords / sentences.

I **don't** like Java .

"I"      "do"      "n't"      "like"      "Java"      ":"

Subwords

- we can use RegEx or libraries like Spacy, Huggingface etc.
- ↓
- include language rules, abbreviation handling etc.

### 3. Removing Punctuation

- Removing punctuation characters ( . , ! ? ; : )
- often using regex.
- it checks for each character's unicode category to determine if it falls in punctuation category.
- Can help in bag of words (Bow) models  
but we can skip this in tasks which require sentence structure.

### 4. Removing Stopwords

- removing words that don't contribute significantly to the semantic meaning of the sentence.
- NLTK & Spacy provide a list of stopwords & it matches each token against these lists & remove.
- goal → Remove noise & Reduce dimensionality.

### 5. Stemming

- Reduces the word to its root form by removing prefixes or suffixes.

"ing".

- rule based & doesn't care about grammar.

root

Relational → Relate

Arguing → Argu   
not grammatically  
correct.

## 6. Lemmatization

→ It uses a vocabulary & morphological analysis to get from the dictionary form of a word ↓  
lemma

→ better  $\xrightarrow{\text{lemma}}$  good

→ It also looks at the POS (Parts of Speech) to convert word into lemma.

e.g. → Lemma

meeting (Noun) → meeting

meeting (Verb) → meet

## 7. Removing Numbers

→ Stripping out digits using regex.

→ context sensitive.

→ It matches numerical characters (unicode) & replaces them with an empty string.

## 8. Removing extra spaces

→ removing extra spaces introduced in the text.

→ mostly prevalent in text obtained from HTML (web) scraping etc.

→ `' '.join(text.split())'`

### 9.) Removing Special Characters

- Removing characters like #, @, &, \* etc.
- we do it using regex.

### 10.) Handling Contractions

- Contractions like don't are expanded to "do not", "you're" → "you are", I'm → "I am".
- use Python libraries → contractions, textSearch.
- internal logic checks each token for expands by looking for pre defined mappings.

### 11.) POS (Parts of Speech) Tagging.

- POS tagging assign grammatical tags
  - ↓  
Parts of speech (Noun, Pron etc)
- In background libraries uses different models to label each word based on features
  - ↓  
Suffix, Surrounding words, Capitalisation etc.

NLTK	→	Perception
Spacy	]	bi-directional LSTM
Stanza	]	Transformer

## 12. Name Entity Recognition (NER)

- NER identifies & classifies entities like people, organisation & location etc.
- under the hood
  - (a) Rule Based → Pattern of POS chunks.
  - (b) Context.
  - (c) Deep learning models → Bi-LSTM, Transformers.

## 13.) Handling missing data.

- Predict next word / missing word
- Replace → mode.

## 14.) spelling correction

→ Pre-trained model on Corpus of words

→ Textblob

word is a 4 lettered word

↓

→ Probabilistic spelling model.

Trained on large corpus of words.

Levenshtein distance.

## 15.) Handling emoji's / emotions

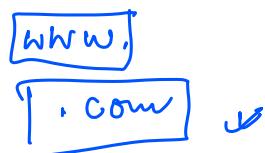
Smiley face → :slightly\_smiling\_face:

- unicode mapping
  - emoji
- ↓  
enjix

## (6.) Handling URL's

→ regex patterns like https://  
http://

Can be detected & removed.



## (7.) Removing HTML Tags

→ HTML is parsed with libraries  
one BeautifulSoup.

→ DOM tree

→ get\_text()

↳ recursively strips tags ↴ X

↳ & left with textual content ✓