

Tellnow

Python, stats, SQL, ML, ML project

Computer Vision

## Natural Language Processing

### Agenda

- Overview and history of NLP.
- Use of NLP
- Text processing, Regex.
- Text Normalization
- Lemmatization & Stemming
- Word embedding
- Word 2 Vec and doc2 vec
- NLTK & Spacy, Textblob, Stanford NLP

\* NLP

↓  
Natural language Processing



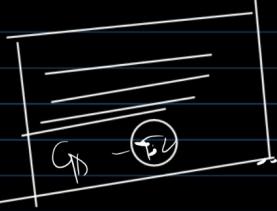
Part of AI that enables  
machines to understand, interpret,  
generate and respond to  
human language

✓ Ram & Shyam is going to temple. He will go movie  
after that.

Animal  
Cat dog

## Motivation

{ Chat Bots  
 auto correction  
 document summarization  
 translation  
 imaging →   
 Search engine  
 Chat GPT  
 Personal Assistant  
 Medical health records  
 Sentiment Analysis

Zoho  
 Recommendation  
 (all about) → factcheck.  


## History of NLP

1950's → during cold war → Machine translation

→ Turing test → Alan Turing → Proposes the idea of machine that can exhibit intelligent behaviors.

1954 → George town IBM Experiment → 60 russian sentence to English.

1980's → POS tagging.

↓  
Statistical NLP

2010 → deep learning

Word embedding, RNN, LSTM, Transformers

(Bert; GPT)

2020 → LLM → GPT3, GPT4, Gptoky

## Few Terminologies / Concepts

① NLP → Natural Language Processing

② Phonetics and phonology →  
Concerned with how sounds function

③ Morphology → It studies the structure of words

unhappiness  
un + happy + ness (suffix)  
(prefix)                            (root word)

④ Semantics

→ literal meaning of words, phrases, sentences.

↳ How meaning is constructed from words.

I went to the bank. (एस्टो)  
\_\_\_\_\_ ↓  
river bank

⑤ Pragmatics

↳ how context affects the meaning.

Can you explain me Normal distribution ??

↳ question, request.

⑥ Corpus → large collection of text documents.

Corpus. → [ "Ajay Devgn is good friend", "This is not geo" ]

⑥ Document → An individual piece of text within corpus

→ "Ajay Devgn and Ravi Kishan is friend"

⑦ Vocabulary

↪ Set of unique words without duplicates.

↪ Ram is going to Delhi. He goes everyday to, is, he  
↪ [Ram, is, going, to, delhi, he, everyday]

⑧ Word (Tokens)

↪ Individual component of sentence.

✓ am ✓ going

task in NLP

① Spell check → detect and correct spelling error

② Keyword based information retrieval →

↓ process of finding relevant information from

③ Topic Modelling a large collection of documents.

“The Government pass a law to regulate salary”

“Nvidia is releasing the most powerful GPU”

“Tesla released GigaRok”.

↪ Topic 1 → Politics → [Government, law, regulate]

↪ 2 → Technology → [Tesla, Nvidia, GPU, GigaRok]

④ Text Classification

Sentiment Analysis → +ve, -ve, Neutral.

Spam detection

News Categorization → Politics, Sports, Tech.

Language detection → French, English

### ⑤ Information Extraction :-

→ Process of automatically extracting structured information like names, dates, relationship, events from Unstructured data.

\* → Named Entity recognition:-

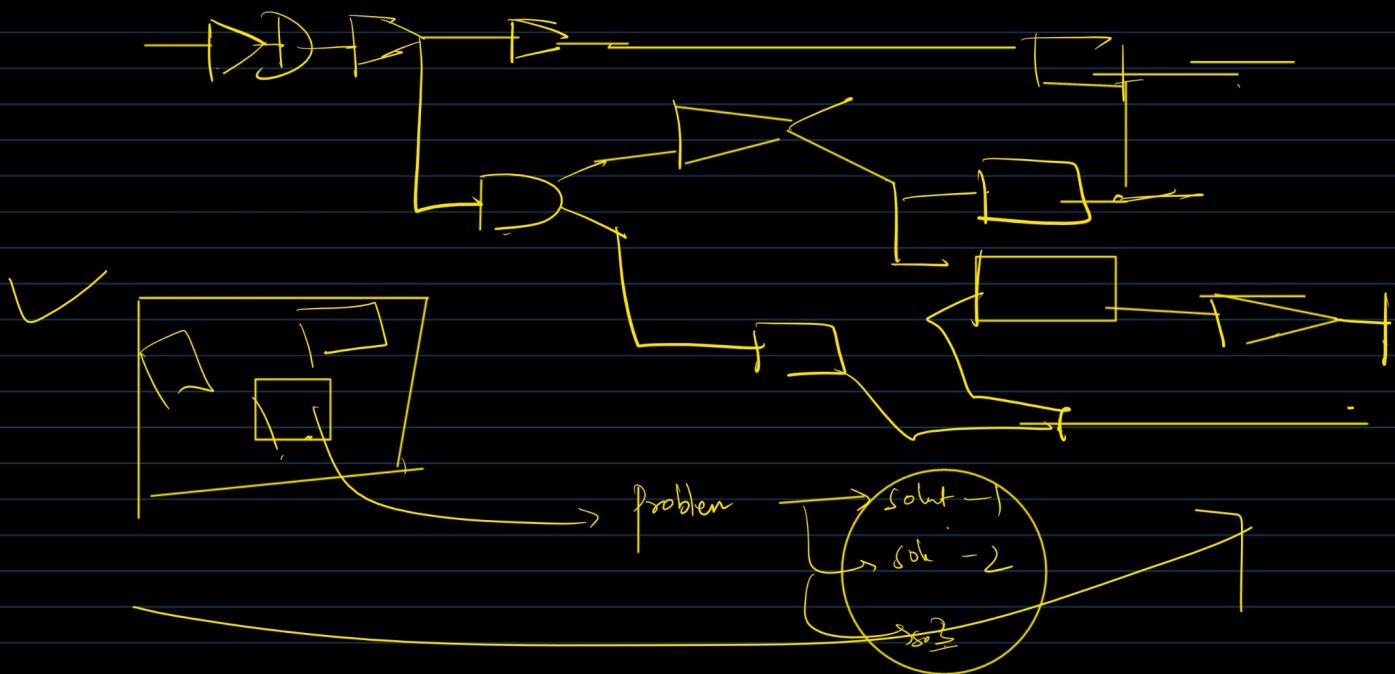
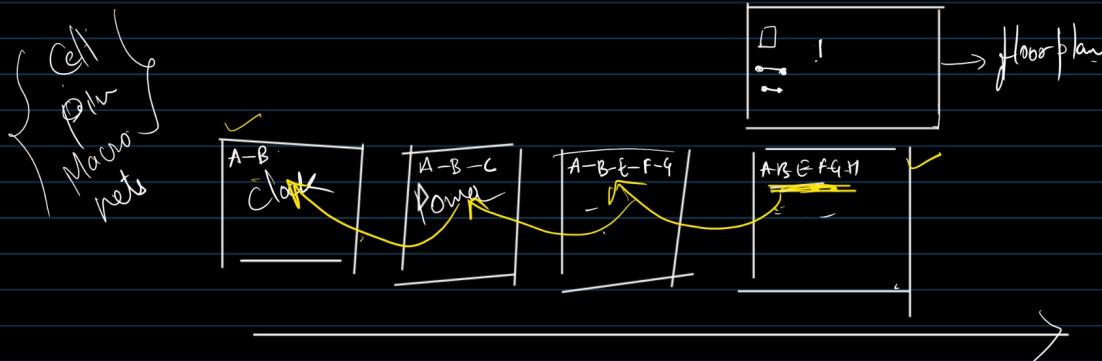
Ajay was born in Bhagatpur

(Person: Ajay

"Location": Bhagatpur)

### ⑥ Closed Domain Conversational Agent

Synopsis →



Banking → What is my account balance? Transfer 500 to UPKAR.

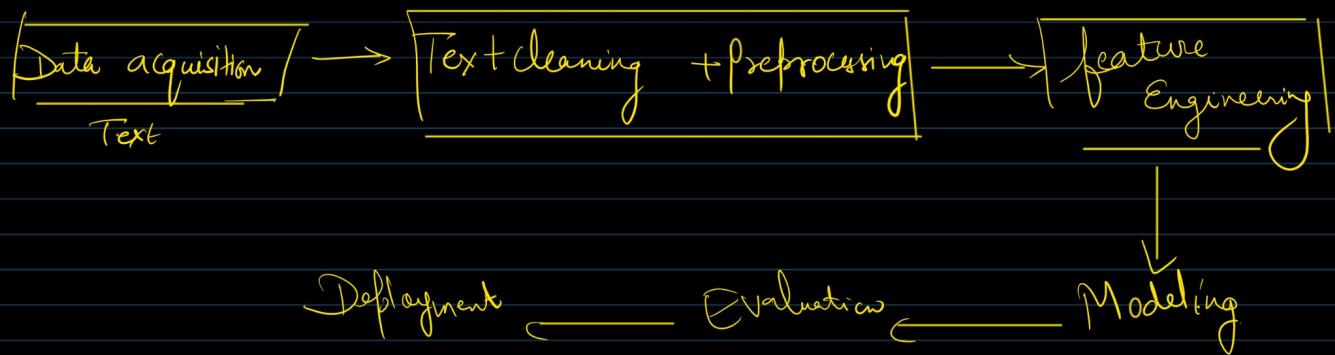
Healthcare → image → diagnosis → disease → solution

⑦ Text Summarization

⑧ Question Answering

⑨ Translation

ML → CRISP Framework



① Data acquisition

↳ Collection of textual data

→ Web Scraping

→ User Input

→ Data Repositories

→ APIs

→ Business Team

② Text Cleaning

"I<sup>①</sup> am going <sup>②</sup> to DELHI TO <sup>③</sup> Morrow <sup>④</sup> after ~oon."

↓      ↓      ↓      ↓

I am. X Caps. Space? in Extraspace - n

→ Once data is acquired, it has noise → (punctuation, extra space, caps, incomplete words, missing value)

### ③ Text processing

After cleaning, the text/data undergoes for pre-processing for modelling.

① Tokenization → (-A small amount is token)

"I am Ajay"  
tokenization of sentence: ("I", "am", "Ajay")

tokenization of word: "Ajay"  
("A", "j", "a", "y")

tokenization of corpus → "I am Ajay. I am learning ds."  
↓  
"I am Ajay", "I am learning"

Tokenization → Tokenization is the process of breaking a text string into smaller meaningful units called tokens.

The tokens can be

- Sentences
- Words
- Subwords
- Characters.

② Removing Punctuation & Special Character

I don't play \$, @, ?

### ③ Stemming & Lemmatization

I love reading → eat  
 I love eating → eat  
 He goes to market daily  
 ↘ 90

good better best  
 ↙ ↘ good

### ④ Stop word removal

↳ eliminating very common words → to, the, a, an and, is

I love going (to) delhi

### ⑤ N-grams

→ A contiguous sequence of n items from given text.

I love NLP → Unigram → "I", "love", "NLP"  
 ↗  
 words character

," → Bi-gram → ("I love", "Love NLP")

," → Tri-gram - ["I love NLP"]

\* In text processing → you would like to understand

① Word Context (which words appear together)

② Phrase pattern (frequent phrases or expression)

③ Probability of sequences (LLMs)

A language model estimates the prob of a sequence of words:

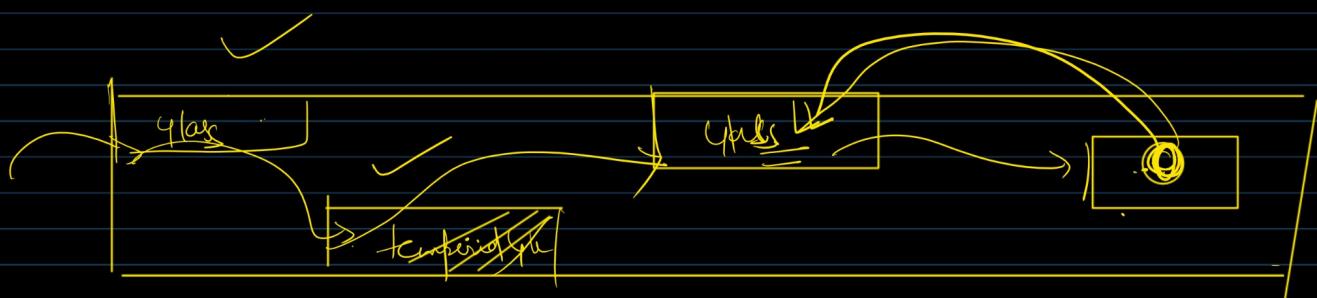
$$P(w_1, w_2, w_3, \dots, w_n) = P(w_1) \cdot P(w_2|w_1) \cdot P(w_3|w_1, w_2) \cdot \dots$$

Markov Assumption

$n-1$  words matter



$$P(w_3|w_1, w_2) \approx P(w_3|w_2)$$



⑥ POS → noun, verb → structure of sentence.

## Feature Engineering

↳ processed text to number

① TF-IDF

② Word Embedding

③ N-Grams → talked above

## modelling

Derivation

