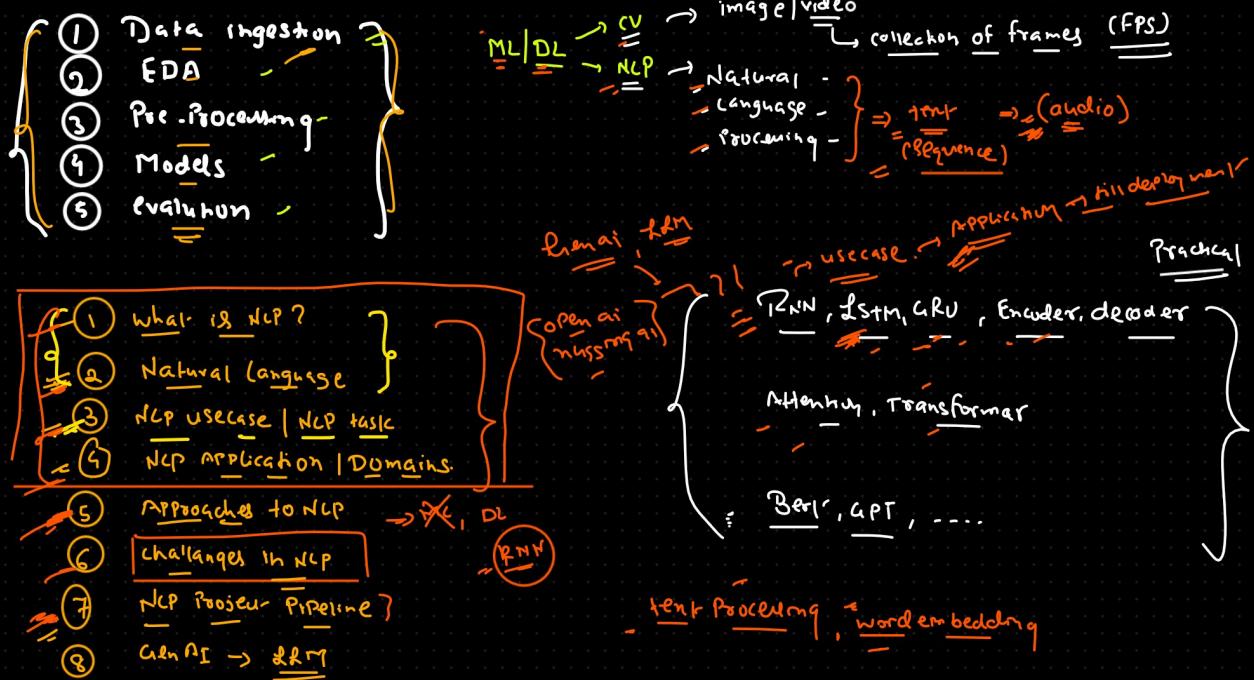


# NLP



## What is a NLP?

(1)

Natural language processing (NLP) is an interdisciplinary subfield of computer science and linguistics. It is primarily concerned with giving computers the ability to support and manipulate human language. It involves processing natural language datasets, such as text corpora or speech corpora, using either rule-based or probabilistic (i.e. statistical and, most recently, neural network-based) machine learning approaches. The goal is a computer capable of "understanding" the contents of documents, including the contextual nuances of the language within them. The technology can then accurately extract information and insights contained in the documents as well as categorize and organize the documents themselves.

Challenges in natural language processing frequently involve speech recognition, natural-language understanding, and natural-language generation.

= Natural + [language] + [processing]

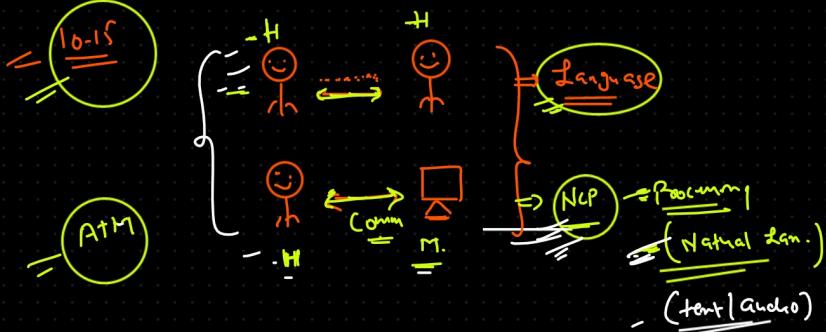
↓  
medium of communication

{Computer Science}

(text) (corpora, corpora)

Primary goal

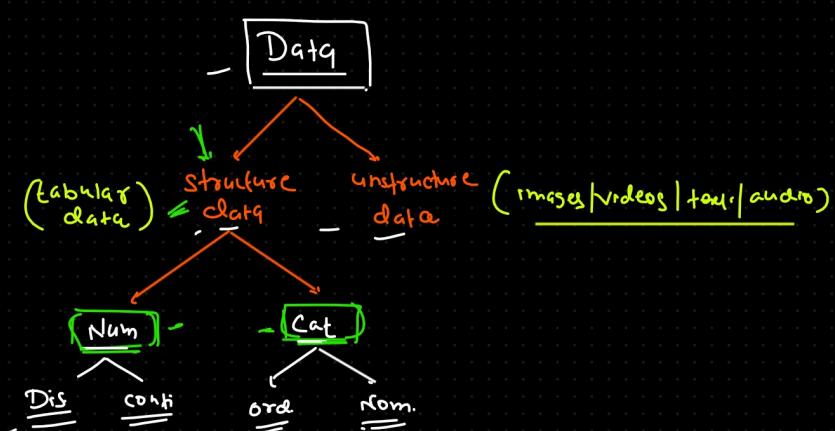
So Machine can understand  
Human language

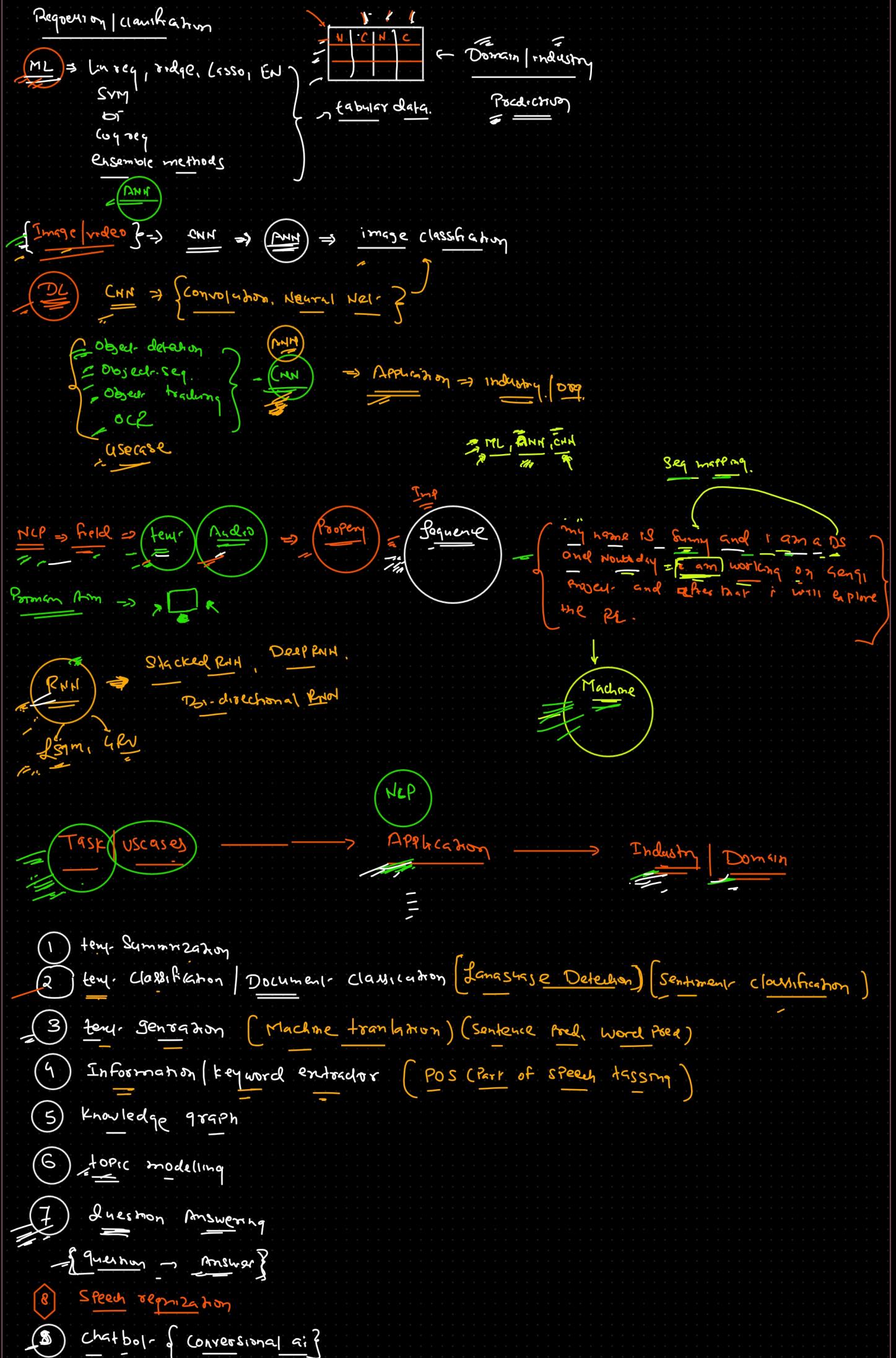


DataScience

ML / DL

{Regression  
Classification}





10 NER

11 Sentence similarity

12 Fill in the blanks / Spell corrector

### Application

1 Spam filtering  $\Rightarrow$  mail  $\xrightarrow{\text{spam}}$

2 near word Pred.  $\Rightarrow$  mail (auto completion)

3 Instructor (text summarization)  $\Rightarrow$  30 word

4 Grammacy (spell corrector) -

5 Paraphrasing (GPTbot) -

6 Google translator (text generation)

7 Google Assistant (speech to text)

8 twiter classification | sentiment analysis

9 Contentful A.I.  $\Rightarrow$  Insta., FB, twitter

$\xrightarrow{\text{text, audio}}$

$\xrightarrow{\text{Process}}$

Content

Model

$\left. \begin{array}{l} \text{Playing Shoe} \\ \text{Nike shoe} \end{array} \right\}$

$\left. \begin{array}{l} \text{Insta.} \\ \text{FB} \end{array} \right\}$

10 ChatPLT (Chatbot)

(Conversational AI)

11 Search engine

$\left. \begin{array}{l} \text{Text} \\ \text{EngP} \end{array} \right\}$

12 Alexa (speech to speech)

$\left. \begin{array}{l} \text{Text} \\ \text{EngP} \end{array} \right\}$

### NLP Pipeline

1 Data ingestion

$(EDA \Rightarrow \text{text})$

2 Feature Eng.

5 Deployment

2 Text Preparation

3 Manual approach

Basic Cleaning  
Basic to advance Processing

4 DL approach

5 Model building

5 Model evaluation

## ① Data ingestion

### Sentiment Analysis

↳ CSV?

SAC | NOSAC

3rd Party Services

huge.



ETL

Kafka

SParc

mysqz

### API

WebScraping [flipkart review]

Image  $\Rightarrow$  text  
Document  $\Rightarrow$  Word, PDF

OCR

## ② EDA $\Rightarrow$ Chart

## ③ Text Preparation

Clean.

NLP  $\Rightarrow$  text  $\Rightarrow$  Numbers

Feature eng

1 manual

Bow  
tf-idf

2 DC  $\Rightarrow$  NN

word2vec  
ELMO  
fasttext

## ⑤ Model building

1 Rule based approach  
(heuristic approach)

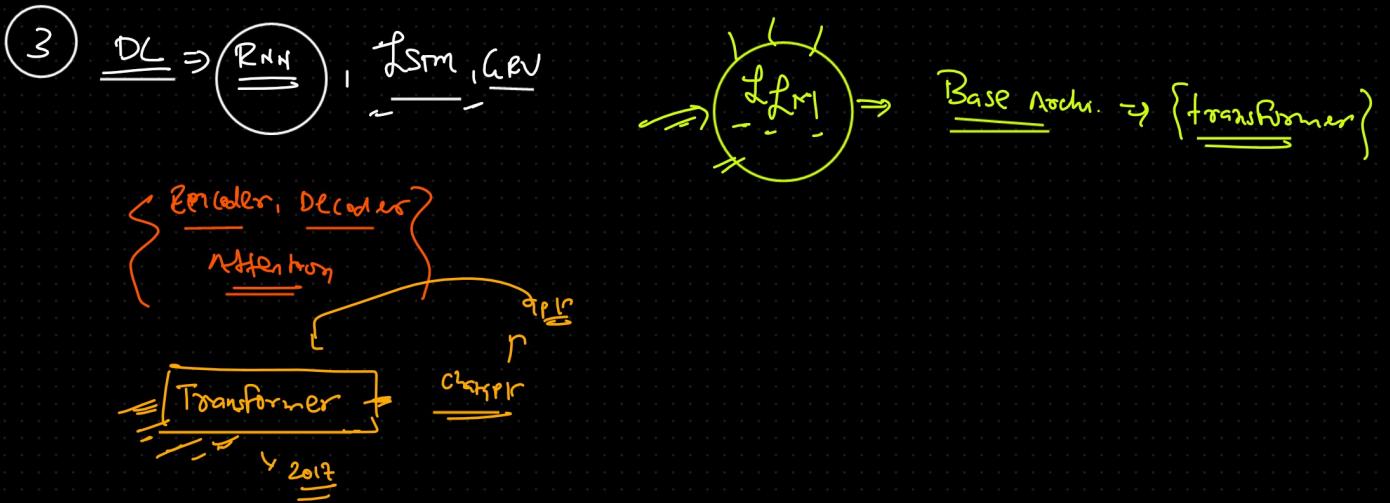
(1950-1960)

NLP  $\Rightarrow$  Seq mapping

Briger sentence  $\Rightarrow$  fail

2 KIC

Naive bayes (Cond. prob)



### 5 Model Evaluation

- 1 Confusion mat (Acc, F1 score, Prec, Recall)
- 2 BLEU
- 3 Perplexity (GPT)  
 $\hookrightarrow \underline{\text{PPL}}$

### Challenges

