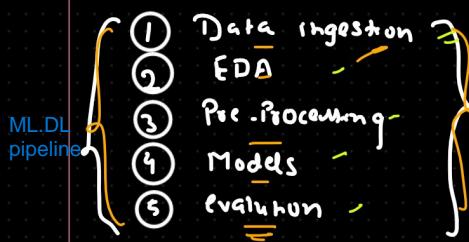


Earlier we talked about NLP wrt to machine learning which was probability based approach. Now we will be learning about NLP using deep learning based approach

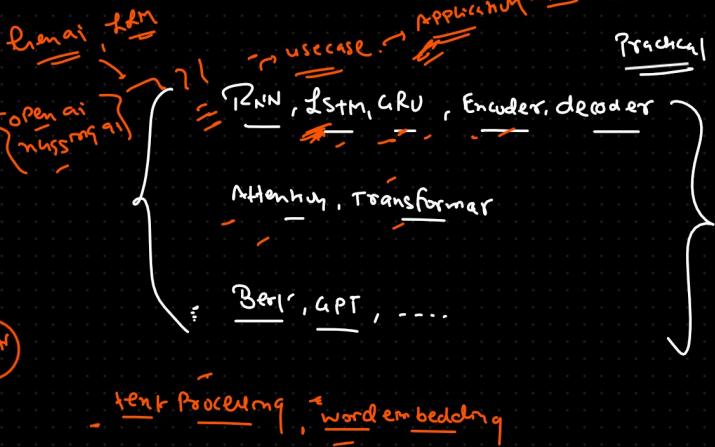
# NLP



In NLP we are going to deal with the dataset having sequence within it like text, audio, timeseries based dataset etc

Agenda:

- 1 What is NLP?
- 2 Natural language
- 3 NLP Usecase / NLP task
- 4 NLP Application / Domains
- 5 Approaches to NLP → DL
- 6 Challenges in NLP
- 7 NLP Project Pipeline
- 8 GPT → LSTM



These all architectures we are going to cover under this NLP

~~What is a NLP?~~

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.

decoding this definition

(1)

= Natural + [language] + [processing]

language is the medium of communication

medium of communication

{Computer Science}

(text) (corpora, corpora)

The way 2 humans communicate with each other using natural language, in the same way main AIM OF NLP is to establish same intelligence when communicating with the machines

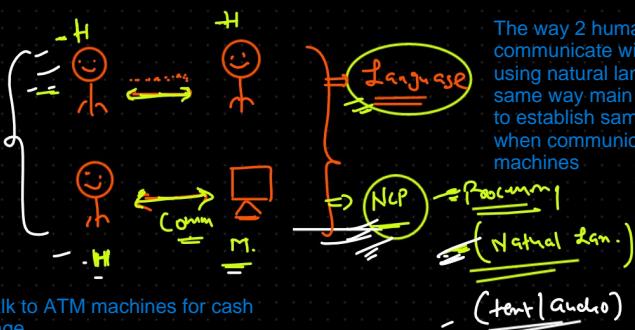
~~Primary goal~~

So Machine can understand  
Human language

Natural language communication b/w humans and machines was not possible 10-15 years back

10-15

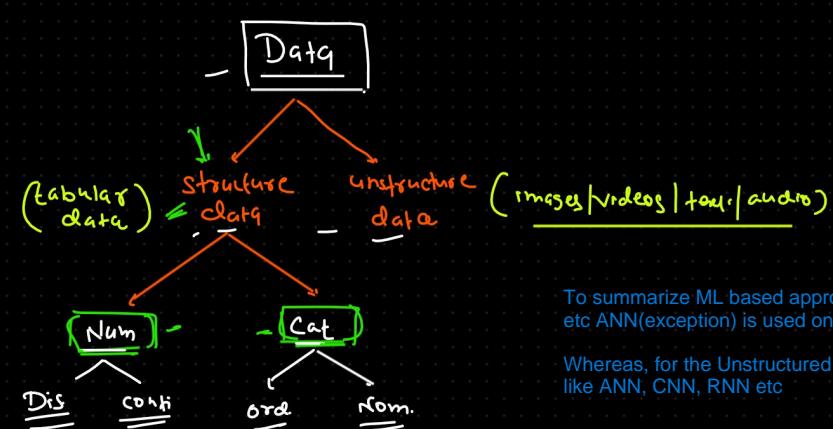
ATM



DataScience

ML / DL

{Regression  
Classification}



To summarize ML based approached like Linear regression, SVM etc ANN(exception) is used on top of Structured data.

Whereas, for the Unstructured data we use DL based approached like ANN, CNN, RNN etc

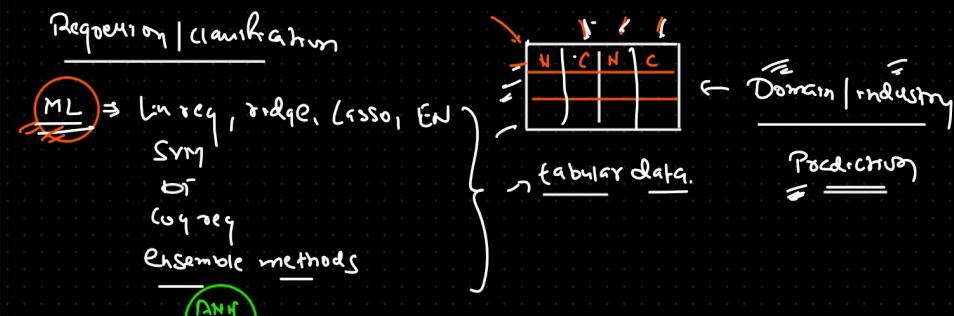
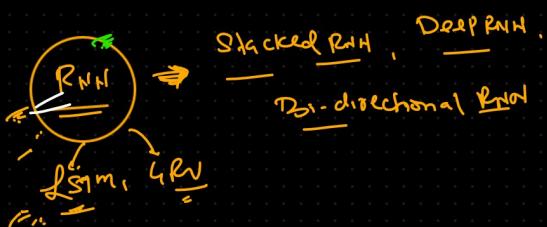
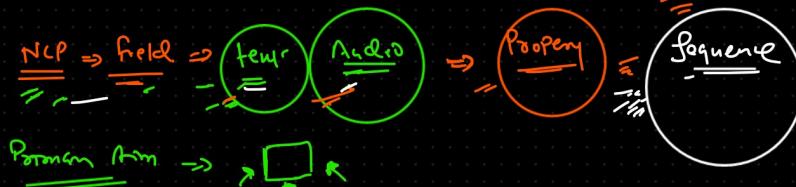
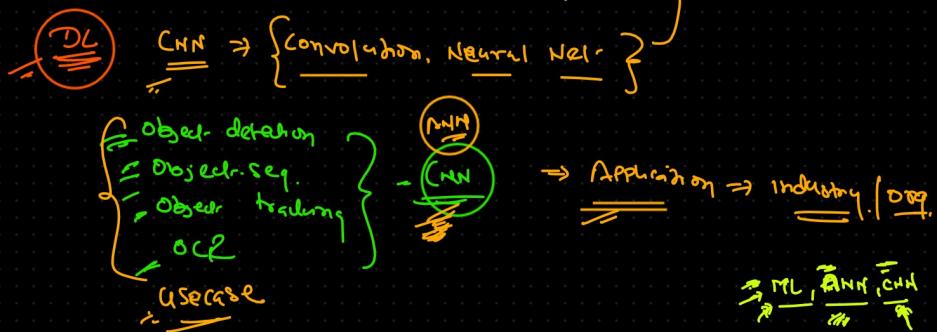


Image | video  $\Rightarrow$  CNN  $\Rightarrow$  ANN  $\Rightarrow$  image classification



# Different use cases of the NLP:

- 1 text summarization
- 2 text classification | Document classification [language detection] [sentiment classification]
- 3 text generation [Machine translation] (Sentence Prod, Word Prod)
- 4 Information | keyword extractor [POS (Part of speech tagging)]  $\gg$  POS tagging
- 5 knowledge graph
- 6 topic modelling (Based on important words present in the sentence or paragraphs we are basically going to give a topic to that sentence)
- 7 question answering {Question  $\rightarrow$  Answer}
- 8 Speech recognition

Both Question Answering and Chatbot(Conversational AI) may seem like similar use case but they are not. Question Answering is done on only some specific topic whereas, Chatbot does conversation over wide variety of topics

10 NER (Named Entity Recognition which is performed using BERT)

11 Sentence similarity

12 Fill in the blanks / Spell corrector

### Application

NLP powered applications:

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

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

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

④ Grammery (spell corrector)

⑤ Paraphrasing (Quillbot)

⑥ Google translator (text generation)

⑦ Google Assistant (speech to text)

⑧ twiter classification | sentiment analysis

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

(text, audio)

Process

Playing Shoe  
nike shoe

Insta FB

Contextual Advertisement: When we are searching for something over google then after then Insta, Fb, Twitter starts showing the advertisements similar to the searched keyword.

10 Chapt (chatbot)

(conversational ai)

11 Search Engine (Sentence matching, Sentence retrieval, Sentence classification)

(Engp)

12 Alexa (speech to speech)

(Engp)

### NLP Pipeline

1 Data ingestion

(Optional)  
EDA  $\Rightarrow$  text

2 Feature Eng.

5 Deployment

2 Text Preparation

{ Basic Cleaning  
Basic to advance Processing }

Manual approach (In machine learning)

DL approach (Automatic feature extraction)

4 Model building

5 Model evaluation

## 1 Data ingestion

### Sentiment Analysis

CSV or JSON

SQL | NoSQL

3rd party Services  
huge.



ETL

Kafka

If data is present with the 3rd party vendor like Azure block etc then data engineers help would be needed who will basically build a data pipeline using kafka and spark and can store the data in the sql database

Spark → MySQL

### API

Web scraping [flipkart review]

Image → text  
Document → Word, PDF  
OCR

## 2 EDA ⇒ Chatter

## 3 Text Preparation

Clean.

NLP → text → Numbers

Feature eng

In this step we are basically going to convert the text into embeddings(numbers)

1 manual

(Machine learning based approaches)

Bow

(Bag of word)

OHE

(One hot encoding)

TFIDF

(TFIDF)

2 DL → NN

(Automatic features extraction using neural network based DL approach)

word2vec

ELMO

fasttext

## 5 Model building

1 Rule based approach  
(heuristic approach)

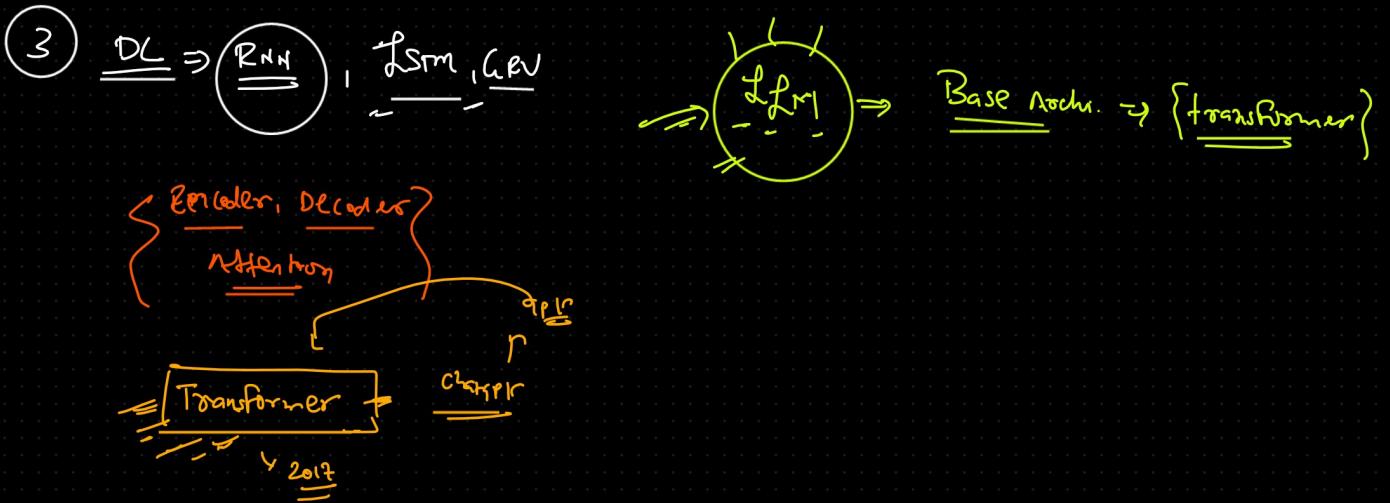
(1950-1960)

NLP ⇒ Seq mapping

2 KIC  
Naive bayes (Conditional Probability)

Bigger sentence ⇒ fail

These two approaches were not able to map the sequences effectively. Meaning that it was getting failed for bigger sentence by not able to retain the contextual understanding effectively.



### 5 Model Evaluation

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

I saw the man with the telescope.

This sentence can be interpreted in 2 ways:  
1. The person saw the man who had a telescope.  
2. The person used the telescope to see a man.

This sentence is ambiguous since there is no clear relationship b/w man and telescope

