# Day 1 – Natural Language Processing (ML & Deep Learning)

## Agenda

① Roadmap of Natural Language Processing ✓
② Why NLP ✓
③ Lot of Examples ✓
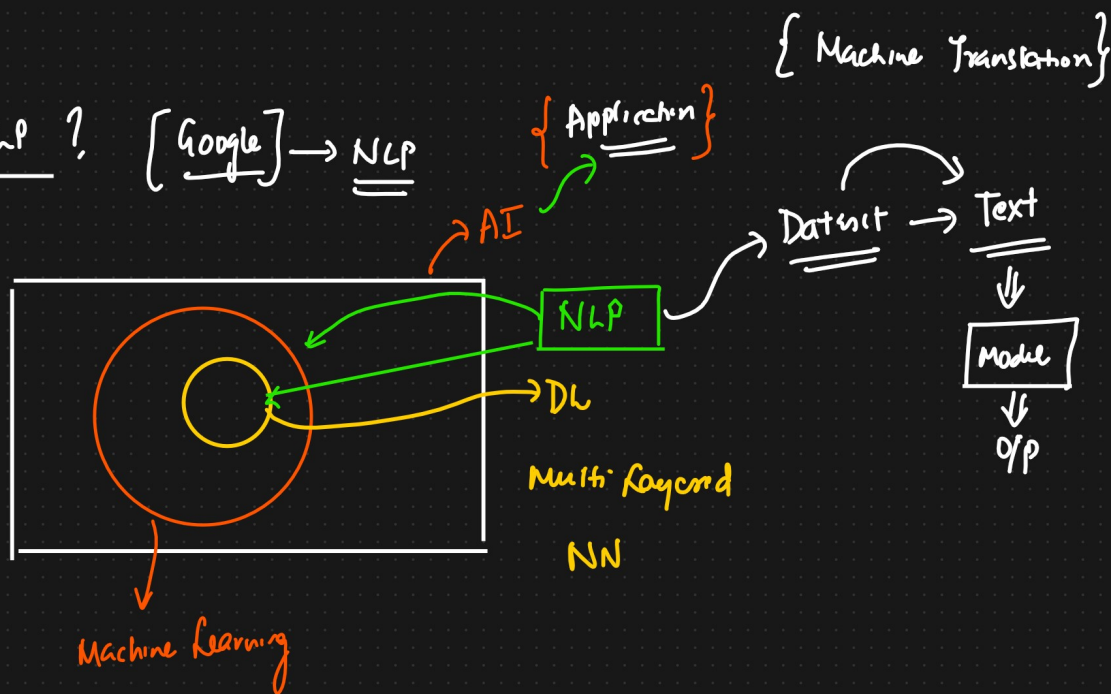④ Tokenization, Stemming, Lemmatization ✓
⑤ Bag of Words

<u>Quiz</u> = 5000 Rs

① 2000 Rs INR
② 1500 JNR
③ 1500 JNR
} { KrishnaiK06 }

## Prerequisites

① Python
② Stats      ⇒ Community Sessions
③ Machine Learning Algo
④ ANN, Optimizers, Loss functions

SPAM CLASSIFICATION

{ Machine Translation }

---

<u>Why NLP ?</u>   [ Google ] → <u>NLP</u>   { Application }

AI

NLP

Dataset → Text → Model → O/P

Multi Layered NN

Machine Learning

DL

# Roadmap of NLP

BERT

TRANSFORMERS

Bidirectional LSTM, Encoders, Decoders, Attention Models

Text Preprocessing → Word Embedding → Word2vec

RNN, LSTM RNN, GRU RNN → Deep Learning

ML Usecases

Text Preprocessing 3 → Word2vec, Avg word2vec

Text Preprocessing 2 → BOW, TFIDF, Unigrams, Bigrams

Text Preprocessing 1 → Tokenization, Lemmatization, Stopwords
Stemming

{ Tensorflow ✓
  Pytorch ✓

{ NLTR
  SPACY
  TextBlob

Huggingface

---

## NLP

① Tokenization

ML Usecase

Mail    Spam Classifier

I/p features = Email body, Email Subject

Spam/ham { Tokenization ①→ Stemming ②
                              ↓
            Stopwords ← Lemmatization ③

| Dataset | f2 | |
|---|---|---|
| f1 Email body | Email Subject | O/p |
| 1) You won 1000000 $$ | Billionaire | Spam |
| 2) Hey KRISH, HOW ARE YOU | Hello | HAM |
| ③ Credit Cards worth | Winner | Spam |

[You] [won] [1000000] [$$]

## Text Preprocessing
① Tokenization :{ Sentence into words}    Sentence → Converts

[Hey buddy I want ~~to~~ go ~~to~~ your house] → [not] ←

↓

Stopwords → yes

③ Stemming { Not have any meaning

Processing of reducing words to their Base

[ . ]

historical
history  } ⇒ histori     Word Stem

↓

Root word
or
Base form

finaly
final
finalized  } ⇒ fina     Meaning is gone ↑

going
goes
gone  } ⇒ go     { Meaningful word } ↑

**Advantages**

① Stemming is really fast

**Disadvantage**

① It is removing the meaning of the word

④ Lemmatization

history
historical  } history

finally
final
finalized  } final

Advantages | Disadvantage
--- | ---

**Advantages**

① meaningful World

**Disadvantage**

① It is slow.

**Usecase**

**Stemming**

① Spam Classification

② Review Classification

**Lemmatization**

① Text Summarization

② Language Translation

③ chatbot

Text Preprocessing

① Tokenization  ② StopWords  ③ Stemming  ④ Lemmatization

Words ⟶ Vectors

① Bag of Words  ② TF-IDF  ③ Word2Vec

⇓

Term Frequency — Inverse Document

Frequency.