# NLP → Natural language Processing

## Sub - sections

NLU
Natural language
understanding

NLG
Natural language
generation
- New text generate

Basic terms of any language :

① Phonemes → Smallest unit of any language
characters, Speech, sound

② Morphemes & lexemes
    ↓                    ↓
  words            Run, Running
                   Swim, swimming

③ Syntax → phrases, Sentences

④ Context → meaning

$\rightarrow$ Combination of Syntaxes

NLP applications:

① Sentiment analysis $\rightarrow$ Text classification

Tweets $\rightarrow$ +ve
  -ve
  Neutral

Movie Reviews $\rightarrow$ +ve
  -ve
  Neutral

| Text | Label |
|------|-------|
| Tweets, Movie Reviews | +ve -ve Neutral |

Product / services $\rightarrow$ % customer satisfied
  Data collection
  - Survey
  - google forms
  - audio files
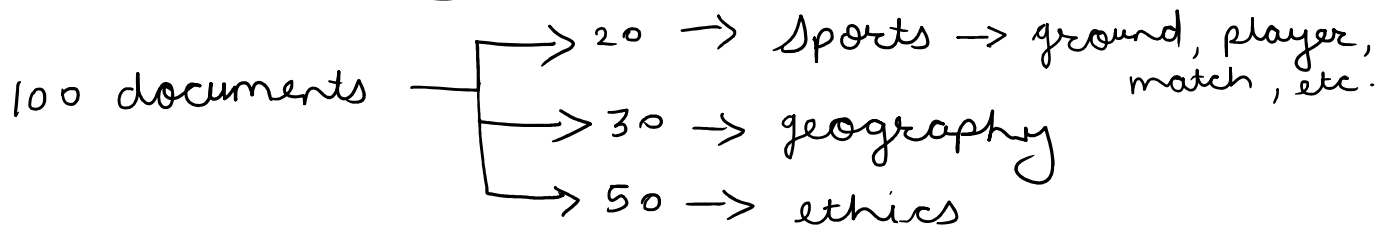
② Document classification

| Text | Label |
|------|-------|
| Document | Document name |

| Document content | Document name |
|---|---|
| | - Aadhar |
| | - PAN |
| | - DL |
| | - VID |

③ Text summarization

Text summarization

extractive T.S.          Abstractive T.S

extracting only          - creating new summary
important lines/              (NLG)
sentences from the
data

④ Topic modelling / identification

100 documents ——⌐→ 20 → Sports → ground, player, match, etc.
                ├→ 30 → geography
                └→ 50 → ethics

- Hidden pattern between text

100 documents ⌐→ 20 —→ 0 →
               ├→ 30 —→ 1 →
               └→ 50 → 2 →

⑤ Chatbot



Hi,
How can I help you?

I am having problem with 4 3 2 1 order no.
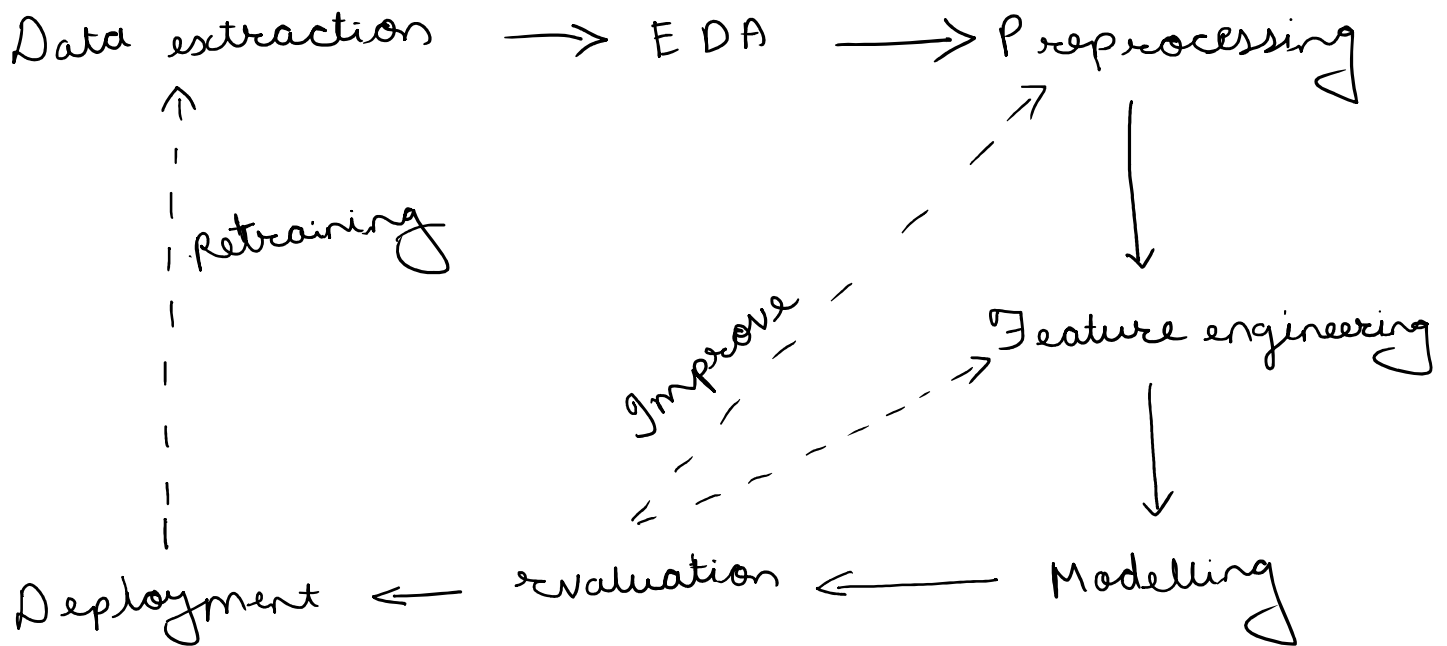– – – – – – – –

– automatic response generate

NLP Pipeline

Data extraction ⟶ EDA ⟶ Preprocessing

Retraining (dashed arrow up to Data extraction)

Improve (dashed arrow from Evaluation to Preprocessing / Feature engineering)

Preprocessing ⟶ Feature engineering ⟶ Modelling

Deployment ⟵ Evaluation ⟵ Modelling

① Data extraction → Data formats → json
txt
csv
images → OCR
↓
Pytesseract

Pytesseract
Amazon textract
Google vision

Data
→ Public
 - easily
available
→ Private
 - this belongs to some
 organisation

Data
Not available → Public options
available → Download

Public options
↓
Web scrapping

Web scrapping → Data ← Download
Data
↓
Quality
↓

what is noise?

" We >>>> loved $ @ / \ # the product _____

" _____ " will definitely _____ recommend."

Quality

Huge noise → extreme data preprocessing

minimal noise → Minimal data preprocessing

extreme data preprocessing → clean data ← Minimal data preprocessing

clean data ↓

Quantity ↓

POC (1 GB)

Problems regarding data:

① Quantity

② Quality

③ Exact data is not available for our use case.

④ Do not have continuous flow of data.

Data cycles → Monthly → 1 ✓ 2 ✗ 3 ✓ 4 ✗

Quarterly ↓

yearly +

EDA :     NLP $\longrightarrow$ ① Ngram

                       ② Word cloud

                       ③ Key phrase extraction

Ngram $\longrightarrow$ ① Unigram ② Bigram ③ Trigram

                ④ Quadragram

"Rajesh is hardworking guy"

Unigram = [ Rajesh, is, hardworking, guy]

Bigram = [ Rajesh is, is hardworking,
                   hardworking guy]

Trigram = [ Rajesh is hardworking,
                is hardworking guy ]
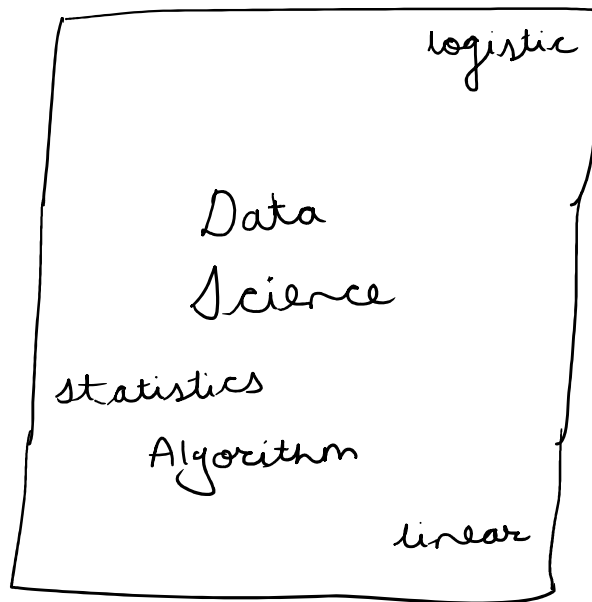
Why Ngrams ?

① To get insights from the data
          words understanding

         Positive reviews $\longrightarrow$ +ve words

         Negative reviews $\longrightarrow$ -ve words

② To get domain specific stopwords

② Word cloud

```
┌─────────────────────────┐
│                logistic │
│                         │
│     Data                │
│     Science             │
│   statistics            │
│        Algorithm        │
│                         │
│                linear   │
└─────────────────────────┘
```

Word frequency ↑
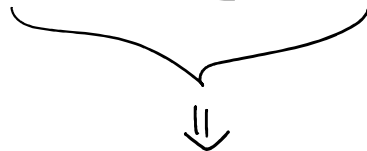word font ↑

Word frequency ↓
word font ↓

③ Keyphrase / Keyword extraction

   To extract → Important keyphrase / Keyword
                    from text
         → R A K E , Y A K E

   We are learning NLP.
         _____⌣_____/
                       ‖
                       ⇓

Preprocessing : →

① Tokenization → ① Sentence tokenization

① Tokenization ⟶ ① Sentence tokenization
② Word tokenization

"We are learning NLP.
NLP is a huge domain."

Sent = [ We are learing NLP., NLP is a huge domain.]

Tokens = [ We, are, learning, NLP, ., NLP, is, a, huge, domain, .]

Sentence tokenization ⟶ Syntax
⟶ Punctuation ⟶ ., ?, !
⟶ Conjunction ⟶ and, but

② Normalization ⟶ GREAT ⟶ Num1
‖                great ⟶ Num 2
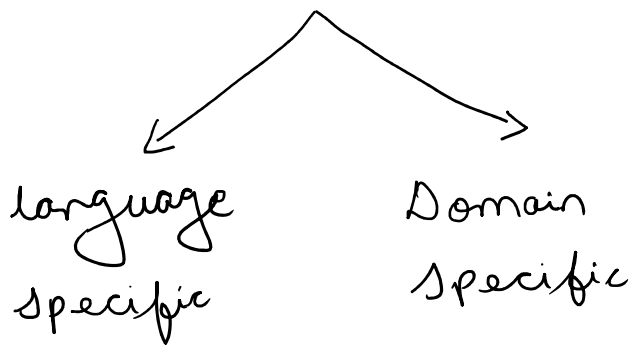⇓

Single case ⟶ lower case
                  upper case

③ Remove punctuation / Symbols

from string import punctuation

. , ; [ / / ] ? $ @ , - - - - - -

④ Remove stopwords

language          Domain
specific          specific
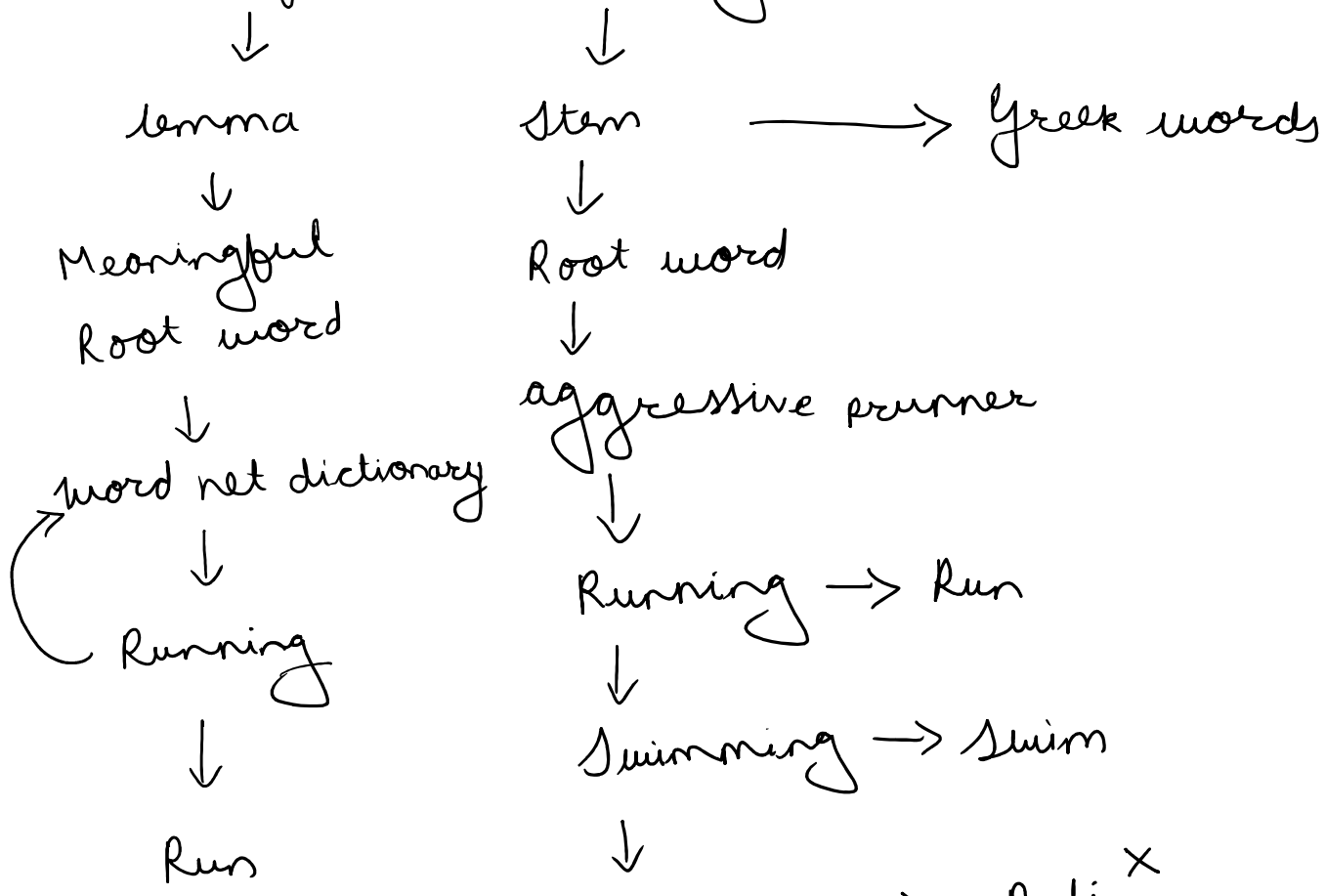
"Rajesh is suffering from cancer. Right now
Doctor Pravin is treating him. We have given
him xyz tablet."

Health domain ⟶ Doctor, tablet, capsule,
                          treatment, etc.

⑤ Lemmatization & stemming ⟹
        ↓                    ↓
     lemma              stem      ⟶   Greek words
        ↓                    ↓
   Meaningful         Root word
   Root word              ↓
        ↓             aggressive pruner
  word net dictionary        ↓
        ↓             Running ⟶ Run
     Running               ↓
        ↓             Swimming ⟶ Swim
      Run                   ↓              ✗

Run        $\downarrow$
           Believe $\longrightarrow$ Beli $^\times$

⑥ Contraction mapping $\longrightarrow$ expanding text

     didn't $\longrightarrow$ did not
     doesn't $\longrightarrow$ does not
     haven't $\longrightarrow$ have not

I didn't like the movie $\longrightarrow$ like movie

I liked the movie      $\longrightarrow$ like movie

I did not like the movie $\longrightarrow$ not like movie

I liked the movie      $\longrightarrow$ like movie

stopword_list = [ I, me, haven't, didn't,
                         no, nor, not, ----]
                       $\times$    $\times$    $\times$

stopword list . remove ("no")

                    ("nor")

                    ("not")

⑦ Handling accented characters $\longrightarrow$ unidecode library

     a $\longrightarrow$ $\hat{a}$, $\bar{a}$, $\acute{a}$, ----
     b $\longrightarrow$ $\hat{\imath}$, $\bar{\imath}$, $\acute{\imath}$,

$a \longrightarrow a, a, \acute{a}, \text{-----}$

$b \longrightarrow \hat{b}, \bar{b}, \dot{b}, \text{------}$

$\hat{a}\,ble \longrightarrow ble$

$a\,ble$

⑧ Autocorrection → correct spellings of words

→ autocorrect library

→ text blob library

Feature engineering

text ⟶ ✓Numerical format / vectors

Word embedding

word frequency ← Frequency Based          Prediction based → Algorithm

① Count Vectorizer ✓          ① Word 2 Vec ✓

② T F I D F ✓          ② Fast text

③ Doc 2 Vec

Modelling ⟶ Data → Numerical format

⇓

Build model

logistic regression , SVM, Random forest, Adaboost, Naive Bayes, Decision tree class, RNN, LSTM

evaluation : Accuracy
             Precision
             Recall

Low accuracy ⟶ Preprocessing
                      ↓
                Feature engineering
                      ↓
                Frequency based
                      ↓
              Prediction based

Good accuracy ⟶ Deployment