# NATURAL LANGUAGE PROCESSING

- Allows computer to understand language naturally, as a person does

## Text preprocessing -
1. Data cleaning
   a. Lowercasing
   - The method lower() converts all uppercase characters into lower case
   b. Punctuation Removal
   - Removing punctuation is a crucial step. We import string. string.punctuation refers all punctuation. Then we can remove using func lambda function.

   c. Stop words removal
   - Words that frequently occur in sentences and carry no significant meaning. Data must be in lowercase
   - We use the NLTK library
2. Spelling correction
   - Improve model accuracy
   - We use textblob library
3. Tokenization
   - Splits text into word (tokenization)

- Sentence tokenize: split a paragraph into meaningful sentence
- Word tokenize: split a sentence into unit, meaningful words. We use NLTK
4. Stemming
   - Convert words into their root word using some sort of rules irrespective of meaning. ek fish, fishes, fishing ⟶ fish. We use NLTK.
5. Lemmatization
   - Convert words into their root word using vocabulary mapping. Done with help of part of speech and its meaning, hua it doesn't generate meaningless root words. We use NLTK : chatbot

## EDA
   - Where explore & understand the data
1. Word frequency in Data
   - Counting unique words in our data gives an idea about our data, most frequent & least frequent
   - We drop the least frequent
   - We use NLTK's FreqDist

# TEXT PREPROCESSING
## II Word Embedding I

### a. Bag of words

- Simplest form of text representation in numbers. Like the term itself, we can represent a sentence a a bag of word vector

E.g.

S1: He is a good boy.
S2: She is a good girl
S3: Boys and girls are good

↓

stop words

S1: good boy
S2: good girl
S3: Boy girl good

then we deno vector, then a histogram
q words & Frequency
good
boy (i)
girl (ii)

Freq for s

|    | good | boy | girl | output |
|----|------|-----|------|--------|
| S1 | 1    | 1   | 0    |        |
| S2 | 1    | 0   | 1    |        |
| S3X|      |     |      |        |

Binary BoW
- ignores order of the words

### b. Term Frequency – Inverse Document Frequency (TF-IDF)

- TF-IDF is a numerical statistic that is intended to reflect how important a word is to a document in a collection

- TF: Measure of how frequently a term $t$ appears in a document $d$

$$tf_{t,d} = \frac{n_{t,d}}{\text{Number of terms in the document}}$$

$n$ = No. of times that term $t$ appears in a document

- IDF: Measure of how important a term is

$$idf_t = \log \frac{\text{No. of document}}{\text{No. of documents with term } t}$$

$$(tf\_idf)_{t,d} = tf_{t,d} * idf_t$$

© N-grams
Contin sequence of words or symbol, or token in a document

```
pip install --upgrade gensim - 2
import gensim
from gensim.models import word2vec
from gensim.models.word2vec import Word2Vec
# gensim is a library, it has API through which
  we can download model


import gensim.downloader as api
print(list(gensim.downloader.info()['models'].keys()))
# prints list of models available in gensim API

                                        dim
wv = api.load('glove-twitter-50')  # loads a pretrained
model
wv = api.load( ... )
wv.save("  path/      ")  # save model in the
local machine


wv['apple']  # vector returns vector rep of the word

from gensim.models import KeyedVectors
wv = KeyedVectors.load("  path ")  # load the
saved model from the local machine

wv.similarity("apple", "mango")  #


        # Embedding matrix
Embedding dim -
embedding_matrix = np.zeros((  input_dim, 100 ))
for word, i in word_index.items():
    if word in wv:
        embedding_matrix[i] = wv[model]


model = Sequential()
model.add(Embedding(input_dim, 100, weights =
  [embedding_matrix], input_len = max_len
```