



Gyan Ganga Institute of Technology and Sciences, Jabalpur
Computer Science & Engineering

Natural Language Processing

(AL – 504 (B))

PRACTICAL FILE

NAME:

ENROLLMENT NUMBER:

BRANCH: CSE-AIML

SESSION: 2023-24

Submitted to:

Dr Saurabh Tewari

Assistant Professor

INDEX

<u>S.No</u>	<u>TITLE</u>	<u>DATE</u>	<u>GRADE</u>	<u>Teachers Sign.</u>
1.	Tutorial for introduction to NLTK toolkit	07-08-2023		
2.	Example using NLTK for preprocessing text	21-08-2023		
3.	Preprocessing of text (Tokenization, Filtration, Script Validation, Stop Word Removal , Stemming	28-08-2023		
4.	Morphological Analysis	04-09-2023		
5.	N-gram Model	11-09-2023		
6.	POS Tagging	18-09-2023		
7.	Chunking	09-10-2023		
8.	Named Entity Recognition	16-10-2023		
9.	Finite State Automata	30-10-2023		
10	Study of PorterStemmer, LancasterStemmer, RegexpStemmer, snowball etc.	06-11-2023		

Experiment 1

Aim:-Tutorial for introduction to NLTK toolkit

Double-click (or enter) to edit

```
1 pip install --user -U nltk
```

```
Requirement already satisfied: nltk in /usr/local/lib/python3.10/dist-packages (3.8.1)
Requirement already satisfied: click in /usr/local/lib/python3.10/dist-packages (from nltk) (8.1.7)
Requirement already satisfied: joblib in /usr/local/lib/python3.10/dist-packages (from nltk) (1.3.2)
Requirement already satisfied: regex>=2021.8.3 in /usr/local/lib/python3.10/dist-packages (from nltk) (2023.6.3)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from nltk) (4.66.1)
```

```
1 import nltk
2 nltk.download()
```



NLTK Downloader

```
-----
d) Download  l) List  u) Update  c) Config  h) Help  q) Quit
-----
Downloader> d

Download which package (l=list; x=cancel)?
Identifier> l
Packages:
[ ] abc..... Australian Broadcasting Commission 2006
[ ] alpino..... Alpino Dutch Treebank
[ ] averaged_perceptron_tagger Averaged Perceptron Tagger
[ ] averaged_perceptron_tagger_ru Averaged Perceptron Tagger (Russian)
[ ] basque_grammars..... Grammars for Basque
[ ] bcp47..... BCP-47 Language Tags
[ ] biocreative_ppi..... BioCreAtIvE (Critical Assessment of Information
Extraction Systems in Biology)
[ ] bllip_wsj_no_aux.... BLLIP Parser: WSJ Model
[ ] book_grammars..... Grammars from NLTK Book
[ ] brown..... Brown Corpus
[ ] brown_tei..... Brown Corpus (TEI XML Version)
[ ] cess_cat..... CESS-CAT Treebank
[ ] cess_esp..... CESS-ESP Treebank
[ ] chat80..... Chat-80 Data Files
[ ] city_database..... City Database
[ ] cmudict..... The Carnegie Mellon Pronouncing Dictionary (0.6)
[ ] comparative_sentences Comparative Sentence Dataset
[ ] comtrans..... ComTrans Corpus Sample
[ ] conll2000..... CONLL 2000 Chunking Corpus
-----
```

```
KeyboardInterrupt                                Traceback (most recent call last)
<ipython-input-7-3b02390c7677> in <cell line: 2>()
      1 import nltk
----> 2 nltk.download()
```

⬆ 6 frames

```
/usr/local/lib/python3.10/dist-packages/ipykernel/kernelbase.py in _input_request(self, prompt, ident, parent,
password)
    893         except KeyboardInterrupt:
    894             # re-raise KeyboardInterrupt, to truncate traceback
--> 895             raise KeyboardInterrupt("Interrupted by user") from None
    896         except Exception as e:
    897             self.log.warning("Invalid Message:", exc_info=True)
```

KeyboardInterrupt: Interrupted by user

SEARCH STACK OVERFLOW

```
1 import nltk
2 pip install --user -U nltk
```

```
1 # import all the resources for Natural Language Processing with Python
2 nltk.download("book")
```

```
[nltk_data] Downloading collection 'book'
[nltk_data] |
[nltk_data] | Downloading package abc to /root/nltk_data...
[nltk_data] | Unzipping corpora/abc.zip.
[nltk_data] | Downloading package brown to /root/nltk_data...
[nltk_data] | Unzipping corpora/brown.zip.
[nltk_data] | Downloading package chat80 to /root/nltk_data...
[nltk_data] | Unzipping corpora/chat80.zip.
[nltk_data] | Downloading package cmudict to /root/nltk_data...
```

```

[nltk_data] | Unzipping corpora/cmudict.zip.
[nltk_data] | Downloading package conll2000 to /root/nltk_data...
[nltk_data] | Unzipping corpora/conll2000.zip.
[nltk_data] | Downloading package conll2002 to /root/nltk_data...
[nltk_data] | Unzipping corpora/conll2002.zip.
[nltk_data] | Downloading package dependency_treebank to
[nltk_data] | /root/nltk_data...
[nltk_data] | Unzipping corpora/dependency_treebank.zip.
[nltk_data] | Downloading package genesis to /root/nltk_data...
[nltk_data] | Unzipping corpora/genesis.zip.
[nltk_data] | Downloading package gutenber to /root/nltk_data...
[nltk_data] | Unzipping corpora/gutenberg.zip.
[nltk_data] | Downloading package ieer to /root/nltk_data...
[nltk_data] | Unzipping corpora/ieer.zip.
[nltk_data] | Downloading package inaugural to /root/nltk_data...
[nltk_data] | Unzipping corpora/inaugural.zip.
[nltk_data] | Downloading package movie_reviews to
[nltk_data] | /root/nltk_data...
[nltk_data] | Unzipping corpora/movie_reviews.zip.
[nltk_data] | Downloading package nps_chat to /root/nltk_data...
[nltk_data] | Unzipping corpora/nps_chat.zip.
[nltk_data] | Downloading package names to /root/nltk_data...
[nltk_data] | Unzipping corpora/names.zip.
[nltk_data] | Downloading package ppattach to /root/nltk_data...
[nltk_data] | Unzipping corpora/ppattach.zip.
[nltk_data] | Downloading package reuters to /root/nltk_data...
[nltk_data] | Downloading package senseval to /root/nltk_data...
[nltk_data] | Unzipping corpora/senseval.zip.
[nltk_data] | Downloading package state_union to /root/nltk_data...
[nltk_data] | Unzipping corpora/state_union.zip.
[nltk_data] | Downloading package stopwords to /root/nltk_data...
[nltk_data] | Unzipping corpora/stopwords.zip.
[nltk_data] | Downloading package swadesh to /root/nltk_data...
[nltk_data] | Unzipping corpora/swadesh.zip.
[nltk_data] | Downloading package timit to /root/nltk_data...
[nltk_data] | Unzipping corpora/timit.zip.
[nltk_data] | Downloading package treebank to /root/nltk_data...
[nltk_data] | Unzipping corpora/treebank.zip.
[nltk_data] | Downloading package toolbox to /root/nltk_data...
[nltk_data] | Unzipping corpora/toolbox.zip.
[nltk_data] | Downloading package udhr to /root/nltk_data...
[nltk_data] | Unzipping corpora/udhr.zip.
[nltk_data] | Downloading package udhr2 to /root/nltk_data...
[nltk_data] | Unzipping corpora/udhr2.zip.
[nltk_data] | Downloading package unicode_samples to
[nltk_data] | /root/nltk_data...
[nltk_data] | Unzipping corpora/unicode_samples.zip.
[nltk_data] | Downloading package webtext to /root/nltk_data...
[nltk_data] | Unzipping corpora/webtext.zip.

```

▼ Take a sentence and tokenize into words. Then apply a part-of-speech tagger

```

1 sentence = """At eight o'clock on Thursday morning
2 Arthur didn't feel very good."""
3 tokens = nltk.word_tokenize(sentence)
4 print(tokens)
5 tagged = nltk.pos_tag(tokens)
6 print(tagged)

['At', 'eight', 'o'clock', 'on', 'Thursday', 'morning', 'Arthur', 'did', 'n't', 'feel', 'very', 'good', '.']
[('At', 'IN'), ('eight', 'CD'), ('o'clock', 'NN'), ('on', 'IN'), ('Thursday', 'NNP'), ('morning', 'NN'), ('Arthur', 'NNP'), ('did', 'VBD'), ('n't', 'RB'), ('feel', 'VB'), ('very', 'RB'), ('good', 'JJ'), ('.', '.')]

```

▼ From the tagged words, identify the proper names.

```

1 entities = nltk.chunk.ne_chunk(tagged)
2 print(entities)

(S
  At/IN
  eight/CD
  o'clock/NN
  on/IN
  Thursday/NNP
  morning/NN
  (PERSON Arthur/NNP)
  did/VBD
  n't/RB
  feel/VB
  very/RB
  good/JJ
  ./.)

```

▼ Get texts for corpus analysis

```
1 %matplotlib inline
2 from nltk.book import *

*** Introductory Examples for the NLTK Book ***
Loading text1, ..., text9 and sent1, ..., sent9
Type the name of the text or sentence to view it.
Type: 'texts()' or 'sents()' to list the materials.
text1: Moby Dick by Herman Melville 1851
text2: Sense and Sensibility by Jane Austen 1811
text3: The Book of Genesis
text4: Inaugural Address Corpus
text5: Chat Corpus
text6: Monty Python and the Holy Grail
text7: Wall Street Journal
text8: Personal Corpus
text9: The Man Who Was Thursday by G . K . Chesterton 1908
```

▼ Generate a key-word in context concordance

```
1 text1.concordance("monstrous")
2 #There are many ways to examine the context of a text apart from simply reading it.
3 #A concordance view shows us every occurrence of a given word, together with some context.
4 #Here we look up the word monstrous in Moby Dick by entering text1 followed by a period, then the term concordance, and t

Displaying 11 of 11 matches:
ong the former , one was of a most monstrous size . ... This came towards us ,
ON OF THE PSALMS . " Touching that monstrous bulk of the whale or ork we have r
ll over with a heathenish array of monstrous clubs and spears . Some were thick
d as you gazed , and wondered what monstrous cannibal and savage could ever hav
that has survived the flood ; most monstrous and most mountainous ! That Himmal
they might scout at Moby Dick as a monstrous fable , or still worse and more de
th of Radney .'" CHAPTER 55 Of the Monstrous Pictures of Whales . I shall ere l
ing Scenes . In connexion with the monstrous pictures of whales , I am strongly
ere to enter upon those still more monstrous stories of them which are to be fo
ght have been rummaged out of this monstrous cabinet there is no telling . But
of Whale - Bones ; for Whales of a monstrous size are oftentimes cast up dead u
```

▼ Find words with similar concordance to a given word

```
1 print(text1)
2 text1.similar("monstrous")
3 print(text2)
4 text2.similar("monstrous")
5 #A concordance permits us to see words in context. For example, we saw that monstrous occurred in contexts such as the __
6 #What other words appear in a similar range of contexts?
7 #We can find out by appending the term similar to the name of the text in question, then inserting the relevant word in p

<Text: Moby Dick by Herman Melville 1851>
true contemptible christian abundant few part mean careful puzzled
mystifying passing curious loving wise doleful gamesome singular
delightfully perilous fearless
<Text: Sense and Sensibility by Jane Austen 1811>
very so exceedingly heartily a as good great extremely remarkably
sweet vast amazingly
```

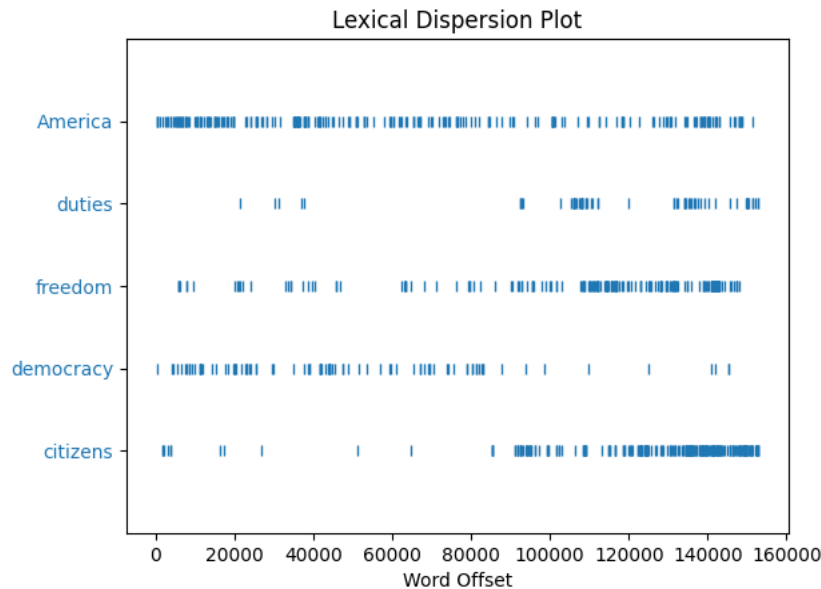
▼ Find contexts which are similar for the given words

```
1 text2.common_contexts(["monstrous", "very"])

am_glad a_pretty a_lucky is_pretty be_glad
```

▼ Plot where in the text certain words appear

```
1 text4.dispersion_plot(["citizens", "democracy", "freedom", "duties", "America"])
2 #Lexical dispersion is a measure of how frequently a word appears across the parts of a corpus.
3 #This plot notes the occurrences of a word and how many words from the beginning of the corpus it appears (word offsets).
4 #This is particularly useful for a corpus that covers a longer time period and for which you want to analyse how specific
```



Print the identity of a text, the length of the text and its vocabulary

```
1 print(text3)
2 print(len(text3))
3 print(sorted(set(text3)))

<Text: The Book of Genesis>
44764
['!', '"', '(', ')', ',', '.', ':', ';', '?', 'A', 'Abel', 'Abelmizraim', 'Abidah', 'Abide', 'Ab
```

```
1 text3.generate()
2 #Let's try generating some random text in the various styles we have just seen. To do this, we type the name of the text

Building ngram index...
laid by her , and said unto Cain , Where art thou , and said , Go to ,
I will not do it for ten ' s sons ; we dreamed each man according to
their generatio the firstborn said unto Laban , Because I said , Nay ,
but Sarah shall her name be . , duke Elah , duke Shobal , and Akan .
and looked upon my affliction . Bashemath Ishmael ' s blood , but Isra
for as a prince hast thou found of all the cattle in the valley , and
the wo The
'laid by her , and said unto Cain , Where art thou , and said , Go to ,\nI wi
ll not do it for ten ' s sons ; we dreamed each man according to\ntheir gener
atio the firstborn said unto Laban , Because I said , Nay ,\nbut Sarah shall
her name be . . duke Elah . duke Shobal . and Akan .\nand looked upon my affl
```

```
1 len(text3)
2 #So Genesis has 44,764 words and punctuation symbols,
3 #or "tokens." A token is the technical name for a sequence of characters – such as hairy, his, or :) – that we want to tr

44764

1 print(sorted(set(text3)))
2 #By wrapping sorted() around the Python expression set(text3) [1],
3 # we obtain a sorted list of vocabulary items, beginning with various punctuation symbols and continuing with words start
4 #Although it has 44,764 tokens, this book has only 2,789 distinct words, or "word types." A word type is the form or spel
5 # that is, the word considered as a unique item of vocabulary. Our count of 2,789 items will include punctuation symbols,
6

['!', '"', '(', ')', ',', '.', ':', ';', '?', 'A', 'Abel', 'Abelmizraim', 'Abidah', 'Abide', 'Ab
```

Print some statistics of word occurrence in the text

```
1 def lexical_diversity(text):
2     return len(set(text)) / len(text)
3 # Lexical richness refers to the range and variety of vocabulary deployed in a text by a speaker/write
4 def percentage(count, total):
5     return 100 * count / total
6 print(lexical_diversity(text3))
7 print(lexical_diversity(text5))
8 print(percentage(text4.count('a'), len(text4)))
9 #The most common approach to measuring lexical richness or diversity is based on the ratio of different words (types) to
```

```
0.06230453042623537
0.13477005109975562
1.457806031353621
```

Experiment 2

AIM:-Example using NLTK for preprocessing text

```

1 # Setup
2 !pip install -q wordcloud
3 import wordcloud

1 import nltk
2 nltk.download('stopwords')
3 nltk.download('wordnet')
4 nltk.download('punkt')
5 nltk.download('averaged_perceptron_tagger')

[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Unzipping corpora/stopwords.zip.
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   /root/nltk_data...
[nltk_data]   Unzipping taggers/averaged_perceptron_tagger.zip.
True

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import io
4 import unicodedata
5 import numpy as np
6 import re
7 import string

1 #Constants
2 # POS (Parts Of Speech) for: nouns, adjectives, verbs and adverbs
3 DI_POS_TYPES = {'NN':'n', 'JJ':'a', 'VB':'v', 'RB':'r'}
4 POS_TYPES = list(DI_POS_TYPES.keys())

1 # Constraints on tokens
2 MIN_STR_LEN = 3
3 RE_VALID = '[a-zA-Z]+'

1 from google.colab import drive
2 drive.mount('/content/drive')

Mounted at /content/drive

1 # Upload from google drive
2 from google.colab import files
3 uploaded = files.upload()
4 print("len(uploaded.keys()):", len(uploaded.keys()))

Choose files quotes.txt
• quotes.txt(text/plain) - 6935 bytes, last modified: 23/11/2023 - 100% done
Saving quotes.txt to quotes.txt
len(uploaded.keys()): 1

1 for fn in uploaded.keys():
2     print('User uploaded file "{name}" with length {length} bytes'.format(name=fn, length=len(uploaded[fn])))

User uploaded file "quotes.txt" with length 6935 bytes

1 # Get list of quotes
2 df_quotes = pd.read_csv(io.StringIO(uploaded['quotes.txt'].decode('utf-8')), sep='\t')

1 # Display
2 print("df_quotes:")
3 print(df_quotes.head().to_string())
4 print(df_quotes.describe())
5

df_quotes:
      Author
0  Agatha Christie  I like living. I have sometimes been wildly, despairingly, acutely miserable, racked with sorrow; bu

```


1	Agatha Christie		
2	Agatha Christie		
3	Agatha Christie		
4	Agatha Christie		
	Author		Quote
count	46		46
unique	5		46
top	Arthur Conan Doyle	I like living. I have sometimes been wildly, d...	
freq	14		1

```

1 # Convert quotes to list
2 li_quotes = df_quotes['Quote'].tolist()
3 print()
4 print("len(li_quotes):", len(li_quotes))

```

```
len(li_quotes): 46
```

Tokenization

```

1 # Get stopwords, stemmer and lemmatizer
2 stopwords = nltk.corpus.stopwords.words('english')
3 stemmer = nltk.stem.PorterStemmer()
4 lemmatizer = nltk.stem.WordNetLemmatizer()

1 # Remove accents function
2 def remove_accents(data):
3     return ''.join(x for x in unicodedata.normalize('NFKD', data) if x in string.ascii_letters or x == " ")
4 # Process all quotes
5 li_tokens = []
6 li_token_lists = []
7 li_lem_strings = []

1

1 for i,text in enumerate(li_quotes):
2 # Tokenize by sentence, then by lowercase word
3     tokens = [word.lower() for sent in nltk.sent_tokenize(text) for word in nltk.word_tokenize(sent)]
4 # Process all tokens per quote
5 li_tokens_quote = []
6 li_tokens_quote_lem = []
7 for token in tokens:
8 # Remove accents
9     t = remove_accents(token)
10    # Remove punctuation
11    t = str(t).translate(string.punctuation)
12    li_tokens_quote.append(t)
13    # Add token that represents "no lemmatization match"
14    li_tokens_quote_lem.append("-") # this token will be removed if a lemmatization match is found below
15    # Process each token
16    if t not in stopwords:
17        if re.search(RE_VALID, t):
18            if len(t) >= MIN_STR_LEN:
19                # Note that the POS (Part Of Speech) is necessary as input to the lemmatizer # (otherwise it assumes the worc
20                pos = nltk.pos_tag([t])[0][1][:2]
21                pos2 = 'n' # set default to noun
22                if pos in DI_POS_TYPES:
23                    pos2 = DI_POS_TYPES[pos]
24
25                stem = stemmer.stem(t)
26                lem = lemmatizer.lemmatize(t, pos=pos2)
27                # lemmatize with the correct POS
28                if pos in POS_TYPES:
29                    li_tokens.append((t, stem, lem, pos))
30                    # Remove the "-" token and append the lemmatization match
31                    li_tokens_quote_lem = li_tokens_quote_lem[:-1]
32                    li_tokens_quote_lem.append(lem)

1 # Build list of token lists from lemmatized tokens
2 li_token_lists.append(li_tokens_quote)
3 # Build list of strings from lemmatized tokens
4 str_li_tokens_quote_lem = ' '.join(li_tokens_quote_lem)
5 li_lem_strings.append(str_li_tokens_quote_lem)

```

```

1 # Build resulting dataframes from lists
2 df_token_lists = pd.DataFrame(li_token_lists)
3 print("df_token_lists.head(5):")
4 print(df_token_lists.head(5).to_string())
5 # Replace None with empty string
6 for c in df_token_lists:
7     if str(df_token_lists[c].dtype) in ('object', 'string_', 'unicode_'):
8         df_token_lists[c].fillna(value='', inplace=True)
9 df_lem_strings = pd.DataFrame(li_lem_strings, columns=['lem quote'])
10 print()
11 print("")
12 print("df_lem_strings.head():")
13 print(df_lem_strings.head().to_string())

df_token_lists.head(5):
   0      1      2      3      4      5      6      7      8      9     10     11     12     13     14     15     16     17     18
0  to  sherlock holmes she is always the woman      i have seldom heard him mention her under any other

df_lem_strings.head():

0 - sherlock holmes - - always - woman - - - seldom heard - mention - - - - name - - - eye - eclipse - predominates - w

```

Process results, and the most popular lemmatized words and group results by Part

```

1 # Add counts
2 print("Group by lemmatized words, add count and sort:")
3 df_all_words = pd.DataFrame(li_tokens, columns=['token', 'stem', 'lem', 'pos'])
4 df_all_words['counts'] = df_all_words.groupby(['lem'])['lem'].transform('count')
5 df_all_words = df_all_words.sort_values(by=['counts', 'lem'], ascending=[False, True]).reset_index()
6 print("Get just the first row in each lemmatized group")
7 df_words = df_all_words.groupby('lem').first().sort_values(by='counts', ascending=False).reset_index()
8 print("df_words.head(10):")
9 print(df_words.head(10))

```

```

Group by lemmatized words, add count and sort:
Get just the first row in each lemmatized group
df_words.head(10):

```

	lem	index	token	stem	pos	counts
0	woman	3	woman	woman	NN	6
1	adler	18	adler	adler	NN	4
2	irene	17	irene	iren	NN	4
3	emotion	14	emotion	emot	NN	4
4	reason	29	reasoning	reason	VB	2
5	observe	30	observing	observ	VB	2
6	particularly	20	particularly	particularli	RB	2
7	perfect	28	perfect	perfect	NN	2
8	precise	23	precise	precis	NN	2
9	predominates	10	predominates	predomin	NN	2

Top 10 words per Part Of Speech (POS)

```

1 df_words = df_words[['lem', 'pos', 'counts']].head(200)
2 for v in POS_TYPES:
3     df_pos = df_words[df_words['pos'] == v]
4     print()
5     print("POS_TYPE:", v)
6     print(df_pos.head(10).to_string())

```

```

POS_TYPE: NN

```

	lem	pos	counts
0	woman	NN	6
1	adler	NN	4
2	irene	NN	4
3	emotion	NN	4
7	perfect	NN	2
8	precise	NN	2
9	predominates	NN	2
11	abhorrent	NN	2
12	name	NN	2
13	seldom	NN	2

```

POS_TYPE: JJ

```

	lem	pos	counts
10	questionable	JJ	2
17	whole	JJ	2
28	dubious	JJ	2

```

POS_TYPE: VB

```

	lem	pos	counts
--	-----	-----	--------

4	reason	VB	2
5	observe	VB	2
16	take	VB	2
19	see	VB	2
26	balance	VB	2

POS_TYPE: RB

	lem	pos	counts
6	particularly	RB	2
23	admirably	RB	2
25	always	RB	2
34	late	RB	2
37	yet	RB	2

Sorted frequency plot for all words

```

1
2 li_token_lists_flat = [y for x in li_token_lists for y in x]
3 # flatten the list of token lists to a single list
4 print("li_token_lists_flat[:10]:", li_token_lists_flat[:10])
5

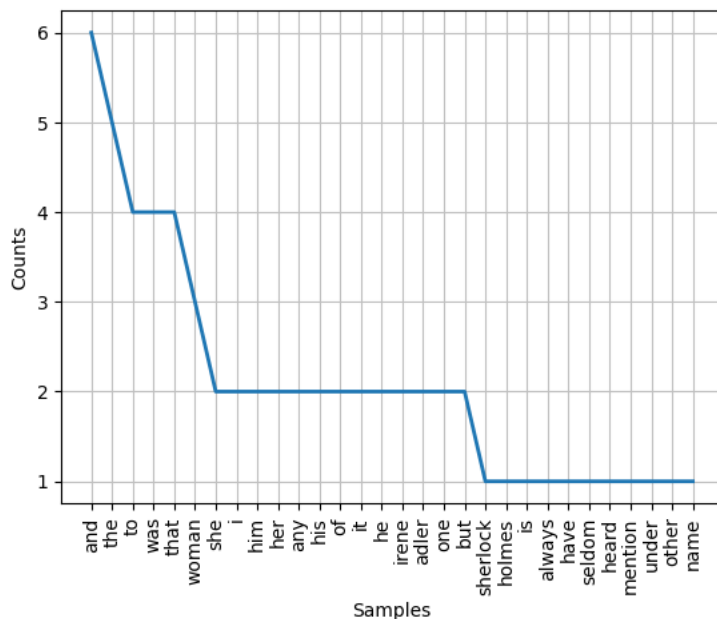
li_token_lists_flat[:10]: ['to', 'sherlock', 'holmes', 'she', 'is', 'always', 'the', 'woman', '', 'i']

```

```

1 di_freq = nltk.FreqDist(li_token_lists_flat)
2 del di_freq['']
3 li_freq_sorted = sorted(di_freq.items(), key=lambda x: x[1], reverse=True)
4 # sorted list print(li_freq_sorted)
5 di_freq.plot(30, cumulative=False)

```



<Axes: xlabel='Samples', ylabel='Counts'>

Sorted frequency plot for Lemmatized words after removing stopwords

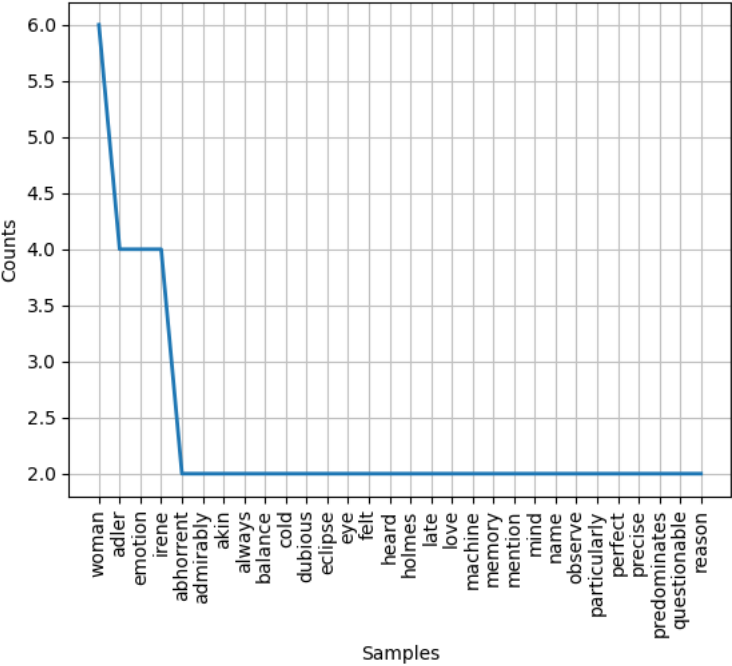
```

1 li_lem_words = df_all_words['lem'].tolist()
2 di_freq2 = nltk.FreqDist(li_lem_words)
3 li_freq_sorted2 = sorted(di_freq2.items(), key=lambda x: x[1], reverse=True)

1 # sorted list print(li_freq_sorted2 SK)
2 di_freq2.plot(30, cumulative=False)

```





<Axes: xlabel='Samples', ylabel='Counts'>

Experiment 3

Aim:-Preprocessing of text (tokenization,Filtration, Script Validation, Stop Word Removal, Stemming)

Code:


String Handling

```
1 print(len("what it is what isnt"))
2 s= ["what", "it", "is", "what", "it","isnt"]
3 print(len(s))
4 x=sorted(s)
5 print(s)
6 print(x)
7 d=x+s
8 print(d)

20
6
['what', 'it', 'is', 'what', 'it', 'isnt']
['is', 'isnt', 'it', 'it', 'what', 'what']
['is', 'isnt', 'it', 'it', 'what', 'what', 'what', 'it', 'is', 'what', 'it', 'isnt']
```

File Handling (tokenization and Filtering)

```
1 from google.colab import files
2 uploaded = files.upload()
```

 No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.
Savinn sk txt to sk txt

```
1 print("Output:-->")
2 for line in open("sk.txt"):
3     for word in line.split():
4         if word.endswith('ing' or 'sk'):
5
6             print(word)
7             print(len(word))
```

```
Output:-->
eating
6
cycling
7
```

Experiment 4

Aim:-To study Morphological Analysis

```
1 #Regular Expression code.
2
3 import re
4 input= "The 5 greatest batsmen in odi cricket 1 Sachin Tendulkar 2 Sir Vivian Richards 3 Brian Lara 4 Virat Kohli 5 Ricky
5 input= input.lower()
6 print(input,'\n')
7
8
9 result=(re.sub(r'\d+', '',input))
10 print(result)
```

```
the 5 greatest batsmen in odi cricket 1 sachin tendulkar 2 sir vivian richards 3 brian lara 4 virat kohli 5 ricky pontin
the greatest batsmen in odi cricket sachin tendulkar sir vivian richards brian lara virat kohli ricky ponting. sk
```

```
1 #StopWord Removal
2 def punctuations(raw_review):
3     text=raw_review
4     text=text.replace("n't",'not')
5     text=text.replace("'s",'is')
6     text=text.replace("'re",'are')
7     text=text.replace("'ve",'have')
8     text=text.replace("'m",'am')
9     text=text.replace("'d",'would')
10    text=text.replace("'ll",'will')
11    text=text.replace("in",'ing')
12    #import re
13    letters_only=re.sub("[^a-z,A-Z]", " ",text)
14    return("".join(letters_only))
15 t="how's my team SK doin , you're supposed to be not loosin"
16 p=punctuations(t)
17 print(p)
```

```
➡ howis my team SK doing , youare supposed to be not loosing
```

```
1
2 #synonym
3 import nltk
4 nltk.download('wordnet')
5
6
[nltk_data] Downloading package wordnet to /root/nltk_data...
[nltk_data] Package wordnet is already up-to-date!
True
```

```
1 from nltk.corpus import wordnet
```

```
1 synonyms=[]
2 for syn in wordnet.synsets("machine"):
3     for lemma in syn.lemmas():
4         synonyms.append(lemma.name())
5     print(synonyms)

['machine']
['machine', 'machine']
['machine', 'machine', 'machine']
['machine', 'machine', 'machine', 'machine']
['machine', 'machine', 'machine', 'machine', 'simple_machine']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto', 'automobil']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto', 'automobil']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto', 'automobil']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto', 'automobil']
['machine', 'machine', 'machine', 'machine', 'simple_machine', 'machine', 'political_machine', 'car', 'auto', 'automobil']
```

```
1 #stemming
2 from nltk.stem import PorterStemmer
3
```

```
1 stemmer=PorterStemmer()  
2 print(stemmer.stem("coding"))  
3 print(stemmer.stem("coded"))
```

```
code  
code
```

Experiment 5

Aim : To study N-gram

▼ CODE:

```
1 import re
2 from nltk.util import ngrams
3 S= "machine learning is an important part of AI""and AI is going to become important for daily functioning"
4 tokens = [token for token in S.split(" ")]
5 output=list(ngrams(tokens,2))
6 print("Generated Output.\n",output,"\n")

Generated Output.
[('machine', 'learning'), ('learning', 'is'), ('is', 'an'), ('an', 'important'), ('important', 'part'), ('part', 'of'),
```


Experiment 6


Aim: To study chunking

CODE:

```
1 import nltk
2 nltk.download('averaged_perceptron_tagger')
3 nltk.download('punkt')

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] /root/nltk_data...
[nltk_data] Package averaged_perceptron_tagger is already up-to-
[nltk_data] date!
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data] Package punkt is already up-to-date!
True
```

```
1 text=nltk.word_tokenize("And now for Everything completely Same")
2 print("Generated Output:\n")
3 nltk.pos_tag(text)
```

 Generated Output:

```
[('And', 'CC'),
 ('now', 'RB'),
 ('for', 'IN'),
 ('Everything', 'VBG'),
 ('completely', 'RB'),
 ('Same', 'JJ')]
```

Experiment 7

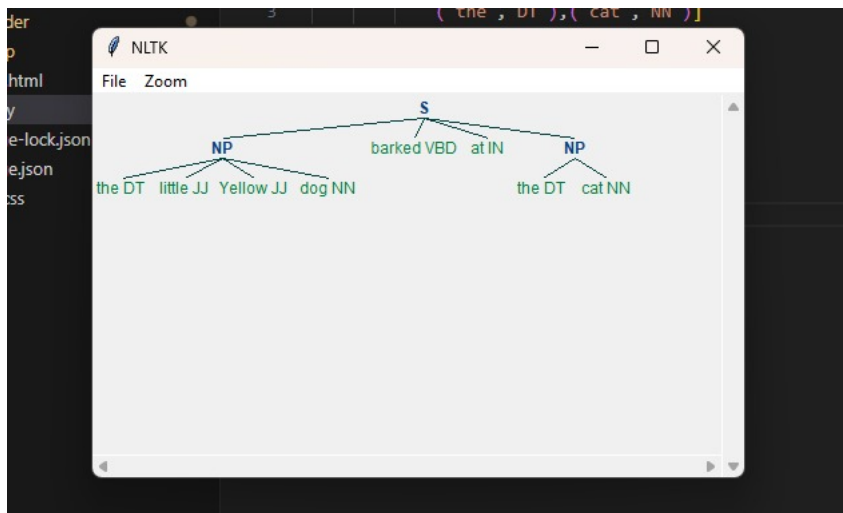
Aim:- To study Chunking

CODE:

```
1 import nltk
2 sentence = [("the","DT"),("little","JJ"),("Yellow","JJ"),("dog","NN"),("barked","VBD"),("at","IN"),
3             ("the","DT"),("cat","NN")]
4 grammar="NP:{<DT>?<JJ>*<NN>}"
5 cp=nltk.RegexpParser(grammar)
6 result= cp.parse(sentence)
7 print("Output :=>")
8 print(result)
```

Output :=>
(S
 (NP the/DT little/JJ Yellow/JJ dog/NN)
 barked/VBD
 at/IN
 (NP the/DT cat/NN))

```
1 result.draw()
```




Experiment 8

Aim:To study Named Entity Recognition

▼ CODE:

```
1 locs=[('Omnicom','IN','New York'),
2       ('DDB Needham','IN','New York'),
3       ('Kalpan Thaler Group','IN','New York'),
4       ('BBDO','IN','Atlanta'),
5       ('Georgia-pacific','IN','Atlanta')]
6 query=[e1 for(e1,rel,e2)in locs if e2=='Atlanta']
7 print("OUTPUT\n",query)
```

 OUTPUT
['BBDO', 'Georgia-pacific']

Experiment 9

Aim: Finite State Automata

a) Define Grammar Using NLTK . Analyze a sentence using the same.

CODE:

```
1 import nltk
2 nltk.download('punkt')
3 from nltk import tokenize
4 grammar1 = nltk.CFG.fromstring("""
5     S -> VP
6     VP -> VP NP
7     NP -> Det NP
8     Det -> 'that'
9     NP -> singular Noun
10    NP -> 'flight'
11    VP -> 'Book'
12    """)
13 sentence= "Book that flight"
14
15 for index in range(len(sentence)):
16     all_tokens=tokenize.word_tokenize(sentence)
17     print(all_tokens)
18
19 parser =nltk.ChartParser(grammar1)
20 for tree in parser.parse(all_tokens):
21     print(tree)
22     tree.draw()
```

[nltk_data] Downloading package punkt to /root/nltk_data...

[nltk_data] Unzipping tokenizers/punkt.zip.

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

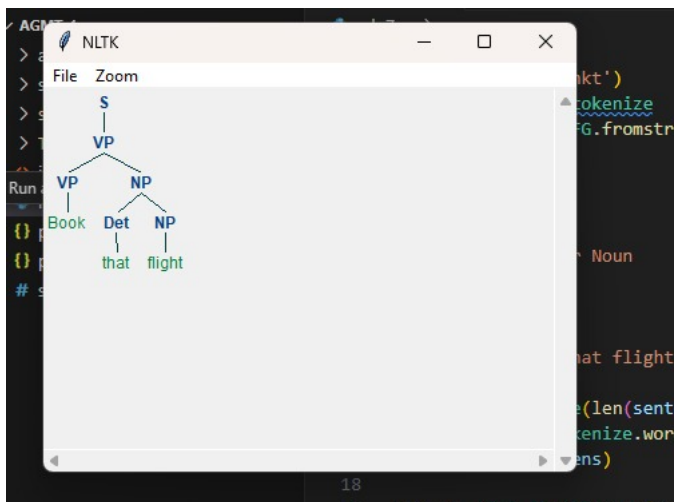
['Book', 'that', 'flight']

['Book', 'that', 'flight']

['Book', 'that', 'flight']

(S (VP (VP Book) (NP (Det that) (NP flight))))

OUTPUT



b) Accept the input string with regular expression of Finite Automata:

CODE:

```

1  def FA(s):
2  #if the lenght is less than 3 it can't be accepted. Therefore program stops.
3      if len(s)<3:
4          return "Rejected"
5  #first 3 characters are fixed then check them using index
6      if s[0]=='1':
7          if s[1]=='0':
8              if s[2]=='1':
9                  for i in range(3,len(s)):
10                     if s[i]!='1':
11                         return "Rejected"
12                     return "Accepted" #if all 4 nested are true
13                 return "Rejected" #else of 3rd if
14             return "Rejected"
15         return "Rejected"
16 inputs=['1','10101','101','10111',
17         '01010','100','',
18         '10111101','1011111']
19 print("output\n")
20 for i in inputs:
21     print(FA(i))
22

```

output

```

Rejected
Rejected
Accepted
Accepted
Rejected
Rejected
Rejected
Rejected
Accepted

```

c) Accept the input string with regular expression pf FA:(a+b)* bba.

CODE:

```

1  def FA(s):
2      size=0
3      for i in s:
4          if i=='a' or i=='b':
5              size+=1
6          else:
7              return "Rejected"
8      if size>=3: #after checking it only contains a and b only.& checking it lenght should be atleast 3
9          if s[size-3]=='b':
10             if s[size-2]=='b':
11                 if s[size-1]=='a':
12                     return "Accepted"
13                 return "Rejected"
14             return "Rejected"
15         return "Rejected"
16     return "Rejected"
17 inputs=['bba','ababbb','abb','baba','bbb',]
18 print("OUTPUT\n")
19 for i in inputs:
20     print(FA(i))

```

OUTPUT

```

Accepted
Accepted
Rejected
Rejected
Rejected

```

d) Implemntation of Deductive chart parsing using context free grammar and a given sentence.

CODE:

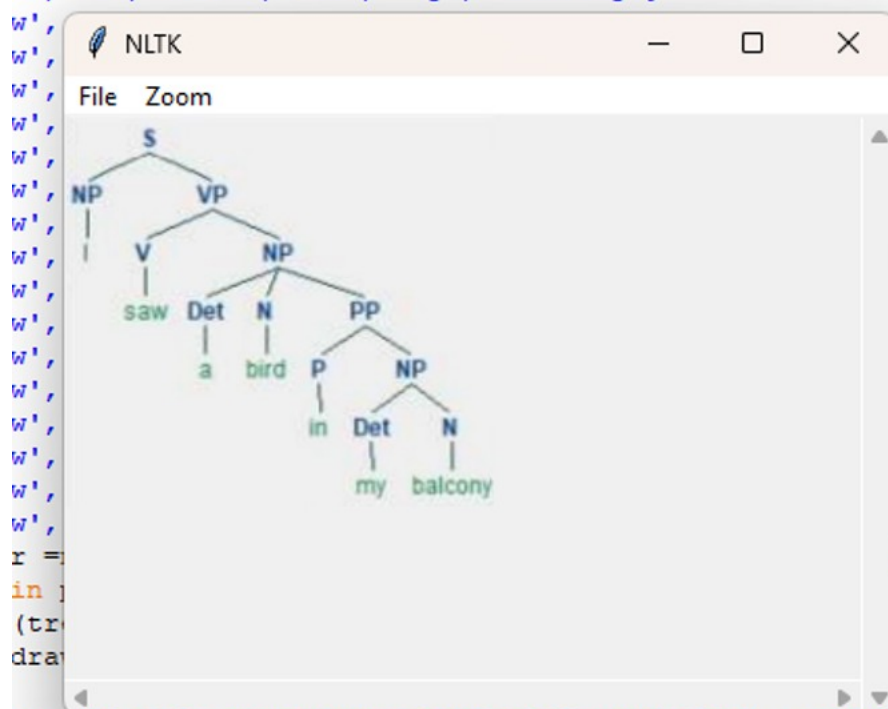
```

1 # import nltk
2 nltk.download('punkt')
3 # from nltk import tokenize
4 grammar2 = nltk.CFG.fromstring("""
5     S -> NPVP
6     PP -> PNP

```

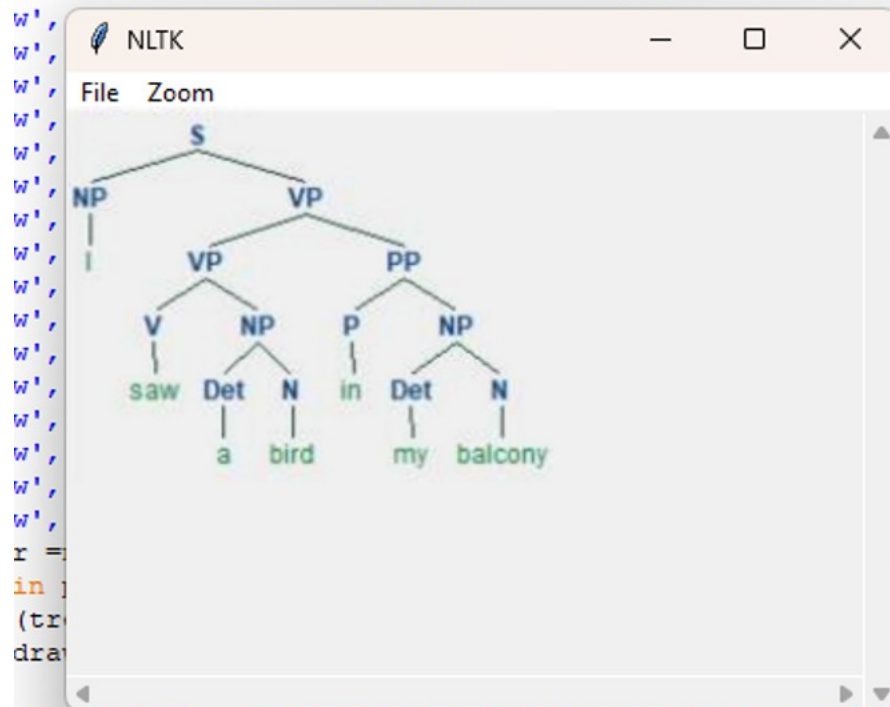
```
7 NP -> DetN|DetNP|'I'  
8 VP -> VNP|VPPP  
9 Det ->'a'|'my'  
10 N -> 'bird'|'balcony'  
11 v ->'saw'  
12 P ->'in'  
13 """)  
14 sentences= "I saw a bird in my balcony"  
15  
16 for index in range(len(sentences)):  
17     all_tokens1=tokenize.word_tokenize(sentences)  
18     print(all_tokens1)  
19  
20 parser1 =nltk.ChartParser(grammar2)  
21 for tree in parser1.parse(all_tokens1):  
22     print(tree)  
23     tree.draw()  
  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']  
[nltk_data] Downloading package punkt to /root/nltk_data...  
[nltk_data] Package punkt is already up-to-date!
```

```
w', 'a', 'bird', 'in', 'my', 'balcony']
```



```
or: multiple statements found while compiling a single statement
r = nltk.ChartParser(grammar1)
ree in parser.parse(all tokens):
```

```
w', 'a', 'bird', 'in', 'my', 'balcony']
```



```
or: multiple statements found while compiling a single statemen
```

```
r = nltk.ChartParser(grammar1)
```

```
ree in parser.parse(all_tokens):
```

Experiment 10

AIM: Study PorterStemmer, LancasterStemmer, RegexpStemmer, SnowballStemmer, WordNetLemmatizer

CODE:

```
1 import nltk
```

```
1 #PorterStemmer
2 from nltk.stem import PorterStemmer
3 word_stemmer=PorterStemmer()
4 print("OUTPUT=>\n")
5 print(word_stemmer.stem('Coding'))
```

OUTPUT=>

code

```
1 #LancasterStemmer
2 from nltk.stem import LancasterStemmer
3 Lanc_stemmer=LancasterStemmer()
4 print("OUTPUT=>\n")
5 print(Lanc_stemmer.stem('Coding'))
```

OUTPUT=>

cod

```
1 #RegexpStemmer
2 from nltk.stem import RegexpStemmer
3 st = RegexpStemmer('ing$|s$|e$|able$', min=4)
4 print("OUTPUT=>")
5 print(st.stem('coding'))
```

OUTPUT=>

cod

```
1 #SnowballStemmer
2 from nltk.stem import SnowballStemmer
3 sk_stemmer=SnowballStemmer('english')
4 print("output=>")
5 print(sk_stemmer.stem('coding'))
```

output=>

code

```
1 nltk.download('wordnet')
```

[nltk_data] Downloading package wordnet to /root/nltk_data...
True

```
1 #WordNetLemmatizer
2 from nltk.stem import WordNetLemmatizer
3 lemmatizer=WordNetLemmatizer()
4 print("word :\t Lemma")
5 print("rocks:", lemmatizer.lemmatize("rocks"))
6 print("Corpora:", lemmatizer.lemmatize("Corpora"))
7 print("better:", lemmatizer.lemmatize("better",pos="a"))
```

word : Lemma
rocks: rock
Corpora: Corpora
better: good