

NATURAL LANGUAGE PROCESSING

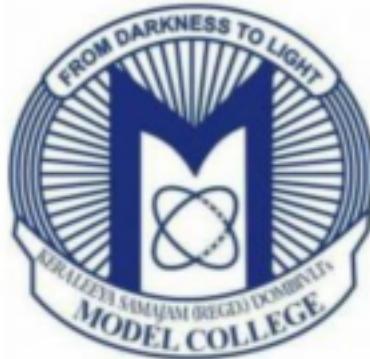
Certified Journal

**Submitted in partial fulfilment of the
Requirements for the award of the Degree of**

**MASTER OF SCIENCE
(INFORMATION TECHNOLOGY)**

By

Anjali Rameshwar Nimje



DEPARTMENT OF INFORMATION TECHNOLOGY

**KERALEYA SAMAJAM (REGD.) DOMBIVLI'S
MODEL COLLEGE (AUTONOMOUS)**
Re-Accredited 'A' Grade by NAAC

(Affiliated to University of Mumbai)

FOR THE YEAR

(2023-24)



Keraleeya Samajam(Regd.) Dombivli's
MODEL COLLEGE
Re-Accredited Grade "A" by NAAC

Kanchan Goan Village, Khambalpada, Thakurli East – 421201
Contact No – 7045682157, 7045682158. www.model-college.edu.in

**DEPARTMENT OF INFORMATION TECHNOLOGY
AND COMPUTER SCIENCE**

CERTIFICATE

This is to certify that Mr. /Miss _____

Studying in Class _____ Seat No. _____

Has completed the prescribed practicals in the subject _____

During the academic year _____

Date : _____

External Examiner

Internal Examiner
M.Sc. Information Technology

INDEX

Sr No	Practical Name	Date	Signature
1A	Install NLTK		
1B	Convert the given text to speech		
1C	Convert audio file Speech to Text.		
2A	Study of various Corpus – Brown, Inaugural, Reuters, udhr with various methods like fields, raw, words, sents, categories,		
2C	Study Conditional frequency distributions		
2D	Study of tagged corpora with methods like tagged_sents, tagged_words.		
2E	Write a program to find the most frequent noun tags.		
2F	Map Words to Properties Using Python Dictionaries		
2G	Study DefaultTagger, Regular expression tagger, UnigramTagger		
3A	Study of Wordnet Dictionary with methods as synsets, definitions, examples,antonyms.		
3B	Study lemmas, hyponyms, hypernyms.		

3C	Write a program using python to find synonym and antonym of word "active" using Wordnet.		
3D	Compare two nouns.		
3E	Handling stopword		
4A	Tokenization using Python's split() function.		
4B	Tokenization using Regular Expressions (RegEx)		
4C	Tokenization using NLTK		
4D	Tokenization using the spaCy library		
4E	Tokenization using Keras		
4F	Tokenization using Gensim		
6A	Illustrate part of speech tagging Part of speech Tagging and chunking of user defined text.		

6B	Named Entity recognition using user defined text.		
6C	Named Entity recognition with diagram using NLTK corpus – treebank.		
7A	Finite state automata Define grammar using nltk. Analyze a sentence using the same.		
7B	Accept the input string with Regular expression of Finite Automaton: 101^+.		
7C	Accept the input string with Regular expression of FA: $(a+b)^*bba$.		
7D	Implementation of Deductive Chart Parsing using context free grammar and a given sentence.		
8	Study PorterStemmer, LancasterStemmer, RegexpStemmer, SnowballStemmer ,Study WordNetLemmatizer		
9	Implement Naive Bayes classifier		

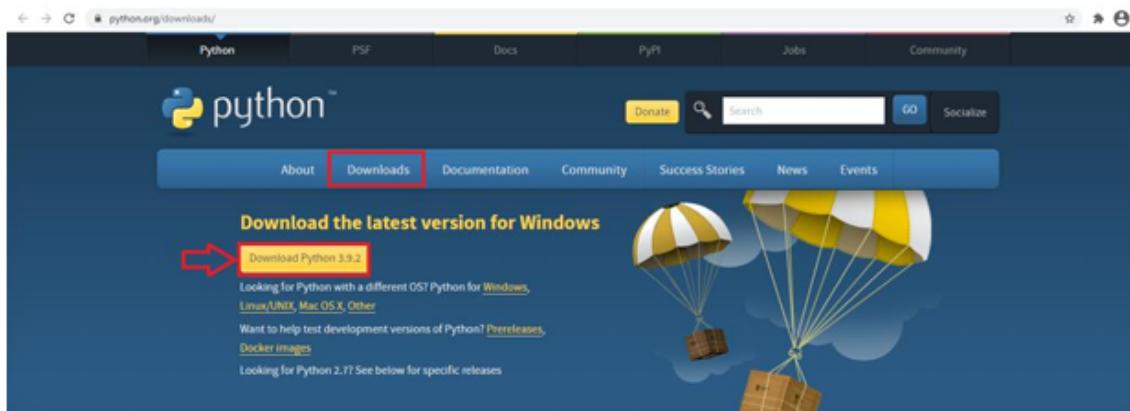
10 A.1	Speech tagging using spacy		
10 A.2	Speech tagging using nltk		
10 B.1	Usage of Give and Gave in the Penn Treebank sample		
10 B.2	Probabilistic parser		
10 C	Malt Parsing		
11 A	Multiword Expressions in NLP		
11 B	Normalized Web Distance and Word Similarity		
11 C	Word Sense Disambiguation		

PRACTICAL 1A

Install NLTK

Python 3.9.2 Installation on Windows

Step 1) Go to link <https://www.python.org/downloads/>, and select the latest version for windows.

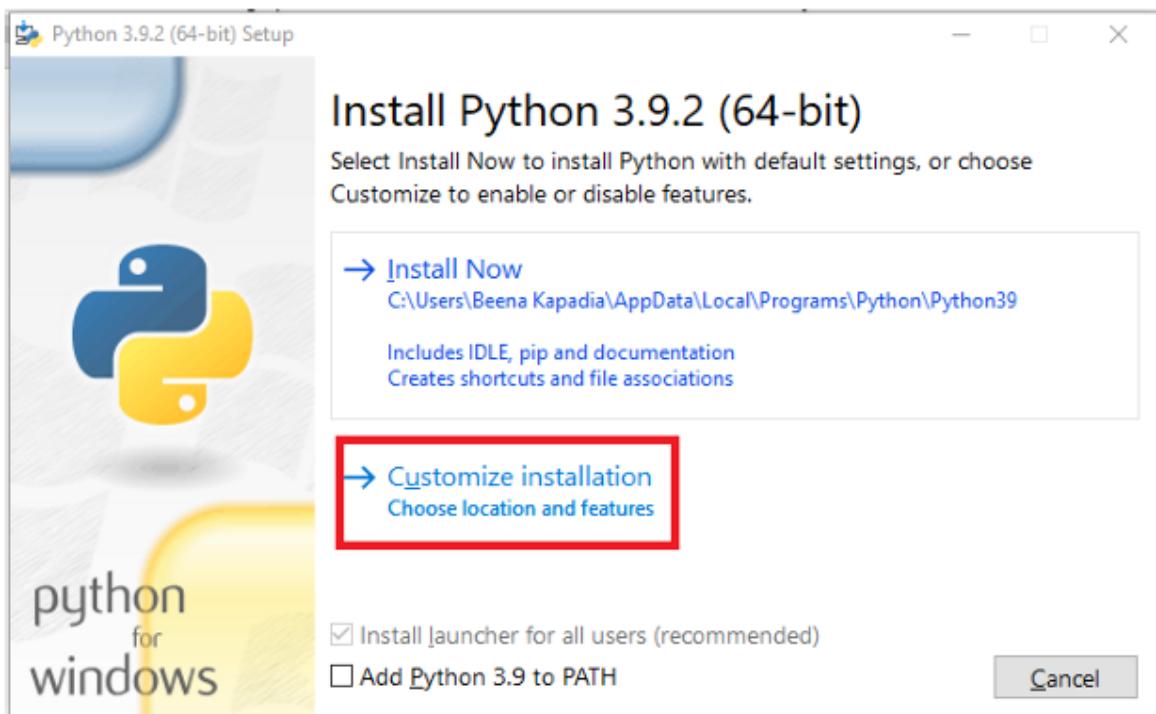


Note: If you don't want to download the latest version, you can visit the download tab and see all releases.

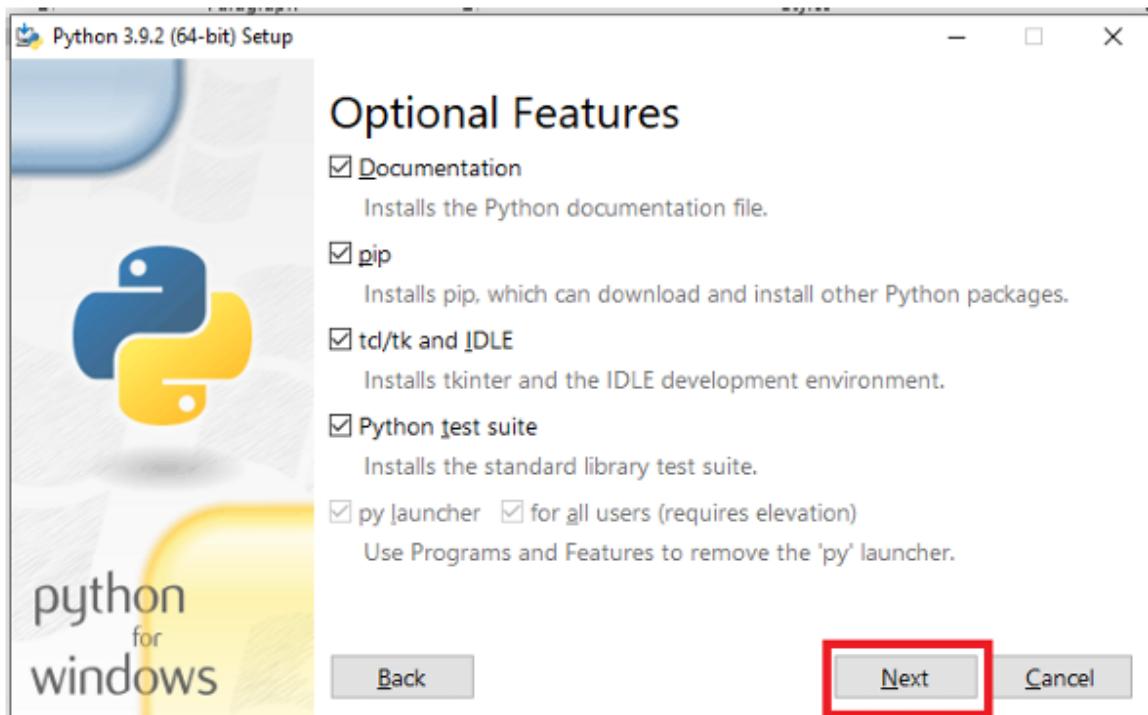
A screenshot of a web browser displaying the Python.org Downloads page for Python 3.9.2. The URL in the address bar is 'python.org/downloads/release/python-392/'. The page title is 'Download the latest version for Windows'. Below the title, there is a table titled 'Files' showing various download options. A red arrow points to the 'Windows Installer (64-bit)' link, which is highlighted with a red box. The table includes columns for Version, Operating System, Description, MD5 Sum, File Size, and GPG. Other links shown in the table include 'Gzipped source tarball', 'XZ compressed source tarball', 'macOS 64-bit Intel installer', 'macOS 64-bit universal2 installer', 'Windows embeddable package (32-bit)', 'Windows embeddable package (64-bit)', and 'Windows help file'. At the bottom of the page, there is a navigation bar with links for About, Downloads, Documentation, Community, Success Stories, News, Applications, Docs, Community Survey, and Python News. A file icon followed by 'python-3.9.2-amd64.exe' is visible at the bottom left, and a 'Show all' link is at the bottom right.

Step 2) Click on the Windows installer (64 bit)

Step 3) Select Customize Installation



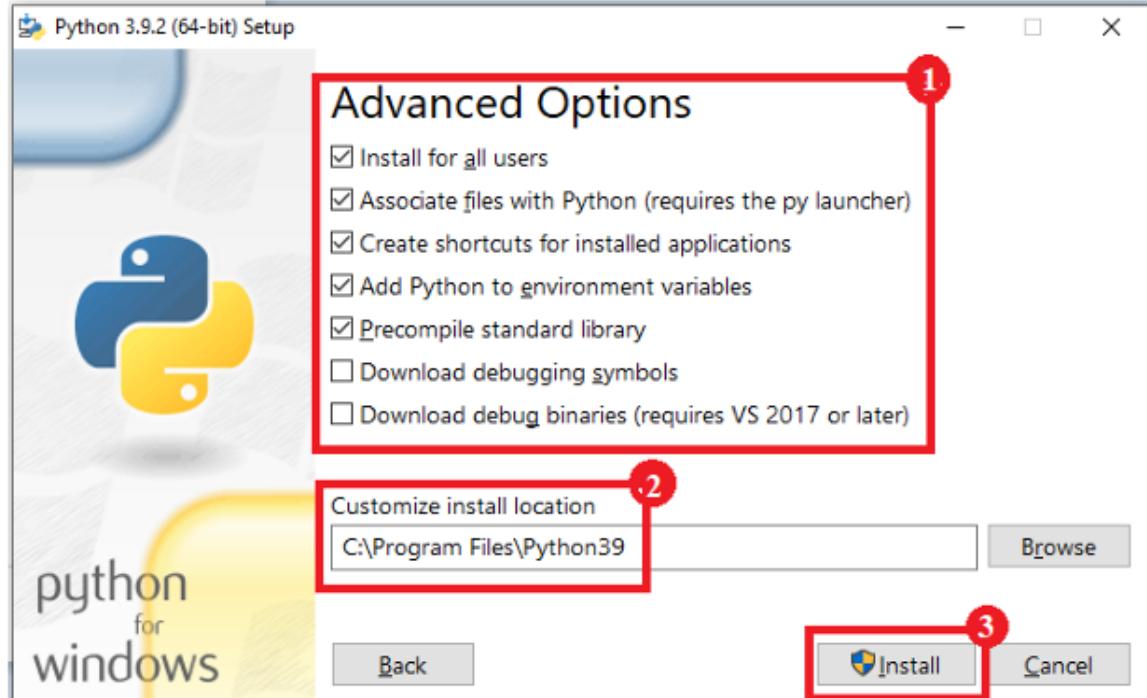
Step 4) Click NEXT



Step 5) In next screen

1. Select the advanced options

2. Give a Custom install location. Keep the default folder as c:\Program files\Python39
3. Click Install



Step 6) Click Close button once install is done.

Step 7) open command prompt window and run the following commands:

```
C:\Users\Beena Kapadia>pip install --upgrade pip
```

```
C:\Users\Beena Kapadia> pip install --user -U nltk
```

```
C:\Users\Beena Kapadia>>pip install --user -U numpy
```

```
C:\Users\Beena Kapadia>python
```

```
>>> import nltk
```

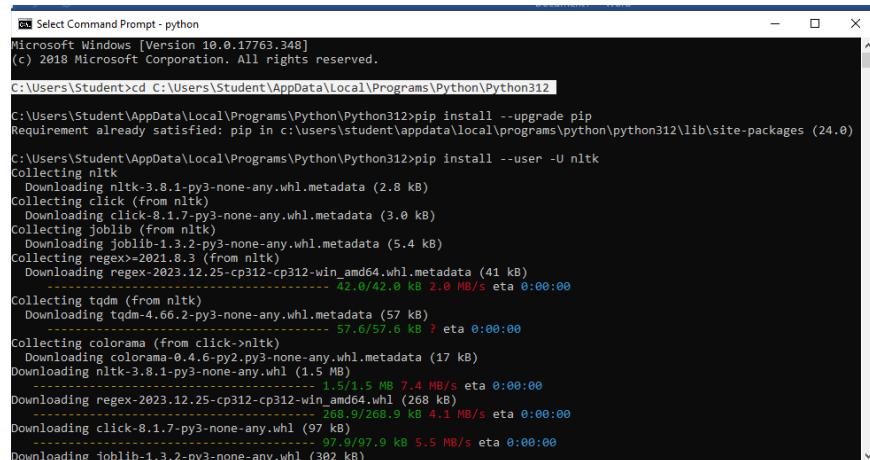
PRACTICAL 1B

Convert the given text to speech.

Code:

```
from playsound import playsound
from gtts import gTTS
mytext="Welcome to Natural Language Processing"
language="en"
myobj=gTTS(text=mytext,lang=language,slow=False)
myobj.save("myfile.mp3")
playsound("myfile.mp3")
```

```
from playsound import playsound
from gtts import gTTS
mytext="Welcome to Natural Language Processing"
language="en"
myobj=gTTS(text=mytext,lang=language,slow=False)
myobj.save("myfile.mp3")
playsound("myfile.mp3")
```



The screenshot shows a Microsoft Windows Command Prompt window titled "Select Command Prompt - python". The window displays the following command and its output:

```
C:\Select Command Prompt - python
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>cd C:\Users\Student\AppData\Local\Programs\Python\Python312>
C:\Users\Student>AppData\Local\Programs\Python\Python312>pip install --upgrade pip
Requirement already satisfied: pip in c:\users\student\appdata\local\programs\python\python312\lib\site-packages (24.0)

C:\Users\Student>AppData\Local\Programs\Python\Python312>pip install --user -U nltk
Collecting nltk
  Downloading nltk-3.8.1-py3-none-any.whl.metadata (2.8 kB)
Collecting click (from nltk)
  Downloading click-8.1.7-py3-none-any.whl.metadata (3.0 kB)
Collecting joblib (from nltk)
  Downloading joblib-1.3.2-py3-none-any.whl.metadata (5.4 kB)
Collecting regex>=2021.8.3 (from nltk)
  Downloading regex-2023.12.25-cp312-cp312-win_amd64.whl.metadata (41 kB)
                                             42.0/42.0 kB 2.0 MB/s eta 0:00:00
Collecting tqdm (from nltk)
  Downloading tqdm-4.66.2-py3-none-any.whl.metadata (57 kB)
                                             57.0/57.6 kB > eta 0:00:00
Collecting colorama (from click->nltk)
  Downloading colorama-0.4.6-py2.py3-none-any.whl.metadata (17 kB)
Downloaded nltk-3.8.1-py3-none-any.whl (1.5 MB)
----- 1.5/1.5 MB 7.4 MB/s eta 0:00:00
Downloaded regex-2023.12.25-cp312-cp312-win_amd64.whl (268 kB)
----- 268.9/268.9 kB 4.1 MB/s eta 0:00:00
Downloaded click-8.1.7-py3-none-any.whl (97 kB)
----- 97.9/97.9 kB 5.5 MB/s eta 0:00:00
Downloaded joblib-1.3.2-py3-none-any.whl (302 kB)
```

```
ca Select Command Prompt - python
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>cd C:\Users\Student\AppData\Local\Programs\Python\Python312>
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --upgrade pip
Requirement already satisfied: pip in c:\users\student\appdata\local\programs\python\python312\lib\site-packages (24.0)

C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --user -U nltk
Collecting nltk
  Downloading nltk-3.8.1-py3-none-any.whl.metadata (2.8 kB)
Collecting click (from nltk)
  Downloading click-8.1.7-py3-none-any.whl.metadata (3.0 kB)
Collecting joblib (from nltk)
  Downloading joblib-1.3.2-py3-none-any.whl.metadata (5.4 kB)
Collecting regex<=2021.8.3 (from nltk)
  Downloading regex-2023.12.25-cp312-cp312-win_amd64.whl.metadata (41 kB)
    42.0/42.0 kB 2.0 MB/s eta 0:00:00

Collecting tqdm (from nltk)
  Downloading tqdm-4.66.2-py3-none-any.whl.metadata (57 kB)
    57.6/57.6 kB ? eta 0:00:00

Collecting colorama (from click->nltk)
  Downloading colorama-0.4.6-py2.py3-none-any.whl.metadata (17 kB)
Downloaded nltk-3.8.1-py3-none-any.whl (1.5 MB)
  1.5/1.5 MB 7.4 MB/s eta 0:00:00
Downloaded regex-2023.12.25-cp312-cp312-win_amd64.whl (268 kB)
  268.9/268.9 kB 4.1 MB/s eta 0:00:00
Downloaded click-8.1.7-py3-none-any.whl (97 kB)
  97.9/97.9 kB 5.5 MB/s eta 0:00:00
Downloaded joblib-1.3.2-py3-none-any.whl (302 kB)
```

```
ca Select Command Prompt - python
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>cd C:\Users\Student\AppData\Local\Programs\Python\Python312>
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --upgrade pip
Requirement already satisfied: pip in c:\users\student\appdata\local\programs\python\python312\lib\site-packages (24.0)

C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --user -U nltk
Collecting nltk
  Downloading nltk-3.8.1-py3-none-any.whl.metadata (2.8 kB)
Collecting click (from nltk)
  Downloading click-8.1.7-py3-none-any.whl.metadata (3.0 kB)
Collecting joblib (from nltk)
  Downloading joblib-1.3.2-py3-none-any.whl.metadata (5.4 kB)
Collecting regex<=2021.8.3 (from nltk)
  Downloading regex-2023.12.25-cp312-cp312-win_amd64.whl.metadata (41 kB)
    42.0/42.0 kB 2.0 MB/s eta 0:00:00

Collecting tqdm (from nltk)
  Downloading tqdm-4.66.2-py3-none-any.whl.metadata (57 kB)
    57.6/57.6 kB ? eta 0:00:00

Collecting colorama (from click->nltk)
  Downloading colorama-0.4.6-py2.py3-none-any.whl.metadata (17 kB)
Downloaded nltk-3.8.1-py3-none-any.whl (1.5 MB)
  1.5/1.5 MB 7.4 MB/s eta 0:00:00
Downloaded regex-2023.12.25-cp312-cp312-win_amd64.whl (268 kB)
  268.9/268.9 kB 4.1 MB/s eta 0:00:00
Downloaded click-8.1.7-py3-none-any.whl (97 kB)
  97.9/97.9 kB 5.5 MB/s eta 0:00:00
Downloaded joblib-1.3.2-py3-none-any.whl (302 kB)
  302.2/302.2 kB 6.2 MB/s eta 0:00:00
Downloaded tqdm-4.66.2-py3-none-any.whl (78 kB)
  78.3/78.3 kB 1.4 MB/s eta 0:00:00
Downloaded colorama-0.4.6-py2.py3-none-any.whl (25 kB)
Installing collected packages: regex, joblib, colorama, tqdm, click, nltk
  WARNING: The script tqdm.exe is installed in 'C:\Users\Student\AppData\Roaming\Python\Python312\Scripts' which is not on PATH.
  Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
  WARNING: The script nltk.exe is installed in 'C:\Users\Student\AppData\Roaming\Python\Python312\Scripts' which is not on PATH.
```

```
  Select Command Prompt - python
Downloaded tqdm-4.66.2-py3-none-any.whl (78 kB)   302.2/302.2 kB 6.2 MB/s eta 0:00:00
Downloaded colorama-0.4.6-py2.py3-none-any.whl (25 kB) 78.3/78.3 kB 1.4 MB/s eta 0:00:00
Installing collected packages: regex, joblib, colorama, tqdm, click, nltk
  WARNING: The script tqdm.exe is installed in 'C:\Users\Student\AppData\Roaming\Python\Python312\Scripts' which is not on PATH.
    Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
  WARNING: The script nltk.exe is installed in 'C:\Users\Student\AppData\Roaming\Python\Python312\Scripts' which is not on PATH.
    Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
Successfully installed click-8.1.7 colorama-0.4.6 joblib-1.3.2 nltk-3.8.1 regex-2023.12.25 tqdm-4.66.2

C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --user -U numpy
Collecting numpy
  Downloading numpy-1.26.4-cp312-cp312-win_amd64.whl.metadata (61 kB)
    ..... 61.0/61.0 kB 406.4 kB/s eta 0:00:00
  Downloading numpy-1.26.4-cp312-cp312-win_amd64.whl (15.5 MB)
    ..... 15.5/15.5 kB 3.9 MB/s eta 0:00:00
Installing collected packages: numpy
  WARNING: The script f2py.exe is installed in 'C:\Users\Student\AppData\Roaming\Python\Python312\Scripts' which is not on PATH.
    Consider adding this directory to PATH or, if you prefer to suppress this warning, use --no-warn-script-location.
Successfully installed numpy-1.26.4

C:\Users\Student\AppData\Local\Programs\Python\Python312>python
Python 3.12.2 (tags/v3.12.2:6abdd9d, Feb 6 2024, 21:26:36) [MSC v.1937 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.

>>> pip install gttcs
  File "<stdin>", line 1
    pip install gttcs
      ^^^^^^

SyntaxError: invalid syntax
>>> exit
Use exit() or Ctrl-Z plus Return to exit
>>>
```

```
PS C:\Select Command Prompt
Type "help", "copyright", "credits" or "license" for more information.
>>> pip install gttts
  File "", line 1
    pip install gttts
      ^^^^^^
SyntaxError: invalid syntax
>>> exit()
KeyboardInterrupt) on Ctrl-Z plus Return to exit
>>> exit()

c:\Users\Student\AppData\Local\Programs\Python\Python312> pip install gttts
Collecting gttts
  Downloading gttts-2.5.1-py3-none-any.whl.metadata (4.1 kB)
Collecting requests<3,>=2.27 (from gttts)
  Downloading requests-2.29.1-py3-none-any.whl.metadata (4.6 kB)
Requirement already satisfied: click<8.2,>=7.3 in c:\users\student\appdata\local\programs\python\python312\site-packages (from gttts) (8.1.7)
Requirement already satisfied: colorama in c:\users\student\appdata\roaming\python\python312\site-packages (from click<8.2,>=7.1>gttts) (0.4.6)
Collecting charset-normalizer<4,>=2. (from requests<3,>=2.27>gttts)
  Downloading charset_normalizer-3.3.2-py312-cp312-win_amd64.whl.metadata (34 kB)
Collecting idna<3.6,>=2.1.1 (from requests<3,>=2.27>gttts)
  Downloading idna-3.6-py3-none-any.whl.metadata (9.9 kB)
Collecting urllib3<3,>=1.21.1 (from requests<3,>=2.27>gttts)
  Downloading urllib3-2.2.0-py3-none-any.whl.metadata (16.6 kB)
Collecting certifi>2017.4.15 (from requests<3,>=2.27>gttts)
  Downloading certifi-2024.2.2-py3-none-any.whl.metadata (2.2 kB)
Collecting requests<2.31.0,>=0-py3-none-any.whl (29 kB)
  Downloading gttts-2.5.1-py3-none-any.whl (29 kB)
  Downloading requests-2.31.0-py3-none-any.whl (62 kB)
  Downloading requests-2.31.0-py3-none-any.whl (163 kB)
----- 163/163.8 kB 10.2 kB/s eta 0:00:00
  Downloading charset_normalizer-3.3.2-py312-win_amd64.whl (106 kB)
----- 106/106.4 kB 5.0 kB/s eta 0:00:00
  Downloading idna-3.6-py3-none-any.whl (61 kB)
----- 61/61.6 kB ? eta 0:00:00
  Downloading urllib3-2.2.0-py3-none-any.whl (121 kB)
----- 121/121.1 kB 6.0 kB/s eta 0:00:00
Installing collected packages: urllib3, idna, charset-normalizer, certifi, requests, gttts
Successfully installed certifi-2024.2.2 charset-normalizer-3.3.2 gttts-2.5.1 idna-3.6 requests-2.31.0 urllib3-2.2.1

C:\Users\Student\AppData\Local\Programs\Python\Python312>
```

```
□ Select Command Prompt
Downloaded urllib3-2.2.1-py3-none-any.whl (121 kB)
Collecting certifi> https://files.pythonhosted.org/packages/3d/2c/certifi-2024.2.2.dist-info/RECORD certifi-2024.2.2.dist-info/METADATA certifi-2024.2.2.dist-info/WHEEL certifi-2024.2.2.dist-info/top_level.txt certifi-2024.2.2.dist-info/requires.txt certifi-2024.2.2.dist-info/namespace_packages.txt certifi-2024.2.2.dist-info/py.typed certifi-2024.2.2.dist-info/manifest.in 1 kB 6.9 MB/s eta 0:00:00
Installing collected packages: urllib3, idna, charset-normalizer, certifi, requests, gttcs
Successfully installed certifi-2024.2.2 charset-normalizer-3.3.2 gttcs-2.5.1 idna-3.6 requests-2.31.0 urllib3-2.2.1

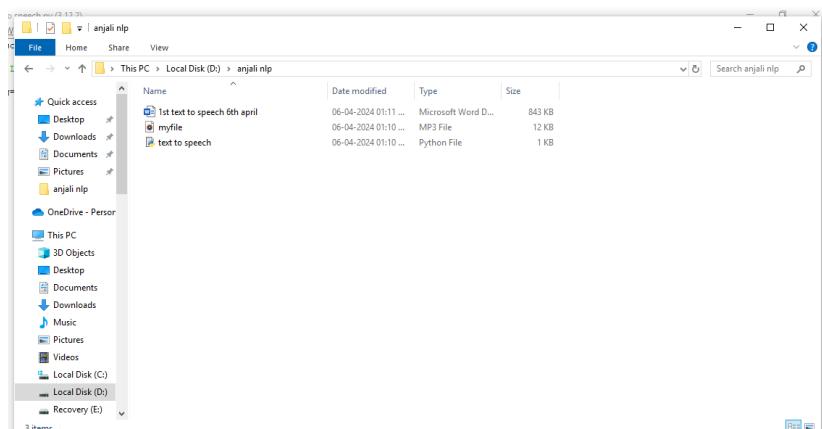
C:\Users\Student\AppData\Local\Programs\Python\Python312\pip install --upgrade setuptools
Collecting setuptools
  Downloading setuptools-69.2.0-py3-none-any.whl.metadata (6.3 kB)
Downloaded setuptools-69.2.0-py3-none-any.whl (821 kB)
Collecting certifi> https://files.pythonhosted.org/packages/3d/2c/certifi-2024.2.2.dist-info/RECORD certifi-2024.2.2.dist-info/METADATA certifi-2024.2.2.dist-info/WHEEL certifi-2024.2.2.dist-info/top_level.txt certifi-2024.2.2.dist-info/requires.txt certifi-2024.2.2.dist-info/namespace_packages.txt certifi-2024.2.2.dist-info/py.typed certifi-2024.2.2.dist-info/manifest.in 1 kB 821.5/821.5 kB 10.4 MB/s eta 0:00:00
Installing collected packages: setuptools
Successfully installed setuptools-69.2.0

C:\Users\Student\AppData\Local\Programs\Python\Python312>
```

```
PS Select Command Prompt
Downloading setuptools-69.2.0-py3-none-any.whl (821 kB)  021.5/821.5 kB 18.4 MB/s eta 0:00:00
Installing collected packages: setuptools
Successfully installed setuptools-69.2.0
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --upgrade wheel
Collecting wheel
  Downloading wheel-0.43.0-py3-none-any.whl.metadata (2.2 kB)
  Downloading wheel-0.43.0-py3-none-any.whl (65 kB)  05.0/65.8 kB 7 eta 0:00:00
Installing collected packages: wheel
Successfully installed wheel-0.43.0
C:\Users\Student\AppData\Local\Programs\Python\Python312>
```

```
PS Select Command Prompt
Installing collected packages: urllib3, idna, charset-normalizer, certifi, requests, gits
Successfully installed certifi-2024.2.2 charset-normalizer-3.3.2 gits-2.5.1 idna-3.6 requests-2.31.0 urllib3-2.2.1
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --upgrade setuptools
Collecting setuptools
  Downloading setuptools-69.2.0-py3-none-any.whl.metadata (6.3 kB)
  Downloading setuptools-69.2.0-py3-none-any.whl (921 kB)  021.5/921.5 kB 18.4 MB/s eta 0:00:00
Installing collected packages: setuptools
Successfully installed setuptools-69.2.0
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install --upgrade wheel
Collecting wheel
  Downloading wheel-0.43.0-py3-none-any.whl.metadata (2.2 kB)
  Downloading wheel-0.43.0-py3-none-any.whl (65 kB)  05.0/65.8 kB 7 eta 0:00:00
Installing collected packages: wheel
Successfully installed wheel-0.43.0
C:\Users\Student\AppData\Local\Programs\Python\Python312>pip install playsound
Collecting playsound
  Downloading playsound-1.3.0.tar.gz (7.7 kB)
  Preparing metadata (setup.py) ... done
Building wheels for collected packages: playsound
  Building wheel for playsound (setup.py) ... done
    Created wheel for playsound: filename=playsound-1.3.0-py3-none-any.whl size=7043 sha256=77f8177abc79d31786d4abc3581d0878ef7587800167b99d06c3e2eb2dfce0
  Stored in directory: C:\Users\Student\AppData\Local\pip\Cache\wheels\cf\42\f8\7c587bae55ecc0b905ca316b258db4daedbfb272a3cbe8907
Successfully built playsound
Installing collected packages: playsound
Successfully installed playsound-1.3.0
C:\Users\Student\AppData\Local\Programs\Python\Python312>
```

Output:



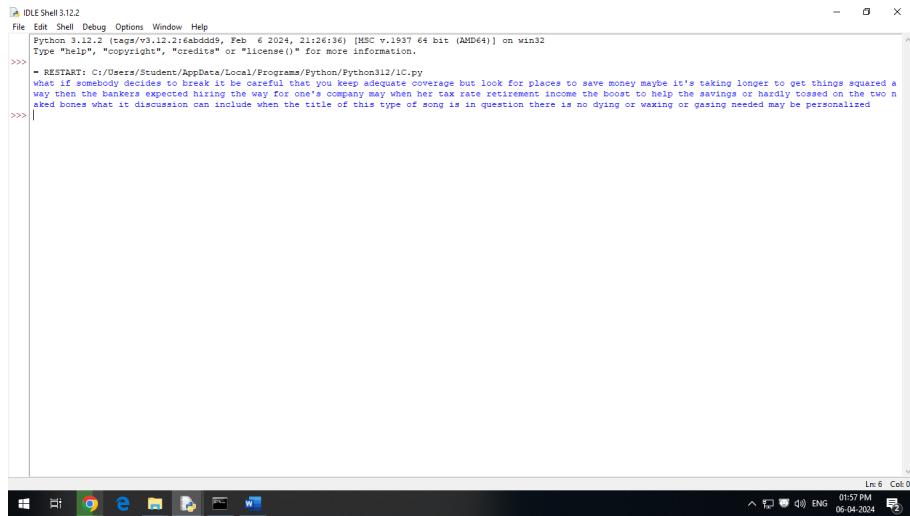
PRACTICAL 1C

Convert audio file Speech to Text.

Code:

```
import speech_recognition as sr
filename="male.wav"
r=sr.Recognizer()
with sr.AudioFile(filename)as source:
    audio_data=r.record(source)
    text=r.recognize_google(audio_data)
    print(text)
```

Output:



The screenshot shows the Python IDLE Shell 3.12.2 interface. The window title is "IDLE Shell 3.12.2". The menu bar includes File, Edit, Shell, Debug, Options, Window, and Help. The status bar at the bottom right shows "Line: 6 Col: 0", "01:57 PM", "ENG", and the date "06-04-2024". The main text area displays the following output:

```
File Edit Shell Debug Options Window Help
Python 3.12.2 (tags/v3.12.2:6abddd9, Feb  6 2024, 21:26:36) [MSC v.1937 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>> - RESTART: C:/Users/Student/AppData/Local/Programs/Python/Python312/1C.py
what if somebody decides to break it be careful that you keep adequate coverage but look for places to save money maybe it's taking longer to get things squared a
way then the bankers expected hiring the way for one's company may when her tax rate retirement income the boost to help the savings or hardly tossed on the two n
aked bones what it discussion can include when the title of this type of song is in question there is no dying or waxing or gassing needed may be personalized
```

PRACTICAL 2A

Study of various Corpus – Brown, Inaugural, Reuters, udhr with various methods like fields, raw, words, sents, categories,

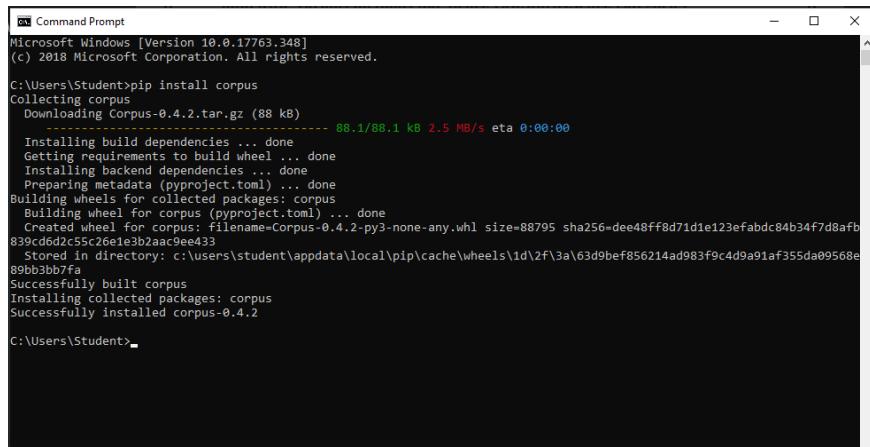
Code:

```
import nltk
from nltk.corpus import brown
print ('File ids of brown corpus\n',brown.fileids())
""Let's pick out the first of these texts — Emma by Jane Austen — and give it a short name, emma, then find out how many words it contains:""
ca01 = brown.words('ca01')
# display first few words
print("\nca01 has following words:\n",ca01)
# total number of words in ca01
print("\nca01 has",len(ca01),'words')

#categories or files
print ('\n\nCategories or file in brown corpus:\n')
print (brown.categories())
""display other information about each text, by looping over all the values of fileid corresponding to the brown file identifiers listed earlier and then computing statistics for each text.""
print ('\n\nStatistics for each text:\n')
print
('AvgWordLen\tAvgSentenceLen\tno.ofTimesEachWordAppearsOnAvg\t\tFileName')
```

```
for fileid in brown.fileids():
```

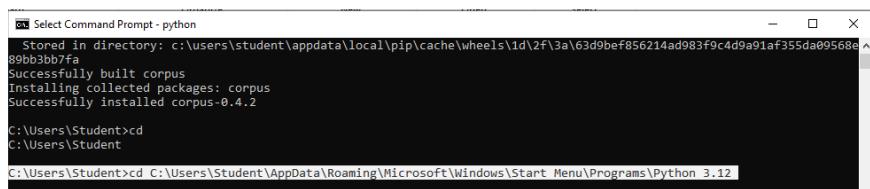
```
    num_chars = len(brown.raw(fileid))
    num_words = len(brown.words(fileid))
    num_sents = len(brown.sents(fileid))
    num_vocab = len(set([w.lower() for w in brown.words(fileid)]))
    print(int(num_chars/num_words),'\t\t\t', int(num_words/num_sents),'\t\t\t',
          int(num_words/num_vocab),'\t\t\t', fileid)
```



```
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>pip install corpus
Collecting corpus
  Downloading Corpus-0.4.2.tar.gz (88 kB)
     88.1/88.1 kB 2.5 MB/s eta 0:00:00
    Installing build dependencies ... done
    Getting requirements to build wheel ... done
    Installing backend dependencies ... done
    Preparing metadata (pyproject.toml) ... done
    Building wheels for collected packages: corpus
      Building wheel for corpus (pyproject.toml) ... done
        Created wheel for corpus: filename=Corpus-0.4.2-py3-none-any.whl size=88795 sha256=dee48ff8d71d1e123efabdc84b34f7d8afb839cd62c5c26e1e3b2aa9ee433
        Stored in directory: c:\users\student\appdata\local\pip\cache\wheels\1d\2f\3a\63d9bef856214ad983f9c4d9a91af355da09568e89bb3bb7fa
    Successfully built corpus
    Installing collected packages: corpus
    Successfully installed corpus-0.4.2

C:\Users\Student>
```



```
Select Command Prompt - python
Stored in directory: c:\users\student\appdata\local\pip\cache\wheels\1d\2f\3a\63d9bef856214ad983f9c4d9a91af355da09568e89bb3bb7fa
Successfully built corpus
Installing collected packages: corpus
Successfully installed corpus-0.4.2

C:\Users\Student>cd C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12
C:\Users\Student>
```

```
C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12>pip install nltk
Collecting nltk
  Downloading nltk-3.8.1-py3-none-any.whl.metadata (2.8 kB)
Collecting click (from nltk)
  Downloading click-8.1.7-py3-none-any.whl.metadata (3.0 kB)
Collecting joblib (from nltk)
  Downloading joblib-1.4.0-py3-none-any.whl.metadata (5.4 kB)
Collecting regex>=2021.8.3 (from nltk)
  Downloading regex-2024.4.16-cp312-cp312-win_amd64.whl.metadata (41 kB) 42.0/42.0 kB 2.1 MB/s eta 0:00:00
Collecting tqdm (from nltk)
  Downloading tqdm-4.66.2-py3-none-any.whl.metadata (57 kB) 57.6/57.6 kB 3.0 MB/s eta 0:00:00
Collecting colorama (from click->nltk)
  Downloading colorama-0.4.6-py2.py3-none-any.whl.metadata (17 kB)
Collecting nltk-3.8.1-py3-none-any.whl (1.5 MB)
  Downloading nltk-3.8.1-py3-none-any.whl (1.5 MB) 1.5/1.5 MB 6.0 MB/s eta 0:00:00
Collecting regex>=2024.4.16-cp312-cp312-win_amd64.whl (268 kB)
```

```
C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12>python
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> import nltk
```

```
AttributeError: module 'nltk' has no attribute 'download'. Did you mean: 'download'?
>>> nltk.download('brown')
[nltk_data] Downloading package brown to
[nltk_data]   C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\brown.zip.
True
>>>
```

Output:

```
IDLE Shell 3.12.3
File Edit Shell Debug Options Window Help
----- RESTART: D:/anjali/dl/2a.py -----
File list of brown corpus:
'c0a0', 'c0a5', 'c0a6', 'c0a7', 'c0a8', 'c0a9', 'c0a10',
'c0a11', 'c0a12', 'c0a13', 'c0a14', 'c0a15', 'c0a16', 'c0a17', 'c0a18', 'c0a19', 'c0a20',
'c0a21', 'c0a22', 'c0a23', 'c0a24', 'c0a25', 'c0a26', 'c0a27', 'c0a28', 'c0a29', 'c0a30',
'c0a31', 'c0a32', 'c0a33', 'c0a34', 'c0a35', 'c0a36', 'c0a37', 'c0a38', 'c0a39', 'c0a40',
'c0a41', 'c0a42', 'c0a43', 'c0a44', 'c0a45', 'c0a46', 'c0a47', 'c0a48', 'c0a49', 'c0a50',
'c0a51', 'c0a52', 'c0a53', 'c0a54', 'c0a55', 'c0a56', 'c0a57', 'c0a58', 'c0a59', 'c0a60',
'c0a61', 'c0a62', 'c0a63', 'c0a64', 'c0a65', 'c0a66', 'c0a67', 'c0a68', 'c0a69', 'c0a70',
'c0a71', 'c0a72', 'c0a73', 'c0a74', 'c0a75', 'c0a76', 'c0a77', 'c0a78', 'c0a79', 'c0a80',
'c0a81', 'c0a82', 'c0a83', 'c0a84', 'c0a85', 'c0a86', 'c0a87', 'c0a88', 'c0a89', 'c0a90',
'c0a91', 'c0a92', 'c0a93', 'c0a94', 'c0a95', 'c0a96', 'c0a97', 'c0a98', 'c0a99', 'c0a100',
'c0a101', 'c0a102', 'c0a103', 'c0a104', 'c0a105', 'c0a106', 'c0a107', 'c0a108', 'c0a109',
'c0a110', 'c0a111', 'c0a112', 'c0a113', 'c0a114', 'c0a115', 'c0a116', 'c0a117', 'c0a118',
'c0a119', 'c0a120', 'c0a121', 'c0a122', 'c0a123', 'c0a124', 'c0a125', 'c0a126', 'c0a127',
'c0a128', 'c0a129', 'c0a130', 'c0a131', 'c0a132', 'c0a133', 'c0a134', 'c0a135', 'c0a136',
'c0a137', 'c0a138', 'c0a139', 'c0a140', 'c0a141', 'c0a142', 'c0a143', 'c0a144', 'c0a145',
'c0a146', 'c0a147', 'c0a148', 'c0a149', 'c0a150', 'c0a151', 'c0a152', 'c0a153', 'c0a154',
'c0a155', 'c0a156', 'c0a157', 'c0a158', 'c0a159', 'c0a160', 'c0a161', 'c0a162', 'c0a163',
'c0a164', 'c0a165', 'c0a166', 'c0a167', 'c0a168', 'c0a169', 'c0a170', 'c0a171', 'c0a172',
'c0a173', 'c0a174', 'c0a175', 'c0a176', 'c0a177', 'c0a178', 'c0a179', 'c0a180', 'c0a181',
'c0a182', 'c0a183', 'c0a184', 'c0a185', 'c0a186', 'c0a187', 'c0a188', 'c0a189', 'c0a190',
'c0a191', 'c0a192', 'c0a193', 'c0a194', 'c0a195', 'c0a196', 'c0a197', 'c0a198', 'c0a199',
'c0a200', 'c0a201', 'c0a202', 'c0a203', 'c0a204', 'c0a205', 'c0a206', 'c0a207', 'c0a208',
'c0a209', 'c0a210', 'c0a211', 'c0a212', 'c0a213', 'c0a214', 'c0a215', 'c0a216', 'c0a217',
'c0a218', 'c0a219', 'c0a220', 'c0a221', 'c0a222', 'c0a223', 'c0a224', 'c0a225', 'c0a226',
'c0a227', 'c0a228', 'c0a229', 'c0a230', 'c0a231', 'c0a232', 'c0a233', 'c0a234', 'c0a235',
'c0a236', 'c0a237', 'c0a238', 'c0a239', 'c0a240', 'c0a241', 'c0a242', 'c0a243', 'c0a244',
'c0a245', 'c0a246', 'c0a247', 'c0a248', 'c0a249', 'c0a250', 'c0a251', 'c0a252', 'c0a253',
'c0a254', 'c0a255', 'c0a256', 'c0a257', 'c0a258', 'c0a259', 'c0a260', 'c0a261', 'c0a262',
'c0a263', 'c0a264', 'c0a265', 'c0a266', 'c0a267', 'c0a268', 'c0a269', 'c0a270', 'c0a271',
'c0a272', 'c0a273', 'c0a274', 'c0a275', 'c0a276', 'c0a277', 'c0a278', 'c0a279', 'c0a280',
'c0a281', 'c0a282', 'c0a283', 'c0a284', 'c0a285', 'c0a286', 'c0a287', 'c0a288', 'c0a289',
'c0a290', 'c0a291', 'c0a292', 'c0a293', 'c0a294', 'c0a295', 'c0a296', 'c0a297', 'c0a298',
'c0a299', 'c0a300', 'c0a301', 'c0a302', 'c0a303', 'c0a304', 'c0a305', 'c0a306', 'c0a307',
'c0a308', 'c0a309', 'c0a310', 'c0a311', 'c0a312', 'c0a313', 'c0a314', 'c0a315', 'c0a316',
'c0a317', 'c0a318', 'c0a319', 'c0a320', 'c0a321', 'c0a322', 'c0a323', 'c0a324', 'c0a325',
'c0a326', 'c0a327', 'c0a328', 'c0a329', 'c0a330', 'c0a331', 'c0a332', 'c0a333', 'c0a334',
'c0a335', 'c0a336', 'c0a337', 'c0a338', 'c0a339', 'c0a340', 'c0a341', 'c0a342', 'c0a343',
'c0a344', 'c0a345', 'c0a346', 'c0a347', 'c0a348', 'c0a349', 'c0a350', 'c0a351', 'c0a352',
'c0a353', 'c0a354', 'c0a355', 'c0a356', 'c0a357', 'c0a358', 'c0a359', 'c0a360', 'c0a361',
'c0a362', 'c0a363', 'c0a364', 'c0a365', 'c0a366', 'c0a367', 'c0a368', 'c0a369', 'c0a370',
'c0a371', 'c0a372', 'c0a373', 'c0a374', 'c0a375', 'c0a376', 'c0a377', 'c0a378', 'c0a379',
'c0a380', 'c0a381', 'c0a382', 'c0a383', 'c0a384', 'c0a385', 'c0a386', 'c0a387', 'c0a388',
'c0a389', 'c0a390', 'c0a391', 'c0a392', 'c0a393', 'c0a394', 'c0a395', 'c0a396', 'c0a397',
'c0a398', 'c0a399', 'c0a400', 'c0a401', 'c0a402', 'c0a403', 'c0a404', 'c0a405', 'c0a406'
```

```
IDE Shef 3.12.3
File Edit Shell Debug Options Window Help

.,'csm1', 'csm2', 'csm3', 'csm4', 'csm5', 'csm6', 'csm7', 'csm8',
.,'csm9', 'csm10', 'csm11', 'csm12', 'csm13', 'csm14', 'csm15', 'csm16',
.,'csm17', 'csm18', 'csm19', 'csm20', 'csm21', 'csm22', 'csm23', 'csm24',
.,'csm25', 'csm26', 'csm27', 'csm28', 'csm29', 'csm30', 'csm31', 'csm32',
.,'csm33', 'csm34', 'csm35', 'csm36', 'csm37', 'csm38', 'csm39', 'csm40',
.,'csm41', 'csm42', 'csm43', 'csm44', 'csm45', 'csm46', 'csm47', 'csm48',
.,'csm49', 'csm50', 'csm51', 'csm52', 'csm53', 'csm54', 'csm55', 'csm56',
.,'csm57', 'csm58', 'csm59', 'csm60', 'csm61', 'csm62', 'csm63', 'csm64',
.,'csm65', 'csm66', 'csm67', 'csm68', 'csm69', 'csm70', 'csm71', 'csm72',
.,'csm73', 'csm74', 'csm75', 'csm76', 'csm77', 'csm78', 'csm79', 'csm80',
.,'csm81', 'csm82', 'csm83', 'csm84', 'csm85', 'csm86', 'csm87', 'csm88',
.,'csm89', 'csm90', 'csm91', 'csm92', 'csm93', 'csm94', 'csm95', 'csm96',
.,'csm97', 'csm98', 'csm99', 'csm100', 'csm101', 'csm102', 'csm103', 'csm104',
.,'csm105', 'csm106', 'csm107', 'csm108', 'csm109', 'csm110', 'csm111', 'csm112',
.,'csm113', 'csm114', 'csm115', 'csm116', 'csm117', 'csm118', 'csm119', 'csm120',
.,'csm121', 'csm122', 'csm123', 'csm124', 'csm125', 'csm126', 'csm127', 'csm128',
.,'csm129', 'csm130', 'csm131', 'csm132', 'csm133', 'csm134', 'csm135', 'csm136',
.,'csm137', 'csm138', 'csm139', 'csm140', 'csm141', 'csm142', 'csm143', 'csm144',
.,'csm145', 'csm146', 'csm147', 'csm148', 'csm149', 'csm150', 'csm151', 'csm152',
.,'csm153', 'csm154', 'csm155', 'csm156', 'csm157', 'csm158', 'csm159', 'csm160',
.,'csm161', 'csm162', 'csm163', 'csm164', 'csm165', 'csm166', 'csm167', 'csm168',
.,'csm169', 'csm170', 'csm171', 'csm172', 'csm173', 'csm174', 'csm175', 'csm176',
.,'csm177', 'csm178', 'csm179', 'csm180', 'csm181', 'csm182', 'csm183', 'csm184',
.,'csm185', 'csm186', 'csm187', 'csm188', 'csm189', 'csm190', 'csm191', 'csm192',
.,'csm193', 'csm194', 'csm195', 'csm196', 'csm197', 'csm198', 'csm199', 'csm200',
.,'csm201', 'csm202', 'csm203', 'csm204', 'csm205', 'csm206', 'csm207', 'csm208',
.,'csm209', 'csm210', 'csm211', 'csm212', 'csm213', 'csm214', 'csm215', 'csm216',
.,'csm217', 'csm218', 'csm219', 'csm220', 'csm221', 'csm222', 'csm223', 'csm224',
.,'csm225', 'csm226', 'csm227', 'csm228', 'csm229', 'csm230', 'csm231', 'csm232',
.,'csm233', 'csm234', 'csm235', 'csm236', 'csm237', 'csm238', 'csm239', 'csm240',
.,'csm241', 'csm242', 'csm243', 'csm244', 'csm245', 'csm246', 'csm247', 'csm248',
.,'csm249', 'csm250', 'csm251', 'csm252', 'csm253', 'csm254', 'csm255', 'csm256',
.,'csm257', 'csm258', 'csm259', 'csm260', 'csm261', 'csm262', 'csm263', 'csm264',
.,'csm265', 'csm266', 'csm267', 'csm268', 'csm269', 'csm270', 'csm271', 'csm272',
.,'csm273', 'csm274', 'csm275', 'csm276', 'csm277', 'csm278', 'csm279', 'csm280',
.,'csm281', 'csm282', 'csm283', 'csm284', 'csm285', 'csm286', 'csm287', 'csm288',
.,'csm289', 'csm290', 'csm291', 'csm292', 'csm293', 'csm294', 'csm295', 'csm296',
.,'csm297', 'csm298', 'csm299', 'csm300', 'csm301', 'csm302', 'csm303', 'csm304',
.,'csm305', 'csm306', 'csm307', 'csm308', 'csm309', 'csm310', 'csm311', 'csm312',
.,'csm313', 'csm314', 'csm315', 'csm316', 'csm317', 'csm318', 'csm319', 'csm320',
.,'csm321', 'csm322', 'csm323', 'csm324', 'csm325', 'csm326', 'csm327', 'csm328',
.,'csm329', 'csm330', 'csm331', 'csm332', 'csm333', 'csm334', 'csm335', 'csm336',
.,'csm337', 'csm338', 'csm339', 'csm340', 'csm341', 'csm342', 'csm343', 'csm344',
.,'csm345', 'csm346', 'csm347', 'csm348', 'csm349', 'csm350', 'csm351', 'csm352',
.,'csm353', 'csm354', 'csm355', 'csm356', 'csm357', 'csm358', 'csm359', 'csm360',
.,'csm361', 'csm362', 'csm363', 'csm364', 'csm365', 'csm366', 'csm367', 'csm368',
.,'csm369', 'csm370', 'csm371', 'csm372', 'csm373', 'csm374', 'csm375', 'csm376',
.,'csm377', 'csm378', 'csm379', 'csm380', 'csm381', 'csm382', 'csm383', 'csm384',
.,'csm385', 'csm386', 'csm387', 'csm388', 'csm389', 'csm390', 'csm391', 'csm392',
.,'csm393', 'csm394', 'csm395', 'csm396', 'csm397', 'csm398', 'csm399', 'csm400',
]

csm01 has following words:
["The", "Fulton", "County", "Grand", "Jury", "said", ...]

csm01 has 2362 words

Categories or file in brown corpus:
['adventure', 'belles_lettres', 'editorial', 'fiction', 'government', 'hobbies',
.'humor', 'learning', 'lore', 'mystery', 'news', 'religion', 'reviews', 'romance',
.'science_fiction']

Statistics for each text:

```

PRACTICAL 2C

Study Conditional frequency distributions

Code:

```
#process a sequence of pairs
text = ['The', 'Fulton', 'County', 'Grand', 'Jury', 'said', ...]
pairs = [('news', 'The'), ('news', 'Fulton'), ('news', 'County'), ...]
import nltk
from nltk.corpus import brown
fd = nltk.ConditionalFreqDist(
    (genre, word)
    for genre in brown.categories()
    for word in brown.words(categories=genre))
genre_word = [(genre, word)
    for genre in ['news', 'romance']
    for word in brown.words(categories=genre)]
print(len(genre_word))
print(genre_word[:4])
print(genre_word[-4:])
cfд = nltk.ConditionalFreqDist(genre_word)
print(cfд)
print(cfд.conditions())
print(cfд['news'])
print(cfд['romance'])
print(list(cfд['romance']))
from nltk.corpus import inaugural
cfд = nltk.ConditionalFreqDist(
    (target, fileid[:4])
    for fileid in inaugural.fileids()
```

```

for w in inaugural.words(fileid)
    for target in ['america', 'citizen']
        if w.lower().startswith(target))
            from nltk.corpus import udhr
            languages = ['Chickasaw', 'English', 'German_Deutsch','Greenlandic_Inuktikut',
            'Hungarian_Magyar', 'Ibibio_Efik']
            cfd = nltk.ConditionalFreqDist(
                (lang, len(word))
                for lang in languages
                for word in udhr.words(lang + '-Latin1'))
            cfd.tabulate(conditions=['English', 'German_Deutsch'],
            samples=range(10), cumulative=True)

```

```

Select Command Prompt - python
Searched in:
- 'C:\\Users\\Student\\nltk_data'
- 'C:\\Users\\Student\\AppData\\Local\\Programs\\Python\\Python312\\nltk_data'
- 'C:\\Users\\Student\\AppData\\Local\\Programs\\Python\\Python312\\share\\nltk_data'
- 'C:\\Users\\Student\\AppData\\Local\\Programs\\Python\\Python312\\lib\\nltk_data'
- 'C:\\Users\\Student\\AppData\\Roaming\\nltk_data'
- 'C:\\nltk_data'
- 'D:\\nltk_data'
- 'E:\\nltk_data'
*****
>> cfd.tabulate(conditions=['English', 'German_Deutsch'],
... samples=range(10), cumulative=True)^Z
File "<stdin>", line 2
    samples=range(10), cumulative=True) ^
SyntaxError: invalid non-printable character U+001A
>> nltk.download('inaugural')
nltk_data] Downloading package inaugural to
nltk_data]     C:\\Users\\Student\\AppData\\Roaming\\nltk_data...
nltk_data]     Unzipping corpora\\inaugural.zip.
rue
>> -

```

```
Command Prompt - python
Searched in:
'C:\Users\Student\nltk_data'
'C:\Users\Student\AppData\Local\Programs\Python\Python312\nltk_data'
'C:\Users\Student\AppData\Local\Programs\Python\Python312\share\nltk_data'
'C:\Users\Student\AppData\Roaming\lib\nltk_data'
'C:\nltk_data'
'D:\nltk_data'
'E:\nltk_data'
*****
>>> cfd.tabulate(conditions=['English', 'German_Deutsch'],
... samples=range(10), cumulative=True)^Z
File "<stdin>", line 2
    samples=range(10), cumulative=True)E
SyntaxError: invalid non-printable character U+001A
>>> nltk.download('inaugural')
[nltk_data] Downloading package inaugural to
[nltk_data]   C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\inaugural.zip.
True
>>> nltk.download('udhr')
[nltk_data] Downloading package udhr to
[nltk_data]   C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\udhr.zip.
True
>>>
>>>
```

Output:

```
IDLE Shell 3.12.3
File Edit Shell Debug Options Window Help
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on Win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: D:/anjali dl/2C.py
170576
[('news', 'The'), ('news', 'Fulton'), ('news', 'County'), ('news', 'Grand')]
[('romance', 'afraid'), ('romance', 'not'), ('romance', "'"), ('romance', '.')]
<ConditionalFreqDist with 2 conditions>
['news', 'romance']
<FreqDist with 14394 samples and 100554 outcomes>
<FreqDist with 8452 samples and 70022 outcomes>
Squeezed text(1147 lines).
      0   1   2   3   4   5   6   7   8   9
English  0  185  525  883  997  1166  1283  1440  1558  1638
German_Deutsch  0  171  263  614  717  894  1013  1110  1213  1275
>>>
```

PRACTICAL 2D

Study of tagged corpora with methods like tagged_sents, tagged_words.

Code:

```
import nltk  
from nltk import tokenize  
nltk.download('punkt')  
nltk.download('words')  
para = "Hello! My name is Anjali Nimje. Today you'll be learning NLTK."  
sents = tokenize.sent_tokenize(para)  
print("\nsentence tokenization\n=====\\n",sents)  
# word tokenization  
print("\nword tokenization\n=====\\n")  
for index in range(len(sents)):  
    words = tokenize.word_tokenize(sents[index])  
    print(words)
```

Output:

```
IDLE Shell 3.12.3  
File Edit Shell Debug Options Window Help  
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr 9 2024, 14:05:25) [MSC v.1938 64 bit  
AMD64] on Win32  
Type "help", "copyright", "credits" or "license()" for more information.  
>>> = RESTART: D:/anjali dl/2D.py  
[nltk_data] Downloading package punkt to  
[nltk_data]     C:/Users/Student/AppData/Roaming/nltk_data...  
[nltk_data] Package punkt is already up-to-date!  
[nltk_data] Downloading package words to  
[nltk_data]     C:/Users/Student/AppData/Roaming/nltk_data...  
[nltk_data] Package words is already up-to-date!  
  
sentence tokenization  
=====  
['Hello!', 'My name is Anjali Nimje.', 'Today you''ll be learning NLTK.'][  
word tokenization  
=====  
['Hello', '!']  
['My', 'name', 'is', 'Anjali', 'Nimje', '.']  
['Today', 'you', 'll', 'be', 'learning', 'NLTK', '.']  
>>>
```

PRACTICAL 2E

Write a program to find the most frequent noun tags.

Code:

```
import nltk
from collections import defaultdict
text = nltk.word_tokenize("Nick likes to play football. Nick does not like to
playcricket.")
tagged = nltk.pos_tag(text)
print(tagged)
# checking if it is a noun or not
addNounWords = []
count=0
for words in tagged:
    val = tagged[count][1]
    if(val == 'NN' or val == 'NNS' or val == 'NNPS' or val == 'NNP'):
        addNounWords.append(tagged[count][0])
    count+=1
print (addNounWords)
temp = defaultdict(int)
# memoizing count
for sub in addNounWords:
    for wrd in sub.split():
        temp[wrd] += 1
# getting max frequency
res = max(temp, key=temp.get)
# printing result
print("Word with maximum frequency : " + str(res))
```

```
Command Prompt - python
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>cd C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12

C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12>python
Python 3.12.3 (tags/v3.12.3:f665bf9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.
>>> nltk.download('averaged_perceptron_tagger')
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'nltk' is not defined
>>> import nltk
>>> nltk.download('averaged_perceptron_tagger')
[nltk_data]  Downloading package averaged_perceptron_tagger to
[nltk_data]    C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]  Unzipping taggers\averaged_perceptron_tagger.zip.
True
>>>

print("Word with maximum frequency : "+str(res))
```

Output:

```
IDLE Shell 3.12.3
File Edit Shell Debug Options Window Help
Python 3.12.3 (tags/v3.12.3:f665bf9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.
>>>
= RESTART: D:/anjali dl/2E.py
[('Nick', 'NNP'), ('likes', 'VBZ'), ('to', 'TO'), ('play', 'VB'), ('football', 'NN'),
 ('.', '.'), ('Nick', 'NNP'), ('does', 'VBZ'), ('not', 'RB'), ('like', 'VB'),
 ('to', 'TO'), ('playcricket', 'VB'), ('.', '.')]

['Nick', 'football', 'Nick']

Word with maximum frequency : Nick
>>>
```

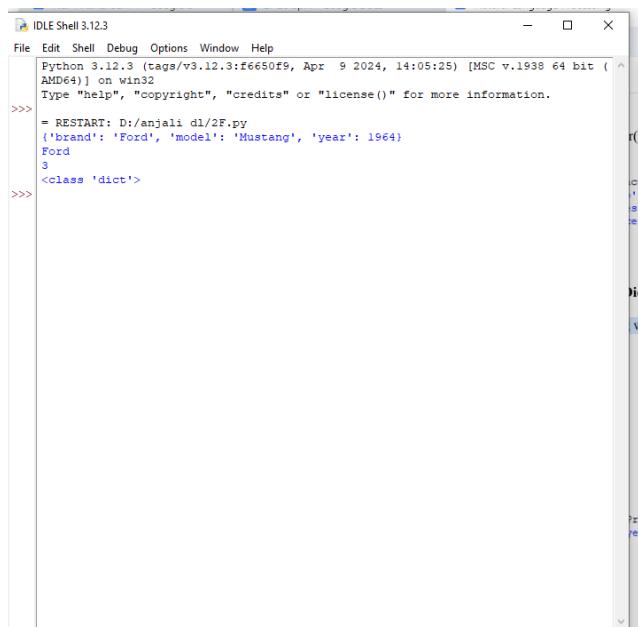
PRACTICAL 2F

Map Words to Properties Using Python Dictionaries

Code:

```
#creating and printing a dictionay by mapping word with its properties
thisdict = {
    "brand": "Ford",
    "model": "Mustang",
    "year": 1964
}
print(thisdict)
print(thisdict["brand"])
print(len(thisdict))
print(type(thisdict))
```

Output



The screenshot shows the Python IDLE Shell interface. The title bar reads "IDLE Shell 3.12.3". The menu bar includes File, Edit, Shell, Debug, Options, Window, and Help. The main window displays the following text output from the shell:

```
File Edit Shell Debug Options Window Help
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.

>>> = RESTART: D:/anjali d1/2F.py
{'brand': 'Ford', 'model': 'Mustang', 'year': 1964}
Ford
3
<class 'dict'>
```

PRACTICAL 2G

Study DefaultTagger, Regular expression tagger, UnigramTagger

A)Default Tagger

Code:

```
import nltk

from nltk.tag import DefaultTagger

exptagger = DefaultTagger('NN')

from nltk.corpus import treebank

testsentences = treebank.tagged_sents()[1000:]

print(exptagger.evaluate(testsentences))

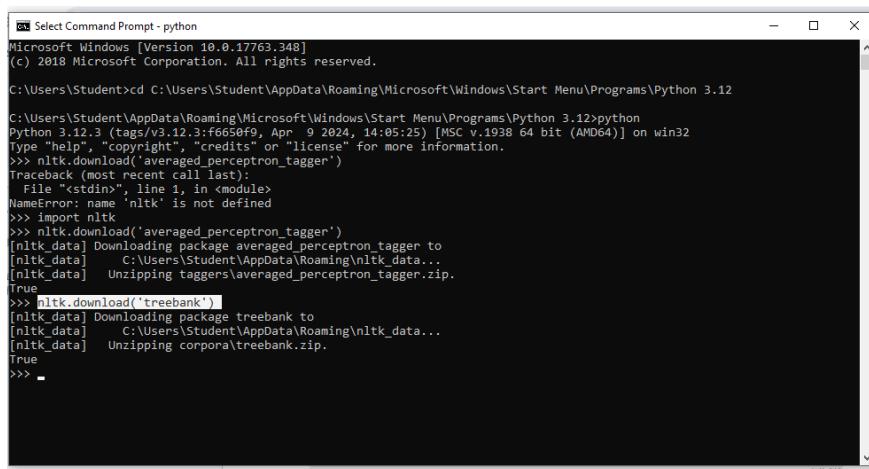
#Tagging a list of sentences

import nltk

from nltk.tag import DefaultTagger

exptagger = DefaultTagger('NN')

print(exptagger.tag_sents([('Hi', ','), ('How', 'are'), ('you', '?')]))
```



```
Microsoft Windows [Version 10.0.17763.348]
(c) 2018 Microsoft Corporation. All rights reserved.

C:\Users\Student>cd C:\Users\Student\AppData\Roaming\Microsoft\Windows\Start Menu\Programs\Python 3.12

C:\Users\Student>python
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license" for more information.

>>> nltk.download('averaged_perceptron_tagger')
Traceback (most recent call last):
  File "<stdin>", line 1, in <module>
NameError: name 'nltk' is not defined
>>> import nltk
>>> nltk.download('averaged_perceptron_tagger')
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]   C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\averaged_perceptron_tagger.zip.
True
>>> nltk.download('treebank')
[nltk_data] Downloading package treebank to
[nltk_data]   C:\Users\Student\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\treebank.zip.
True
>>>
```

Output:

The screenshot shows the Python IDLE Shell interface. The menu bar includes File, Edit, Shell, Debug, Options, Window, and Help. The status bar at the bottom right indicates L=14 Col 0, 09:46 PM, ENG, 20-04-2024.

```
idle Shell 3.12.2
File Edit Shell Debug Options Window Help
Python 3.12.3 (tags/v3.12.3-95605fb9, Apr 9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.

>>> = RESTART: D:/anjali/dl/20.1.py

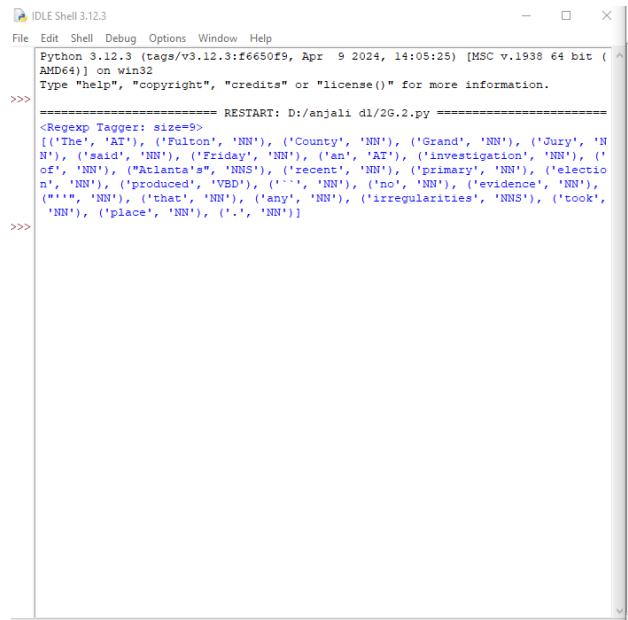
Warning (from warnings module):
  File "D:/anjali/dl/20.1.py", line 6
    from sympy import symbols, evaluate, tesselate
DeprecationWarning:
  Function evaluate() has been deprecated. Use accuracy(gold)
0.13196749534374715
[('Hi', 'NN'), ('.', 'NN'), ('are', 'NN'), ('you', 'NN'), ('?', 'NN')]
```

B)Regular Expressions

Code:

```
from nltk.corpus import brown  
from nltk.tag import RegexpTagger  
test_sent = brown.sents(categories='news')[0]  
regexp_tagger = RegexpTagger(  
    [(r'^-?[0-9]+([0-9]+)?$', 'CD'), # cardinal numbers  
     (r'(The|the|A|a|An|an)$', 'AT'), # articles  
     (r'.*able$', 'JJ'), # adjectives  
     (r'.*ness$', 'NN'), # nouns formed from adjectives  
     (r'.*ly$', 'RB'), # adverbs  
     (r'.*s$', 'NNS'), # plural nouns  
     (r'.*ing$', 'VBG'), # gerunds  
     (r'.*ed$', 'VBD'), # past tense verbs  
     (r'.*', 'NN') # nouns (default)  
])  
print(regexp_tagger)  
print(regexp_tagger.tag(test_sent))
```

Output:



The screenshot shows the Python IDLE Shell interface with the title bar "IDLE Shell 3.12.3". The menu bar includes File, Edit, Shell, Debug, Options, Window, and Help. The main window displays the following text:

```
File Edit Shell Debug Options Window Help
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr 9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.

>>> ===== RESTART: D:/anjali/dl/2G.2.py =====
<Regexp Tagger: size=9>
([('The', 'AT'), ('Fulton', 'NN'), ('County', 'NN'), ('Grand', 'NN'), ('Jury', 'N
N'), ('said', 'NN'), ('Friday', 'NN'), ('an', 'AT'), ('Investigation', 'NN'), ('o
f', 'NN'), ('Atlanta', 'NNS'), ('recess', 'NN'), ('primary', 'NN'), ('electio
n', 'NN'), ('produced', 'VBD'), ('a', 'NN'), ('no', 'NN'), ('evidence', 'NN'),
('in', 'NN'), ('that', 'NN'), ('any', 'NN'), ('irregularities', 'NNS'), ('took',
'NN'), ('place', 'NN'), ('.', 'NN')]
```

C) Unigram Tagger

Code:

```
# Loading Libraries
from nltk.tag import UnigramTagger
from nltk.corpus import treebank

# Training using first 10 tagged sentences of the treebank corpus as data.

# Using data
train_sents = treebank.tagged_sents()[:10]

# Initializing
tagger = UnigramTagger(train_sents)

# Lets see the first sentence
# (of the treebank corpus) as list
print(treebank.sents()[0])
print("\n",tagger.tag(treebank.sents()[0]))

#Finding the tagged results after training.
tagger.tag(treebank.sents()[0])

#Overriding the context model
tagger = UnigramTagger(model ={'Pierre': 'NN'})
print("\n",tagger.tag(treebank.sents()[0]))
```

Output:

```
IDLE Shell 3.12.3
Edit  Shell  Debug  Options  Window  Help
Python 3.12.3 (tags/v3.12.3:f6650f9, Apr  9 2024, 14:05:25) [MSC v.1938 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.

===== RESTART: D:/anjali/d1/26_3.py =====
[('Pierre', 'Vinken', '.', '61', 'years', 'old', '.', 'will', 'join', 'the', 'board', 'as', 'a', 'nonexecutive', 'director', 'Nov.', '29', '.')]
[('Pierre', 'NNP'), ('Vinken', 'NNP'), ('.', 'DT'), ('61', 'CD'), ('years', 'NNS'), ('.', 'CD'), ('.', 'DT'), ('.', 'DT'), ('will', 'MD'), ('join', 'VB'), ('the', 'DT'), ('board', 'NN'), ('as', 'IN'), ('.', 'DT'), ('nonexecutive', 'JJ'), ('director', 'NN'), ('Nov.', 'NNP'), ('29', 'CD'), ('.', 'DT')]

[('Pierre', 'NN'), ('Vinken', None), ('.', None), ('61', None), ('years', None), ('.', None), ('.', None), ('will', None), ('join', None), ('the', None), ('board', None), ('as', None), ('.', None), ('nonexecutive', None), ('director', None), ('Nov.', None), ('29', None), ('.', None)]
```

PRACTICAL 3A

Study of Wordnet Dictionary with methods as synsets, definitions, examples,antonyms.

Code:

```
import nltk
nltk.download('wordnet')
from nltk.corpus import wordnet
print(wordnet.synsets("computer"))
print(wordnet.synset("computer.n.01").definition())
print("Examples:", wordnet.synset("computer.n.01").examples())
print(wordnet.lemma('buy.v.01.buy').antonyms())
```

Output:

```
-----[REDACTED]-----
Edit Shell Debug Options Window Help
Python 3.12.4 (tags/v3.12.4:48e8a4ba, Jun 6 2024, 19:30:16) [MSC v.1940 64 bit (AMD64)] on win32
Type "help", "copyright", "credits" or "license()" for more information.

= RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/3a.py
[nltk_data] Downloading package wordnet to
[nltk_data]   C:\Users\Admin\AppData\Roaming\nltk_data...
[Synset('computer.n.01'), Synset('calculator.n.01')]
a machine for performing calculations automatically
Examples: []
[Lemma('sell.v.01.sell')]
```

PRACTICAL 3B

Study lemmas, hyponyms, hypernyms.

Code:

```
import nltk
from nltk.corpus import wordnet
print(wordnet.synsets("computer"))
print(wordnet.synset("computer.n.01").lemma_names())
for e in wordnet.synsets("computer"):
    print(f'{e} --> {e.lemma_names()}'')
print(wordnet.synset("computer.n.01").lemmas())
print(wordnet.lemma('computer.n.01.computing_device').synset())
print(wordnet.lemma('computer.n.01.computing_device').name())
syn = wordnet.synset('computer.n.01')
print(syn.hyponyms())
print([lemma.name() for synset in syn.hyponyms() for lemma in
synset.lemmas()])
vehicle = wordnet.synset('vehicle.n.01')car = wordnet.synset('car.n.01')
print(car.lowest_common_hypernyms(vehicle))
```

Output:

```
>>> RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/3b.py =====
[Synset('computer.n.01'), Synset('calculator.n.01')]
['computer', 'computing_machine', 'computing_device', 'data_processor', 'electronic_computer', 'information_processing_system']
Synset('computer.n.01') --> ['computer', 'computing_machine', 'computing_device', 'data_processor', 'electronic_computer', 'information_processing_system']
Synset('calculator.n.01') --> ['calculator', 'reckoner', 'figurer', 'estimator', 'computer']
[Lemma('computer.n.01.computer'), Lemma('computer.n.01.computing_machine'), Lemma('computer.n.01.computing_device'), Lemma('computer.n.01.data_processor'), Lemma('computer.n.01.electronic_computer'), Lemma('computer.n.01.information_processing_system')]
Synset('computer.n.01')
computing_device
<bound method _WordNetObject.hyponyms of Synset('computer.n.01')>
['analog_computer', 'analogue_computer', 'digital_computer', 'home_computer', 'node', 'client', 'guest', 'number_cruncher', 'pari-mutuel_machine', 'totalizer', 'totaliser', 'totalizator', 'totalisator', 'predictor', 'server', 'host', 'Turing_machine', 'web_site', 'website', 'internet_site', 'site']
[Synset('vehicle.n.01')]
>>>
```

PRACTICAL 3C

Write a program using python to find synonym and antonym of word "active" using Wordnet.

Code:

```
from nltk.corpus import wordnet  
  
print(wordnet.synsets("active"))  
  
print(wordnet.lemma('active.a.01.active').antonyms())
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/3c.py =====  
[Synset('active_agent.n.01'), Synset('active_voice.n.01'), Synset('active.n.03'), Synset('active.a.01'), Synset('active.s.02'), Synset('active.a.03'), Synset('acti  
ve.s.04'), Synset('active.a.05'), Synset('active.a.06'), Synset('active.a.07'), Synset('active.s.08'), Synset('active.a.09'), Synset('active.a.10'), Synset('active  
.a.11'), Synset('active.a.12'), Synset('active.a.13'), Synset('active.a.14')]  
[Lemma('inactive.a.02.inactive')]  
>>> |
```

PRACTICAL 3D

Compare two nouns.

Code:

```
import nltk
from nltk.corpus import wordnet
syn1 = wordnet.synsets('football')
syn2 = wordnet.synsets('soccer')
for s1 in syn1:
    for s2 in syn2:
        print("path similarity of:")
        print(s1,'(',s1.pos(),'',[',s1.definition(),']')
        print(s2,'(',s2.pos(),'',[',s2.definition(),']')
        print("is", s1.path_similarity(s2))
print()
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/3d.py =====
path similarity of:
Synset('football.n.01') ( n ) [ any of various games played with a ball (round or oval) in which two teams try to kick or carry or propel the ball into each other's goal ]
Synset('soccer.n.01') ( n ) [ a football game in which two teams of 11 players try to kick or head a ball into the opponents' goal ]
is 0.5

path similarity of:
Synset('football.n.02') ( n ) [ the inflated oblong ball used in playing American football ]
Synset('soccer.n.01') ( n ) [ a football game in which two teams of 11 players try to kick or head a ball into the opponents' goal ]
is 0.05
```

PRACTICAL 3E

Handling stopword

A) Using nltk Adding or Removing Stop Words in NLTK's Default Stop Word List

Code:

```
import nltk  
from nltk.corpus import stopwords  
nltk.download('stopwords')  
  
from nltk.tokenize import word_tokenize  
  
text = "Yashesh like to play football, however he is not too fond of tennis."  
  
text_tokens = word_tokenize(text)  
  
tokens_without_sw = [word for word in text_tokens if not word in  
stopwords.words()]  
  
print(tokens_without_sw)  
  
all_stopwords = stopwords.words('english')  
all_stopwords.append('play')  
  
text_tokens = word_tokenize(text)  
  
tokens_without_sw = [word for word in text_tokens if not word in  
all_stopwords]  
  
print(tokens_without_sw)  
  
all_stopwords.remove('not')  
  
text_tokens = word_tokenize(text)  
  
tokens_without_sw = [word for word in text_tokens if not word in  
all_stopwords]  
  
print(tokens_without_sw)
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/3e1.py ======  
[nltk_data] Downloading package stopwords to  
[nltk_data]   C:\Users\Admin\AppData\Roaming\nltk_data...  
[nltk_data] Package stopwords is already up-to-date!  
['Yashesh', 'play', 'football', ',', 'fond', 'tennis', '.']  
['Yashesh', 'like', 'football', ',', 'however', 'fond', 'tennis', '.']  
['Yashesh', 'like', 'football', ',', 'however', 'not', 'fond', 'tennis', '.']  
>>> | Ln: 5
```

B) Using Gensim Adding and Removing Stop Words in Default Gensim Stop Words List

Code:

```
import gensim  
from gensim.parsing.preprocessing import remove_stopwords  
from nltk.tokenize import word_tokenize  
text = "Yashesh likes to play football, however he is not too fond of tennis."  
filtered_sentense = remove_stopwords(text)  
print(filtered_sentense)  
all_stpwords = gensim.parsing.preprocessing.STOPWORDS  
print(all_stpwords)  
from gensim.parsing.preprocessing import STOPWORDS  
all_stpwords_genism = STOPWORDS.union(set(['likes','play']))  
text = "Yashesh likes to play football, however he is not too fond of tennis."  
text_tokens = word_tokenize(text)  
tokens_without_sw = [word for word in text_tokens if not word in  
all_stpwords]print(tokens_without_sw)  
print("====")  
all_stpwords_genism = STOPWORDS  
sw_list = {"not"}  
all_stpwords_genism = STOPWORDS.difference(sw_list)  
text = "Yashesh likes to play football, however he is not too fond of tennis."  
text_tokens = word_tokenize(text)  
tokens_without_sw = [word for word in text_tokens if not word in  
all_stpwords_genism]  
print(all_stpwords_genism)  
print(tokens_without_sw)
```

```
C:\Windows\system32\cmd.exe

C:\Users\Admin>pip install gensim
Collecting gensim
  Downloading gensim-4.3.3-cp312-cp312-win_amd64.whl.metadata (8.2 kB)
Collecting numpy<2.0,>=1.18.5 (from gensim)
  Downloading numpy-1.26.4-cp312-cp312-win_amd64.whl.metadata (61 kB)
    100% |██████████| 61.0/61.0 kB 651.3 kB/s eta 0:00:00
Collecting wrapt<1.14.0,>=1.7.0 (from gensim)
  Downloading wrapt-1.13.1-cp312-cp312-win_amd64.whl.metadata (60 kB)
    100% |██████████| 60.0/60.0 kB 799.0 kB/s eta 0:00:00

Collecting smart-open<4.8.1 (from gensim)
  Downloading smart_open-7.0.4-py3-none-any.whl.metadata (23 kB)
Collecting wrapt (<from smart-open>1.8.1->gensim)
  Downloading wrapt-1.16.0-cp312-cp312-win_amd64.whl.metadata (6.8 kB)
  Downloading wrapt-1.16.0-cp312-cp312-win_amd64.whl (24.0 kB)
Collecting gensim<4.3.3-cp312-cp312-win_amd64.whl (45.5 MB)
  Downloading gensim-4.3.3-cp312-cp312-win_amd64.whl (45.5 MB)
    100% |██████████| 45.5/45.5 kB 681.6 kB/s eta 0:00:00
  Downloading scipy-1.13.1-cp312-cp312-win_amd64.whl (45.9 kB)
    100% |██████████| 45.9/45.9 kB 681.6 kB/s eta 0:00:00

Collecting smart-open_7.0.4-py3-none-any.whl (61 kB)
  Downloading smart-open_7.0.4-py3-none-any.whl (61 kB)
    100% |██████████| 61.2/61.2 kB 542.3 kB/s eta 0:00:00

Installing collected packages: wrapt, numpy, smart-open, scipy, gensim
  Attempting uninstall: numpy
    Found existing installation: numpy 2.0.1
    Uninstalling numpy-2.0.1...
      Successfully uninstalled numpy-2.0.1
WARNING: Failed to remove contents in a temporary directory 'C:\Users\Admin\AppData\Roaming\Python\Python312\site-packages\~numpy.libs'.
You can safely remove it manually.
WARNING: Failed to remove contents in a temporary directory 'C:\Users\Admin\AppData\Roaming\Python\Python312\site-packages\~scipy'.
You can safely remove it manually.
Successfully installed gensim-4.3.3 numpy-1.26.4 scipy-1.13.1 smart-open-7.0.4 wrapt-1.16.0

C:\Users\Admin>
```

Output:

C) Using Spacy Adding and Removing Stop Words in Default Spacy Stop Words List.

Code:

```
import spacy
import nltk
from nltk.tokenize import word_tokenize
sp = spacy.load('en_core_web_sm')
all_stopwords = sp.Defaults.stop_words
all_stopwords.add('play')
text = "Yashesh like to play football, however he is not too fond of tennis."
text_tokens = word_tokenize(text)
tokens_without_sw = [word for word in text_tokens if not word in
all_stopwords]
print(tokens_without_sw)
all_stopwords.remove('not')
tokens_without_sw = [word for word in text_tokens if not word in
all_stopwords]
print(tokens_without_sw)
```

Commands:

```
#pip install spacy
#python -m spacy download en_core_web_sm
#python -m spacy download en
```

```
C:\Windows\system32\cmd.exe - pip install spacy
Collecting spacy
  Downloading spacy-3.7.5-cp312-cp312-win_amd64.whl.metadata (27 kB)
Collecting spacy-legacy<3.7.5>,>=3.0.11 (from spacy)
  Downloading spacy_legacy-3.0.12-py3.py3-none-any.whl.metadata (2.8 kB)
Collecting spacy_loggers<1.0.0>,>=0.28.0 (from spacy)
  Downloading spacy_loggers-0.28.0-py3-none-any.whl.metadata (23 kB)
Collecting murmurhash<1.1.0,>=0.28.0 (from spacy)
  Downloading murmurhash-0.28.0-py3-none-any.whl.metadata (2.0 kB)
Collecting cymem<2.1.0,>=2.0.2 (from spacy)
  Downloading cymem-2.0.8-cp312-cp312-win_amd64.whl.metadata (8.6 kB)
Collecting preshed<3.0.0,>=2.0.0 (from spacy)
  Downloading preshed-3.0.9-cp312-cp312-win_amd64.whl.metadata (2.2 kB)
Collecting thinc<3.8.0,>=3.8.2 (from spacy)
  Downloading thinc-3.8.2-cp312-cp312-win_amd64.whl.metadata (15 kB)
Collecting wasabi<2.0.0,>=0.9.3 (from spacy)
  Downloading wasabi-1.1.3-py3-none-any.whl.metadata (28 kB)
Collecting srilay<2.4.8,>=2.4.8 (from spacy)
  Downloading srilay-2.4.8-cp312-cp312-win_amd64.whl.metadata (20 kB)
Collecting cataloger<2.0.0,>=0.6 (from spacy)
  Downloading cataloger-0.6.0-py3-none-any.whl.metadata (14 kB)
Collecting weasel<0.5.0,>=0.1.8 (from spacy)
  Downloading weasel-0.4.1-py3-none-any.whl.metadata (4.6 kB)
Collecting typing_extensions<4.0.0,>=3.0.0 (from spacy)
  Downloading typing_extensions-3.1.0-py3-none-any.whl.metadata (15 kB)
Requirement already satisfied: tghc<0.1.0> in C:\Users\admin\AppData\Roaming\Python\Python312\site-packages (from spacy) (4.66.4)
Requirement already satisfied: typing<4.0.0,>=3.0.0 (from spacy)
Collecting pydantic<1.8.1,>=1.8.1 (from spacy)
  Downloading pydantic-2.8.2-py3-none-any.whl.metadata (12 kB)
  Downloading pydantic-2.8.2-py3-none-any.whl.metadata (12 kB)  1.1 MB/s eta 0:00:00
Collecting jinja2 (from spacy)
  Downloading jinja2-3.1.2-py3-none-any.whl.metadata (2.6 kB)
Requirement already satisfied: setuptools in C:\Users\admin\AppData\Local\Programs\Python\Python312\lib\site-packages (from spacy) (71.1.0)
Collecting packaging<20.0> (from spacy)
  Downloading packaging-20.0-py3-none-any.whl.metadata (3.2 kB)
Collecting language-data<4.0.0,>=3.0.2 (from spacy)
  Downloading language_data-3.2.0-py3-none-any.whl.metadata (4.3 kB)
Collecting annotated-types<0.7.0,>=0.6.0 (from spacy)
  Downloading annotated_types-0.6.0-py3-none-any.whl.metadata (15 kB)
Collecting pydantic-core<2.20.1,>=2.20.1 (from pydantic<1.8.1,>=1.8.1,3.0.0,>=3.0.0,4->spacy)
  Downloading pydantic_core-2.20.1-cp312-none-win_amd64.whl.metadata (6.7 kB)
Requirement already satisfied: typing_extensions<4.0.0,>=3.0.0 (in C:\Users\admin\AppData\Local\Programs\Python\Python312\lib\site-packages (from pydantic<1.8.1,>=1.8.1,<3.0.0,>))
```

Output:

```
IDLE Shell 3.12.4
File Edit Shell Options Window Help
Python 3.12.4 (tags/v3.12.4:5e69fb0, Jun  6 2024, 19:20:16) [MSC v.1940 64 bit (AMD64)] on win32
Type "help()", "copyright()", "credits()" or "license()" for more information.
>>> = RESTART: C:\Users\admin\AppData\Local\Programs\Python\Python312\3e3.py
['Tashash', 'like', 'Football', '.', 'Rom', 'Tennis', '.']
['Tashash', 'like', 'Football', '.', 'not', 'Tennis', '.']
```

PRACTICAL 4A

Tokenization using Python's split() function.

Code:

```
text = "Monism is a thesis about oneness: that only one thing exists in a certain sense. The denial of monism is pluralism, the thesis that, in a certain sense, more than one thing exists.[7] There are manyforms of monism and pluralism, but in relation to the world as a whole, two are of special interest:existence monism/pluralism and priority monism/pluralism."
```

```
data = text.split('.')
for i in data:
    print(i)
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/4a.py =====
Monism is a thesis about oneness: that only one thing exists in a certain sense
The denial of monism is pluralism, the thesis that, in a certain sense, more than one thing exists
[7] There are manyforms of monism and pluralism, but in relation to the world as a whole, two are of special interest:existence monism/pluralism and priority monis
m/pluralism
>>> |
```

PRACTICAL 4B

Tokenization using Regular Expressions (RegEx)

Code:

```
import nltk  
from nltk.tokenize import RegexpTokenizer  
tk = RegexpTokenizer('\s+', gaps=True)  
str = "I love to study Natrual Language processing in python"  
tokkens = tk.tokenize(str)  
print(tokkens)
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/4b.py =====  
['I', 'love', 'to', 'study', 'Natrual', 'Language', 'processing', 'in', 'python']
```

PRACTICAL 4C

Tokenization using NLTK

Code:

```
import nltk  
from nltk.tokenize import word_tokenize  
str = "I love to study Natrual Language processing in python"  
print(word_tokenize(str))
```

Output

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/4c.py ======  
['I', 'love', 'to', 'study', 'Natrual', 'Language', 'processing', 'in', 'python']
```

PRACTICAL 4D

Tokenization using the spaCy library

Code:

```
import spacy  
nlp = spacy.blank("en")  
str = "I love to study Natrual Language processing in python"  
doc = nlp(str)  
words = [ word.text for word in doc]  
print(words)
```

Output

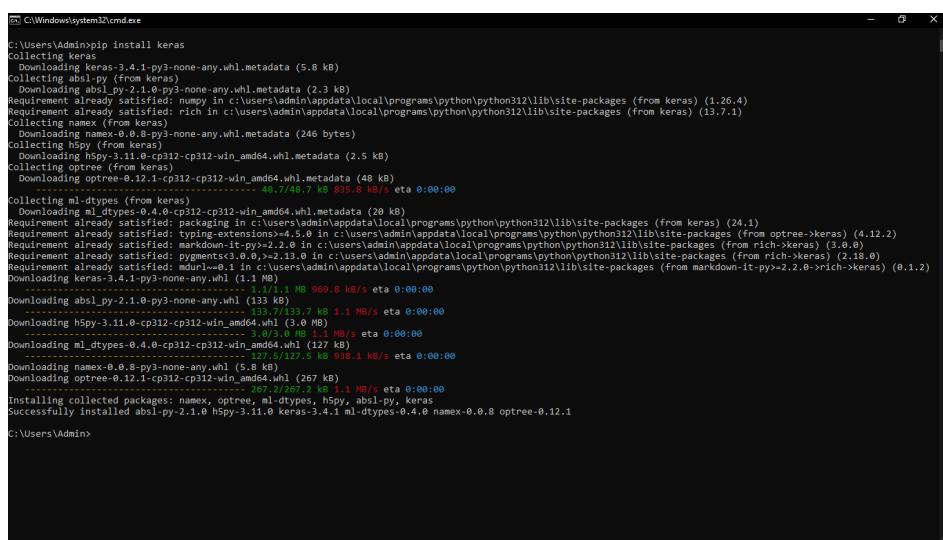
```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/4d.py =====  
>>> | ['I', 'love', 'to', 'study', 'Natrual', 'Language', 'processing', 'in', 'python']
```

PRACTICAL 4E

Tokenization using Keras

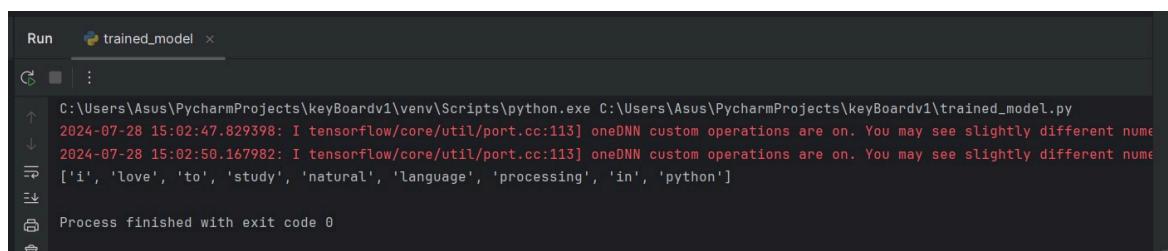
Code:

```
import tensorflow
import keras
from tensorflow.keras.preprocessing.text import text_to_word_sequence
# Create a string input
str = "I love to study Natural Language Processing in Python"
# tokenizing the text
tokens = text_to_word_sequence(str)
print(tokens)
```



```
C:\Windows\system32\cmd.exe
C:\Users\Admin>pip install keras
Collecting keras
  Downloading keras-2.4.1-py3-none-any.whl.metadata (5.8 kB)
Collecting absl-py (from keras)
  Downloading absl-py-2.1.0-py3-none-any.whl.metadata (2.3 kB)
Requirement already satisfied: numpy in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from keras) (1.26.4)
Requirement already satisfied: rich in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from keras) (13.7.1)
Collecting optree (from keras)
  Downloading optree-0.12.1-cp312-cp312-win_amd64.whl.metadata (48 kB)
Collecting h5py (from keras)
  Downloading h5py-3.11.0-cp312-cp312-win_amd64.whl.metadata (2.5 kB)
Collecting optree (from keras)
  Downloading optree-0.12.1-cp312-cp312-win_amd64.whl.metadata (48 kB)
Collecting ml-dtypes (from keras)
  Downloading ml-dtypes-4.0-cp312-cp312-win_amd64.whl.metadata (28 kB)
Requirement already satisfied: packaging in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from keras) (24.1)
Requirement already satisfied: typing-extensions<4.5.0 in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from optree->keras) (4.12.2)
Requirement already satisfied: markdown-it-py>=2.2.0 in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from rich->keras) (3.0.0)
Requirement already satisfied: pygments<3.0.0,>=2.13.0 in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from rich->keras) (2.18.0)
Requirement already satisfied: mdurl=>0.1 in c:\users\admin\appdata\local\programs\python\python312\lib\site-packages (from markdown-it-py>=2.2.0->rich->keras) (0.1.2)
Collecting keras
  Downloading keras-3.4.1-py3-none-any.whl (1:1/1.1 MB 969.8 kB/s eta 0:00:00
Collecting ml-dtypes (from keras)
  Downloading ml-dtypes-4.0-cp312-cp312-win_amd64.whl (133 kB)
  133/133.7 kB 1.1 MB/s eta 0:00:00
Collecting h5py-3.11.0-cp312-cp312-win_amd64.whl (3.8 MB)
  3.8 MB/3.8 MB 1.1 MB/s eta 0:00:00
Collecting ml-dtypes-0.4.0-cp312-cp312-win_amd64.whl (127 kB)
  127.5/127.5 kB 938.1 kB/s eta 0:00:00
Collecting namex-0.0.8-py3-none-any.whl (5.8 kB)
  Downloading namex-0.0.8-py3-none-any.whl (267 kB)
  267.3/267.3 kB 1.1 MB/s eta 0:00:00
Installing collected packages: namex, optree, ml-dtypes, h5py, absl-py, keras
Successfully installed absl-py-2.1.0 h5py-3.11.0 keras-3.4.1 ml-dtypes-0.4.0 namex-0.0.8 optree-0.12.1
C:\Users\Admin>
```

Output:



```
Run   trained_model x
G: [ ] :
C:\Users\Asus\PycharmProjects\keyBoardv1\venv\Scripts\python.exe C:\Users\Asus\PycharmProjects\keyBoardv1\trained_model.py
2024-07-28 15:02:47.829398: I tensorflow/core/util/port.cc:113] oneDNN custom operations are on. You may see slightly different nume
[ 'i', 'love', 'to', 'study', 'natural', 'language', 'processing', 'in', 'python' ]
Process finished with exit code 0
```

PRACTICAL 4F

Tokenization using Gensim

Code:

```
from gensim.utils import tokenize  
str = "I love to study Natrual Language processing in python"  
print(list(tokenize(str)))
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/4f.py =====  
>>> ['I', 'love', 'to', 'study', 'Natrual', 'Language', 'processing', 'in', 'python']
```

PRACTICAL 6A

Illustrate part of speech tagging

Part of speech Tagging and chunking of user defined text.

Code:

```
import nltk  
from nltk import tokenize  
nltk.download('punkt')  
from nltk import tag  
from nltk import chunk  
nltk.download('averaged_perceptron_tagger')  
nltk.download('maxent_ne_chunker')  
nltk.download('words')
```

```
para = "Hello! My name is anjali nimje . Today you'll be learning NLTK."  
sents = tokenize.sent_tokenize(para)  
print("\nsentence tokenization\n=====\\n",sents)
```

```
# word tokenization  
print("\nword tokenization\\n=====\\n")  
for index in range(len(sents)):  
    words = tokenize.word_tokenize(sents[index])  
    print(words)
```

```
# POS Tagging  
tagged_words = []  
for index in range(len(sents)):  
    tagged_words.append(tag.pos_tag(words))
```

```
print("\nPOS Tagging\n=====\\n",tagged_words)
```

chunking

```
tree = []
```

```
for index in range(len(sents)):
```

```
tree.append(chunk.ne_chunk(tagged_words[index]))
```

```
print("\nchunking\n=====\\n")
```

print(tree)

Output:

PRACTICAL 6B

Named Entity recognition using user defined text.

Code:

```
import spacy
nlp = spacy.load("en_core_web_sm")
text = ("when seabastian Thrun started working on self-driving cars at"
"Google in 2007, few people oustside of the company took him"
"Seroiusly,I can tell you very senior CEO's of major American"
"car companies would shake my hand and turn away beacause I wasn't"
"worth talking to,said Thurn , in an interview with recorder earlier")doc =
nlp(text)
print("nouns:\n",[chunk.text for chunk in doc.noun_chunks])
print("verbs",[token.lemma_ for token in doc if token.pos_ == "VERB"])
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/6b.py =====
Noun phrases: ['Sebastian Thrun', 'self-driving cars', 'Google', 'few people', 'the company', 'him', 'I', 'you', 'very senior CEOs', 'major American car companies'
, 'my hand', 'I', 'Thrun', 'an interview', 'Recode']
Verbs: ['start', 'work', 'drive', 'take', 'tell', 'shake', 'turn', 'talk', 'say']
>>>
```

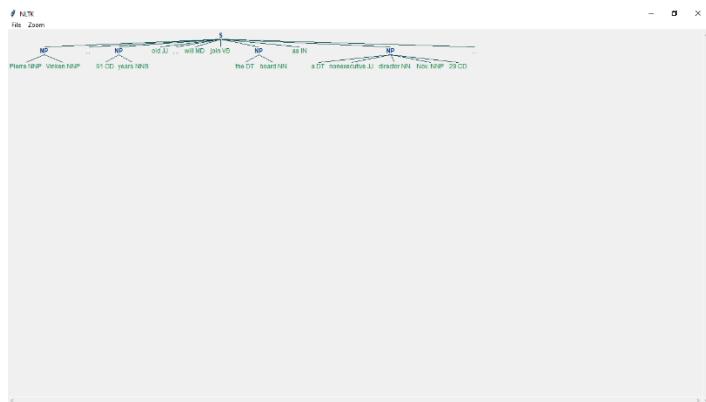
PRACTICAL 6C

Named Entity recognition with diagram using NLTK corpus – treebank.

Code:

```
import nltk  
nltk.download('treebank')  
from nltk.corpus import treebank_chunk  
treebank_chunk.tagged_sents()[0]  
treebank_chunk.chunked_sents()[0]  
treebank_chunk.chunked_sents()[0].draw()
```

Output:



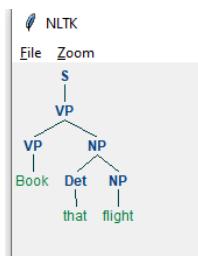
PRACTICAL 7A

Define grammar using nltk. Analyze a sentence using the same.

Code:

```
import nltk  
from nltk import tokenize  
grammar1 = nltk.CFG.fromstring(""  
    S -> VP  
    VP -> VP NP  
    NP -> Det NP  
    Det -> 'that'  
    NP -> singular Noun  
    NP -> 'flight'  
    VP -> 'Book'  
    "")  
  
sentence = "Book that flight"  
for index in range(len(sentence)):  
    all_tokens = tokenize.word_tokenize(sentence)  
    print(all_tokens)  
parser = nltk.ChartParser(grammar1)  
for tree in parser.parse(all_tokens):  
    print(tree)  
    tree.draw()
```

Output:



PRACTICAL 7B

Accept the input string with Regular expression of Finite Automaton: 101

Code:

```
def FA(s):
    #if the length is less than 3 then it can't be accepted, Therefore end the process.
    if len(s)<3:
        return "Rejected"
    #first three characters are fixed. Therefore, checking them using index
    if s[0]=='1':
        if s[1]=='0':
            if s[2]=='1':
                # After index 2 only "1" can appear. Therefore break the process if any other
                # character is detected
                for i in range(3,len(s)):
                    if s[i]!='1':
                        return "Rejected"
                return "Accepted" # if all 4 nested if true
            return "Rejected" # else of 3rd if
        return "Rejected" # else of 2nd if
    return "Rejected" # else of 1st if
inputs=['1','10101','101','10111','01010','100','','1011101','1011111']
for i in inputs:
    print(FA(i))
```

Output:

```
>>> ===== RESTART: D:/nlp/7b.py =====
Rejected
Rejected
Rejected
Accepted
None
Rejected
Rejected
Accepted
Accepted
>>> |
```

PRACTICAL 7C

Accept the input string with Regular expression of FA: $(a+b)^*bba$.

Code:

```
def FA(s):
    size=0
    #scan complete string and make sure that it contains only 'a' & 'b'
    for i in s:
        if i=='a' or i=='b':
            size+=1
        else:
            return "Rejected"
    #After checking that it contains only 'a' & 'b'
    #check it's length it should be 3 atleast
    if size>=3:
        #check the last 3 elements
        if s[size-3]=='b':
            if s[size-2]=='b':
                if s[size-1]=='a':
                    return "Accepted" # if all 4 if true
                return "Rejected" # else of 4th if
            return "Rejected" # else of 3rd if
        return "Rejected" # else of 2nd if
    return "Rejected" # else of 1st if
inputs=['bba', 'ababbba', 'abba','abb', 'baba','bbb',"]
for i in inputs:
    print(FA(i))
```

Output:

```
>>> ===== RESTART: D:/nlp/7c.py =====
Accepted
Accepted
Accepted
Rejected
Rejected
Rejected
Rejected
>>> |
```

PRACTICAL 7D

Implementation of Deductive Chart Parsing using context free grammar and a given sentence.

Code:

```
import nltk
from nltk import tokenize
grammar1 = nltk.CFG.fromstring("""
    S -> NP VP
    PP -> P NP
    NP -> Det N | Det N PP | 'T'
    VP -> V NP | VP PP
    Det -> 'a' | 'my'
    N -> 'bird' | 'balcony'
    V -> 'saw'
    P -> 'in'
    """)

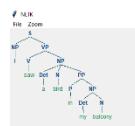
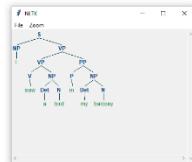
sentence = "I saw a bird in my balcony"

for index in range(len(sentence)):
    all_tokens = tokenize.word_tokenize(sentence)
    print(all_tokens)

# all_tokens = ['I', 'saw', 'a', 'bird', 'in', 'my', 'balcony']
parser = nltk.ChartParser(grammar1)
for tree in parser.parse(all_tokens):
    print(tree)
    tree.draw()

Output:
```

```
[1]: In [1]: Help Cell Kernel Window Help
Python 3.7.4 (v3.7.4:e09359112e, Jun 6 2018, 14:49:54) [MSC v.1940 64 bit (AMD64)] in win32
Type "help()" for help.
>>> print("Hello, world!")
Hello, world!
>>> print("Hello, world!", "Hello, world!")
Hello, world! Hello, world!
>>> print("Hello, world!", "Hello, world!", sep=" ")
Hello, world. Hello, world.
>>> print("Hello, world!", "Hello, world!", sep=" ", end=" ")
Hello, world!Hello, world!
>>> print("Hello, world!", "Hello, world!", sep=" ", end=" ")
Hello, world!Hello, world!
```



PRACTICAL 8

Study PorterStemmer, LancasterStemmer, RegexpStemmer, SnowballStemmer ,Study WordNetLemmatizer

Code:

```
import nltk
from nltk.stem import PorterStemmer
word_stemm = PorterStemmer()
print(word_stemm.stem('writing'))

import nltk
from nltk.stem import LancasterStemmer
lanc_stemm = LancasterStemmer()
print(lanc_stemm.stem('writing'))

import nltk
from nltk.stem import RegexpStemmer
Reg_stemm = RegexpStemmer('ing\$|s\$|e\$|able$', min=4)
print(Reg_stemm.stem('writing'))

import nltk
from nltk.stem import SnowballStemmer
english_stemm = SnowballStemmer('english')
print(english_stemm.stem('writing'))

import nltk
from nltk.stem import WordNetLemmatizer
lemetizer = WordNetLemmatizer()
print("word:\t lemma")
print("rocks:", lemetizer.lemmatize("rocks"))
print("corpa:", lemetizer.lemmatize("corpora"))
```

Output:

```
----- RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/8.py -----
PorterStemmer:
write

LancasterStemmer:
writ

RegexpStemmer:
writ

SnowballStemmer:
write

WordNetLemmatizer:
word : lemma
rocks : rock
corpora : corpus
better : good
>>> |
```

PRACTICAL 9

Implement Naive Bayes classifier

Code:

```
#pip install pandas
#pip install sklearn
import pandas as pd
import numpy as np
sms_data =
pd.read_csv('C:\\\\Users\\\\Student\\\\Desktop\\\\Mayuri\\\\spam.csv',encoding='latin-1')
import re
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords
from nltk.stem.porter import PorterStemmer
stemming = PorterStemmer()
corpus = []
for i in range (0,len(sms_data)):
    s1 = re.sub('[^a-zA-Z]',repl="",string=sms_data['v2'][i])
    s1.lower()
    s1 = s1.split()
    s1 = [stemming.stem(word) for word in s1 if word not in
set(stopwords.words('english'))]
    s1 = ' '.join(s1)
    corpus.append(s1)

from sklearn.feature_extraction.text import CountVectorizer
countvectorizer =CountVectorizer()
x = countvectorizer.fit_transform(corpus).toarray()
```

```
print(x)
y = sms_data['v1'].values
print(y)

from sklearn.model_selection import train_test_split
x_train,x_test,y_train,y_test = train_test_split(x,y,test_size = 0.3,
stratify=y,random_state=2)

#Multinomial Naïve Bayes.
from sklearn.naive_bayes import MultinomialNB
multinomialnb = MultinomialNB()
multinomialnb.fit(x_train,y_train)

# Predicting on test data:
y_pred = multinomialnb.predict(x_test)
print(y_pred)

#Results of our Models
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.metrics import accuracy_score
print(classification_report(y_test,y_pred))
print("accuracy_score: ",accuracy_score(y_test,y_pred))
```

Output:

```
[nltk_data]  Downloading package stopwords to
[nltk_data]    C:\Users\Admin\AppData\Roaming\nltk_data...
[nltk_data]      Package stopwords is already up-to-date!
[[0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 ...
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
 [0 0 0 ... 0 0 0]
['ham' 'ham' 'spam' ... 'ham' 'ham' 'ham']
['ham' 'ham' 'ham' ... 'ham' 'ham' 'ham']
          precision    recall   f1-score   support
              ham        0.87       1.00     0.93    1448
              spam        1.00       0.02     0.04     224

          accuracy         0.87    1672
      macro avg        0.93       0.51     0.49    1672
  weighted avg        0.89       0.87     0.81    1672

accuracy_score: 0.8690191387558809
>>>
```

PRACTICAL 10 A.1

Speech tagging using spacy

Code:

```
import spacy  
sp = spacy.load('en_core_web_sm')  
sen = sp(u"I like to play football. I hated it in my childhood though")  
print(sen.text)  
print(sen[7].pos_)  
print(sen[7].tag_)  
print(spacy.explain(sen[7].tag_))  
for word in sen:  
    print(f'{word.text}:{word.pos_}:{word.tag_}'  
          {spacy.explain(word.tag_)}')
```

```
sen = sp(u'Can you google it?')  
word = sen[2]  
print(f'{word.text}:{word.pos_}:{word.tag_}'  
      {spacy.explain(word.tag_)}')  
sen = sp(u'Can you search it on google?')  
word = sen[5]  
print(f'{word.text}:{word.pos_}:{word.tag_}'  
      {spacy.explain(word.tag_)}')  
  
#Finding the Number of POS Tags  
sen = sp(u"I like to play football. I hated it in my childhood though")  
num_pos = sen.count_by(spacy.attrs.POS)  
num_pos  
for k,v in sorted(num_pos.items()):  
    print(f'{k}: {v}')  
  
#Visualizing Parts of Speech Tags
```

```
from spacy import displacy
```

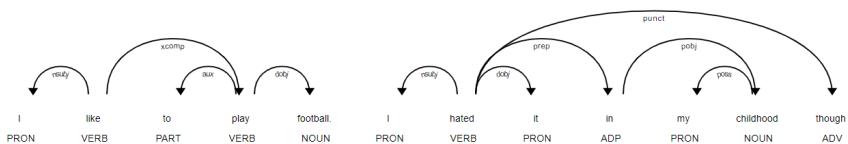
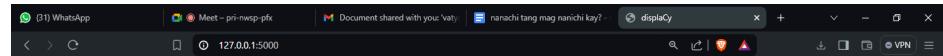
```
sen = sp(u"I like to play football. I hated it in my childhood though")
```

```
displacy.serve(sen, style='dep', options={'distance': 120})
```

Output:

```
===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/10a1.py =====
I like to play football. I hated it in my childhood though
VERB
VBD
verb, past tense
PRON PRP pronoun, personal
like VERB VBP verb, non-3rd person singular present
to PART TO infinitival "to"
play VERB VB verb, base form
Football NOUN NN noun, singular or mass
PUNCT .
punctuation mark, sentence closer
I PRON PRP pronoun, personal
hated VERB VBD verb, past tense
it PRON PRP pronoun, personal
in ADP IN conjunction, subordinating or preposition
my PRON PRPS pronoun, possessive
childhood NOUN NN noun, singular or mass
though ADV RB adverb
google VERB VB verb, base form
google PROPN NNP noun, proper singular
85. ADP : 1
86. ADV : 1
92. NOUN : 2
94. PART : 1
95. PRON : 4
97. PUNCT : 1
100. VERB : 3

Using the 'dep' visualizer
Serving on http://0.0.0.0:5000 ...
```



PRACTICAL 10 A.2

Speech tagging using nltk

Code:

```
import nltk  
nltk.download('state_union')  
nltk.download('punkt')  
nltk.download('averaged_perceptron_tagger')  
from nltk.corpus import state_union  
from nltk.tokenize import PunktSentenceTokenizer
```

```
#create our training and testing data:
```

```
train_text = state_union.raw("2005-GWBush.txt")  
sample_text = state_union.raw("2006-GWBush.txt")
```

```
#train the Punkt tokenizer like:
```

```
custom_sent_tokenizer = PunktSentenceTokenizer(train_text)  
# tokenize:  
tokenized = custom_sent_tokenizer.tokenize(sample_text)  
  
def process_content():  
    try:  
        for i in tokenized[:2]:  
            words = nltk.word_tokenize(i)  
            tagged = nltk.pos_tag(words)  
            print(tagged)  
    except Exception as e:  
        print(str(e))  
  
process_content()
```

Output:

```
===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/10a2.py =====
[nltk_data]  Downloading package state_union to C:/Users/Admin/AppData/Roaming/nltk_data...
[nltk_data]    Unzipping corpora/state_union.zip.
[nltk_data]  Downloading package punkt to C:/Users/Admin/AppData/Roaming/nltk_data...
[nltk_data]    Package punkt is already up-to-date!
[nltk_data]  Downloading package averaged_perceptron_tagger to C:/Users/Admin/AppData/Roaming/nltk_data...
[nltk_data]    Package averaged_perceptron_tagger is already up-to-date.
>>> 
```

PRACTICAL 10 B.1

Usage of Give and Gave in the Penn Treebank sample

Code:

```
import nltk  
  
import nltk.parse.viterbi  
  
import nltk.parse.pchart  
  
def give(t):  
  
    return t.label() == 'VP' and len(t) > 2 and t[1].label() == 'NP'\  
    and (t[2].label() == 'PP-DTV' or t[2].label() == 'NP')\  
    and ('give' in t[0].leaves() or 'gave' in t[0].leaves())  
  
def sent(t):  
  
    return ' '.join(token for token in t.leaves() if token[0] not in '*-0')  
  
def print_node(t, width):  
  
    output = "%s %s: %s / %s: %s" %\  
        (sent(t[0]), t[1].label(), sent(t[1]), t[2].label(), sent(t[2]))  
  
    if len(output) > width:  
  
        output = output[:width] + "..."  
  
    print (output)  
  
for tree in nltk.corpus.treebank.parsed_sents():  
  
    for t in tree.subtrees(give):  
  
        print_node(t, 72)
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/1obj.py =====  
give NP: the charts / NP: a standing ovation  
give NP: advertisers / NP: discontents for maintaining or increasing ad sp...  
give NP: it / PP-DTV: to the politicians  
give NP: them / NP: similar help  
give NP: them / NP:  
give NP: federal judges / NP: a raise  
give NP: consumers / NP: the straight scoop on the U.S. waste crisis  
give NP: Mitsui / NP: access to a high-tech medical product  
give NP: Mitsubishi / NP: a window on the U.S. glass industry  
give NP: the chairman / NP: the right to sue stockholders , nor to ...  
give NP: your Foster Savings Institution / NP: the gift of hope and free...  
give NP: market operators / NP: the authority to suspend trading in futu...  
give NP: quick approval / PP-DTV: to $ 3.18 billion in supplemental appr...  
give NP: the International Development / NP: up to 50 days to review any...  
give NP: the President / NP: such power  
give NP: me / NP: the heebie-jeebies  
give NP: holders / NP: the right , but not the obligation , to buy a cal...  
give NP: Mr. Thomas / NP: only a `` qualified '' rating , rather than '...  
give NP: the president / NP: line-item veto power
```

PRACTICAL 10 B.2

Probabilistic parser

Code:

```
import nltk
from nltk import PCFG
grammar = PCFG.fromstring("""
NP -> NNS [0.5] | JJ NNS [0.3] | NP CC NP [0.2]
NNS -> "men" [0.1] | "women" [0.2] | "children" [0.3] | NNS CC NNS [0.4]
JJ -> "old" [0.4] | "young" [0.6]
CC -> "and" [0.9] | "or" [0.1]
")
print(grammar)
viterbi_parser = nltk.ViterbiParser(grammar)
token = "old men and women".split()
obj = viterbi_parser.parse(token)
print("Output: ")
for x in obj:
    print(x)
```

Output:

```
>>> ===== RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/10b2.py =====
Grammar with 11 productions (start state = NP)
NP -> NNS [0.5]
NP -> JJ NNS [0.3]
NP -> NP CC NP [0.2]
NNS -> "men" [0.1]
NNS -> "women" [0.2]
NNS -> "children" [0.3]
NNS -> NNS CC NNS [0.4]
JJ -> "old" [0.4]
JJ -> "young" [0.6]
CC -> "and" [0.9]
CC -> "or" [0.1]
Output:
(NP (JJ old) (NNS (NNS men) (CC and) (NNS women))) (p=0.000864)
>>> |
```

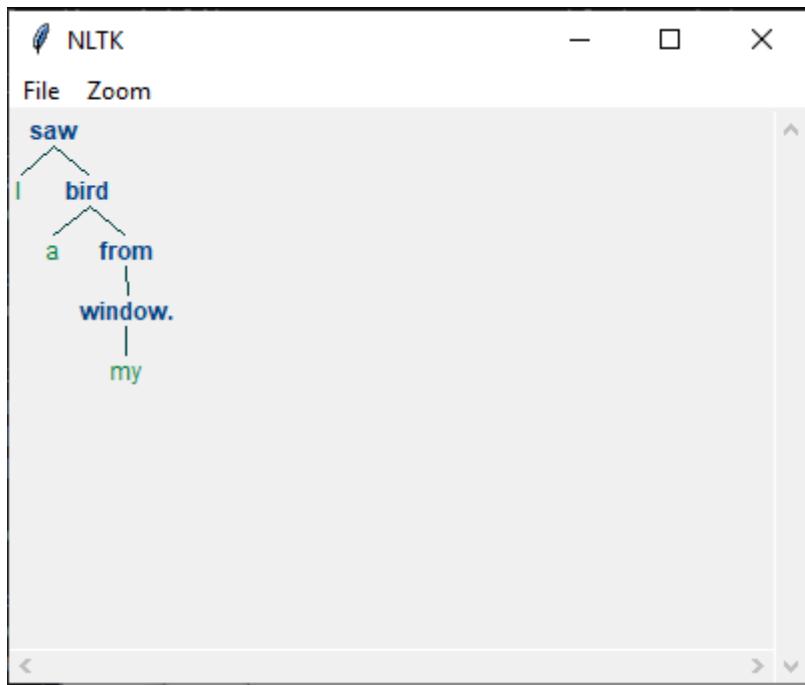
PRACTICAL 10C

Malt Parsing

Code:

```
from nltk.parse import malt  
  
mp = malt.MaltParser('maltparser-1.9.2', 'engmalt.linear-1.7.mco')#file  
  
t = mp.parse_one('I saw a bird from my window.'.split()).tree()  
  
print(t)  
  
t.draw()
```

Output:



PRACTICAL 11A

Multiword Expressions in NLP

Code:

```
from nltk.tokenize import MWETokenizer  
  
from nltk import sent_tokenize, word_tokenize  
  
s = "Good cake cost Rs.1500\kg in Mumbai. Please buy me one of  
them.\n\nThanks."  
  
mwe = MWETokenizer([('New', 'York'), ('Hong', 'Kong')], separator='_')  
  
for sent in sent_tokenize(s):  
  
    print(mwe.tokenize(word_tokenize(sent)))
```

Output:

```
>>> |  
| = RESTART: C:/Users/Admin/AppData/Local/Programs/Python/Python312/11a.py  
| ['Good', 'cake', 'cost', 'Rs.1500\\kg', 'in', 'Mumbai', '.']  
| ['Please', 'buy', 'me', 'one', 'of', 'them', '.']  
| ['Thanks', '.']  
>>> |
```

PRACTICAL 11B

Normalized Web Distance and Word Similarity

Code:

```
import numpy as np
import re
import textrdistance # pip install textrdistance
# we will need scikit-learn>=0.21
import sklearn #pip install sklearn
from sklearn.cluster import AgglomerativeClustering
texts = [
'Reliance supermarket', 'Reliance hypermarket', 'Reliance', 'Reliance',
'Reliancedowntown', 'Relianc market',
'Mumbai', 'Mumbai Hyper', 'Mumbai dxb', 'mumbai airport',
'k.m trading', 'KM Trading', 'KM trade', 'K.M. Trading', 'KM.Trading'
]
def normalize(text):
    """ Keep only lower-cased text and numbers"""
    return re.sub('[^a-zA-Z0-9]+', ' ', text.lower())
def group_texts(texts, threshold=0.4):
    """ Replace each text with the representative of its cluster"""
    normalized_texts = np.array([normalize(text) for text in texts])
    distances = 1 - np.array([
        [textrdistance.jaro_winkler(one, another) for one in normalized_texts]
        for another in normalized_texts
    ])
    clustering = AgglomerativeClustering(
        distance_threshold=threshold, # this parameter needs to be tuned
        carefully
```

```
metric="precomputed", linkage="complete", n_clusters=None  
).fit(distances)  
  
centers = dict()  
  
for cluster_id in set(clustering.labels_):  
    index = clustering.labels_ == cluster_id  
    centrality = distances[:, index][index].sum(axis=1)  
    centers[cluster_id] = normalized_texts[index][centrality.argmin()]  
  
    return [centers[i] for i in clustering.labels_]  
  
print(group_texts(texts))
```

Output:



PRACTICAL 11C

Word Sense Disambiguation

Code:

```
from nltk.corpus import wordnet as wn

def get_first_sense(word, pos=None):
    if pos:
        synsets = wn.synsets(word,pos)
    else:
        synsets = wn.synsets(word)
    return synsets[0]

best_synset = get_first_sense('bank')
print ('%s: %s' % (best_synset.name, best_synset.definition))
best_synset = get_first_sense('set','n')
print ('%s: %s' % (best_synset.name, best_synset.definition))
best_synset = get_first_sense('#set','v')
print ('%s: %s' % (best_synset.name, best_synset.definition))
```

Output:

```
>>> ===== RESTART: C:\Users\Admin\AppData\Local\Programs\Python\Python312\l1c.py =====
<bound method Synset.name of Synset('bank.n.01')>: <bound method Synset.definition of Synset('bank.n.01')>
<bound method Synset.name of Synset('set.n.01')>: <bound method Synset.definition of Synset('set.n.01')>
<bound method Synset.name of Synset('put.v.01')>: <bound method Synset.definition of Synset('put.v.01')>
>>>
```

