

**Deccan Education Society's**  
**Kirti M. Doongursee College of Arts, Science and Commerce**  
**[AUTONOMOUS]**



**M.Sc. [Computer Science]**  
**Practical Journal**

**Course Name: Natural Language Processing**

**Seat Number [    ]**

**(Academic Year 2024-2025)**  
**Department of Computer Science and Information Technology**

**Department of Computer Science and Information  
Technology Deccan Education Society's  
Kirti M. Doongursee College of Arts, Science and Commerce  
[AUTONOMOUS]**

**C E R T I F I C A T E**

This is to certify that Mr./Miss \_\_\_\_\_ of M.Sc. (C.S.) with  
Seat No. \_\_\_\_ has complete \_\_\_\_\_ Practical of Paper( Course Name- **Natural  
Language Processing**) under my supervision in this College during the Second  
Semester of academic year 2024-2025

Prof. Jaymala Deshpande  
Lecturer-In-Charge

Dr. Apurva Yadav  
H.O.D.  
Department of  
Computer Science & IT

Date:     /     /2024

Date:

**Examined by:**

**Remarks:**

Date

\_\_\_\_\_

## Index

Sr.No	Practical Names	Sign
1	A. Convert the given text to speech.	
	B. Convert audio file to Text	
2	A. Study of Brown Corpus with various methods like fields, raw, words, sents ,categories.	
	B.Create and use your own corpora (plaintext, categorical)	
	C. Study of tagged corpora with method like tagged_sents , tagged_words.	
	D. Map Words to Properties Using Python Dictionaries	
3	A.Study Default tagger	
	B. Study Unigram Tagger	
4	A. Study of Wordnet Dictionary with methods as synsets, definitions, examples, antonyms	
	B.Write a program using python to find synonym and antonym of word "active" using Wordnet.	
	C. Compare two nouns	
5	A. Tokenization using Python's split() function	
	B. Tokenization using Regular Expressions (RegEx)	
	C. Tokenization using nltk	
	D. Tokenization using Spacy	
	E. Tokenization using Keras	
	F.Tokenization using Gensim	
6	A. Named Entity recognition using user defined text.	
	B. Named Entity recognition with diagram using NLTK corpus – treebank. Program:	
7	A. Define grammar using nltk. Analyze a sentence using the same	
	B. Implementation of Deductive Chart Parsing using context free grammar and a given sentence	
8	Study PorterStemmer, LancasterStemmer, RegexpStemmer, SnowballStemmer Study WordNetLemmatizer	
9	Finite Automata	

# Practical - 1(a)

**Aim:** Convert the given text to speech.

**Requirements:**

```
pip install nltk
pip install gtts
pip install --upgrade wheel
pip install playsound
```

**Program:**

```
from gtts import gTTS
from playsound import playsound

mytext = "Hello World"
lang = "en"

myobj = gTTS(text = mytext, lang = lang, slow = False)

myobj.save("./myFile.mp3")
playsound("./myFile.mp3")
```

**Output:**

myFile.mp3 audio file is getting created and it plays the file with playsound() method, while running the program

# Practical - 1(b)

**Aim: Convert audio file to Text.**

**Requirements:**

```
pip install SpeechRecognition
```

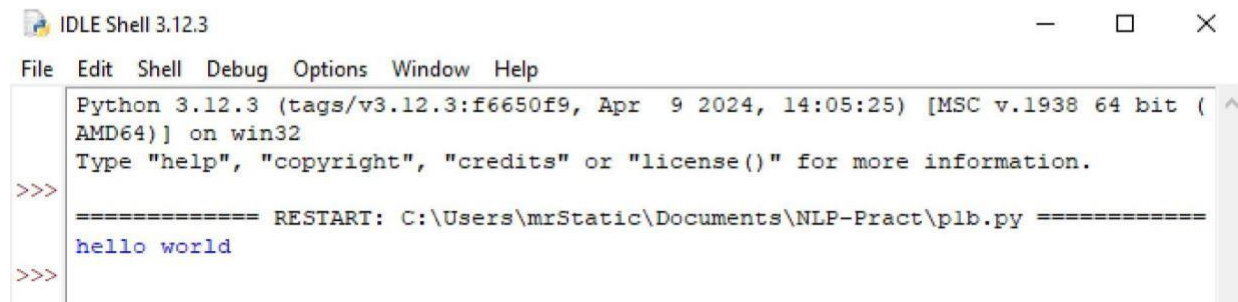
**Program:**

```
import speech_recognition as sr

filename = "./myFile.wav"
r = sr.Recognizer()

with sr.AudioFile(filename) as source:
    audio_data = r.record(source)
    text = r.recognize_google(audio_data)
    print(text)
```

**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: C:\Users\mrStatic\Documents\NLP-Pract\plb.py =====
hello world
>>>
```

## Practical - 2(a)

**Aim: Study of Brown Corpus with various methods like filelds, raw, words, sents, categories.**

### Requirements:

```
pip install nltk
nltk.download("brown")
```

### Program:

```
import nltk
nltk.download("brown")
from nltk.corpus import brown

print("File ids of brown corpus\n", brown.fileids())

ca01 = brown.words("ca01")
print("\nca01 has following words:\n", ca01)
print("\nca01 has ", len(ca01), " words")

print("\n\nCategories or file in brown corpus: \n")
print(brown.categories())

print("\n\nStatistics for each text:\n")
print("Avg-Word-Len\tAvg-Sentence-Len\tNo.-of-Times-Each-Word-Appears-On-Avg\t\tFile-Name")
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", num_words/num_sents,
          "\t\t\t", int(num_words/num_vocab), "\t\t\t", fileid)
```

### Output:

---

```

===== RESTART: C:/Users/mrStatic/Document/NLP-Pract/p3a.py =====
[nltk_data] Downloading package brown to
C:\Users\mrStatic\AppData\Roaming\nltk_data...
[nltk_data]   Unzipping corpora\brown.zip.
File ids of brown corpus
['ca01', 'ca02', 'ca03', 'ca04', 'ca05', 'ca06', 'ca07', 'ca08', 'ca09', 'ca10', 'ca11', 'ca12', 'ca13', 'ca14', 'ca15', 'ca16', 'ca17', 'ca18', 'ca19', 'ca20', 'ca21', 'ca22', 'ca23', 'ca24',
'ca25', 'ca26', 'ca27', 'ca28', 'ca29', 'ca30', 'ca31', 'ca32', 'ca33', 'ca34', 'ca35', 'ca36', 'ca37', 'ca38', 'ca39', 'ca40', 'ca41', 'ca42', 'ca43', 'ca44', 'ca45', 'ca46', 'ca47', 'ca48',
'ca49', 'ca50', 'ca51', 'ca52', 'ca53', 'ca54', 'ca55', 'ca56', 'ca57', 'ca58', 'ca59', 'ca60', 'ca61', 'ca62', 'ca63', 'ca64', 'ca65', 'ca66', 'ca67', 'ca68', 'ca69', 'ca70', 'ca71', 'ca72', 'ca73',
'ca74', 'ca75', 'ca76', 'ca77', 'ca78', 'ca79', 'ca80', 'ca81', 'ca82', 'ca83', 'ca84', 'ca85', 'ca86', 'ca87', 'ca88', 'ca89', 'ca90', 'ca91', 'ca92', 'ca93', 'ca94', 'ca95', 'ca96', 'ca97',
'ca98', 'ca99', 'ca100', 'ca101', 'ca102', 'ca103', 'ca104', 'ca105', 'ca106', 'ca107', 'ca108', 'ca109', 'ca110', 'ca111', 'ca112', 'ca113', 'ca114', 'ca115', 'ca116', 'ca117', 'ca118', 'ca119',
'ca120', 'ca121', 'ca122', 'ca123', 'ca124', 'ca125', 'ca126', 'ca127', 'ca128', 'ca129', 'ca130', 'ca131', 'ca132', 'ca133', 'ca134', 'ca135', 'ca136', 'ca137', 'ca138', 'ca139', 'ca140',
'ca141', 'ca142', 'ca143', 'ca144', 'ca145', 'ca146', 'ca147', 'ca148', 'ca149', 'ca150', 'ca151', 'ca152', 'ca153', 'ca154', 'ca155', 'ca156', 'ca157', 'ca158', 'ca159', 'ca160', 'ca161', 'ca162',
'ca163', 'ca164', 'ca165', 'ca166', 'ca167', 'ca168', 'ca169', 'ca170', 'ca171', 'ca172', 'ca173', 'ca174', 'ca175', 'ca176', 'ca177', 'ca178', 'ca179', 'ca180', 'ca181', 'ca182', 'ca183', 'ca184',
'ca185', 'ca186', 'ca187', 'ca188', 'ca189', 'ca190', 'ca191', 'ca192', 'ca193', 'ca194', 'ca195', 'ca196', 'ca197', 'ca198', 'ca199', 'ca200', 'ca201', 'ca202', 'ca203', 'ca204', 'ca205',
'ca206', 'ca207', 'ca208', 'ca209', 'ca210', 'ca211', 'ca212', 'ca213', 'ca214', 'ca215', 'ca216', 'ca217', 'ca218', 'ca219', 'ca220', 'ca221', 'ca222', 'ca223', 'ca224', 'ca225', 'ca226', 'ca227',
'ca228', 'ca229', 'ca230', 'ca231', 'ca232', 'ca233', 'ca234', 'ca235', 'ca236', 'ca237', 'ca238', 'ca239', 'ca240', 'ca241', 'ca242', 'ca243', 'ca244', 'ca245', 'ca246', 'ca247', 'ca248', 'ca249',
'ca250', 'ca251', 'ca252', 'ca253', 'ca254', 'ca255', 'ca256', 'ca257', 'ca258', 'ca259', 'ca260', 'ca261', 'ca262', 'ca263', 'ca264', 'ca265', 'ca266', 'ca267', 'ca268', 'ca269', 'ca270', 'ca271',
'ca272', 'ca273', 'ca274', 'ca275', 'ca276', 'ca277', 'ca278', 'ca279', 'ca280', 'ca281', 'ca282', 'ca283', 'ca284', 'ca285', 'ca286', 'ca287', 'ca288', 'ca289', 'ca290', 'ca291', 'ca292', 'ca293',
'ca294', 'ca295', 'ca296', 'ca297', 'ca298', 'ca299', 'ca300', 'ca301', 'ca302', 'ca303', 'ca304', 'ca305', 'ca306', 'ca307', 'ca308', 'ca309', 'ca310', 'ca311', 'ca312', 'ca313', 'ca314', 'ca315',
'ca316', 'ca317', 'ca318', 'ca319', 'ca320', 'ca321', 'ca322', 'ca323', 'ca324', 'ca325', 'ca326', 'ca327', 'ca328', 'ca329', 'ca330', 'ca331', 'ca332', 'ca333', 'ca334', 'ca335', 'ca336', 'ca337',
'ca338', 'ca339', 'ca340', 'ca341', 'ca342', 'ca343', 'ca344', 'ca345', 'ca346', 'ca347', 'ca348', 'ca349', 'ca350', 'ca351', 'ca352', 'ca353', 'ca354', 'ca355', 'ca356', 'ca357', 'ca358', 'ca359',
'ca360', 'ca361', 'ca362', 'ca363', 'ca364', 'ca365', 'ca366', 'ca367', 'ca368', 'ca369', 'ca370', 'ca371', 'ca372', 'ca373', 'ca374', 'ca375', 'ca376', 'ca377', 'ca378', 'ca379', 'ca380', 'ca381',
'ca382', 'ca383', 'ca384', 'ca385', 'ca386', 'ca387', 'ca388', 'ca389', 'ca390', 'ca391', 'ca392', 'ca393', 'ca394', 'ca395', 'ca396', 'ca397', 'ca398', 'ca399', 'ca400', 'ca401', 'ca402', 'ca403',
'ca404', 'ca405', 'ca406', 'ca407', 'ca408', 'ca409', 'ca410', 'ca411', 'ca412', 'ca413', 'ca414', 'ca415', 'ca416', 'ca417', 'ca418', 'ca419', 'ca420', 'ca421', 'ca422', 'ca423', 'ca424', 'ca425',
'ca426', 'ca427', 'ca428', 'ca429', 'ca430', 'ca431', 'ca432', 'ca433', 'ca434', 'ca435', 'ca436', 'ca437', 'ca438', 'ca439', 'ca440', 'ca441', 'ca442', 'ca443', 'ca444', 'ca445', 'ca446', 'ca447',
'ca448', 'ca449', 'ca450', 'ca451', 'ca452', 'ca453', 'ca454', 'ca455', 'ca456', 'ca457', 'ca458', 'ca459', 'ca460', 'ca461', 'ca462', 'ca463', 'ca464', 'ca465', 'ca466', 'ca467', 'ca468', 'ca469',
'ca470', 'ca471', 'ca472', 'ca473', 'ca474', 'ca475', 'ca476', 'ca477', 'ca478', 'ca479', 'ca480', 'ca481', 'ca482', 'ca483', 'ca484', 'ca485', 'ca486', 'ca487', 'ca488', 'ca489', 'ca490', 'ca491',
'ca492', 'ca493', 'ca494', 'ca495', 'ca496', 'ca497', 'ca498', 'ca499', 'ca500', 'ca501', 'ca502', 'ca503', 'ca504', 'ca505', 'ca506', 'ca507', 'ca508', 'ca509', 'ca510', 'ca511', 'ca512', 'ca513',
'ca514', 'ca515', 'ca516', 'ca517', 'ca518', 'ca519', 'ca520', 'ca521', 'ca522', 'ca523', 'ca524', 'ca525', 'ca526', 'ca527', 'ca528', 'ca529', 'ca530', 'ca531', 'ca532', 'ca533', 'ca534', 'ca535',
'ca536', 'ca537', 'ca538', 'ca539', 'ca540', 'ca541', 'ca542', 'ca543', 'ca544', 'ca545', 'ca546', 'ca547', 'ca548', 'ca549', 'ca550', 'ca551', 'ca552', 'ca553', 'ca554', 'ca555', 'ca556', 'ca557',
'ca558', 'ca559', 'ca560', 'ca561', 'ca562', 'ca563', 'ca564', 'ca565', 'ca566', 'ca567', 'ca568', 'ca569', 'ca570', 'ca571', 'ca572', 'ca573', 'ca574', 'ca575', 'ca576', 'ca577', 'ca578', 'ca579',
'ca580', 'ca581', 'ca582', 'ca583', 'ca584', 'ca585', 'ca586', 'ca587', 'ca588', 'ca589', 'ca590', 'ca591', 'ca592', 'ca593', 'ca594', 'ca595', 'ca596', 'ca597', 'ca598', 'ca599', 'ca600', 'ca601',
'ca602', 'ca603', 'ca604', 'ca605', 'ca606', 'ca607', 'ca608', 'ca609', 'ca610', 'ca611', 'ca612', 'ca613', 'ca614', 'ca615', 'ca616', 'ca617', 'ca618', 'ca619', 'ca620', 'ca621', 'ca622', 'ca623',
'ca624', 'ca625', 'ca626', 'ca627', 'ca628', 'ca629', 'ca630', 'ca631', 'ca632', 'ca633', 'ca634', 'ca635', 'ca636', 'ca637', 'ca638', 'ca639', 'ca640', 'ca641', 'ca642', 'ca643', 'ca644', 'ca645',
'ca646', 'ca647', 'ca648', 'ca649', 'ca650', 'ca651', 'ca652', 'ca653', 'ca654', 'ca655', 'ca656', 'ca657', 'ca658', 'ca659', 'ca660', 'ca661', 'ca662', 'ca663', 'ca664', 'ca665', 'ca666', 'ca667', 'ca668',
'ca669', 'ca670', 'ca671', 'ca672', 'ca673', 'ca674', 'ca675', 'ca676', 'ca677', 'ca67
```



## Practical - 2(b)

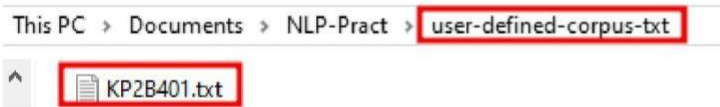
**Aim: Create and use your own corpora (plaintext, categorical)**

**Requirements:**

```
pip install nltk
nltk.download("punkt")
```

### Creating user defined corpus

1. In same program directory create new folder with name `user-defined-corpus-txt`
2. Open `user-defined-corpus-txt` and create new `txt` file with any name here I have given `KP2B401.txt`
3. Open the `txt` file and enter any sentence here I have given `The Quick Brown Fox Jumped Over The Lazy Dog.`



**Program:**

```
import nltk
nltk.download("punkt")
from nltk.corpus import PlaintextCorpusReader
corpus_root = "./user-defined-corpus-txt"
fileList = PlaintextCorpusReader(corpus_root, ".*")

print("\n File List \n")
print(fileList.fileids())
print(fileList.root)
print("\n\nStatistics for each text:\n")
print("Avg_Word_Len\tAvg_Sentence_Len\tno._of_Times_Each_Word_Appears_On_Avg\tFile_Name")

for fileid in fileList.fileids():
    num_chars = len(fileList.raw(fileid))
    num_words = len(fileList.words(fileid))
    num_sents = len(fileList.sents(fileid))
    num_vocab = len(set([w.lower() for w in fileList.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)
```



## Output:

```
===== RESTART: C:/Users/mrStatic/Documents/NLP-Pract/p2b.py =====  
[nltk_data] Downloading package punkt to  
[nltk_data] C:\Users\mrStatic\AppData\Roaming\nltk_data...  
[nltk_data] Package punkt is already up-to-date!
```

File List

```
['KP2B401.txt']  
C:\Users\mrStatic\Documents\NLP-Pract\user-defined-corpus-txt
```

Statistics for each text:

Avg_Word_Len	Avg_Sentence_Len	no._of_Times_Each_Word_Appears_On_Avg	File_Name
5	9	1	KP2B401.txt

---

## Practical - 2(c)

**Aim:** Study of tagged corpora with methods like `tagged_sents`, `tagged_words`.

### Requirements:

```
pip install nltk
import nltk
nltk.download('punkt')
nltk.download('words')
```

### Program:

```
import nltk
nltk.download('punkt')
nltk.download('words')
from nltk import tokenize

para = "Hello World! From kirti college. Today we will be learning NLTK."
sents = tokenize.sent_tokenize(para)
print("\nsentence tokenization\n=====\n", sents)

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

### Output:

```
= RESTART: /home/mrportable/Documents/MiniHacker2/MS-CSP-2/NLP/Pract/Pract-2  
/pract-2c.py
```

```
[nltk_data] Downloading package punkt to /home/mrportable/nltk_data...  
[nltk_data]   Unzipping tokenizers/punkt.zip.  
[nltk_data] Downloading package words to /home/mrportable/nltk_data...  
[nltk_data]   Unzipping corpora/words.zip.
```

```
sentence tokenization
```

```
=====
```

```
['Hello World!', 'From kirti college.', 'Today we will be learning NLTK.']
```

```
word tokenization
```

```
=====
```

```
['Hello', 'World', '!']
```

```
['From', 'kirti', 'college', '.']
```

```
['Today', 'we', 'will', 'be', 'learning', 'NLTK', '.']
```

## Practical - 2(d)

**Aim: Map Words to Properties Using Python Dictionaries**

**Program:**

```
thisdict={
    "brand": "Porsche",
    "model": "911 gt3 rs",
    "year": 2003
}

print(thisdict)
print(thisdict["brand"])
print(len(thisdict))
print(type(thisdict))
```

**Output:**

```
= RESTART: /home/mrportable/Documents/MiniHacker2/MS-C-2/NLP/Pract/Pract-2
/pract-2d.py
{'brand': 'Porsche', 'model': '911 gt3 rs', 'year': 2003}
Porsche
3
<class 'dict'>
```

---

# Practical - 3(a)

## Aim: Study Default Tagger

### Requirements:

```
pip install nltk
import nltk
nltk.download('treebank')
```

### Program:

```
import nltk
nltk.download('treebank')
from nltk.tag import DefaultTagger
from nltk.corpus import treebank
exptagger = DefaultTagger("NN")
testsentences=treebank.tagged_sents()[1000:]
print(exptagger.accuracy(testsentences))
print(exptagger.tag_sents([[ 'Hey', ',', ''], [ 'How', 'are', 'you', '?' ]]))
```

### Output:

```
= RESTART: /home/mrportable/Documents/MiniHacker2/MS-CSP-2/NLP/Pract/Pract-3
/pract-3a.py
[nltk_data] Downloading package treebank to
[nltk_data]   /home/mrportable/nltk_data...
[nltk_data]   Package treebank is already up-to-date!
0.13198749536374715
[[('Hey', 'NN'), (',', 'NN')], [('How', 'NN'), ('are', 'NN'), ('you', 'NN'), (
'?', 'NN')]]
```

# Practical - 3(b)

## Aim: Study Unigram Tagger

### Requirements:

```
pip install nltk
import nltk
nltk.download('treebank')
```

### Program:

```
from nltk.tag import UnigramTagger
from nltk.corpus import treebank

train_sents = treebank.tagged_sents()[10]
tagger = UnigramTagger(train_sents)
print(treebank.sents()[0])
print('\n', tagger.tag(treebank.sents()[0]))

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

### Output:

```
= RESTART: /home/mrportable/Documents/MiniHacker2/MS-CSP-2/NLP/Pract/Pract-3
/pract-3b.py
['Pierre', 'Vinken', ',', '61', 'years', 'old', ',', 'will', 'join', 'the', 'b
oard', 'as', 'a', 'nonexecutive', 'director', 'Nov.', '29', '.']

[('Pierre', 'NNP'), ('Vinken', 'NNP'), (',', ','), ('61', 'CD'), ('years', 'N
NS'), ('old', 'JJ'), (',', ','), ('will', 'MD'), ('join', 'VB'), ('the', 'DT')
, ('board', 'NN'), ('as', 'IN'), ('a', 'DT'), ('nonexecutive', 'JJ'), ('direct
or', 'NN'), ('Nov.', 'NNP'), ('29', 'CD'), ('.', '.')]

[('Pierre', 'NN'), ('Vinken', None), (',', None), ('61', None), ('years', Non
e), ('old', None), (',', None), ('will', None), ('join', None), ('the', None),
('board', None), ('as', None), ('a', None), ('nonexecutive', None), ('directo
r', None), ('Nov.', None), ('29', None), ('.', None)]
```

## Practical - 4(a)

**Aim: Study of Wordnet Dictionary with methods as synsets, definitions, examples, antonyms**

### Requirements:

```
pip install nltk
nltk.download("wordnet")
```

### Program:

```
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:

```
===== RESTART: E:/Prat-4/a/pract-4a.py =====
[nltk_data] Downloading package wordnet to C:\Users\IT LAB PC NO
[nltk_data] 29\AppData\Roaming\nltk_data...
[nltk_data] Package wordnet is already up-to-date!
[Synset('computer.n.01'), Synset('calculator.n.01')]
a machine for performing calculations automatically
Examples: []
[Lemma('sell.v.01.sell')]
```



## Practical - 4(b)

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

**Requirements:**

```
pip install nltk
nltk.download("wordnet")
```

**Program:**

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

**Output:**

```
===== RESTART: E:/Pract-4/b/pract-4b.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('a
ctive.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 - 4(c)

**Aim: Compare two nouns**

**Requirements:**

```
pip install nltk
nltk.download("wordnet")
```

**Program:**

```
import nltk
from nltk.corpus import wordnet
syn1 = wordnet.synsets("cricket")
syn2 = wordnet.synsets("hokey")
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))
```

**Output:**

```
===== RESTART: E:/Pract-4/c/pract4c.py =====
Path similarity of:
Synset('football.n.01') ( n ) [ any of various games played with a ball (round o
r oval) in which two teams try to kick or carry or propel the ball into each oth
er's goal ]
Synset('football.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 America
n football ]
Synset('football.n.02') ( 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 - 5(a)

## Aim: Tokenization using Python's split() function

### Program:

```
text = """ This tool is an a beta stage. Alexa developers can use Get  
Metrics API to seamlessly analyse  
metric. It also supports custom skill model, prebuilt Flash Briefing  
model, and the Smart Home Skill API.  
You can use this tool for creation of monitors, alarms, and dashboards  
that spotlight changes. The  
release of these three tools will enable developers to create visual rich  
skills for Alexa devices with  
screens. Amazon describes these tools as the collection of tech and tools  
for creating visually rich and  
interactive voice experiences. """  
data = text.split('.')  
for i in data:  
    print (i)
```

### Output:

```
===== RESTART: E:/Pract-5/a/pract-5a.py =====  
    This tool is an a beta stage  
    Alexa developers can use Get Metrics API to seamlessly analyse  
metric  
    It also supports custom skill model, prebuilt Flash Briefing model, and the Smart  
Home Skill API  
  
You can use this tool for creation of monitors, alarms, and dashboards that spotli  
ght changes  
    The  
release of these three tools will enable developers to create visual rich skills f  
or Alexa devices with  
screens  
    Amazon describes these tools as the collection of tech and tools for creating vis  
ually rich and  
interactive voice experiences
```

---

# Practical - 5(b)

## Aim: Tokenization using Regular Expressions (RegEx)

### Requirements:

```
pip install nltk
```

### Program:

```
import nltk
from nltk.tokenize import RegexpTokenizer
tk = RegexpTokenizer('\s+',gaps = True)
str = "Let's use RegexpTokenizer to split"
tokens = tk.tokenize(str)
print(tokens)
```

### Output:

```
===== RESTART: E:/Pract-5/b/pract-5b.py =====
["Let's", 'use', 'RegexpTokenizer', 'to', 'split']
```

# Practical - 5(c)

## Aim: Tokenization using NLTK

### Requirements:

```
pip install nltk
```

### Program:

```
import nltk
from nltk.tokenize import word_tokenize
str = "Let's use word_tokenize to split"
print(word_tokenize(str))
```

### Output:

```
===== RESTART: E:/Pract-5/e/Pract-5e.py =====
['Let', "'", 's', 'use', 'word_tokenize', 'to', 'split']
|
```

# Practical - 5(d)

**Aim:**

**Requirements:**

```
pip install spacy
```

**Program:**

```
import spacy
nlp = spacy.blank("en")
str = "Let's use spacy to split"
doc = nlp(str)
words = [word.text for word in doc]
print(words)
```

**Output:**

```
===== RESTART: E:/Pract-5/d/pract-5d.py =====
['Let', "'s", 'use', 'spacy', 'to', 'split']
```

---

# Practical - 5(e)

## Aim: Tokenization using Keras

### Requirements:

```
pip install tensorflow
pip install keras
pip install Keras-Preprocessing
```

### Program:

```
import keras
from keras.preprocessing.text import text_to_word_sequence
str = "Let's use keras to split"
tokens = text_to_word_sequence(str)
print(tokens)
```

### Output:

```
➡ ["let's", 'use', 'keras', 'to', 'split']
```



# Practical - 5(f)

## Aim: Tokenization using Gensim

### Requirements:

```
pip install gensim
```

```
gensim requires C++
```

### Program:

```
from gensim.utils import tokenize  
str = "I love to study Natural Language Processing in Python"  
list(tokenize(str))
```

### Output:

```
⇒ ['I',  
   'love',  
   'to',  
   'study',  
   'Natural',  
   'Language',  
   'Processing',  
   'in',  
   'Python']
```

# Practical - 6(a)

**Aim:** Named Entity recognition using user defined text.

**Requirements:**

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

**Program:**

```
import spacy
nlp = spacy.load("en_core_web_sm")
text = ("When Sebastian Thrun started working on self-driving cars at "
"Google in 2007, few people outside of the company took him "
"seriously. "I can tell you very senior CEOs of major American "
"car companies would shake my hand and turn away because I wasn't "
"worth talking to," said Thrun, in an interview with Recode earlier "
"this week.")

doc = nlp(text)
print("Noun phrases: ", [chunk.text for chunk in doc.noun_chunks])
print("Verbs: ", [token.lemma_ for token in doc if token.pos_=="VERB"])
```

**Output:**

```
= RESTART: /home/mrportable/Documents/MiniHacker2/MS-CSP-2/NLP/Pract/Pract-6/p
ract-6a.py
Noun phrases: ['Sebastian Thrun', 'self-driving cars', 'Google', 'few people',
'the company', 'him', 'I', 'you', 'very senior CEOs', 'major American car compan
ies', 'my hand', 'I', 'Thrun', 'an interview', 'Recode']
Verbs: ['start', 'work', 'drive', 'take', 'tell', 'shake', 'turn', 'talk', 'say
']
```

## Practical - 6(b)

**Aim: Named Entity recognition with diagram using NLTK corpus – treebank.**

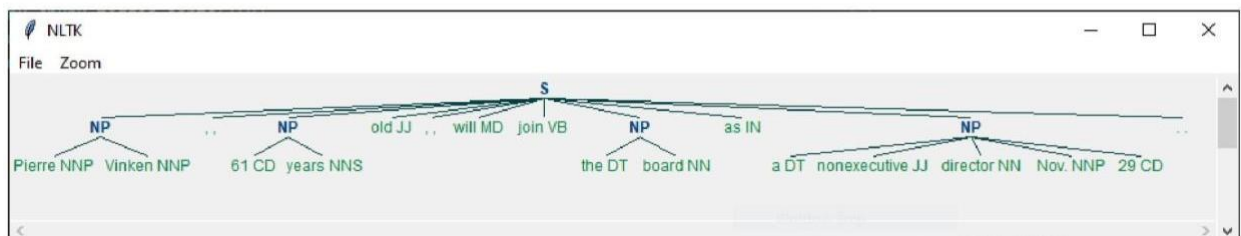
### Requirements:

```
pip install nltk  
nltk.download('treebank')
```

### Program:

```
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 - 7(a)

**Aim:** Define grammar using nltk. Analyze a sentence using the same

### Requirements:

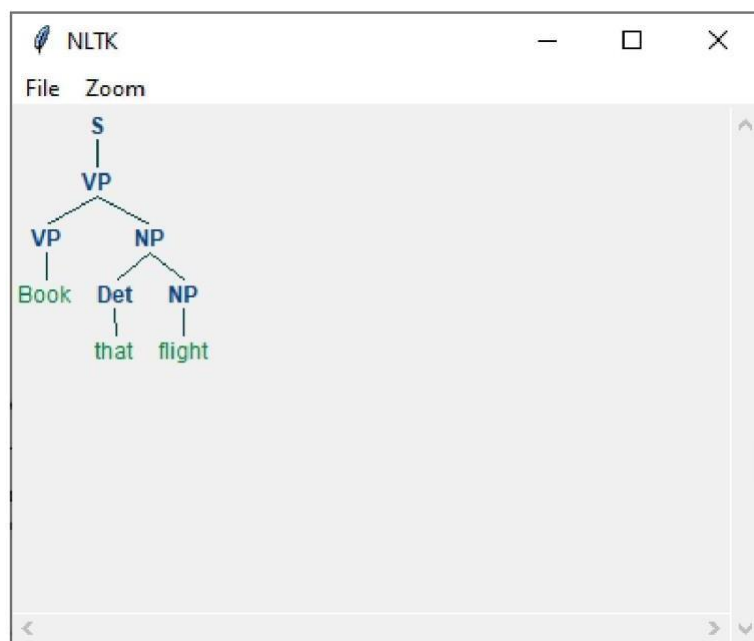
```
pip install nltk
import nltk
nltk.download("punkt")
nltk.download("treebank")
```

### Program:

```
import nltk
from nltk import tokenize
nltk.download("punkt")
nltk.download("treebank")
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 - 7(b)

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

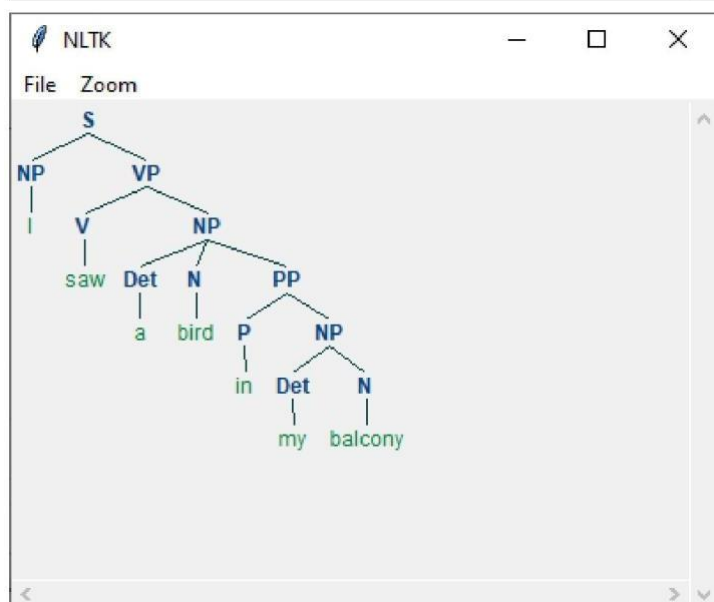
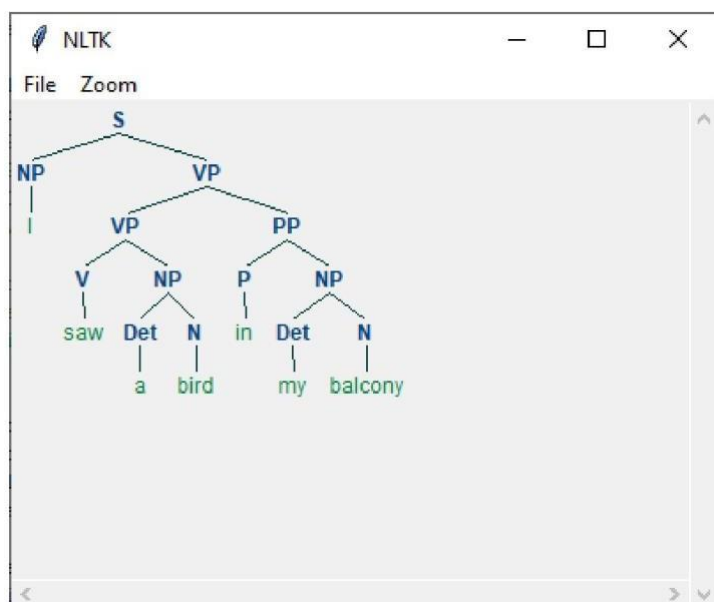
**Requirements:**

```
pip install nltk
```

**Program:**

```
import nltk
from nltk import tokenize
grammar1 = nltk.CFG.fromstring("""
S -> NP VP
PP -> P NP
NP -> Det N | Det N PP | 'I'
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)
parser = nltk.ChartParser(grammar1)
for tree in parser.parse(all_tokens):
    print(tree)
    tree.draw()
```

**Output:**





# Practical - 8

**Aim: Study PorterStemmer, LancasterStemmer, RegexpStemmer, SnowballStemmer Study WordNetLemmatizer**

## Requirements:

```
pip install nltk
import nltk
nltk.download("wordnet")
```

## Program:

```
import nltk
nltk.download("wordnet")
word = "running"
print("PorterStemmer")
from nltk.stem import PorterStemmer
word_stemmer = PorterStemmer()
print(word_stemmer.stem(word))

print("-----")
print("Lancaster Stemmer")
from nltk.stem import LancasterStemmer
Lanc_stemmer = LancasterStemmer()
print(Lanc_stemmer.stem(word))

print("-----")
print('RegexpStemmer')
import nltk
from nltk.stem import RegexpStemmer
Reg_stemmer = RegexpStemmer('ing|$s|$e|$able$', min=4)
print(Reg_stemmer.stem(word))

print("-----")
print('SnowballStemmer')
from nltk.stem import SnowballStemmer
english_stemmer = SnowballStemmer('english')
print(english_stemmer.stem(word))

print("-----")
print('WordNetLemmatizer')
```

```

from nltk.stem import WordNetLemmatizer
lemmatizer = WordNetLemmatizer()
print("word :\tlemma")
print("rocks :", lemmatizer.lemmatize("rocks"))
print("corpora :", lemmatizer.lemmatize("corpora"))

print("better :", lemmatizer.lemmatize("better", pos ="a"))

```

## Output:

```

===== RESTART: E:/Pract-8/Pract-8.py =====
[nltk_data] Downloading package wordnet to C:\Users\IT LAB PC NO
[nltk_data]      29\AppData\Roaming\nltk_data...
[nltk_data]   Package wordnet is already up-to-date!
PoterStemmer
run
-----
Lancaster Stemmer
run
-----
RegexpStemmer
runn
-----
SnowballStemmer
run
-----
WordNetLemmatizer
word : lemma
rocks : rock
corpora : corpus
better : good

```

# Practical - 9

## Aim: Finite Automata

### Program:

```
def FA(s):
    if len(s) < 3:
        return "Rejected"

    if s[0] == '1':
        if s[1] == '0':
            if s[2] == '1':
                for i in range(3, len(s)):
                    if s[i] != "1":
                        return "Rejected"
                return "Accepted"
            return "Rejected"
        return "Rejected"
    return "Rejected"

inputs=['1','10101','101','10111','01010','100','','10111101','1011111']
for i in inputs:
    print(i, FA(i))
```

### Output:

```
===== RESTART: E:/Pract-9/pract-9.py =====
1 Rejected
10101 Rejected
101 Accepted
10111 Accepted
01010 Rejected
100 Rejected
Rejected
10111101 Rejected
1011111 Accepted
```

---