

```
1 #!/usr/bin/env python
2 # coding: utf-8
3
4 # In[6]:
5
6
7 #task 5: chinese segmentation of jieba
8 #!pip install jieba
9 import jieba
10 seg=jieba.cut("把句子中所有的可以成词的词语都扫描出来",cut_all=True)
11 print(" ".join(seg))
12
13
14 # In[7]:
15
16
17 #TASK 1: BASIC TEXT PROCESSING PIPELINE
18 import nltk
19 texts="""When I dare to be powerful - to use my strength in the service of my vision,
20 then it becomes less and less important whether I am afraid"""
21 for text in texts:
22     sentences=nltk.sent_tokenize(texts)
23     for sentence in sentences:
24         words=nltk.word_tokenize(sentence)
25         tagged=nltk.pos_tag(words)
26         print(tagged)
27
28 for text in texts:
29     words=nltk.word_tokenize(texts)
30     for word in words:
31         word
32
33 # In[9]:
34
35
36 #task 2: tweettokenizer
37 import nltk
38 from nltk.tokenize import TweetTokenizer
39 text='The party was soooo fun: D #superfun'
40 twtkn=TweetTokenizer()
41 twtkn.tokenize(text)
42
43
44 # In[12]:
45
46
47 #Task 3: Scrapping Data from Web
48 from urllib import request
49 url="http://www.gutenberg.org/files/2554/2554-0.txt"
50 response=request.urlopen(url)
51 raw=response.read().decode('utf8')
52 type(raw)
53 from nltk.tokenize import word_tokenize
54 tokens=word_tokenize(raw)
55 type(tokens)
56
57 #HOMEWORK
58
```

```
59 #HTML => ASCII => TOKEN => TEXT => VOCAB
60
61
62 # In[18]:
63
64
65 from urllib import request
66 url = "https://timesofindia.indiatimes.com/"
67 html = request.urlopen(url).read().decode('utf8')
68 html[:60]
69
70 print(html)
71
72
73 # In[51]:
74
75
76 from nltk.tokenize import word_tokenize
77 from bs4 import BeautifulSoup
78 raw = BeautifulSoup(html, 'html.parser').get_text()
79 tokens = word_tokenize(raw)
80 print(len(tokens))
81 print(tokens)
82
83
84 # In[53]:
85
86
87 import nltk
88 tokens = tokens[:10615]
89 text = nltk.Text(tokens)
90 text.concordance('News')
91
92
93 # In[ ]:
94
95
96
97
98
99
```