

WEB MINING

by Sharadindu Adhikari, 19BCE2105

Experiment 1:

1(a).Stop Words removal from a paragraph.

Solution:-

```
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
ex = """Web Mining is the process of Data Mining techniques to automatically
discover and extract information from Web documents and services. The main purpose
of web mining is discovering useful information from the World-Wide Web and its
usage patterns."""
stop_words = set(stopwords.words('english'))
token = word_tokenize(ex)
fs = [w for w in word_tokens if not w in stop_words]
fs = []
for w in word_tokens:
    if w not in stop_words:
        fs.append(w)
print("Before Removal=>")
print(token)
print("After Removal=>")
print(fs)
```

Output:-

```
Command Prompt
Microsoft Windows [Version 10.0.19042.804]
(c) 2020 Microsoft Corporation. All rights reserved.

C:\Users\rishi>cd documents

C:\Users\rishi\Documents>python yolo.txt
Before Removal=>
['Web', 'Mining', 'is', 'the', 'process', 'of', 'Data', 'Mining', 'techniques', 'to', 'automatically', 'discover', 'and',
, 'extract', 'information', 'from', 'Web', 'documents', 'and', 'services', '.', 'The', 'main', 'purpose', 'of', 'web', 'mining',
, 'is', 'discovering', 'useful', 'information', 'from', 'the', 'World-Wide', 'Web', 'and', 'its', 'usage', 'patterns', '.']
After Removal=>
['Web', 'Mining', 'process', 'Data', 'Mining', 'techniques', 'automatically', 'discover', 'extract', 'information', 'Web',
, 'documents', 'services', '.', 'The', 'main', 'purpose', 'web', 'mining', 'discovering', 'useful', 'information', 'World-Wide', 'Web', 'usage', 'patterns', '.']

C:\Users\rishi\Documents>
```

1(b). Create a set of stop words given below and print the output:-

```
stop_words['.', ',', 'a', 'they', 'the', 'his', 'so', 'and', 'were', 'from',
'that', 'of', 'in', 'only', 'with', 'to']
```

Solution: -

```
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

ex = """. , a they the his so and were from that of in only with to"""
stop_words = set(stopwords.words('english'))
token = word_tokenize(example_sent)
fs = [w for w in word_tokens if not w in stop_words]
fs = []

for w in word_tokens:
    if w not in stop_words:
        fs.append(w)

print("Before Removal=>")
print(token)
print("After Removal=>")
print(fs)
```

Output:

```
Command Prompt
C:\Users\rishi\Documents>python yolo.txt
Before Removal=>
['.', ',', 'a', 'they', 'the', 'his', 'so', 'and', 'were', 'from', 'that', 'of', 'in', 'only', 'with', 'to']
After Removal=>
['.', ',']

C:\Users\rishi\Documents>
C:\Users\rishi\Documents>
```

2(a,b). Write a program to tokenise a sentence(using nltk).

```
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
```

```
word_data = "It originated from the idea that there are readers who prefer learning
new skills from the comforts of their drawing rooms"
nltk_tokens = word_tokenize(word_data)
print (nltk_tokens)
```

```
Administrator: Command Prompt
Microsoft Windows [Version 10.0.19042.804]
(c) 2020 Microsoft Corporation. All rights reserved.

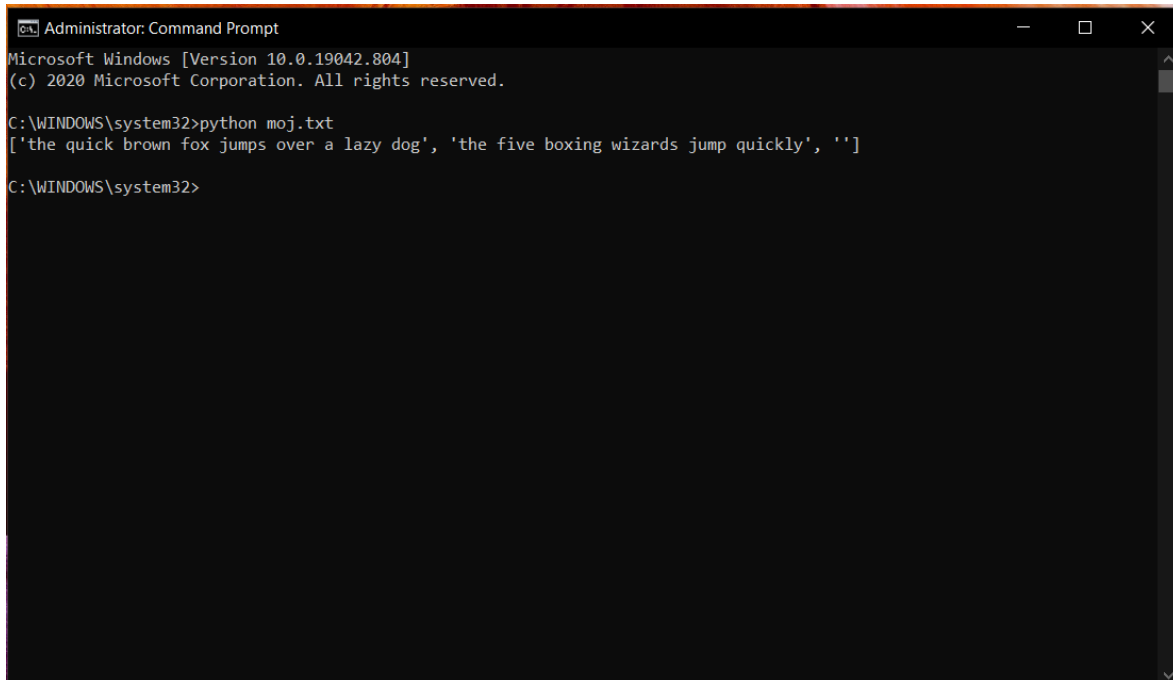
C:\WINDOWS\system32>python moj.txt
['the quick brown fox jumps over a lazy dog', 'the five boxing wizards jump quickly', '']

C:\WINDOWS\system32>python lan.txt
['It', 'originated', 'from', 'the', 'idea', 'that', 'there', 'are', 'readers', 'who', 'prefer', 'learning', 'new', 'skills', 'from', 'the', 'comforts', 'of', 'their', 'drawing', 'rooms']

C:\WINDOWS\system32>
```

2(c). Write a program to tokenize multiple sentences without using nltk.

```
text = """the quick brown fox jumps over a lazy dog.the five boxing wizards jump  
quickly"""  
# Splits at '.'  
print(text.split('. '))
```



The screenshot shows a Windows Command Prompt window titled "Administrator: Command Prompt". The window displays the following text:

```
Microsoft Windows [Version 10.0.19042.804]  
(c) 2020 Microsoft Corporation. All rights reserved.  
  
C:\WINDOWS\system32>python moj.txt  
['the quick brown fox jumps over a lazy dog', 'the five boxing wizards jump quickly', '']  
  
C:\WINDOWS\system32>
```