1) **Write a program to remove the stopwords for any given paragraph. Create a set of stopwords given below and print the output.**
   **stop_words =**
   **['.',',','a','they','the','his','so','and','were','from',that','of','in','only','with','to']**

**PROCEDURE**
Filereader = open(File, read)
Text[]= []
Stop_words[]=['.',',','a','they','the','his','so','and','were','from','that','of','in','only','with','to']
Textarray[] = Filereader.read().split()
For i in (0 to lengthof(Textarray))
        If TextArray[i] is not in Stop_words
                Text.append(TextArray[i])
Print (Text)

**CODE**
```
fr = open("../SampleText.txt")
stop_words = ['.',',','a','they','the','his','so','and','were','from','that','of','in','only','with','to']
text_arr = fr.read().split()
text_without_sw=[]
for i in range(0,len(text_arr)):
    if text_arr[i] not in stop_words:
        text_without_sw.append(text_arr[i])
print (text_without_sw)
fw = open("../TextWithoutStopwords.txt","w")
for j in range(0,len(text_without_sw)):
    fw.write(text_without_sw[j])
    fw.write(" ")
fw.close()
print(fr.read())
fr.close()
```
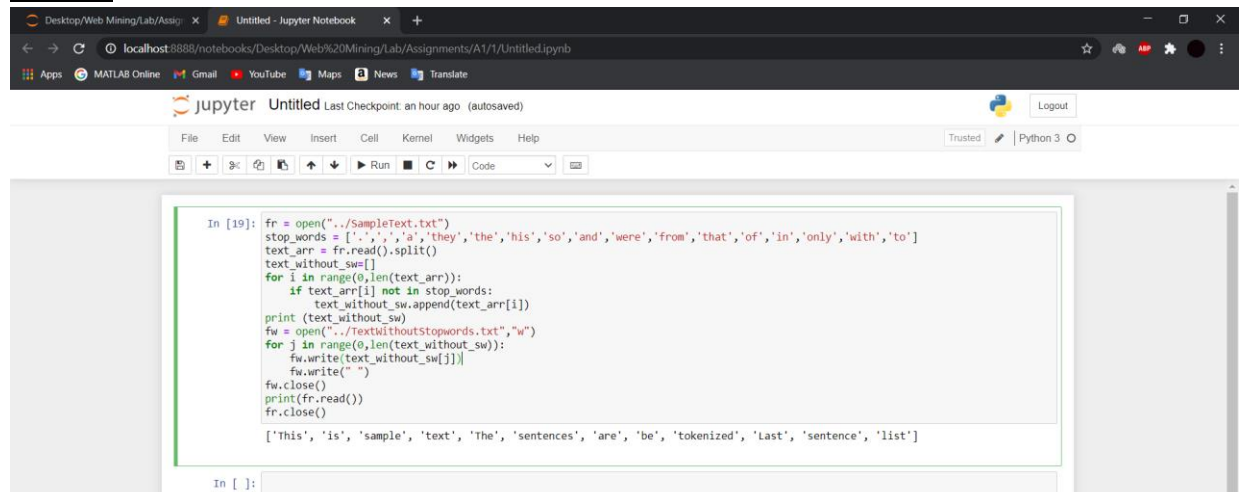
**OUTPUT**

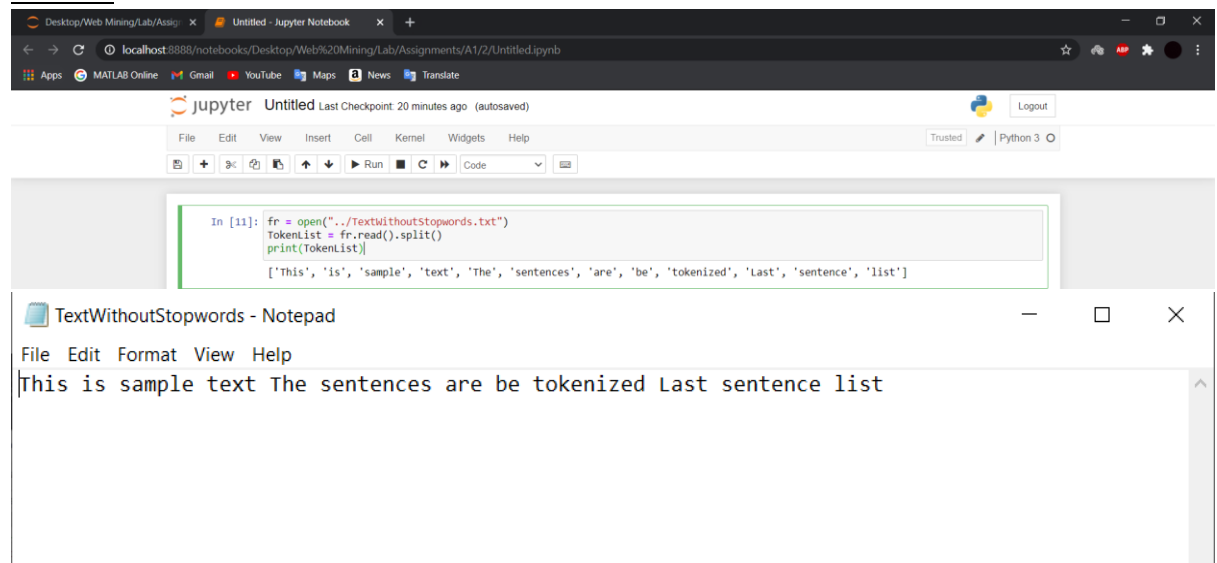## 2) Write a program to tokenize (without Nltk)

### a) A sentence

**PROCEDURE**

FileReader = open (File, read)
Tokens= FileReader.read().split()
Print(Tokens)

**CODE**

```
fr = open("../TextWithoutStopwords.txt")
TokenList = fr.read().split()
print(TokenList)
```
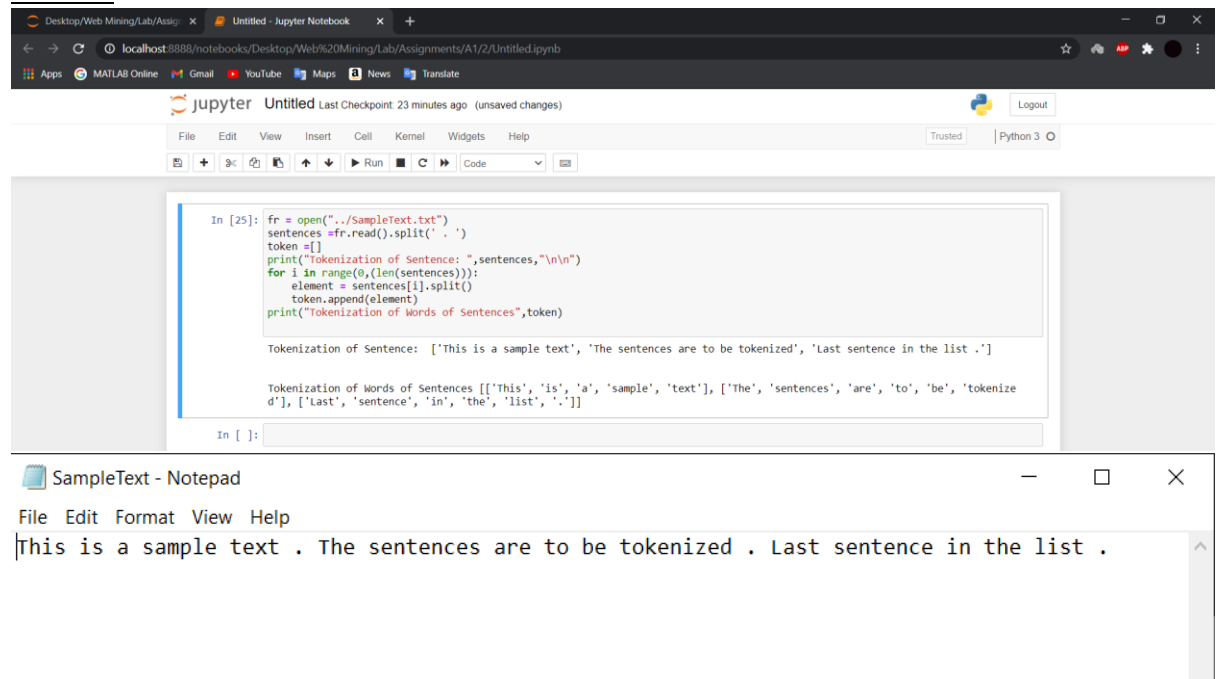
**OUTPUT**

**b) Multiple Sentences**

**PROCEDURE**

FileReader = open (File, read)

Sentences= FileReader.read().split(' . ')

Print(Sentences)

For i in (0 to lengthof(Sentences))

    Arr[] = Sentences[i].split()

    Text.append(Arr)

Print (Text)

**CODE**

```
fr = open("../SampleText.txt")
sentences =fr.read().split(' . ')
token =[]
print("Tokenization of Sentence: ",sentences,"\n\n")
for i in range(0,(len(sentences))):
    element = sentences[i].split()
    token.append(element)
print("Tokenization of Words of Sentences",token)
```

**OUTPUT**

3) **Write a program (using nltk toolkit in python environment) to tokenize**
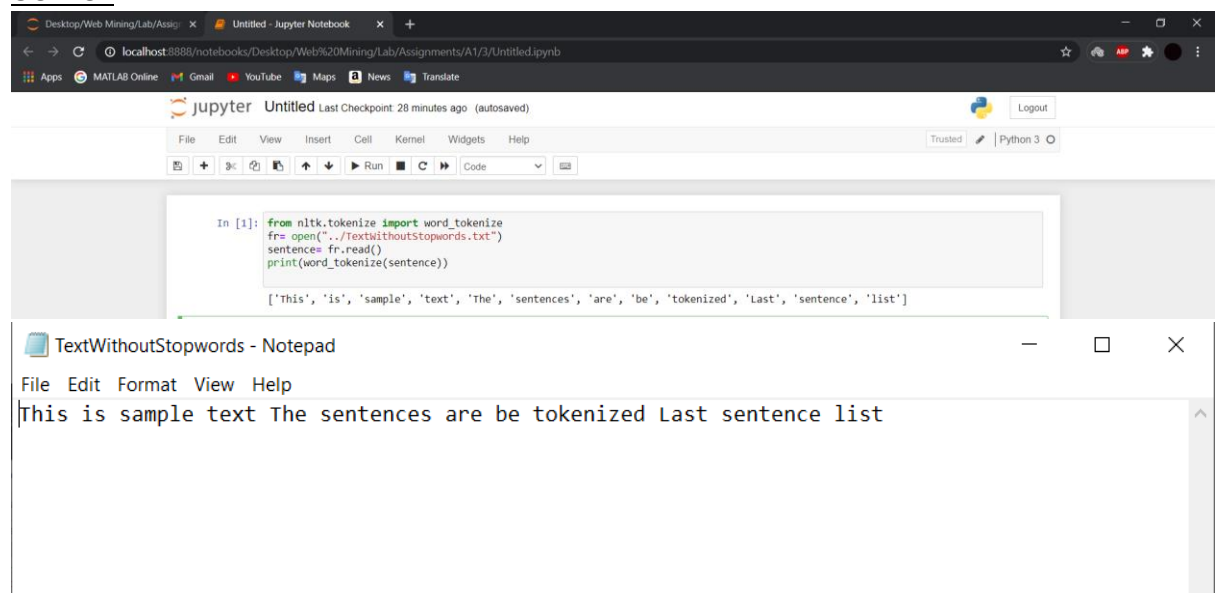   a) **A sentence**

   **PROCEDURE**
   //Importing NLTK library to use the word_tokenize function
   Import nltk.tokenize.word_tokenize
   FileReader=open(File, read)
   Sentence=FileReader.read()
   Print(word_tokenize(Sentence))

   **CODE**
   ```
   from nltk.tokenize import word_tokenize
   fr= open("../TextWithoutStopwords.txt")
   sentence= fr.read()
   print(word_tokenize(sentence))
   ```

   **OUTPUT**

**b) A paragraph**

**PROCEDURE**

//Importing NLTK library to use the sent_tokenize function
Import nltk.tokenize.sent_tokenize
FileReader=open(File, read)
Paragraph=FileReader.read()
Print(sent_tokenize(Paragraph))

**CODE**

```
from nltk.tokenize import sent_tokenize
fr= open("../SampleText.txt")
para= fr.read()
print(sent_tokenize(para))
```

**OUTPUT**