

EXPERIMENT NO 4

A. CALCULATE WORD COUNT OF A GIVEN SPECIFIC DOCUMENT AND SHOW TOP 10 FREQUENT WORDS WITH THEIR FREQUENCY

IMPORTING LIBRARIES

```
import pandas as pd
import matplotlib.pyplot as plt
import nltk as ntk
%matplotlib inline
```

```
# Read input file, note the encoding
```

```
file = open('/DOCUMENTS/COLLEGE/CLASSES/EXPERIMENT_NO_4/EXPERIMENT4.txt',
            encoding="utf8")
a = file.read()

In [3]: # stopwords
```

CALCULATE THE WORD COUNTING

```
# Instantiate a dictionary, and for every
# Add to the dictionary if it doesn't ex
```

```
# To eliminate duplicates, splitting by punctuation, and use case demiliters.
for word in a.lower().split():
    word = word.replace(".", "")
    word = word.replace(", ", "")
    word = word.replace(":", "")
    word = word.replace("\\""", "")
    word = word.replace("!", "")
    word = word.replace("'", "")
    if word not in stopwords:
        if word not in wordcount:
            wordcount[word] = 1
        else:
            wordcount[word] += 1
```

```
# Print most common word
```

```
word_counts = collections.Counter(word_counts)
for word, count in word_counts.most_common(n_print):
    print(word, ":", count)

# Close the file
file.close()

# Create a data frame of the most common words
lst = word_counts.most_common(n_print)
df = pd.DataFrame(lst, columns=['Word', 'Count'])
```

```
How many most common words to print: 10
```

OK. The 10 most common words are as follows

```

amice : 6224
sit : 6173
ut : 5258
id : 5247
eget : 5068
et : 4702
nunc : 4648
vitae : 4565
enim : 4075

```

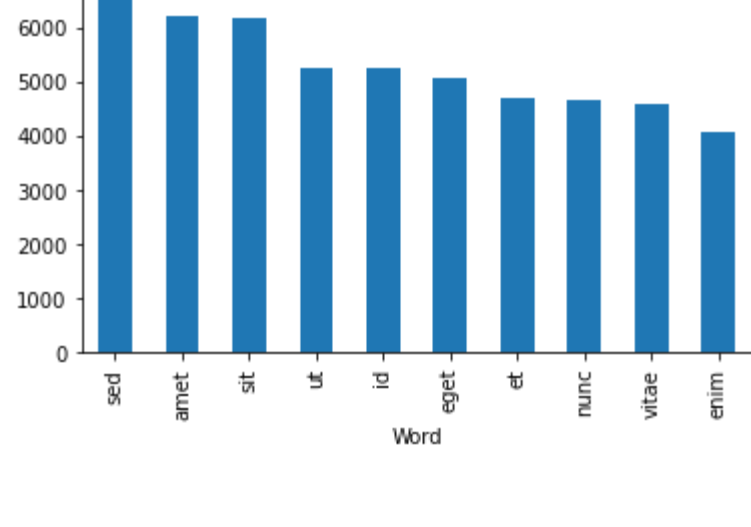
0	sed	7650
1	amet	6224
2	sit	6173
3	ut	5258
4	id	5247

VISUALIZING THE DATASET

```
df.plot.bar(x='Word', y='Count')

<AxesSubplot: xlabel='Word'>
```

Category	Number of people
Do not use the Internet	7500



INSTALLING THE WORD CLOUD

```
# install wordcloud
```

```
!pip install wordcloud

# import package and its set of stopwords

print('Wordcloud is installed and imported!')
```

Requirement already satisfied: wordcloud in c:
Requirement already satisfied: numpy>=1.6.1 in c:
3)

```
Requirement already satisfied: pillow in c:\users\hp\anaconda3\lib\site-packages (from wordcloud) (8.4.0)
Requirement already satisfied: kiwisolver>=1.0.1 in c:\users\hp\anaconda3\lib\site-packages (from matplotlib->wordcloud) (1.3.1)
Requirement already satisfied: cyclor>=0.10 in c:\users\hp\anaconda3\lib\site-packages (from matplotlib->wordcloud) (0.10.0)
Requirement already satisfied: python-dateutil>=2.7 in c:\users\hp\anaconda3\lib\site-packages (from matplotlib->wordcloud) (2.8.2)
Requirement already satisfied: pyparsing>=2.2.1 in c:\users\hp\anaconda3\lib\site-packages (from matplotlib->wordcloud) (3.0.4)
Requirement already satisfied: six in c:\users\hp\anaconda3\lib\site-packages (from cyclor>=0.10->matplotlib->wordcloud) (1.16.0)
Wordcloud is installed and imported!
```

READING THE DOCUMENT FILE

```
# Read input file, note the encoding is specified here
# It may be different in your text file
```

```
encoding="utf8")
```

CREATING THE WORD CLOUD

```
In [9]: # use the function set to remove any redundant stopwords.  
stopwords = set(STOPWORDS)
```

```
[10]: # Create a word cloud object and generate a word cloud using only the first 2000 words.
      # instantiate a word cloud object
      alice_wc = WordCloud(background_color='white',
                           max_words=2000,
                           stopwords=stopwords)
```

```
# generate the word cloud
alice_wc.generate(a)
```

```
Out[10]: \wordcloud.wordcloud.wordcloud at 0x2c424092340>
```

```
# Now that the word cloud is created, let's visualize it.
# display the word cloud
fig = plt.figure(figsize=(14, 18))
```

```
plt.axis('off')
plt.show()
```

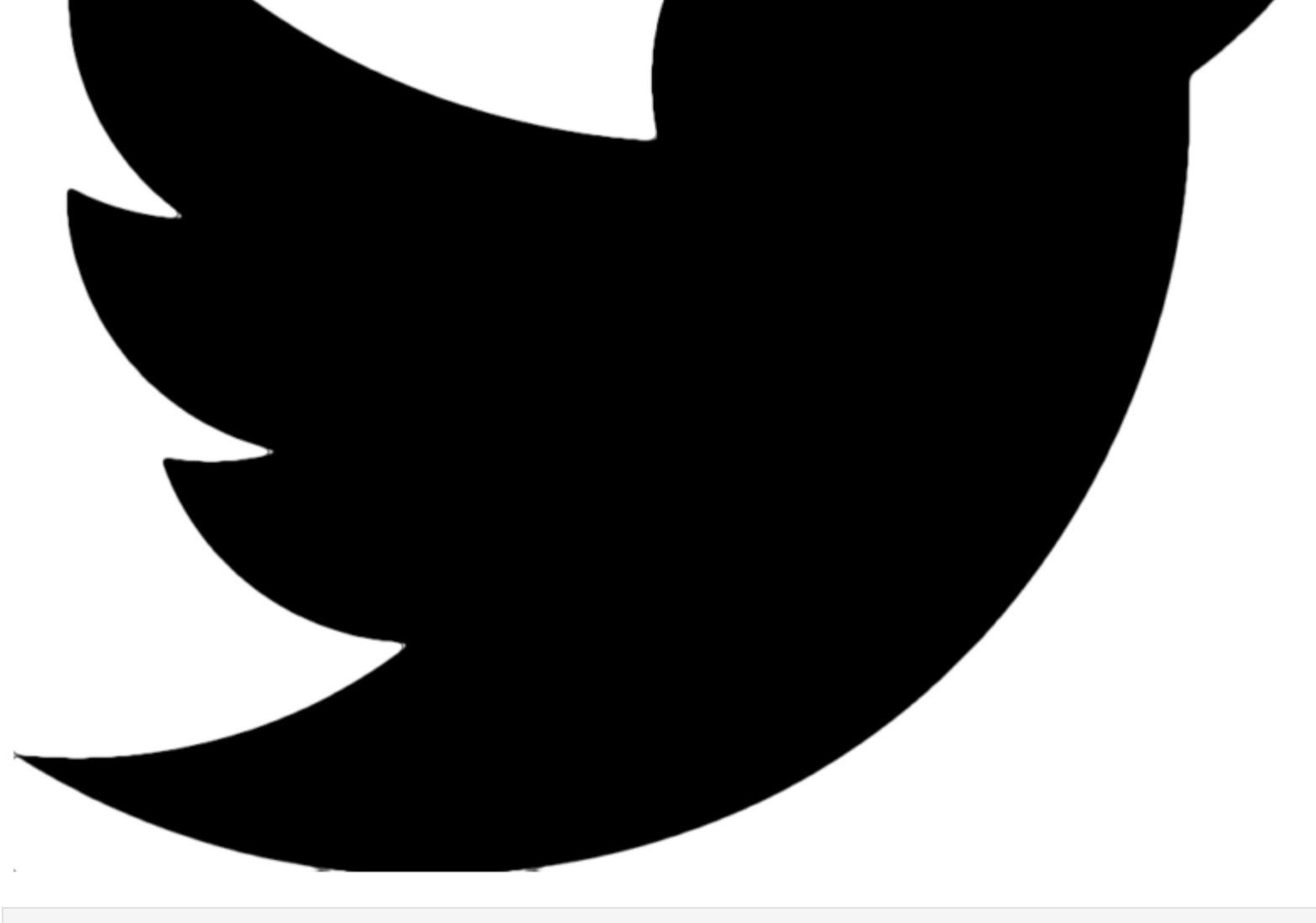


```
import urllib
```

```
in [13]: # masked image
mask_img = np
```

```
fig = plt.figure(figsize=(14, 18))

plt.imshow(mask_img, cmap=plt.cm.gray, interpolation='bilinear')
plt.axis('off')
plt.show()
```



```
alice_wc = WordCloud(background_color='white',
                     max_words=2000,
                     mask=mask_img,
```

```
# generate the word cloud
alice_wc.generate(a)

# display the word cloud
fig = plt.figure(figsize=(14, 18))
```

```
plt.imshow(alice_wc, interpolation='bilinear')
plt.axis('off')
plt.show()
```

et malesu pellentesque habitant

[illegible]

Word cloud visualization of the most frequent words in the corpus. The words are arranged in a circular pattern, with the most frequent words in the center and less frequent words on the periphery. The words are color-coded by frequency, with red indicating the highest frequency and blue indicating the lowest frequency.

senectus net lacus sed morbi tincidunt et tortor cras fermentum toror capere ut lectualligam dui ut egestas agestas quoniam ut no viverra arci nulla pharetra aris vna

-----X-----X-----X-----

