```python
import pandas as pd
```

```python
import os
```

```python
#!pip install nltk
```

```python
import nltk
```

```python
from nltk import FreqDist
```

```python
#!pip install matplotlib
```

```python
#!pip install wordcloud
```

```python
import matplotlib.pyplot as plt
```

```python
folder_path = r'C:/Users/fpate/Desktop/AIT-580'
os.chdir(folder_path)
```

```python
os.getcwd()
```

```
'C:\\Users\\fpate\\Desktop\\AIT-580'
```

```python
text = open("C:\\Users\\fpate\\Desktop\\AIT-580\\WABA-crash copy.txt","r
mytext = text.read()
```

```python
mytext
```

```
',6810,,,,Crash Description,"Turned right from 38th and Nebraska into a g
ap in traffic. Car approached from behind and tried to pass too fast and
too close, and car\'s wing mirror hit my handlebars. Lost control of the
bike and crashed on the pavement, hitting my head but sustaining no other
injuries. Driver stopped, another witness stopped, both offered aid. Driv
er offered his name and phone number but I failed to ask for insurance in
fo. Didn\'t get witness\'s name or phone number. Also, didn\'t call the p
olice as I didn\'t think I was seriously injured at the time. Only later
discovered I had a concussion and went to the ER. Filed a police report a
fter the fact. Now stuck in an insurance fight to cover the medical treat
ment - I didn\'t realize that car insurance is assumed to be the payer of
first resort, not sure if my medical insurance will ultimately cover the
treatment. Tried to contact the driver on multiple occasions but thus far
getting no response. Would prefer not to go the lawyer route, but I\'m no
t sure how else to find out his car insurance info. Any WABA advice welco
me!"\n11310384,6810,WASHINGTON,DC,20009,Crash Description,"While recoveri
ng from a turn, hit a parallel rut in the road and went down."\n11310789,
6810,Washington,DC,20001,Crash Description,"I was southbound on 3rd Stree
t and saw a northbound motorist with his left signal on. I knew he didn
```

In [95]: ▶ 1 #1.Calculate the word frequency counts (eliminating stop words)

In [96]: ▶ 1 tokenize_wd = nltk.word_tokenize(mytext)
          2 print("count of words :- ",len(tokenize_wd))
          3 tokenize_wd[:10]

count of words :-  119452

Out[96]: [',6810', ',', ',', ',', ',Crash', 'Description', ',', "''", 'Turned', 'rig
          ht']

In [97]: ▶ 1 #Removing extra special characters like , space etc.
          2 tokenizer = nltk.RegexpTokenizer(r"\w+")
          3 remove_s_word = tokenizer.tokenize(mytext)
          4 remove_s_word[:10]

Out[97]: ['6810',
          'Crash',
          'Description',
          'Turned',
          'right',
          'from',
          '38th',
          'and',
          'Nebraska',
          'into']

In [98]: ▶ 1 print("count of word :- ",len(remove_s_word))

count of word :-  106313

In [99]: ▶ 1 #converting text into the Lower character for NLP procedure:

In [100]: ▶ 1 from nltk.tokenize import RegexpTokenizer

```
In [101]:  ▶|    1  tokenizer = RegexpTokenizer(r'\w+')
                 2  Token_lower = tokenizer.tokenize(mytext.lower())
                 3  print(Token_lower)
```

['6810', 'crash', 'description', 'turned', 'right', 'from', '38th', 'an
d', 'nebraska', 'into', 'a', 'gap', 'in', 'traffic', 'car', 'approached',
'from', 'behind', 'and', 'tried', 'to', 'pass', 'too', 'fast', 'and', 'to
o', 'close', 'and', 'car', 's', 'wing', 'mirror', 'hit', 'my', 'handlebar
s', 'lost', 'control', 'of', 'the', 'bike', 'and', 'crashed', 'on', 'th
e', 'pavement', 'hitting', 'my', 'head', 'but', 'sustaining', 'no', 'othe
r', 'injuries', 'driver', 'stopped', 'another', 'witness', 'stopped', 'bo
th', 'offered', 'aid', 'driver', 'offered', 'his', 'name', 'and', 'phon
e', 'number', 'but', 'i', 'failed', 'to', 'ask', 'for', 'insurance', 'inf
o', 'didn', 't', 'get', 'witness', 's', 'name', 'or', 'phone', 'number',
'also', 'didn', 't', 'call', 'the', 'police', 'as', 'i', 'didn', 't', 'th
ink', 'i', 'was', 'seriously', 'injured', 'at', 'the', 'time', 'only', 'l
ater', 'discovered', 'i', 'had', 'a', 'concussion', 'and', 'went', 'to',
'the', 'er', 'filed', 'a', 'police', 'report', 'after', 'the', 'fact', 'n
ow', 'stuck', 'in', 'an', 'insurance', 'fight', 'to', 'cover', 'the', 'me
dical', 'treatment', 'i', 'didn', 't', 'realize', 'that', 'car', 'insuran
ce', 'is', 'assumed', 'to', 'be', 'the', 'payer', 'of', 'first', 'resor
t', 'not', 'sure', 'if', 'my', 'medical', 'insurance', 'will', 'ultimatel
y', 'cover', 'the', 'treatment', 'tried', 'to', 'contact', 'the', 'drive

```
In [102]:  ▶|    1  remove_word_lower = list(map(lambda word: word.lower(),remove_s_word))
```

```
In [103]:  ▶|    1  from nltk.corpus import stopwords
```

```
In [104]:  ▶|    1  #Using Filter for tokenzing word by removing stopwords i.e. one of the i
                 2  filtered_words = [word for word in remove_word_lower if word not in stopw
                 3  #filtered_words[:10]
                 4  print(filtered_words)
```

['6810', 'crash', 'description', 'turned', 'right', '38th', 'nebraska',
'gap', 'traffic', 'car', 'approached', 'behind', 'tried', 'pass', 'fast',
'close', 'car', 'wing', 'mirror', 'hit', 'handlebars', 'lost', 'control',
'bike', 'crashed', 'pavement', 'hitting', 'head', 'sustaining', 'injurie
s', 'driver', 'stopped', 'another', 'witness', 'stopped', 'offered', 'ai
d', 'driver', 'offered', 'name', 'phone', 'number', 'failed', 'ask', 'ins
urance', 'info', 'get', 'witness', 'name', 'phone', 'number', 'also', 'ca
ll', 'police', 'think', 'seriously', 'injured', 'time', 'later', 'discove
red', 'concussion', 'went', 'er', 'filed', 'police', 'report', 'fact', 's
tuck', 'insurance', 'fight', 'cover', 'medical', 'treatment', 'realize',
'car', 'insurance', 'assumed', 'payer', 'first', 'resort', 'sure', 'medic
al', 'insurance', 'ultimately', 'cover', 'treatment', 'tried', 'contact',
'driver', 'multiple', 'occasions', 'thus', 'far', 'getting', 'response',
'would', 'prefer', 'go', 'lawyer', 'route', 'sure', 'else', 'find', 'ca
r', 'insurance', 'info', 'waba', 'advice', 'welcome', '11310384', '6810',
'washington', 'dc', '20009', 'crash', 'description', 'recovering', 'tur
n', 'hit', 'parallel', 'rut', 'road', 'went', '11310789', '6810', 'washin
gton', 'dc', '20001', 'crash', 'description', 'southbound', '3rd', 'stree
t', 'saw', 'northbound', 'motorist', 'left', 'signal', 'knew', 'time', 'm
```

```
In [105]:  ▶|    1  print("count of word :- ",len(filtered_words))
```

count of word :-  55689

```
In [106]:  ▶|    1  from nltk.stem import WordNetLemmatizer
```

```
In [107]:  ▶|    1  #lemmatizer is a process of grouping same meaning words together.
                 2  lemmatizer = WordNetLemmatizer()
                 3  rsp_word_lower_lem = list(map(lambda word:lemmatizer.lemmatize(word),fil
                 4  rsp_word_lower_lem[:10]
```

Out[107]:  ['6810',
           'crash',
           'description',
           'turned',
           'right',
           '38th',
           'nebraska',
           'gap',
           'traffic',
           'car']

```
In [108]:  ▶|    1  frequency_dist_words = FreqDist(rsp_word_lower_lem)
                 2  frequency_dist_words.most_common(10)
```

Out[108]:  [('car', 1191),
           ('bike', 1074),
           ('crash', 949),
           ('lane', 877),
           ('description', 847),
           ('6810', 841),
           ('right', 790),
           ('left', 711),
           ('driver', 680),
           ('dc', 602)]

```
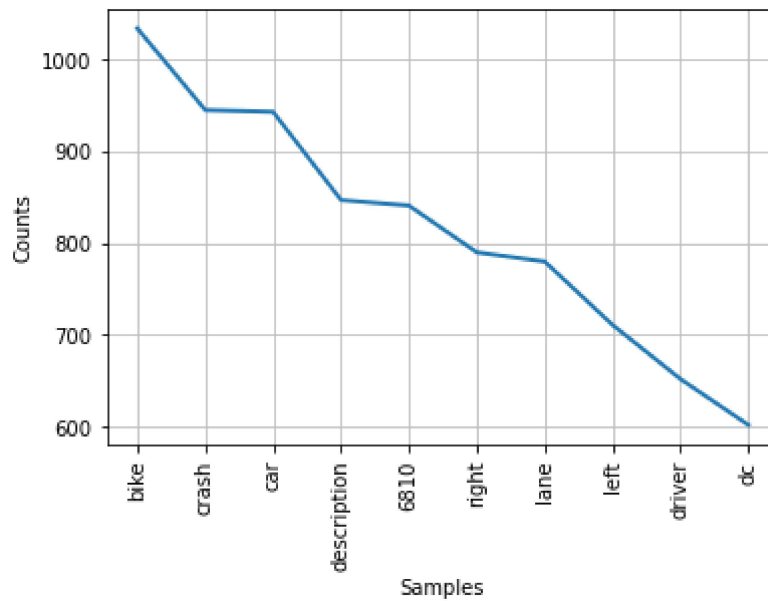In [109]:  ▶|    1  #Plot top 10-word frequency (eliminating stop words)
```

```
In [110]:  ▶|    1  import matplotlib.pyplot as plt
```

```
1  freq_dist = nltk.FreqDist(filtered_words)
2  freq_dist
3
4  print(freq_dist)
5  print(freq_dist.most_common(10))
6  freq_dist.plot(10)
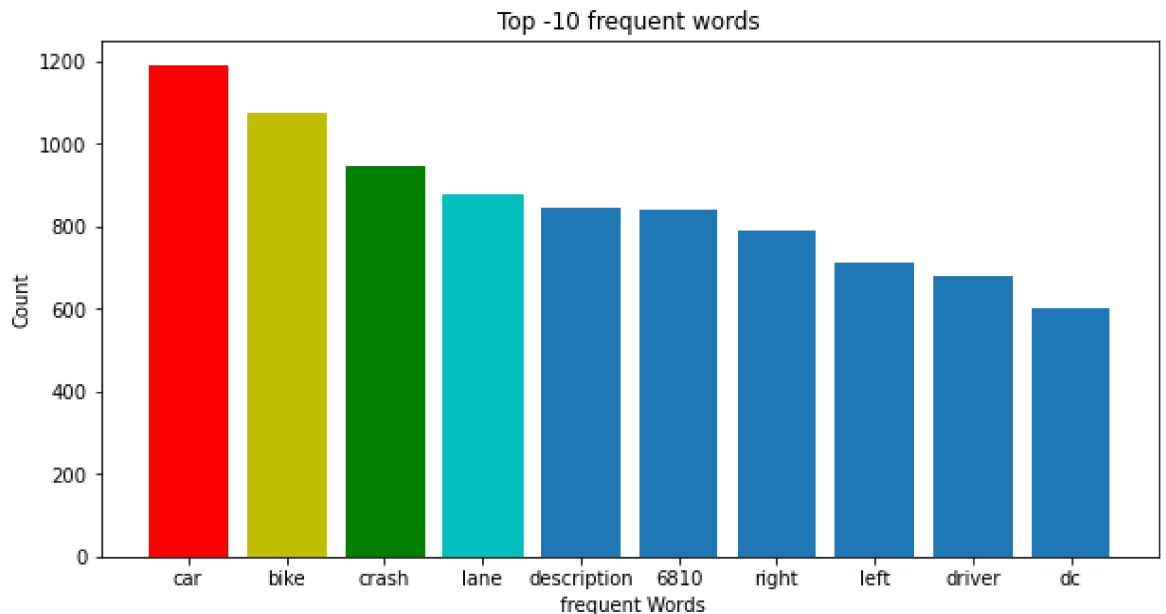```

```
<FreqDist with 6347 samples and 55689 outcomes>
[('bike', 1034), ('crash', 945), ('car', 943), ('description', 847), ('681
0', 841), ('right', 790), ('lane', 780), ('left', 711), ('driver', 652),
('dc', 602)]
```



`<AxesSubplot:xlabel='Samples', ylabel='Counts'>`

```
In [112]: ▶| 1  plt.figure(figsize=(10,5))
              2  barlist=plt.bar(*zip(*frequency_dist_words.most_common(10)))
              3  barlist[0].set_color('r')
              4  barlist[1].set_color('y')
              5  barlist[2].set_color('g')
              6  barlist[3].set_color('c')
              7  plt.title("Top -10 frequent words")
              8  plt.xlabel(" frequent Words")
              9  plt.ylabel("Count")
             10  plt.show()
```

Top -10 frequent words

```
In [113]: ▶| 1  #Interpret the result; what can you learn from the word frequency analys
```

```
In [114]: ▶| 1  #car, bike, crash, lane, 6810 are the most frequently used words in text
              2  #cyran color in the bar graph.
              3
              4  # Through this word frequency analysis we can get the worthy data from th
              5  # For example, when some organization wants to get some specific kind of
              6  # frequency analysis.
              7
              8  #So, basically the frequency analysis means getting the proper insights
```

```
In [115]: ▶| 1  #1.(b)  Extract and display the sentences that include the word 'police'
```

```
In [116]:  ▶|    1  sentences = nltk.sent_tokenize(mytext)
```

```
In [117]:  ▶|    1  police=[]
                2  for s in sentences:
                3      if "police" in s:
                4          police.append(s)
                5  print(police)
```

["Also, didn't call the police as I didn't think I was seriously injured
at the time.", 'Filed a police report after the fact.', 'Since this occur
red just feet from the FBI parking garage entrance, I was immediately sur
rounded by police (I believe FBI police) who must have contacted police a
nd an ambulance.', 'It took about 15 minutes for EMS to arrive—police fol
lowed shortly afterwards.', '(I did give a statement to the police.)"',
'I got a scrape on my arm but nothing significant and my bike wasn\'t tou
ched, so I continued home and then filed a police report."', 'The driver
remained at the scene and appeared to give a statement to the police.',
'As I was still just realizing what all was happening I noticed a police
car stopped ahead at the Vermont intersection, the officer was out of the
car but I was able to yell to him as we stopped at the light.', 'A police
officer was near by but not close enough to see the entire accident.', 'H
e can file a traffic ticket but that would require that I and a witness g
o to police department to complete it.', 'I asked if I need to contact hi
m and he basically said that any police officer can do it.', 'Following m
y returning home from hospital, I call the Montgomery PD to see which pol
ice depart would be most convenient for me and the witness that I had con
tacted to go to.', 'I described to him what happen at the accident and th

```
In [118]:  ▶|    1  police_sentence= len(police)
```

```
In [119]:  ▶|    1  freq_dist = nltk.FreqDist(filtered_words)
                2  freq_dist
```

Out[119]:  FreqDist({'bike': 1034, 'crash': 945, 'car': 943, 'description': 847, '681
           0': 841, 'right': 790, 'lane': 780, 'left': 711, 'driver': 652, 'dc': 602,
           ...})

```
In [120]:  ▶|    1  police_word=freq_dist['police']
```

```
In [121]:  ▶|    1  print("count :",police_sentence)
```

count : 284

```
In [122]:  ▶|    1  print("count :",police_word)
```

count : 339

```
In [123]:  ▶|    1  #3.Determine and interpret the general sentiment of the comments about th
```

```
In [124]:  ▶|    1  from nltk.sentiment import SentimentIntensityAnalyzer
```

```python
In [125]:  ▶|    1  sia = SEIA()
```

```python
In [126]:  ▶|    1  mydata = sia.polarity_scores('.'.join(police))
```

```python
In [127]:  ▶|    1  mydata
```

Out[127]:  {'neg': 0.083, 'neu': 0.879, 'pos': 0.038, 'compound': -0.9998}

```python
In [128]:  ▶|    1  # Through the general sentiment of the word police we get to know that th
                 2  #rather than the police in the data i .e. on adding upto 1 it says negat
                 3  # of 0.879/1. And the compound scores means to calculate all the lexicon
                 4  # case the compound value -0.998 is less than -0.05 so it is negative ser
                 5
                 6  # Hence, it is NEGATIVE SENTIMENT ANALYSIS.
```