# Week 3 - Text Analysis & Working with Text Files

(Be sure to copy to drive)

Text data is a bit different from numeric data. We can easily find the average of a series of numbers and things like the highest and lowest values in a range to get some ideas on what we are dealing with. We can't really do that with text. We'll focus on some tools that you can use to actually analyze text. We'll start with a library called [TextBlob](#).

```python
1  #Load up our libraries
2  from textblob import TextBlob
3  from google.colab import drive
4
5  #these should look familar
6  import pandas as pd
7  import matplotlib.pyplot as plt
8  %matplotlib inline
9  import requests
10
11 #Some extra libraries we'll need for text analysis
12 import nltk
13 nltk.download('punkt')
14 nltk.download('brown')
15 nltk.download('punkt_tab')
16
17
18 #Connect to Gdrive
19 drive.mount('/content/gdrive')
20
21 print("Libraries and Drive Ready!")
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt.zip.
[nltk_data] Downloading package brown to /root/nltk_data...
[nltk_data]   Unzipping corpora/brown.zip.
[nltk_data] Downloading package punkt_tab to /root/nltk_data...
[nltk_data]   Unzipping tokenizers/punkt_tab.zip.
Mounted at /content/gdrive
Libraries and Drive Ready!
```
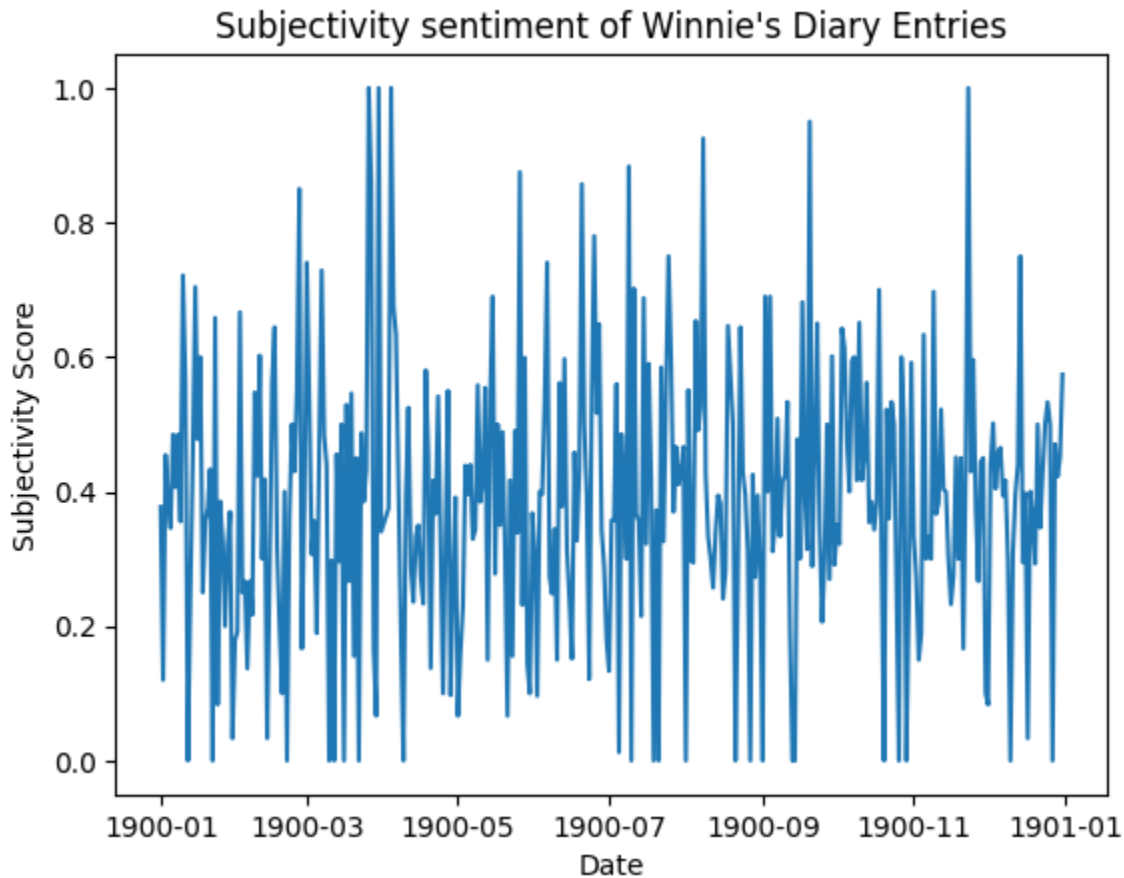
**Q2** Modify the following code cell to create a line graph of the *subjectivity* of her diary entries for the year.

```python
1  plt.plot(winnie_corpus["date"],winnie_corpus["subjectivity"])
2  plt.title("Subjectivity sentiment of Winnie's Diary Entries")
```

```
3 plt.ylabel("Subjectivity Score")
4 plt.xlabel("Date")
5 plt.show()
```



Subjectivity sentiment of Winnie's Diary Entries

## Noun Phrases for the Diary

Now let's generate the noun phrases for January's entries

```
 1 #we use some pandas work to just grab the January entries
 2 #stretch back into your memory to think about conditionals again
 3 jan_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-01-01') & (winnie_
 4
 5 jan_phrases = dict()
 6
 7 for entry in jan_corpus.entry:
 8
 9     tb = TextBlob(entry)
10     #we create a dictionary that will hold the noun phrases
11     #if it is the first time we see this np we put it in the dictionary
12     #if not, we must have a count already, so we increase that by one
13     for np in tb.noun_phrases:
14         if np in jan_phrases:
```

```
15              jan_phrases[np] += 1
16          else:
17              jan_phrases[np] = 1
18
19 #Print the top 10 things she mentioned in January
20
21 for np in sorted(jan_phrases, key=jan_phrases.get, reverse=True)[0:10]:
22      print(np, jan_phrases[np])
23
24
```

```
went 23
mamma 19
papa 16
beatrice 8
took 8
trusty 5
music lesson 4
mr read 4
annie gardiner 4
got 4
```

**Q5** Modify the next series of cells to generate noun phrases for the next 5 months of the year.

```
1 #February Entries
2 feb_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-02-01') & (winnie_
3
4 feb_phrases = dict()
5
6 for entry in feb_corpus.entry:
7      tb = TextBlob(entry)
8      for np in tb.noun_phrases:
9          if np in feb_phrases:
10              feb_phrases[np] += 1
11          else:
12              feb_phrases[np] = 1
13
14 #Print the top 10 things she mentioned in February
15
16 for np in sorted(feb_phrases, key=feb_phrases.get, reverse=True)[0:10]:
17      print(np, feb_phrases[np])
```

```
mamma 20
went 18
papa 12
beatrice 7
pay 6
was 5
took 5
got 5
music lesson 4
```

musical 4


```
 1 #March Entries
 2 mar_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-03-01') & (winnie_
 3
 4
 5 mar_phrases = dict()
 6
 7 for entry in mar_corpus.entry:
 8     tb = TextBlob(entry)
 9     for np in tb.noun_phrases:
10         if np in mar_phrases:
11             mar_phrases[np] += 1
12         else:
13             mar_phrases[np] = 1
14
15 #Print the top 10 things she mentioned in March
16
17 for np in sorted(mar_phrases, key=mar_phrases.get, reverse=True)[0:10]:
18     print(np, mar_phrases[np])
```

    mamma 21
    took 9
    papa 9
    went 8
    got 7
    pay 7
    dwyer 7
    mr perry 6
    godard 5
    beatrice 5


```
 1 #April Entries
 2 april_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-04-01') & (winni
 3
 4 april_phrases = dict()
 5
 6 for entry in april_corpus.entry:
 7     tb = TextBlob(entry)
 8     for np in tb.noun_phrases:
 9         if np in april_phrases:
10             april_phrases[np] += 1
11         else:
12             april_phrases[np] = 1
13
14 #Print the top 10 things she mentioned in April
15
16 for np in sorted(april_phrases, key=april_phrases.get, reverse=True)[0:10]:
17     print(np, april_phrases[np])
```

    mamma 24

```
   mamma 24
   papa 20
   went 11
   got 9
   pay 8
   took 7
   beatrice 7
   sullivan 7
   christ 6
   lizzie 5
```

```
 1 #May Entries
 2 may_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-05-01') & (winnie_
 3
 4 may_phrases = dict()
 5
 6 for entry in may_corpus.entry:
 7     tb = TextBlob(entry)
 8     for np in tb.noun_phrases:
 9         if np in may_phrases:
10             may_phrases[np] += 1
11         else:
12             may_phrases[np] = 1
13
14 #Print the top 10 things she mentioned in may
15
16 for np in sorted(may_phrases, key=may_phrases.get, reverse=True)[0:10]:
17     print(np, may_phrases[np])
```

```
   mamma 30
   papa 21
   went 17
   got 11
   ella 7
   lizzie 6
   rode 5
   helen 5
   pay 4
   dwyer 4
```

```
 1 #June Entries
 2 june_corpus = winnie_corpus[(winnie_corpus['date'] >= '1900-06-01') & (winnie_
 3
 4 june_phrases = dict()
 5
 6 for entry in june_corpus.entry:
 7     tb = TextBlob(entry)
 8     for np in tb.noun_phrases:
 9         if np in june_phrases:
10             june_phrases[np] += 1
11         else:
```

```
12            june_phrases[np] = 1
13
14 #Print the top 10 things she mentioned in june
15
16 for np in sorted(june_phrases, key=june_phrases.get, reverse=True)[0:10]:
17     print(np, june_phrases[np])
```

```
mamma 21
papa 16
got 13
went 13
took 7
lizzie 7
was 4
mrs. sullivan 4
ella 4
aunt lillie 4
```