## ⌄ Arsh misra

```
1 from google.colab import drive
2 import nltk
3 drive.mount('/content/drive')
4 nltk.download('punkt')
```

```
⤷  Drive already mounted at /content/drive; to attempt to forcibly remount, call drive.mount("/content/drive", force_remount=True).
    [nltk_data] Downloading package punkt to /root/nltk_data...
    [nltk_data]   Unzipping tokenizers/punkt.zip.
    True
```

I decided to use google drive on colab because spyder does not work for me. The content of the files is the same aside from the input text, which I will list in full in this notebook.

```
1 import os
2
3 # Path to the directory
4 path = '/content/drive/My Drive/NLPHWPOSNEG/'
5
6 # List files
7 files = os.listdir(path)
8 print(files)
```

```
⤷  ['positive-words.txt', 'input-text2.txt', 'negative-words.txt', 'example-text.txt']
```

## ⌄ 1. Create a function called gen_senti that tokenizes arbritrary text and compares each token with the positive and negative lexicons of each dictionary and ouputs the sentiment score S. Positive and negative words, pw and nw, count as a score of 1 and -1, respectively, for each word matched. The total count for pw and nw are pc and nc, respectively. Each message sentiment, S, is normalized between -1 and 1. Any text that does not have any positive or negative words would have been ignored, and not scored.

I used a different example since the question did not specify that we needed to use the original example. The input text is listed as:

"Karl Peterson has been living the cruel inverse of the American dream. His rent keeps getting higher, but his apartments keep getting smaller. Peterson left the Midwest nine years ago for the epicenter of an economic boom, only to gradually learn that endless sunshine and desert views are increasingly among the few bargains left in Arizona. Peterson married his wife and they struggled to save for a home, moving through four apartments as their rent nearly tripled from $625 to 1,800$ a month"

This is an article taken from yahoo finance: https://finance.yahoo.com/news/american-despair-arizona-high-home-151927862.html

```
 1 from nltk.tokenize import word_tokenize
 2
 3 with open('/content/drive/My Drive/NLPHWPOSNEG/input-text2.txt', 'r', encoding='ISO-8859-1') as f:
 4     input_text = f.read()
 5
 6 def gen_senti(input_text, path="/content/drive/My Drive/NLPHWPOSNEG/"):
 7     # Load positive words
 8     with open(path + 'positive-words.txt', 'r') as f:
 9         positive_words = set(f.read().splitlines())
10
11     # Load negative words
12     with open(path + 'negative-words.txt', 'r') as f:
13         negative_words = set(f.read().splitlines())
14
15     tokens = word_tokenize(input_text.lower())
16
17     # counts and lists to store matched words
18     pc = 0  # Positive count
```

```
19      nc = 0  # Negative count
20      matched_positive_words = []  # To store matched positive words
21      matched_negative_words = []  # To store matched negative words
22
23      # positive and negative tokens
24      for token in tokens:
25          if token in positive_words:
26              pc += 1
27              matched_positive_words.append(token)  # Add to matched positive words
28          elif token in negative_words:
29              nc += 1
30              matched_negative_words.append(token)  # Add to matched negative words
31
32      # sentiment score
33      total_count = pc + nc
34      if total_count > 0:
35          S = (pc - nc) / total_count  # sentiment score between -1 and 1
36      else:
37          S = 0  # Neutral sentiment if no matched tokens
38
39      return S, pc, nc, matched_positive_words, matched_negative_words
40
41 # function with the loaded input text
42 score, positive_count, negative_count, matched_pos, matched_neg = gen_senti(input_text)
43 print(f"Sentiment Score: {score}, Positive Count: {positive_count}, Negative Count: {negative_count}")
44 print(f"Matched Positive Words: {matched_pos}")
45 print(f"Matched Negative Words: {matched_neg}")
46
```

```
Sentiment Score: -0.5, Positive Count: 1, Negative Count: 3
Matched Positive Words: ['boom']
Matched Negative Words: ['cruel', 'desert', 'struggled']
```

## 2. Using the dataframe from the lecture, the_data, column body, apply this function to each corpus and add a column called simple_senti.

```
1 nltk.download('stopwords')
2 import pandas as pd
3
4 the_path = "/content/drive/My Drive/NLP Text Data/"
5 # Some preliminary cleaning functions
6 def clean_txt(var_in):
7     import re
8     # Modify the regex to preserve punctuation and capitalization
9     tmp_t = re.sub("[^A-Za-z0-9!@#\$%\^&\*\(\)\-_\+=\.,;:!?'\"]+", " ", var_in).strip()
10
11     return tmp_t
12
13 def read_file(full_path_in):
14     with open(full_path_in, "r", encoding="UTF-8") as f_t:
15         text_t = f_t.read() # Reads the entire file
16         text_t = clean_txt(text_t)
17     return text_t
18
19 def file_crawler(path_in):
20     import os
21     import pandas as pd
22     my_pd_t = pd.DataFrame()
23     for root, dirs, files in os.walk(path_in, topdown=False):
24         for name in files:
25             try:
26                 txt_t = read_file(root + "/" + name)
27                 if len(txt_t) > 0:
28                     the_lab = root.split("/")[-1]
29                     tmp_pd = pd.DataFrame({"body": txt_t, "label": the_lab}, index=[0])
30                     my_pd_t = pd.concat([my_pd_t, tmp_pd], ignore_index=True)
31             except Exception as e:
32                 print(f"Error with file {root}/{name}: {e}")
33                 pass
34     return my_pd_t
35
36 df = file_crawler(the_path)
37 print("Sample of the DataFrame:")
```

```
38 print(df.head())  # Print the first few rows of the DataFrame
39
```

```
[nltk_data] Downloading package stopwords to /root/nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK MongoBD report Nov122018.xls: 'utf-8' codec can't decode byte 0xd0 in p
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK segment count Nov122018.xlsx: 'utf-8' codec can't decode bytes in posit
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK vendor count Nov122108 .xlsx: 'utf-8' codec can't decode bytes in posit
Sample of the DataFrame:
                                                body          label
0  Machine Learning Total 239.99 Computer Science...  machinelearning
1  Rendezvous Server to the Rescue: Dealing with ...  machinelearning
2  The 10 Algorithms Machine Learning Engineers N...  machinelearning
3  Find a Job in Artificial Intelligence or Machi...  machinelearning
4  xkcd: Machine Learning Archive What If? Blag S...  machinelearning
```

```
 1
 2 # Load positive and negative words
 3 def load_word_lists(path):
 4     with open(path + 'positive-words.txt', 'r') as f:
 5         positive_words = set(f.read().splitlines())
 6     with open(path + 'negative-words.txt', 'r') as f:
 7         negative_words = set(f.read().splitlines())
 8     return positive_words, negative_words
 9
10 # sentiment analysis function
11 def gen_senti(input_text, positive_words, negative_words):
12     # Tokenize the input text
13     tokens = word_tokenize(input_text.lower())
14
15     # Initialize counts
16     pc = 0  # Positive count
17     nc = 0  # Negative count
18
19     # Count positive and negative tokens
20     for token in tokens:
21         if token in positive_words:
22             pc += 1
23         elif token in negative_words:
24             nc += 1
25
26     # Calculate sentiment score
27     total_count = pc + nc
28     if total_count > 0:
29         S = (pc - nc) / total_count  # Normalized sentiment score between -1 and 1
30     else:
31         S = 0  # Neutral sentiment if no matched tokens
32
33     return S
34
35
36 path = "/content/drive/My Drive/NLPHWPOSNEG/"
37 positive_words, negative_words = load_word_lists(path)
38
39 the_path = "/content/drive/My Drive/NLP Text Data/"
40 df = file_crawler(the_path)
41
42 # Apply sentiment analysis to the 'body' column
43 df['simple_senti'] = df['body'].apply(lambda x: gen_senti(x, positive_words, negative_words))
44
45 # Display the updated DataFrame
46 print("Updated DataFrame with Sentiment Scores:")
47 print(df[['body', 'simple_senti']].head())  # Display the 'body' and 'simple_senti' columns
48
```

```
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK MongoBD report Nov122018.xls: 'utf-8' codec can't decode byte 0xd0 in p
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK segment count Nov122018.xlsx: 'utf-8' codec can't decode bytes in posit
Error with file /content/drive/My Drive/NLP Text Data/fishing/UK vendor count Nov122108 .xlsx: 'utf-8' codec can't decode bytes in posit
Updated DataFrame with Sentiment Scores:
                                                body  simple_senti
0  Machine Learning Total 239.99 Computer Science...      0.170732
1  Rendezvous Server to the Rescue: Dealing with ...      0.629630
2  The 10 Algorithms Machine Learning Engineers N...      0.411765
3  Find a Job in Artificial Intelligence or Machi...      0.730769
4  xkcd: Machine Learning Archive What If? Blag S...      0.000000
```

## 3. Using vaderSentiment, apply the compound value of sentiment for each corpus in column body on a new column of the_data called vader.

```
1 # Install and import necessary libraries
2 import nltk
3 from nltk.sentiment.vader import SentimentIntensityAnalyzer
4 import pandas as pd
5
6 # VADER lexicon
7 nltk.download('vader_lexicon')
8
9 #VADER sentiment analyzer
10 sia = SentimentIntensityAnalyzer()
11
12
13 # Function to compute the compound score
14 def compute_vader_sentiment(text):
15     score = sia.polarity_scores(text)
16     return score['compound']  # Return the compound score
17
18 # Apply the function to the 'body' column and create a new 'vader' column
19 df['vader'] = df['body'].apply(compute_vader_sentiment)
20
21 # Display the updated DataFrame
22 print("Updated DataFrame with VADER Sentiment Scores:")
23 print(df[['body', 'vader']].head())  # Display the 'body' and 'vader' columns
24
```

```
[nltk_data] Downloading package vader_lexicon to /root/nltk_data...
Updated DataFrame with VADER Sentiment Scores:
                                        body    vader
0  Machine Learning Total 239.99 Computer Science...  0.9941
1  Rendezvous Server to the Rescue: Dealing with ...  0.9968
2  The 10 Algorithms Machine Learning Engineers N...  0.9473
3  Find a Job in Artificial Intelligence or Machi...  0.9992
4  xkcd: Machine Learning Archive What If? Blag S...  0.0000
```

## 4. Compute the mean, median, and standard deviations of both sentiment measures, simple_senti and vader

```
1 import pandas as pd
2
3 # mean, median, and standard deviation for both sentiment measures
4 statistics = {
5     'simple_senti': {
6         'mean': df['simple_senti'].mean(),
7         'median': df['simple_senti'].median(),
8         'std_dev': df['simple_senti'].std()
9     },
10     'vader': {
11         'mean': df['vader'].mean(),
12         'median': df['vader'].median(),
13         'std_dev': df['vader'].std()
14     }
15 }
16
17 # results
18 print("Sentiment Statistics:")
19 for measure, stats in statistics.items():
20     print(f"{measure} - Mean: {stats['mean']:.4f}, Median: {stats['median']:.4f}, Standard Deviation: {stats['std_dev']:.4f}")
21
```

```
Sentiment Statistics:
simple_senti - Mean: 0.3784, Median: 0.4286, Standard Deviation: 0.4338
vader - Mean: 0.8758, Median: 0.9851, Standard Deviation: 0.2641
```

we see that simple_senti and vader differ significantly across all measurements.This is mostly because vader is a lot more flexible and smarter than a simple matching function like simple_senti.

Homework 2 concluded

####################

####################

+ Code    + Text