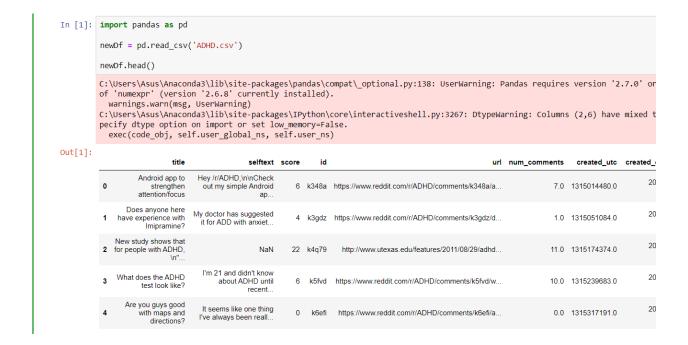
CS 8920 - Independent Study Project Report

K M Tawsik Jawad

U01034189

Domain: Finding out commonly used word sequences in the reddit posts of patients suffering from ADHD. The 2nd dataset is collected from the channels like r/anxiety, r/depression etc. from reddit posts. Same technique in patients suffering from different mental health issues in the 2nd dataset was applied to find out the most common word sequences in both of them. Performing word tokenization and padding, one hot encoding to feed into LSTM to predict next words from commonly occurring sequences of words.

Data Pre-Processing: Since only the posts are required for analysis, other columns of the dataset were dropped along with null columns. Then necessary pre-processing was performed as follows:



```
In [3]: newDf = newDf.drop(['selftext','url','num_comments','created_utc','score', 'id', 'created_datetime'],axis='columns')
           newDf
Out[3]:
             4815
                     My doctor lost his licence a few days from nex...
            78504
                     Besides not being neurotypical, is anyone else.
            56091
                                             Casinos, yes or no?
            77259 Anybody here has had success with the Getting .
           332202
           296548
                        Hardcore Avoidance Recently- anyone else?
            73617
                   Drugs made my life much better than i thought ...
           125490
                                                Ritalin euphoria?
            34990
                    I know this is ridiculous but can someone plea...
          6721 rows × 1 columns
```

Next, the unnecessary urls are removed from the corpus. The entire corpus is converted to lowercase and removal of extra white spaces was done for easier analysis.

```
In [128]: newDf['title'] = newDf['title'].str.lower()
    newDf['title']
  Out[128]: 43598
                      when you feel a new hyper fixation coming... (...
             11093
                      i can't seem to keep friends for very long and...
             10600
                                             trouble processing auditory
             11294
                                      what's your favorite adhd podcast?
                                                     energy? what's that?
             8913
             23356
                                           constant cold hands and feet?
             7450
                                                             mv angel dog
                              jumbly brain and expressing ideas/thoughts
             33704
             4599
                      how i am gamifying my life to increase product...
                      2:34am and i've just finished this! i meant to...
             Name: title, Length: 888, dtype: object
  In [129]: def remove_punctuation(text):
    no_punctuation = "".join(
                                     ".join([word for word in text if word not in PUNCTUATION])
                 return no_punctuation
             newDf['title'] = newDf['title'].apply(remove_punctuation)
             newDf['title'].head()
  Out[129]: 43598
                      when you feel a new hyper fixation coming post...
             11093
                      i can't seem to keep friends for very long and...
             10600
                                             trouble processing auditory
                                       what's your favorite adhd podcast
             11294
                                                       energy what's that
             8913
             Name: title, dtype: object
  In [130]: def remove_extra_white_spaces(text):
                  single_char_pattern = r'\s+[a-zA-Z]\s+'
                 without_sc = re.sub(pattern=single_char_pattern, repl=" ", string = text)
                 return without_sc
             newDf['title'] = newDf['title'].apply(remove_extra_white_spaces)
```

Next, we performed word lemmatization which basically reduces a word to its root form. The reason for choosing lemmatization in stead of stemming is to keep a better root mapping of the word. If the word extracted by stemming from 'better' is 'bet', for lemmatization it would be 'good'.

```
In [133]: from nltk.stem import WordNetLemmatizer
           def lemmatize text(text):
               lemmatizer = WordNetLemmatizer()
               tokens = word_tokenize(text)
               for i in range(len(tokens)):
                   lemma word = lemmatizer.lemmatize(tokens[i])
                   tokens[i] = lemma_word
               return " ".join(tokens)
           newDf['title'] = newDf['title'].apply(lemmatize_text)
newDf['title'].head()
                                   feel new hyper fixation coming post
           11093
                    ' seem keep friend long ' thing make different...
           10600
                                            trouble processing auditory
           11294
                                                ' favorite adhd podcast
           8913
           Name: title, dtype: object
```

Moving on to the last part of the pre-processing is the generation of N-Grams. Experiments were performed with unigrams, bigrams, trigrams. Here the most commonly used bigrams in the entire corpus is shown:

```
def generate_N_grams(text,ngram=1):
                         text = re.sub(r'[^a-zA-Z0-9\s]', '', text)
tokens = [token for token in text.split(" "
output = list(ngrams(tokens, ngram))
return output
                                                                                                               ") if token != ""]
                          return output
In [21]: lst = []
                  for X in newDf['title']:
                           y = generate_N_grams(X,2)
                           1st.append(y)
                   def flatten(input):
                          new_list = []
                           for i in input:
                                  for i in i:
                                          new_list.append(j)
                           return new_list
                   lst1 = flatten(lst)
                   finalList1 = []
                   from collections import Counter
                   c = Counter(lst1)
                  M = c.most_common(30)
                   for key, val in M:
                           finalList1.append(key)
                   print(finalList1)
                  [('anyone', 'else'), ('feel', 'like'), ('diagnosed', 'adhd'), ('adhd', 'med'), ('adderall', 'xr'), ('adhd', 'medication'), (
rst', 'time'), ('people', 'adhd'), ('need', 'advice'), ('need', 'help'), ('adult', 'adhd'), ('side', 'effect'), ('dont', 'kr
w'), ('got', 'diagnosed'), ('think', 'adhd'), ('recently', 'diagnosed'), ('year', 'old'), ('anybody', 'else'), ('taking', '&
rall'), ('adhd', 'symptom'), ('like', 'im'), ('please', 'help'), ('someone', 'adhd'), ('take', 'med'), ('long', 'term'), ('&
ing', 'diagnosed'), ('finally', 'got'), ('first', 'day'), ('adhd', 'diagnosis'), ('else', 'feel')]
```

Next, we perform the similar pre-processing steps of data cleaning in the 2nd dataset of mental health issues and generate N-Grams to find out what common N-Grams occur in the 2nd dataset.

```
In [15]: import pandas as pd
            df = pd.read_csv('anxiety_2018_features_tfidf_256.csv')
Out[15]:
                                                               post automated_readability_index coleman_liau_index flesch_kincaid_grade_level flesch_reading_ease gulpease
                                                             anyone
else like
taking
long
walks
                   arxiety Watch_Me_Get_ 2018/01/01
                                                                                         5.007692
                                                                                                                                            3.283462
                                                          Meditation
                                                                                                                                                                  95.231515
                                                           Rant
about
anxiety
meds -
I've been
self
                                                                                         6.582939
                                                                                                                6.555885
                                                                                                                                            6.839898
                                                                                                                                                                                   65.2
                            coffeegreentea 2018/01/01
                                                                                         5.459158
                                                                                                                6.349314
                                                                                                                                            6.922802
                                                                                                                                                                  72.443031
```

```
df['post'] = df['post'].apply(remove_punctuation)
df['post'] = df['post'].apply(remove_extra_white_spaces)
df['post'] = df['post'].apply(remove_stopwords)
df['post'] = df['post'].apply(drop_duplicates)
df['post'] = df['post'].apply(lemmatize_text)
lst = []
for X in df['post']:
      y = generate_N_grams(X,2)
      lst.append(y)
def flatten(input):
     new_list = []
      for i in input:
           for j in i:
                new_list.append(j)
      return new list
lst1 = flatten(lst)
finalList2 = []
from collections import Counter
c = Counter(lst1)
M = c.most_common(30)
for key, val in M:
      finalList2.append(key)
print(finalList2)
```

[('feel', 'like'), ('panic', 'attack'), ('dont', 'know'), ('anyone', 'else'), ('dont', 'want'), ('like', 'im'), ('im', 'goin g'), ('anxiety', 'attack'), ('even', 'though'), ('social', 'anxiety'), ('first', 'time'), ('year', 'ago'), ('make', 'feel'), ('year', 'old'), ('felt', 'like'), ('every', 'time'), ('feeling', 'like'), ('gon', 'na'), ('im', 'sure'), ('anxiety', 'depressi on'), ('really', 'bad'), ('im', 'scared'), ('every', 'dav'), ('month', 'ago'), ('side', 'effect'), ('high', 'school'), ('last', 'like'), ('gon', 'na'), ('month', 'ago'), ('side', 'effect'), ('high', 'school'), ('last', 'month', 'ago'), ('side', 'effect'), ('last', 'month', 'ago'), ('side', 'effect'), ('last', 'month', 'last', 'month', 'month',

Finally, we use the set operator in the two lists of that contain the most common sequences to find out exactly which word sequences are there in each of the dataset:

```
In [25]: Y = set(finalList2) & set(finalList1)
print(Y)
{('like', 'im'), ('dont', 'know'), ('side', 'effect'), ('feel', 'like'), ('year', 'old'), ('first', 'time'), ('anyone', 'els e')}
```

Next, we analyze the n-gram sequences using trigrams. First lets see what sequences are most common using trigrams in the ADHD dataset:

```
In [26]: lst = []
                         for X in newDf['title']:
                                   y = generate_N_grams(X,3)
lst.append(y)
                        def flatten(input):
                                   new_list = []
for i in input:
                                             for j in i:
                                                      new_list.append(j)
                                   return new_list
                        lst1 = flatten(lst)
                        finalList1 = []
                        from collections import Counter
                         c = Counter(lst1)
                        M = c.most_common(30)
                         for key,val in M:
                                   finalList1.append(key)
                        print(finalList1)
                       [('feel', 'like', 'im'), ('anyone', 'else', 'feel'), ('ever', 'feel', 'like'), ('anyone', 'else', 'experience'), ('recently', 'diagnosed', 'adhd'), ('else', 'feel', 'like'), ('anyone', 'else', 'find'), ('got', 'diagnosed', 'adhd'), ('rejection', 'sensit ive', 'dysphoria'), ('anyone', 'else', 'issue'), ('anyone', 'else', 'struggle'), ('kick', 'start', 'sunday'), ('anyone', 'else', 'problem'), ('anyone', 'else', 'get'), ('med', 'first', 'time'), ('make', 'feel', 'like'), ('think', 'might', 'adhd'), ('anyone', 'else', 'rreally'), ('finally', 'got', 'diagnosed'), ('get', 'stuff', 'done'), ('short', 'term', 'memory'), ('started', 'taking', 'adderall'), ('could', 'use', 'advice'), ('feel', 'like', 'failure'), ('diagnosed', 'adult', 'adhd'), ('get', 'thing', 'done'), ('taking', 'adhd', 'medication'), ('might', 'adhd', 'i
```

Next, we do the similar approach in the 2nd dataset of mental health:

```
for X in df['post']:
           y = generate_N_grams(X,3)
           lst.append(y)
 def flatten(input):
           new_list = []
           for i in input:
                   for j in i:
                            new list.append(j)
           return new_list
 lst1 = flatten(lst)
  finalList2 = []
 from collections import Counter
  c = Counter(lst1)
 M = c.most_common(30)
 for key, val in M:
           finalList2.append(key)
 print(finalList2)
[('feel', 'like', 'im'), ('anyone', 'else', 'experience'), ('anyone', 'else', 'get'), ('anyone', 'else', 'feel'), ('make', 'fee l', 'like'), ('anxiety', 'panic', 'attack'), ('im', 'gon', 'na'), ('like', 'im', 'going'), ('long', 'story', 'short'), ('dont', 'know', 'im"), ('first', 'panic', 'attack'), ('dont', 'even', 'know'), ('feel', 'like', 'cant'), ('feel', 'like', 'going'), ('a nxiety', 'feel', 'like', 'ive'), ('always', 'feel', 'like'), ('feel', 'like', 'ive'), ('always', 'feel', 'like'), ('feel', 'like', 'ive'), ('always', 'feel', 'like'), ('every', 'single', 'day'), ('anyone', 'else', 'experienced'), ('feel', 'like', 'everyone'), ('first', 'time', 'postin g'), ('even', 'though', 'know'), ('feel', 'like', 'dont'), ('feeling', 'like', 'im')]
```

And find out the common trigrams in the two datasets in the following:

```
In [28]: Y = set(finalList2) & set(finalList1)
print(Y)|
{('make', 'feel', 'like'), ('anyone', 'else', 'feel'), ('anyone', 'else', 'get'), ('feel', 'like', 'im'), ('anyone', 'else', 'e
xperience')}
```

Let's analyze the common sequences in bigrams and trigrams putting the results together in the following two snapshots. First one is the bigram and second is trigram:

Although the bigram produces a greater number of common sequences than trigrams, the bigrams don't really put too much meaning in the sequences. If we look at trigrams, they are shown to provide more meaning with more words. So, intuition about the next words in the sentence can be better developed using tri-grams rather than bigrams for this case.

Next, women suffering from ADHD also put up posts in reddit. Analyzing this "adhdwomen.csv" dataset after performing the similar pre-processing as the previous datasets, they were prepared to be used for a word prediction task. Here it can be seen in the picture below how word tokenization is performed on the dataset to assign an id to each word in the corpus.



```
In [29]: import pandas as pd
          import os
          import numpy as np
          import tensorflow as tf
          from tensorflow.keras.preprocessing.sequence import pad_sequences
          from tensorflow.keras.layers import Embedding, LSTM, Dense, Bidirectional
          from tensorflow.keras.preprocessing.text import Tokenizer
          from tensorflow, keras, models import Sequential
          from tensorflow.keras.optimizers import Adam
          tokenizer = Tokenizer(oov_token='<oov>')
          tokenizer.fit_on_texts(newDf['title'])
          total words = len(tokenizer.word index) + 1
          print("Total number of words: ", total_words)
          print("Word: ID")
print("adhd: ", tokenizer.word_index['adhd'])
print("anybody: ", tokenizer.word_index['anybody'])
          Total number of words: 1942
          Word: ID
          adhd: 2
          anybody: 158
```

Next, I convert the texts in the corpus into number sequences. Then, normalize the dataset by taking the max length sentence and padding the rest of the sentences with the exact same length in the following way:

Next, we split the dependent and independent variables. Since we are formulating a word prediction problem, the final word id according to the input sequences is taken as the dependent variable and the ids before that are the independent variables. Analyzing the x and y it can be clear:

Next, I perform one hot encoding on the labels variable to prepare the dependent set to be fed to the LSTM network. Since I wanted to perform the analysis on a random sample of the dataset and LSTM training is a convoluted network, iteration number was kept to 20 where in each step loss value and accuracy metric is seen. For updating weight networks, adam optimizer was used for its good performance in deep learning NLP tasks.

```
In [21]: model = Sequential()
         model.add(Embedding(total_words, 100, input_length=max_sequence_len-1))
         model.add(Bidirectional(LSTM(150)))
         model.add(Dense(total_words, activation='softmax'))
         adam = Adam(1r=0.01)
         model.compile(loss='categorical_crossentropy', optimizer=adam, metrics=['accuracy'])
         history = model.fit(xs, ys, epochs=20, verbose=1)
         #print model.summary()
         print(model)
         WARNING:tensorflow:From C:\Users\Asus\Anaconda3\lib\site-packages\tensorflow\python\ops\resource_variable_ops.py:435: colocate_
         with (from tensorflow.python.framework.ops) is deprecated and will be removed in a future version.
         Instructions for updating:
         Colocations handled automatically by placer.
         WARNING:tensorflow:From C:\Users\Asus\Anaconda3\lib\site-packages\tensorflow\python\ops\math_ops.py:3066: to_int32 (from tensor
         flow.python.ops.math_ops) is deprecated and will be removed in a future version.
         Instructions for updating:
         Use tf.cast instead.
         Epoch 1/20
         Epoch 2/20
         3999/3999 [============== ] - 9s 2ms/sample - loss: 6.7427 - acc: 0.0413
         Epoch 3/20
         3999/3999 [=
                  Epoch 4/20
         Epoch 5/20
         3999/3999 [=
                 Epoch 6/20
         3999/3999 [============= ] - 11s 3ms/sample - loss: 1.8161 - acc: 0.5934
         Epoch 7/20
         3999/3999 [=
                   3999/3999 [=
                  Fnoch 9/20
         3999/3999 [=:
```

Finally, I performed the text prediction with predicting next 20 words after training the sequences with LSTM. For the input, I used one of the common trigrams "anyone else experience" which was found in both the trigrams of ADHD dataset and the mental health dataset. The result is as follows:

anyone else experience overwhelm avalanche feel defeated general like suicidey would share something helpful im overwhelmed try ing organize prioritize figure daily task night

```
In [35]: seed_text = "feel like im"
    next_words = 20

for _ in range(next_words):
    token_list = tokenizer.texts_to_sequences([seed_text])[0]
    token_list = pad_sequences([token_list], maxlen=max_sequence_len-1, padding='pre')
    predicted = model.predict_classes(token_list, verbose=0)
    output_word = ""
    for word, index in tokenizer.word_index.items():
        if index == predicted:
            output_word = word
            break
        seed_text += " " + output_word
        print(seed_text)
```

feel like im really productive day better failure planner planner love swear erin condren seems month start missing deadline relatable read sentence immediately

Conclusion: So, analyzing the output sentences it can be inferred that the women suffering from ADHD have a hard time adjusting to the daily functions and many of them have higher degree of depressive symptoms. It can also be inferred that many of them are trying to organize their day-to-day activities, to focus more on their priorities and make a productive outcome keeping the ADHD symptoms aside. As usual with any kind of sickness, they are hoping to get help from people suffering from similar symptoms. Lastly, it can be observed that Adderall is used commonly when medication is provided as input where some common side effects of the medicine is seen.