```python
import numpy as np
import pandas as pd
import warnings
import math
import torch
import re
import torch.nn as nn
import torch.optim as optim
from torch import tensor as tt
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from sklearn.metrics import f1_score
import matplotlib.pyplot as plt
warnings.filterwarnings("ignore")
```

```python
enc_type = 'utf-8'
```

```python
In [3]:   1   def data_clean_tweets(txt_file):
          2
          3       with open(txt_file,'r',encoding=enc_type) as f :
          4           data = f.readlines()
          5
          6       data = [dt.strip() for dt in data]
          7       for text1 in data :
          8           # Remove unnecessary punctuations,spaces brackets
          9           text1 = re.sub(r'\[\d+\]', '', text1)
         10           # Remove underscores and hyphens from the sides of words
         11           text1 = re.sub(r'[-_]', '', text1)
         12           # Remove numbers used for points
         13           text1 = re.sub(r'\d+\.', '', text1)
         14           # remove "
         15           text1 = re.sub(r'"', '', text1)
         16           # remove '
         17           text1 = re.sub(r"'", '', text1)
         18           # remove special characters
         19           text1 = re.sub(r'[+-.,!@#$%^&<>?/\{}()*_=:;|]', '', text1)
         20           # remove multiple spaces
         21           text1 = re.sub(r'\n\s*\n', '\n', text1)
         22           # remove numbers
         23           text1 = re.sub(r'\s\d+\s', ' ', text1)
         24           # remove new line, tabs
         25           text1 = re.sub(r'\n|\t', ' ', text1)
         26           # Remove bullets
         27           text1 = re.sub(r'^[\s\u2022\u2023\u25E6\u2043]*', '', text1, flags=re.MULTILINE)
         28           # remove multiple spaces
         29           text1 = re.sub(r' +', ' ', text1)
         30           # convert the entire text lo lower case
         31           text1 = text1.lower()
         32       return data
```

```python
In [4]:   1   def read_label_data(label_data):
          2       with open(label_data,'r',encoding = enc_type) as f :
          3           data = f.readlines()
          4       data = [dt.strip() for dt in data]
          5       return data
```

```
In [5]:  1  train_tweets = data_clean_tweets('sentiment/train_text.txt')
         2  train_labels = read_label_data('sentiment/train_labels.txt')
         3  test_tweets = data_clean_tweets('sentiment/test_text.txt')
         4  test_labels = read_label_data('sentiment/test_labels.txt')
         5  val_tweets = data_clean_tweets('sentiment/val_text.txt')
         6  val_labels = read_label_data('sentiment/val_labels.txt')
```

```
In [6]:  1  train_tweets
```

```
Out[6]:  ['"QT @user In the original draft of the 7th book, Remus Lupin survived the Battle of Hogwarts. #HappyBir
         thdayRemusLupin"',
          '"Ben Smith / Smith (concussion) remains out of the lineup Thursday, Curtis #NHL #SJ"',
          'Sorry bout the stream last night I crashed out but will be on tonight for sure. Then back to Minecraft
         in pc tomorrow night.',
          "Chase Headley's RBI double in the 8th inning off David Price snapped a Yankees streak of 33 consecutive
         scoreless innings against Blue Jays",
          '@user Alciato: Bee will invest 150 million in January, another 200 in the Summer and plans to bring Mes
         si by 2017"',
          "@user LIT MY MUM 'Kerry the louboutins I wonder how many Willam owns!!! Look Kerry Warner Wednesday!'",
          '"\\"""" SOUL TRAIN\\"""" OCT 27 HALLOWEEN SPECIAL ft T.dot FINEST rocking the mic...CRAZY CACTUS NIGHT
         CLUB ..ADV ticket $10 wt out costume $15..."',
          'So disappointed in wwe summerslam! I want to see john cena wins his 16th title',
          '"This is the last Sunday w/o football .....,NFL is back baby"',
          "@user @user CENA & AJ sitting in a tree K-I-S-S-I-N-G 1st goes AJ's  job then John's cred then goes Vic
         ki with the GM position.",
          '@user Well said on HMW. Can you now address why Texans fans file out of the stadium midway through the
         4th qtr of every game?',
          "Just said hello to Dennis Kucinich as he walked casually through campus with his #hotwife. He's on 22nd
```

```
In [7]:  1  train_labels
```

```
Out[7]:  ['2',
          '1',
          '1',
          '1',
          '2',
          '2',
          '2',
          '0',
          '2',
          '1',
          '1',
          '1',
          '2',
          '0',
          '1',
          '1',
          '2',
          '2',
          '0',
```

```
In [8]:   1  # define a function to read the lexicon files
          2  def read_lexicons_files(lex_files):
          3      with open(lex_files, "r", encoding = enc_type) as ff:
          4          data = ff.readlines()
          5  #      basic cleaning of data
          6      data = [dt.strip().split('\t') for dt in data]
          7      lexicon_dict = {}
          8      for v in data:
          9          if len(v) == 3:
         10              key = v[0]
         11              value = {'-ve': float(v[1]), '+ve': float(v[2])}
         12              lexicon_dict[key] = value
         13
         14      return lexicon_dict
```

```
In [9]:    1  all_lexicon_files = ['3DS.tsv', '4chan.tsv', '2007scape.tsv', 'ACTrade.tsv',
           2                       'amiugly.tsv', 'BabyBumps.tsv', 'baseball.tsv', 'canada.tsv',
           3                       'CasualConversation.tsv', 'DarknetMarkets.tsv', 'darksouls.tsv', 'elderscrollsonline.tsv
           4                       'Eve.tsv', 'Fallout.tsv', 'fantasyfootball.tsv', 'GameDeals.tsv', 'gamegrumps.tsv', 'hal
           5                       'Homebrewing.tsv', 'IAmA.tsv', 'india.tsv', 'jailbreak.tsv', 'Jokes.tsv', 'KerbalSpacePr
           6                       'Keto.tsv', 'leagueoflegends.tsv', 'Libertarian.tsv', 'magicTCG.tsv', 'MakeupAddiction.t
           7                       'Naruto.tsv', 'nba.tsv', 'oculus.tsv', 'OkCupid.tsv', 'Parenting.tsv', 'pathofexile.tsv'
           8                       'raisedbynarcissists.tsv', 'Random_Acts_Of_Amazon.tsv', 'science.tsv', 'Seattle.tsv',
           9                       'TalesFromRetail.tsv', 'talesfromtechsupport.tsv', 'ultrahardcore.tsv', 'videos.tsv',
          10                       'Warthunder.tsv', 'whowouldwin.tsv', 'xboxone.tsv', 'yugioh.tsv']
```

```
In [10]:   1  adj = "adjectives/2000.tsv"
           2  freq = "adjectives/2000.tsv"
```
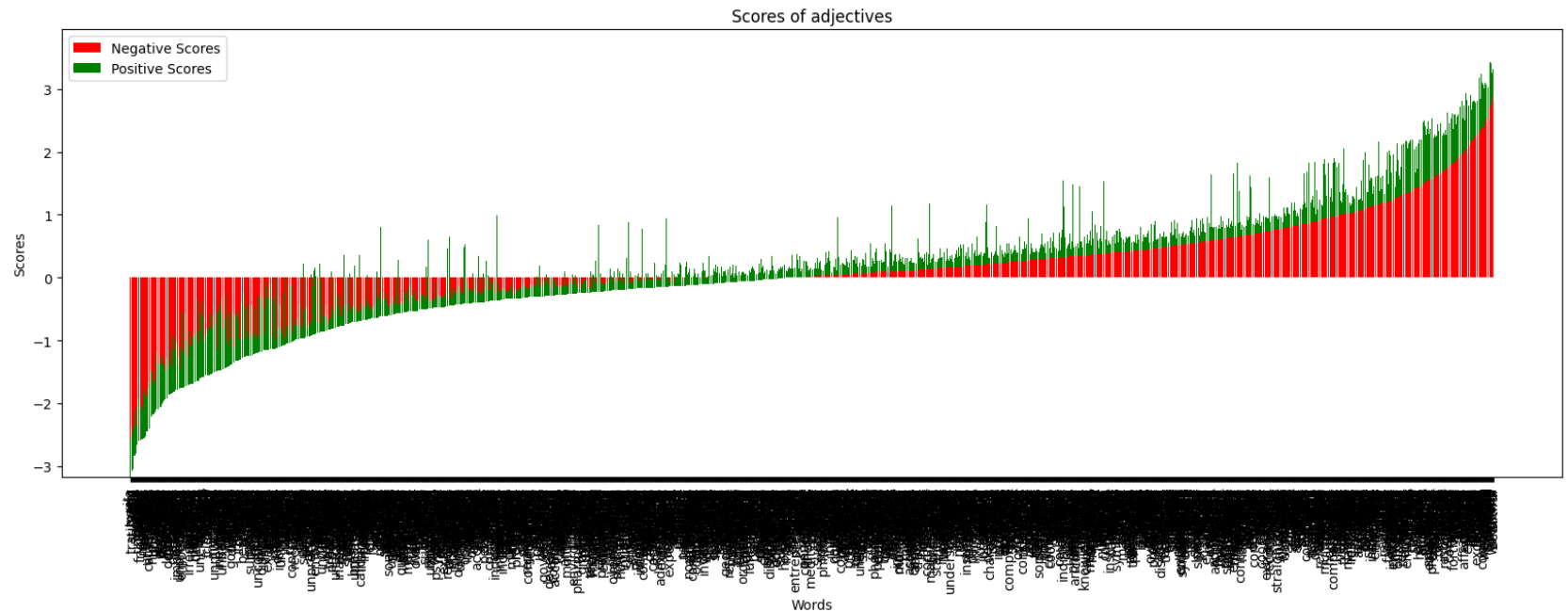
```
In [11]:   1  # reading the adjectives file and frequency of the adjective files
           2  adjectives = read_lexicons_files(adj)
           3  frequency = read_lexicons_files(freq)
```

```
In [12]:   1  negative_scores = [adjectives[word]['-ve'] for word in adjectives]
           2  positive_scores = [adjectives[word]['+ve'] for word in adjectives]
```

```
In [13]:   1  x_axis = list(adjectives.keys())
```

```
In [14]:    1  # plotting a graph for words in adjectives file
            2  plt.figure(figsize=(20, 6))
            3  plt.bar(x_axis, negative_scores, color='red', label='Negative Scores')
            4  plt.bar(x_axis, positive_scores, color='green', bottom=negative_scores, label='Positive Scores')
            5  plt.xticks(rotation=90)
            6  plt.title('Scores of adjectives')
            7  plt.xlabel('Words')
            8  plt.ylabel('Scores')
            9  plt.legend()
           10  plt.show()
```
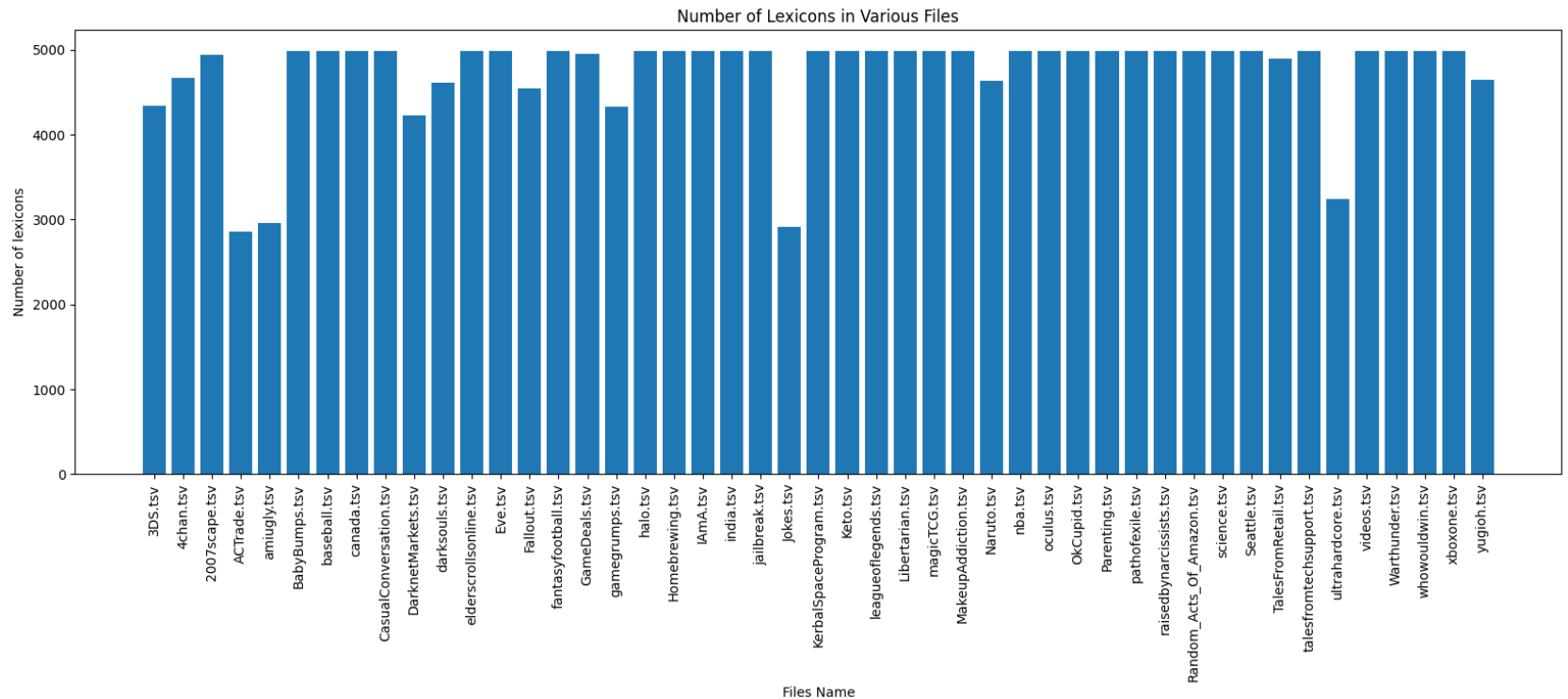


```
In [15]:    1  # Adding all the values to the lexicon_values variable which is a list
            2  lexicon_values = []
            3  xz = {}
            4  for i in all_lexicon_files:
            5      z = read_lexicons_files('subreddits/'+i)
            6      xz[i] = len(z)
            7      lexicon_values.append(z)
```

```
In [16]:    1  names = list(xz.keys())
            2  values = list(xz.values())
```

```
In [17]:    1  plt.figure(figsize=(20, 6))
            2  plt.bar(names, values)
            3  plt.xticks(rotation=90)
            4  # for i, v in enumerate(values):
            5  #     plt.annotate(str(v), xy=(i, v), ha='center', va='bottom',rotation=90)
            6  plt.xlabel('Files Name')
            7  plt.ylabel('Number of lexicons')
            8  plt.title('Number of Lexicons in Various Files')
            9
           10  # Display the chart
           11  plt.show()
```

```python
In [18]:    1  # adding the lexicons of adjectives and lexicons from other files
            2  combined = [adjectives, frequency] + lexicon_values
```

```python
In [19]:
def get_feature(tweets, combined):
    # divide into list of words
    wordings = tweets.split()
    # Count words in the tweet
    total = len(wordings)
    # Finding the longest word
    longest = max(wordings, key=len)

    # set 12 features to the list
    feature_set = [0] * 12

    for i, lex_dict in enumerate(combined[:9]):
        score = 0
        for word in wordings:
            sentiment_dict = lex_dict.get(word, {'-ve': 0, '+ve': 0})
            score += sentiment_dict['-ve'] + sentiment_dict['+ve']
        feature_set[i] = score


    # log of the word count for the tweet
    if total > 0:
        feature_set[9] = math.log(total)
    else:
        feature_set[9] = 0


    # log of length of longest word
    if longest:
        feature_set[10] = math.log(len(longest))
    else:
        feature_set[10] = 0


    # Count of words that have 5 characters or more
    long_word_count = 0
    for word in wordings:
        if len(word) >= 5:
            long_word_count += 1

    # log of count of long words
    if long_word_count > 0:
        feature_set[11] = math.log(long_word_count)
    else:
```

```
44          feature_set[11] = 0
45
46
47      return feature_set
```

In [20]:
```
1  train_features = [get_feature(tweet, combined) for tweet in train_tweets]
2  validation_features = [get_feature(tweet, combined) for tweet in val_tweets]
3  testing_feature = [get_feature(tweet, combined) for tweet in test_tweets]
```

In [21]:
```
1  # making as a input using PyTorch for training set
2  dtype = torch.float32
3  label_dtype = torch.float32
4  X_train = tt(train_features, dtype=dtype)
5  y_train = tt(list(map(int, train_labels)), dtype=label_dtype).unsqueeze(1)
6
7  # making as a input using PyTorch for validation set
8  X_val = tt(validation_features, dtype=dtype)
9  y_val = tt(list(map(int, val_labels)), dtype=label_dtype).unsqueeze(1)
10
11 # making as a input using PyTorch for test set
12 X_test = tt(testing_features, dtype=dtype)
13 y_test = tt(list(map(int, test_labels)), dtype=label_dtype).unsqueeze(1)
14
```

```python
class LogisticRegressionDef(nn.Module):
    def __init__(self, input_size):
        super(LogisticRegressionDef, self).__init__()
        self.linear = nn.Linear(input_size, 1)
#         using sigmoid function from torch library
        self.sigmoid = nn.Sigmoid()

    def forward(self, x):
        nex = self.linear(x)
        nex = self.sigmoid(nex)
        return nex

    def train(self, X_train, y_train, lr=0.01, epochs=100):
        optimizer = optim.SGD(self.parameters(), lr=lr)
        loss_fn = nn.BCELoss()

        for i in range(epochs):
            y_pred = self(X_train)
            loss = loss_fn(y_pred, y_train)
            loss.backward()
            optimizer.step()
            optimizer.zero_grad()

    def predict(self, X):
        with torch.no_grad():
            y_predict = self(X)
            y_predict = (y_predict >= 0.5).float()
        return y_predict

    def evaluate(self, X, y):
        y_predict = self.predict(X)
        find_accuracy = accuracy_score(y, y_predict)
        find_f1 = f1_score(y, y_predict, average='weighted')
        return find_accuracy, find_f1
```

```python
# creating the object for the class
lr1 = LogisticRegressionDef(12)
# calling the train method
lr1.train(X_train, y_train,lr=0.01, epochs=200)
# getting the accuracy and f1-score for the test sets
accuracy, f1_score = lr1.evaluate(X_test, y_test)
print("Accuracy : ",accuracy)
print("F1 Score : ",f1_score)
```

```
Accuracy :   0.48363725170954086
F1 Score :   0.31582110771558486
```