

## Project-2

### Introduction

This project uses a dataset of labeled tweets to train a logistic regression model in PyTorch for sentiment analysis. The model takes 12 different features, including sentiment scores and word count, and outputs a probability of the tweet being positive. The model is trained on a training set and evaluated using accuracy and F1-score on a test set.

### Discussion

```
1 import numpy as np
2 import pandas as pd
3 import warnings
4 import math
5 import torch
6 import re
7 import torch.nn as nn
8 import torch.optim as optim
9 import torch.tensor as tt
10 from sklearn.feature_extraction.text import CountVectorizer
11 from sklearn.linear_model import LogisticRegression
12 from sklearn.metrics import accuracy_score
13 from sklearn.metrics import f1_score
14 import matplotlib.pyplot as plt
15 warnings.filterwarnings("ignore")
```

Importing the required libraries to run the code and analyze it.

```
1 enc_type = 'utf-8'
```

Setting the encoding type to the variable. Defining this as this is used by all the tsv files while loading them on to the system.

## Project-2

```
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 to lower case
31        text1 = text1.lower()
32    return data
```

This function or method that is defined is used to clean the data using the regular expressions. This function clears all the data like extra spaces, emojis, multiple spaces, square brackets, special symbols and punctuation marks. This also converts the data to lowercase.

```
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
```

This method is used to read the data from the file using the r mode and encoding format as utf-8.

```
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')
```

## Project-2

Reading the text files and label files.

```
train_tweets
np The Way You Make Me Feel - 2012 Remaster by Michael Jackson from the album Bad 25th Anniversary.',
I'm going to see Paper Towns. Saturday, 22 August 2015 at 17:40 in Leigh #Cineworld"',
user Yepo I came in Milan it was my 1st concert and it was so good you guys did an amazing job! You are f
FECT! ',
davs #dool Tuesday Hone has to pick up Ciara. Rafe wants a real case not just publicity like Justin'.
```

Train tweets file sample data.

```
|: 1 train_labels
|: ['2',
   '1',
   '1',
   '1',
   '2',
   '2',
   '2',
   '2',
   '0',
   '2',
   '1',
```

Train labels sample data

```
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
```

Here in this we have defined the function to read the lexicon files and converting it to a dictionary and returning it to the calling function.

## Project-2

```

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', 'halo.tsv',
5                     'Homebrewing.tsv', 'IAMA.tsv', 'india.tsv', 'jailbreak.tsv', 'Jokes.tsv', 'KerbalSpaceProgram.tsv',
6                     'Keto.tsv', 'leagueoflegends.tsv', 'Libertarian.tsv', 'magicTCG.tsv', 'MakeupAddiction.tsv',
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']

```

These are all the lexicon files that are used for the project.

```

1 negative_scores = [adjectives[word]['-ve'] for word in adjectives]
2 positive_scores = [adjectives[word]['+ve'] for word in adjectives]

```

Assigning all the positive and negative values to the lists from the adjectives lexicon file which is named as 2000.tsv.

```

1 x_axis = list(adjectives.keys())

```

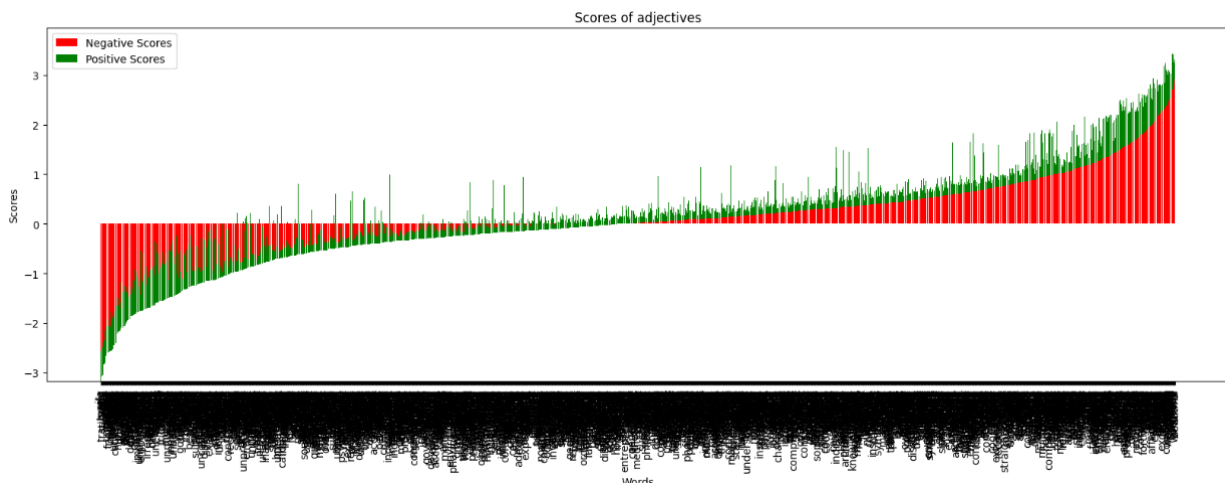
Getting the names of all the words in the 2000.tsv file.

```

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()

```

Plotting a bar graph to show the positive and -ve scores for the given word.



The wording at the bottom is not clear due to the screen size that is present.

## Project-2

```
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)
```

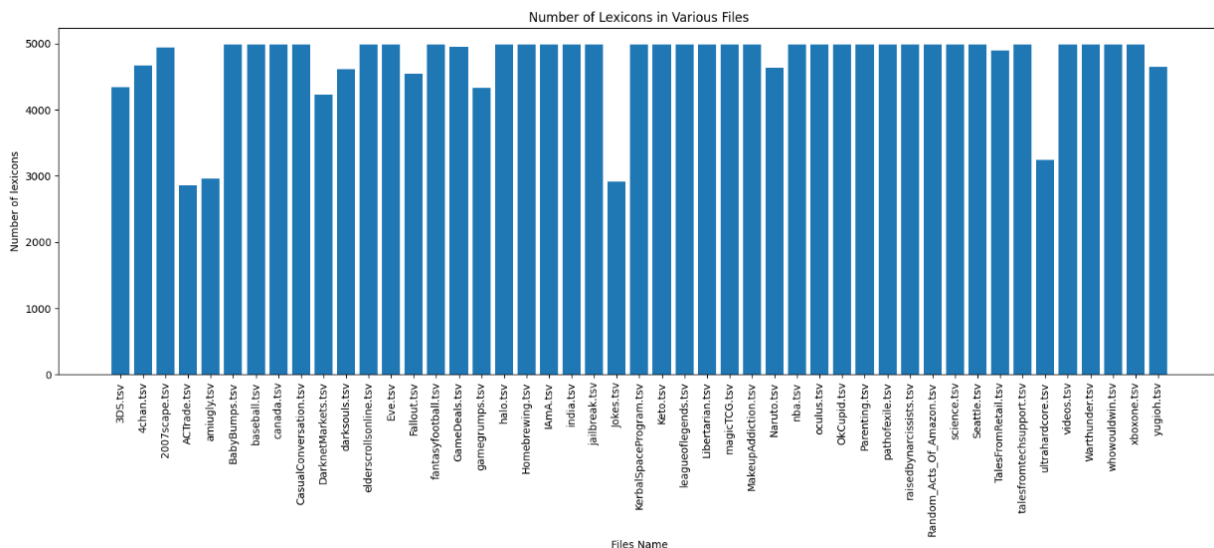
Adding all the lexicon values to the list.

```
: 1 names = list(xz.keys())
   2 values = list(xz.values())
```

This stores the list of file names and contents length of all the files.

```
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()
```

This prints the below graph



## Project-2

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

This is all the lexicons combined from all the files.

```
1 def get_feature(tweets, combined):
2     # divide into list of words
3     wordings = tweets.split()
4     # Count words in the tweet
5     total = len(wordings)
6     # Finding the longest word
7     longest = max(wordings, key=len)
8
9     # set 12 features to the list
10    feature_set = [0] * 12
11
12    for i, lex_dict in enumerate(combined[:9]):
13        score = 0
14        for word in wordings:
15            sentiment_dict = lex_dict.get(word, {'-ve': 0, '+ve': 0})
16            score += sentiment_dict['-ve'] + sentiment_dict['+ve']
17        feature_set[i] = score
18
19
20    # log of the word count for the tweet
21    if total > 0:
22        feature_set[9] = math.log(total)
23    else:
24        feature_set[9] = 0
25
```

## Project-2

```
# 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:
    feature_set[11] = 0

return feature_set
```

The above code will store the features on feature\_set variable. The first nine are taken from the lexicon files and the other three features manually added by doing some actions on dataset.

```
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_features = [get_feature(tweet, combined) for tweet in test_tweets]
```

Loading the data on the variables.

## Project-2

```
# making as a input using PyTorch for training set
dtype = torch.float32
label_dtype = torch.float32
X_train = tt(train_features, dtype=dtype)
y_train = tt(list(map(int, train_labels)), dtype=label_dtype).unsqueeze(1)

# making as a input using PyTorch for validation set
X_val = tt(validation_features, dtype=dtype)
y_val = tt(list(map(int, val_labels)), dtype=label_dtype).unsqueeze(1)

# making as a input using PyTorch for test set
X_test = tt(testing_features, dtype=dtype)
y_test = tt(list(map(int, test_labels)), dtype=label_dtype).unsqueeze(1)
```

Preparing the framework using pytorch to create the data to provide data as input.



Project-2

```
1 class LogisticRegressionDef(nn.Module):
2     def __init__(self, input_size):
3         super(LogisticRegressionDef, self).__init__()
4         self.linear = nn.Linear(input_size, 1)
5         # using sigmoid function from torch library
6         self.sigmoid = nn.Sigmoid()
7
8     def forward(self, x):
9         nex = self.linear(x)
10        nex = self.sigmoid(nex)
11        return nex
12
13    def train(self, X_train, y_train, lr=0.01, epochs=100):
14        optimizer = optim.SGD(self.parameters(), lr=lr)
15        loss_fn = nn.BCELoss()
16
17        for i in range(epochs):
18            y_pred = self(X_train)
19            loss = loss_fn(y_pred, y_train)
20            loss.backward()
21            optimizer.step()
22            optimizer.zero_grad()
23
24    def predict(self, X):
25        with torch.no_grad():
26            y_predict = self(X)
27            y_predict = (y_predict >= 0.5).float()
28        return y_predict
29
30    def evaluate(self, X, y):
31        y_predict = self.predict(X)
32        find_accuracy = accuracy_score(y, y_predict)
33        find_f1 = f1_score(y, y_predict, average='weighted')
34        return find_accuracy, find_f1
35
```

The above class is defined to provide a base for performing the logistic regression for the tweets.

## Project-2

```
# 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)
```

This provides a reference for the class that is defined above. The object that is generated is calling the logistic regression functions and performing the training and test to find the accuracy and f1\_score.

### Output Analysis:

---

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

The Accuracy of the model is 48.3% where as the f1 score is 31.5%

## References

<https://numpy.org/>

<https://pandas.pydata.org/>

<https://matplotlib.org/>

<https://scikit-learn.org/stable/>

<https://pytorch.org/get-started/locally/>

[https://scikit-learn.org/stable/modules/generated/sklearn.feature\\_extraction.text.CountVectorizer.html](https://scikit-learn.org/stable/modules/generated/sklearn.feature_extraction.text.CountVectorizer.html)

<https://stackoverflow.com/questions/34093264/python-logistic-regression-beginner>

<https://pytorch.org/docs/stable/optim.html>

<https://stackoverflow.com/questions/22540449/how-can-i-rotate-a-matplotlib-plot-through-90-degrees>

<https://stackoverflow.com/questions/53975717/pytorch-connection-between-loss-backward-and-optimizer-step>

<https://pytorch.org/docs/stable/generated/torch.nn.BCELoss.html>

[https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1\\_score.html](https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html)

<https://stackoverflow.com/questions/42471082/how-to-find-out-the-accuracy>

## Project-2