## Introduction

This project uses a dataset of labeled tweets to train a Neural Network 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.

**Implementation:**

```python
1   import math
2   import re
3   import torch
4   import torch.nn as nn
5   import torch.optim as optim
6   from torch.utils.data import Dataset as DT
7   from torch.utils.data import DataLoader as DL
8   from sklearn.metrics import confusion_matrix, f1_score
9   from sklearn.preprocessing import LabelEncoder
10  from sklearn.model_selection import KFold
11  import warnings
12  import numpy as np
13  import nltk
14  from nltk.tokenize import word_tokenize
15  import string
16  from nltk.corpus import stopwords
17  from nltk.stem import PorterStemmer
18  from gensim.models import Word2Vec
19  import matplotlib.pyplot as plt
20
21  warnings.filterwarnings('ignore')
```

Importing the required libraries for running the code and analyzing it.

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

This function definition is used to read the lexicon files using the encoding utf-8, this code reads the data ,cleans the data and assigns positive and negative values to the key.

After parsing all the values the function returns the dictionary of lexicons.

```python
1  def load_data(tfile, lfile):
2      with open(tfile, 'r',encoding = 'utf-8') as n:
3          tweetsandtext = n.readlines()
4
5      with open(lfile, 'r',encoding = 'utf-8') as n:
6          labelsforencoding = n.readlines()
7
8      return tweetsandtext, labelsforencoding
```

In the above code we are taking 2 input files, one is the text file and the other is the label file which is related to the text file.

```python
1  def data_clean_tweets(t):
2
3      t = re.sub(r'http\S+|www\S+|https\S+|\S+@\S+', '', t)
4       # Remove numbers
5      t = re.sub(r'\d+', '', t)
6      t = re.sub('[^a-zA-Z0-9\s]', '', t)
7
8      tz = word_tokenize(t)
9      tz = [x.lower() for x in tz if x not in string.punctuation]
10     sws = set(stopwords.words('english'))
11     tz = [x for x in tz if x not in sws]
12     st = PorterStemmer()
13     tz = [st.stem(x) for x in tz]
14     cleaned_text = ' '.join(tz)
15     return cleaned_text
```

In the above function, the data is cleaned  by using regular expressions, by using stop-words from nltk module. We are also using porter stemmer which applies stemming to each word.

```python
 1  def get_feature(tweet, lexicons):
 2      # Split tweet into words
 3      words = tweet.split()
 4
 5      # Count words in the tweet
 6      total_words = len(words)
 7
 8      # Finding the longest word
 9      longest_word = max(words, key=len)
10
11      # Set 12 features to the list
12      feature_set = [0] * 12
13
14      # Calculate lexicon scores for each word in the tweet
15      for i, lex_dict in enumerate(lexicons[:9]):
16          score = 0
17          for word in words:
18              sentiment_dict = lex_dict.get(word, {'-ve': 0, '+ve': 0})
19              score += sentiment_dict['-ve'] + sentiment_dict['+ve']
20          feature_set[i] = score
21
22      # Log of the word count for the tweet
23      if total_words > 0:
24          feature_set[9] = math.log(total_words)
25      else:
26          feature_set[9] = 0
27
28      # Log of length of longest word
29      if longest_word:
30          feature_set[10] = math.log(len(longest_word))
31      else:
32          feature_set[10] = 0
33
34      # Count of words that have 5 characters or more
35      long_word_count = sum([1 for word in words if len(word) >= 5])
36
37      # Log of count of long words
38      if long_word_count > 0:
39          feature_set[11] = math.log(long_word_count)
40      else:
41          feature_set[11] = 0
42
43      return feature_set
```

The above function creates the features for passing. This is being used from project 2. This generates 12 features and returns the features list.

```
 1  class SentimentDataset(DT):
 2      def __init__(self, feat, l):
 3          # Convert  to PyTorch tensors
 4          if not isinstance(features, tt):
 5              feat = tt(feat)
 6          if not isinstance(labels, tt):
 7              l = tt(l)
 8
 9          self.feat = feat
10          self.l = l
11
12      def __len__(self):
13          # Return length of dataset
14          return len(self.feat)
15
16      def __getitem__(self, index):
17          # Checking if the index is in range
18          assert index < len(self), "Index out of range"
19          feat = self.feat[index]
20          l = self.l[index]
21
22          return feat, l
```

Sentiment data set is a class that is extended from Data set of torch library. Here in this we convert the features to the tensors.

The length returns the length of feat variable which is the features that is converted to tensors. Get Item returns the feat and label.

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

```
 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']
```

Setting the encoding type as utf-8 and assigning all the names of lexicon files to the list to read.

```
1  adj = "adjectives/2000.tsv"
2  freq = "adjectives/2000.tsv"
3  adjectives = read_lexicons_files(adj)
4  frequency = read_lexicons_files(freq)
```

This reads the adjective file.

```
1  lexicon_values = []
2  xz = {}
3  for i in all_lexicon_files:
4      z = read_lexicons_files('subreddits/'+i)
5      xz[i] = len(z)
6      lexicon_values.append(z)
```
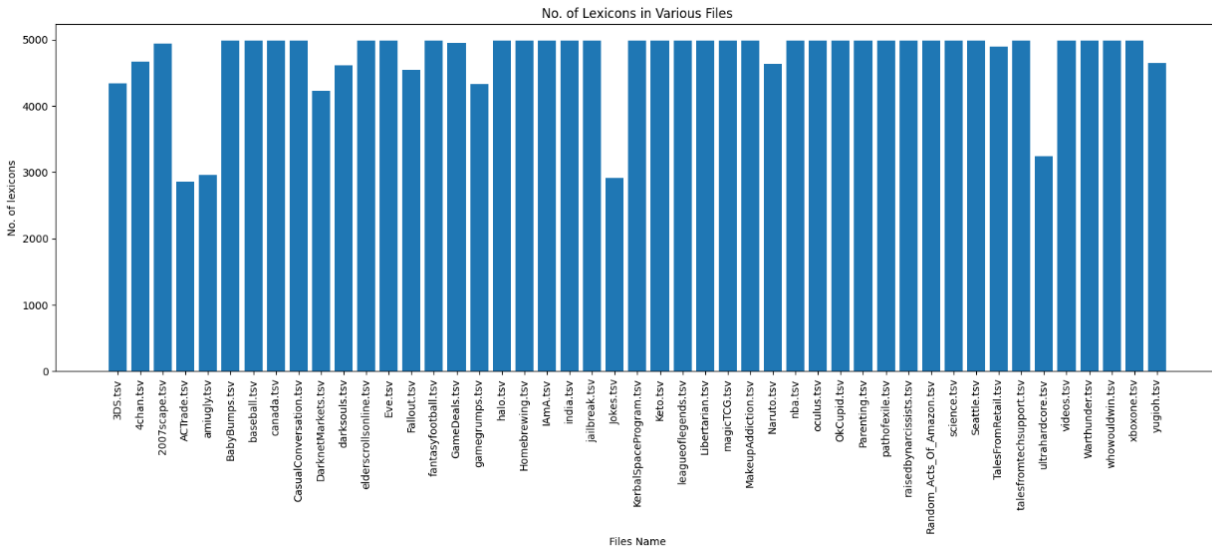
Reading all the lexicon files using a loop and storing in lexicon_values variable.

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

Stores keys and values in the variables name ad values.

These are used in visualization.

```
1  plt.figure(figsize=(20, 6))
2  plt.bar(names, values)
3  plt.xticks(rotation=90)|
4  plt.xlabel('Files Name')
5  plt.ylabel('No. of lexicons')
6  plt.title('No. of Lexicons in Various Files')
7  plt.show()
```

This graph shows how many lexicons are there in which file.

```
1  combined = [adjectives, frequency] + lexicon_values
```

Combining lexicon and adjective and frequency values.

```
1  train_text, train_labels = load_data('sentiment/train_text.txt', 'sentiment/train_labels.txt')
2  val_text, val_labels = load_data('sentiment/val_text.txt', 'sentiment/val_labels.txt')
3  test_text, test_labels = load_data('sentiment/test_text.txt', 'sentiment/test_labels.txt')
```

Reading the values from the files.

```
1  train_text = [data_clean_tweets(tweet) for tweet in train_text]
2  val_text = [data_clean_tweets(tweet) for tweet in val_text]
3  test_text = [data_clean_tweets(tweet) for tweet in test_text]
```

Cleaning the data.

```
1  def get_embedding(tweet, model, size):
2      words = tweet.split()
3      vec = np.zeros(size)
4      count = 0
5      for word in words:
6          try:
7              vec += model.wv[word]
8              count += 1
9          except KeyError:
10             pass
11     if count != 0:
12         vec /= count
13     return vec
14 tt = torch.tensor
```

This is used to get the embedding vector by passing the tweet , model and the size.

```
1  tokenized_tweets = [tweet.split() for tweet in train_text + val_text + test_text]
2
3  #Word2Vec model
4  word2vec_model = Word2Vec(tokenized_tweets, vector_size=300, window=5, min_count=1, workers=4)
```

Creating tokens by combining all the data.

Then passing all the tokenized to the word 2 vector model.

```
1  # Encode labels
2  dts = torch.long
3  encoder = LabelEncoder()
4  train_labels = encoder.fit_transform(train_labels)
5  val_labels = encoder.transform(val_labels)
6  test_labels = encoder.transform(test_labels)
7
8  y_train = tt(train_labels, dtype=dts)
9  y_val = tt(val_labels, dtype=dts)
10 y_test = tt(test_labels, dtype=dts)
11
```

Encoding the labels of the files. The labels are encoded to long variable of torch type.

```
X_train_embed = np.array([get_embedding(tweet, word2vec_model, 300) for tweet in train_text])
X_val_embed = np.array([get_embedding(tweet, word2vec_model, 300) for tweet in val_text])
X_test_embed = np.array([get_embedding(tweet, word2vec_model, 300) for tweet in test_text])

# Combine train, validation, and test data
X_all = np.concatenate((X_train_embed, X_val_embed, X_test_embed), axis=0)
y_all = np.concatenate((y_train, y_val, y_test), axis=0)
```

Getting the train validation, testing set of size 300 . After that adding all the values to the x_all variable as a numpy library.

```
class SentimentDataset(DT):
    def __init__(self, data, targets):
        # Checking if data and targets are of the same length
        assert len(data) == len(targets), "Data and targets must have the same length"

        # Convert data to tensors
        self.data = tt(data)
        self.targets = tt(targets)

    def __len__(self):
        return len(self.data)

    def __getitem__(self, index):
        # finding data and target at the given index
        data_point = self.data[index]
        target = self.targets[index]
        return data_point, target
```

This is extended from the Dataset of the torch module. This checks the input values are of same length or not. It converts the data into torch tensors for running the code.

```
1  class SentimentModel(nn.Module):
2      def __init__(self, indim, hdim, odim):
3          super(SentimentModel, self).__init__()
4          self.leyar1 = nn.Linear(indim, hdim)
5          self.activation = nn.ReLU()
6          self.l2 = nn.Linear(hdim, odim)
7          self.output_layer = nn.Softmax(dim=1)
8
9      def forward(self, z):
10         z = self.leyar1(z)
11         z = self.activation(z)
12         z = self.l2(z)
13         z = self.output_layer(z)
14         return z
```

This is extended from the model nn module of torch library. This model uses the linear function in nn module(layers) , relu activation, softmax. The forward function returns the output of the output layer after performing the activation functions on the data.

```
1  input_size = 300
2  dty = torch.float32
3
4  data_train = SentimentDataset(tt(X_train_embed, dtype=dty), y_train)
5  data_val = SentimentDataset(tt(X_val_embed, dtype=dty), y_val)
6  data_test = SentimentDataset(tt(X_test_embed, dtype=dty), y_test)
7
8  load_train = DL(data_train, batch_size=64, shuffle=True)
9  load_val = DL(data_val, batch_size=64, shuffle=False)
10 load_test = DL(data_test, batch_size=64, shuffle=False)
11
```

This has variables which holds the data, inputsize. These variables holds the data that is passed to the DataLoader of the torch library and holds those values in the provided variables.

```
1  hidden_size = 64
2  num_classes = 3
3  batch_size = 32
4  model = SentimentModel(input_size, hidden_size, num_classes)
```

Setting the hidden size of layers as 64, classes division as positive, negative , neutral, batch size as input to neural network as 32 and model as sentimentModel. This is the class object created for the sentiment Model.

```
1  criterion = nn.CrossEntropyLoss()
2  optimizer = optim.Adam(model.parameters(), lr=0.001)
```

using the crossEntropyLoss as the criterion.

Using the adam optimizer for the model.

```
1  def train(model, X, y, criterion, optimizer, num_epochs):
2      ds = torch.utils.data.TensorDataset(X, y)
3      dataloader = DL(ds, batch_size=batch_size, shuffle=True)
4
5      train_losses = []
6      train_accuracies = []
7
8      for e in range(num_epochs):
9          running_loss = 0.0
10         running_corrects = 0
11         total_samples = 0
12
13         for inp, lab in dataloader:
14             optimizer.zero_grad()
15             out = model(inp)
16             loss = criterion(out, lab)
17             loss.backward()
18             optimizer.step()
19
20             _, preds = torch.max(out, 1)
21             running_loss += loss.item() * inp.size(0)
22             running_corrects += torch.sum(preds == lab.data)
23             total_samples += inp.size(0)
24
25         epoch_loss = running_loss / total_samples
26         epoch_accuracy = running_corrects.double() / total_samples
27         if (e + 1) % 25 == 0:
28             print(f"Epoch {e+1}/{num_epochs} Loss: {epoch_loss:.6f} Accuracy: {epoch_accuracy:.6f}")
29         train_losses.append(epoch_loss)
30         train_accuracies.append(epoch_accuracy.item())
31
32     return train_losses, train_accuracies
33
```

This is the train model that calculates the losses and stores the values for every fold and prints the loss for every epoch. Finally it returns the appened values of losses and accuracies

```
1  def validate(model, dataloader, criterion):
2      model.eval()
3      runLoss = 0.0
4
5      with torch.no_grad():
6          for ipnputs, lebles in dataloader:
7              oput = model(ipnputs)
8              lose = criterion(oput, lebles)
9              runLoss += lose.item()
10
11     return runLoss / len(dataloader)
```

This function sets the model in evaluation This turns off training-specific layers like dropout and batch normalization which should not be active during evaluation.

Runloss is used to keep track of the overall loss. This uses the criterion. This returns the average loss of the batch that is sent.

```python
1  def predict(model, X):
2      model.eval()
3      ds = torch.utils.data.TensorDataset(X)
4      dloader = DL(ds, batch_size=batch_size, shuffle=False)
5
6      pre = []
7      with torch.no_grad():
8          for inputs, in dloader:
9              outputs = model(inputs)
10             _, preds = torch.max(outputs, 1)
11             pre.extend(preds.tolist())
12     return pre
```

This method predict takes the model and tensor. After all the work this function sends the predicted values from the pytorch model.

```python
1  num_epochs = 100
```

Setting number of epochs for each fold as 100.

```
1  from sklearn.model_selection import KFold
2  from sklearn.metrics import accuracy_score
3  from sklearn.metrics import f1_score
4  from sklearn.metrics import confusion_matrix
5
6
7  foldsCount = 5
8  kf = KFold(n_splits=foldsCount, shuffle=True, random_state=42)
9  fold_accuracies = []
10 fold_f1s = []
11
12 all_train_losses = []
13 all_train_accuracies = []
14
15 for fold, (ti, tei) in enumerate(kf.split(X_all, y_all)):
16     X_train_fold = tt(X_all[ti], dtype=dty)
17     y_train_fold = tt(y_all[ti], dtype=torch.long)
18     X_test_fold = tt(X_all[tei], dtype=dty)
19     y_test_fold = tt(y_all[tei], dtype=torch.long)
20
21     # cross-validation folds data
22     model_w2v = SentimentModel(input_size, hidden_size, num_classes)
23     opti = optim.Adam(model_w2v.parameters(), lr=0.001)
24     train_losses, train_accuracies = train(model_w2v, X_train_fold, y_train_fold, criterion, opti, num_epochs)
25     all_train_losses.append(train_losses)
26     all_train_accuracies.append(train_accuracies)
27
28     # testing using testdata
29     y_pred_fold = predict(model_w2v, X_test_fold)
30     fold_confusion_matrix = confusion_matrix(y_test_fold, y_pred_fold)
31     print()
32     print(f"Confusion Matrix for Fold {fold+1}:\n{fold_confusion_matrix}\n")
33     # Finding accuracy
34     fold_accuracy = accuracy_score(y_test_fold, y_pred_fold)
35     fold_f1 = f1_score(y_test_fold, y_pred_fold, average='weighted')
36     fold_accuracies.append(fold_accuracy)
37     fold_f1s.append(fold_f1)
38     print(f"Fold {fold+1} Accuracy: {fold_accuracy:.6f}, F1 Score: {fold_f1:.6f}")
39     print()
40     print()
41
42 average_accuracy = sum(fold_accuracies) / foldsCount
43 print()
44 print()
45 print(f"Average Accuracy: ",average_accuracy*100)
46 average_f1 = sum(fold_f1s) / foldsCount
47 print("Average F1-score:",average_f1*100)
```

This is the main cell block which calls the train model and calculates the losses and accuracies of the model for each epoch. In this we hold all the losses and accuracies for graph visualization. This prints the visualizations and confusion matrix , accuracy , f1 score.

```
Epoch 25/100 Loss: 0.948883 Accuracy: 0.578622
Epoch 50/100 Loss: 0.941048 Accuracy: 0.589891
Epoch 75/100 Loss: 0.935682 Accuracy: 0.596924
Epoch 100/100 Loss: 0.931273 Accuracy: 0.602663

Confusion Matrix for Fold 1:
[[ 623 1218  428]
 [ 451 3749 1280]
 [ 115 1526 2590]]

Fold 1 Accuracy: 0.581135, F1 Score: 0.569202


Epoch 25/100 Loss: 0.951017 Accuracy: 0.577474
Epoch 50/100 Loss: 0.943405 Accuracy: 0.585864
Epoch 75/100 Loss: 0.936826 Accuracy: 0.594691
Epoch 100/100 Loss: 0.933241 Accuracy: 0.600367

Confusion Matrix for Fold 2:
[[ 610 1280  374]
 [ 434 3968 1089]
 [ 102 1711 2412]]

Fold 2 Accuracy: 0.583472, F1 Score: 0.569812
```
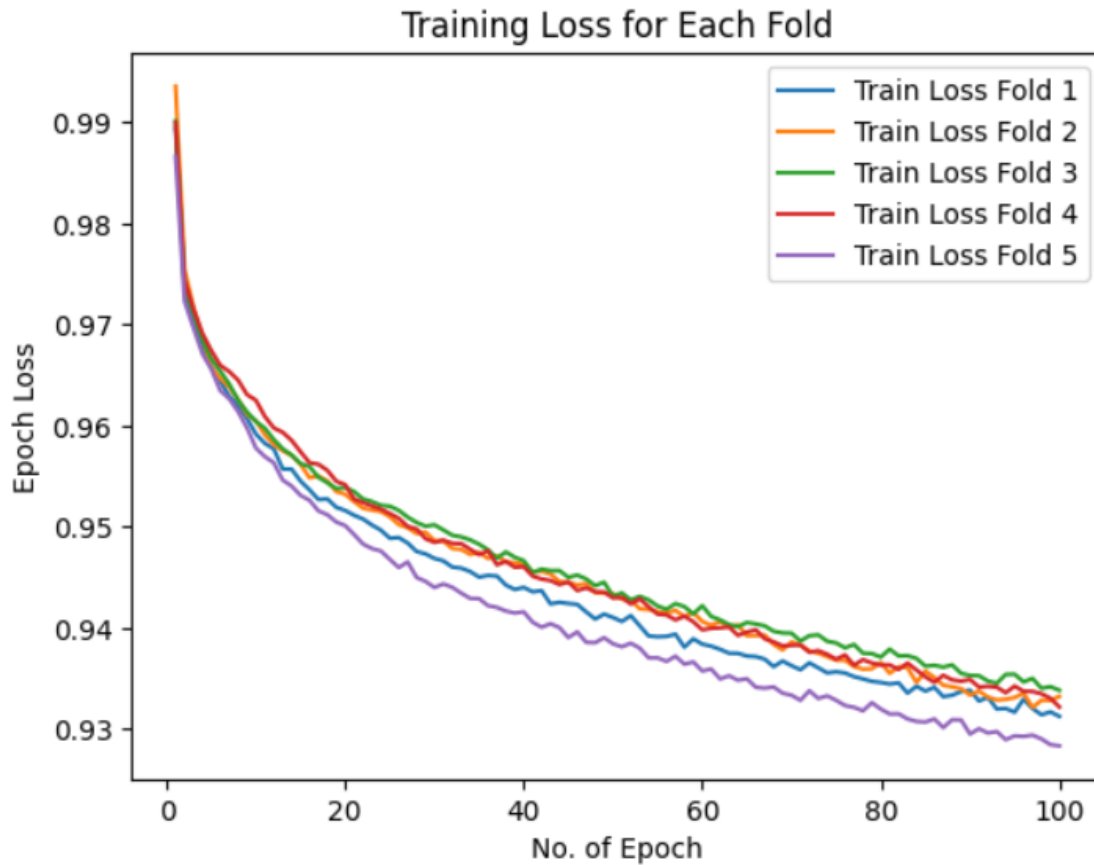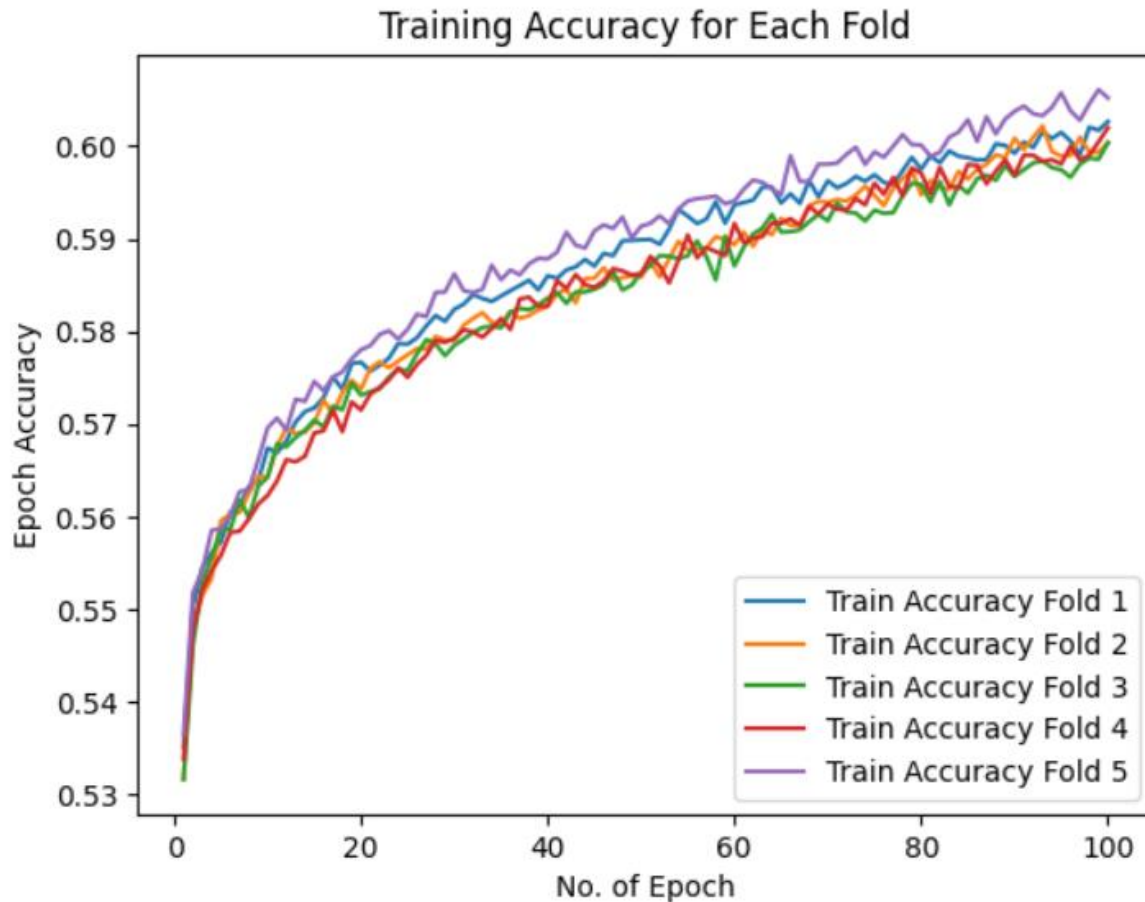
The above is the sample of the output.

```
1  for fold, train_losses in enumerate(all_train_losses):
2      plt.plot(range(1, num_epochs+1), train_losses, label=f'Train Loss Fold {fold+1}')
3
4  plt.title('Training Loss for Each Fold')
5  plt.xlabel('No. of Epoch')
6  plt.ylabel('Epoch Loss')
7  plt.legend()
8  plt.show()
```

The above code shows graph for the epoch and epoch loss for every fold.

Training Loss for Each Fold

```
1  for fold, train_accuracies in enumerate(all_train_accuracies):
2      plt.plot(range(1, num_epochs+1), train_accuracies, label=f'Train Accuracy Fold {fold+1}')
3
4  plt.title('Training Accuracy for Each Fold')
5  plt.xlabel('No. of Epoch')
6  plt.ylabel('Epoch Accuracy')
7  plt.legend()
8  plt.show()
```

The above code shows the epoch number and accuracy of the  model.

Training Accuracy for Each Fold

**References:**

https://pypi.org/project/nets/0.0.3.1/?
https://discuss.pytorch.org/t/batch-normalization-and-weight-initialization-in-seq2seq/54934
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://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://www.deeplearningwizard.com/deep_learning/practical_pytorch/pytorch_feedforward_neuralnetwork/

Project 3
Harshavardhan Reddy Mallannagari                                                          11603389

https://www.deeplearningwizard.com/deep_learning/boosting_models_pytorch/weight_initialization_activation_functions/