Pipeline

1) Removing HTML tagas and URLs, Punctiation*, Replacing emoticons *.
2) Tokenization
3) Removing Stop Words
4) Splitting data: Training, Validation, Test
5) TF-IDF Calculation

**In previous code I created another pipeline, it happens that if I apply TF-IDF I can avoid a many other preprocessing steps.**

# Note about the code:

Note that the dataset takes two paths, one towards the sentiment analysis
and the other towards the analysis of some statistics related to the words/text

In [1]:
```python
# Multilayer perceptron working environment.
# Getting ready the work environment. Importing libraries and modules:
import time
import pandas as pd
import re
import nltk
import torch
import torch.nn as nn
import numpy as np
import string
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.feature_extraction.text import CountVectorizer, TfidfVecto
from sklearn.model_selection import train_test_split
from sklearn.metrics import precision_recall_fscore_support, classifica
from collections import Counter
from bs4 import BeautifulSoup
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize

#===========         Extra tools for the statistic analysis
from nltk.stem import WordNetLemmatizer

lemmatizer = WordNetLemmatizer()
#------------------------------------------------------------------

stop_words = stopwords.words('english')
tfidf_vectorizer = TfidfVectorizer()
vectorizer = CountVectorizer()
```

KeyboardInterrupt

In [ ]:
```python
# Support Vector Machine working environment.
import re
import nltk
import torch
import torch.nn as nn
import numpy as np
import string

from sklearn.feature_extraction.text import CountVectorizer, TfidfVecto
from sklearn.svm import SVC
from sklearn.metrics import accuracy_score, classification_report
from sklearn.model_selection import train_test_split, GridSearchCV
from collections import Counter
from bs4 import BeautifulSoup


tfidf_vectorizer = TfidfVectorizer(stop_words='english')
vectorizer = CountVectorizer()
```

# 1) Importing data set

```
In [ ]:  1  #Importing dataset
         2  imdb_path = 'IMDB.csv'
         3  imdb = pd.read_csv(imdb_path)
         4
         5  #Convert sentiment column to binary class
         6  imdb['sentiment'] = imdb['sentiment'].map({'positive': 1, 'negative': 0
         7
         8  #Checking data and columns
         9  print(imdb.head())
```

## 2) Reducing size

```
In [ ]:  1  #During the applycation of the model in earlier stages, the machine pro
         2  #Hence data will be cut off to 5000 rows:
         3
         4  #Firstly, let's segregate the sentiment column:
         5  positive_reviews = imdb[imdb['sentiment'] == 1]
         6  negative_reviews = imdb[imdb['sentiment'] == 0]
         7
         8  #Secondly, sampling randomly 2500 reviews from each (+/-)
         9  positive_sample = positive_reviews.sample(n=2500, random_state=42)
        10  negative_sample = negative_reviews.sample(n=2500, random_state=42)
        11
        12  #Putting them together again
        13  imdb_reduced = pd.concat([positive_sample, negative_sample])
        14
        15  #Suffling the new dataset
        16  imdb_reduced = imdb_reduced.sample(frac=1, random_state=42).reset_index
        17
        18
        19  #Sources:
        20  # https://stackoverflow.com/questions/71758460/effect-of-pandas-datafra
        21  # https://stackoverflow.com/questions/57300260/how-to-drop-added-column
        22  # https://docs.python.org/3/library/fractions.html
        23  # https://datascience.stanford.edu/news/splitting-data-randomly-can-rui
        24  # https://stats.stackexchange.com/questions/484000/how-to-appropriately
```

## 3) Preprocessing

### 3.1) Removing HTML tags and URLs, lower

Note: I used a function to get rid of the punctuations however the
dataset became massive and my machine was unable to manage. That's
why I am avoiding it.

```python
#Function to remove HTML tags:
def remove_html(text):
    soup = BeautifulSoup(text, "html.parser")
    return soup.get_text()

#Sources:
#https://stackoverflow.com/questions/328356/extracting-text-from-html-f
#https://beautiful-soup-4.readthedocs.io/en/latest/
#https://www.datacamp.com/tutorial/web-scraping-using-python
#https://www.geeksforgeeks.org/how-to-write-the-output-to-html-file-wit
```

```python
#Function to remove URLs characteres:
def remove_urls(text):
    return re.sub(r'https?://\S+|www\.\S+', '', text)

#Source: https://www.geeksforgeeks.org/remove-urls-from-string-in-pytho
```

```python
#Function to put together all the previous functions:
def preprocess_1(text):
    text = remove_html(text)
    text = remove_urls(text)
    text = text.lower()
    #text = remove_punctuation(text)///\\\\Initially used a function to
    return text
```

```python
#Running the function to make the first preprocessing step.
imdb_reduced['review'] = imdb_reduced['review'].apply(preprocess_1)
imdb['review_preprocess_1'] = imdb['review'].apply(preprocess_1)
```

## 3.2) Tokenization and stopwords elimination.

```python
#Function to tokenize and convert to lower case the text in review colu
def tokenize(text):
    tokens = re.findall(r'\b\w+\b', text)
    return tokens

#Tokenization
imdb_reduced['Token'] = imdb_reduced['review'].apply(tokenize)
imdb['Token'] = imdb['review_preprocess_1'].apply(tokenize)
```

```python
#Function to remove stop words from the tokenized review column
def remove_stopwords(tokens):
    filtered_tokens = [word for word in tokens if word not in stop_word
    return filtered_tokens

#Remove stopwords
imdb_reduced['Token'] = imdb_reduced['Token'].apply(remove_stopwords)
imdb['Token'] = imdb['Token'].apply(remove_stopwords)
```

## ** Statistic Analysis **

# ......................Starts

## A) Checking text

```
In [ ]:    1  print(imdb.head())
```

## B) Avg Words Positive Vs Negative:

```
In [ ]:    1  #Calculating the total tokens for each review
           2  imdb['token_count'] = imdb['Token'].apply(lambda x: len(x) if isinstanc
           3
           4  #Dispersion and central tendency measurements
           5  statistics = imdb.groupby('sentiment')['token_count'].agg(['min', 'max'
           6
           7  #Avg words per review:
           8  avg_words = imdb['Token'].apply(len).mean()
           9
          10  #Print the statistics
          11  print("Statistics by Sentiment: ")
          12  print('\n')
          13  print(statistics)
          14  print('\n')
          15  print('\n')
          16  print('Average Words: ', f"{avg_words:.0f}")
          17
          18  #Resources:
          19  #https://www.geeksforgeeks.org/pandas-groupby-one-column-and-get-mean-m
          20  #https://www.kaggle.com/code/akshaysehgal/ultimate-guide-to-pandas-grou
```

## C) Word Frequency

```
In [ ]:    1  #Iterating through the list of lists(each row) to create a new list wit
           2  def word_freq(list_of_list):
           3      single_list = [item for sublist in list_of_list for item in sublist
           4      token_freq = Counter(single_list)
           5      return token_freq
           6
           7  #Counting the frequency for each word.
           8  word_frequency = word_freq(imdb['Token'])
           9  print(word_frequency)
          10
          11  #Sources: https://www.datacamp.com/tutorial/pandas-apply
```

## D) Unique Words

```
In [ ]:    1  unique_words = len(word_frequency.keys())
           2  print('Unique_words: ',f'{unique_words}')
```

## E) Most common words

```python
1  positive_words = Counter()
2  negative_words = Counter()
3
4  for index, row in imdb.iterrows():
5      words = row['Token']
6      sentiment = row['sentiment']
7      if sentiment == 1:
8          positive_words.update(words)
9      else:
10         negative_words.update(words)
11
12 #Resources:
13 #https://pandas.pydata.org/docs/reference/api/pandas.DataFrame.iterrows
14 #https://www.kaggle.com/code/juicykn/imdb-movie-list-analysis-in-python
```

```python
1  top_positive_words = positive_words.most_common(10)
2  top_negative_words = negative_words.most_common(10)
3
4  print('Positive: ', top_positive_words)
5  print('\n')
6  print('Negative: ', top_negative_words)
```

```python
1  #Splitting the tupple we got earlier
2  positive_words, positive_counts = zip(*top_positive_words)
3  negative_words, negative_counts = zip(*top_negative_words)
4
5  #Charts-----------------------------------------------------
6  fig, axs = plt.subplots(2,1,figsize=(10,8))
7
8  #Positive words plot
9  axs[0].bar(positive_words, positive_counts, color='green')
10 axs[0].set_title('Most Frequent Positive Words')
11 axs[0].set_ylabel('Frequency')
12
13 #Negative words plot
14 axs[1].bar(negative_words, negative_counts, color='red')
15 axs[1].set_title('Most Frequent Negative Words')
16 axs[1].set_ylabel('Frequency')
17
18 #Space between charts
19 plt.tight_layout(pad=4.0)
20 plt.show()
21
22 #Resources:
23 # https://realpython.com/python-zip-function/#using-zip-in-python
24 # https://matplotlib.org/stable/index.html
```

## ** Statistic Analysis **

## ...................................Finishes

# 4) Splitting Data

```
In [ ]:    1  #Splitting data into train 70%, validation 15%, test 15%
           2
           3  X_train_val, X_test, y_train_val, y_test = train_test_split(imdb_reduce
           4
           5  X_train, X_val, y_train, y_val = train_test_split(X_train_val, y_train_
```

# 5) TF-IDF Calculation

```
In [ ]:    1  #Feature extraction: Transforming data into TF-IDF features.
           2  X_train = tfidf_vectorizer.fit_transform(X_train)
           3  X_val = tfidf_vectorizer.transform(X_val)
           4  X_test = tfidf_vectorizer.transform(X_test)
```

```
In [ ]:    1  print(X_train.shape)   # Should output (number_of_samples, 33154)
           2  print(X_test.shape)    # Should also output (number_of_samples, 33154)
           3  print(X_val.shape)
```

# 6) Format conversion

```
In [ ]:    1  #Turning sparse matrix into dense
           2  X_train = X_train.toarray()
           3  X_val = X_val.toarray()
           4  X_test = X_test.toarray()
           5
           6  #Turning into PyTorch tensors
           7  X_train = torch.tensor(X_train, dtype=torch.float32)
           8  X_val = torch.tensor(X_val, dtype=torch.float32)
           9  X_test = torch.tensor(X_test, dtype=torch.float32)
          10  y_train = torch.tensor(y_train.values, dtype=torch.float32)
          11  y_val = torch.tensor(y_val.values, dtype=torch.float32)
          12  y_test = torch.tensor(y_test.values, dtype=torch.float32)
          13
          14  #Resources:
          15  # https://pytorch.org/docs/stable/tensors.html
```

```
In [ ]:    1  #We need to know the shape of the input vector to set the in put dimens
           2  input_dim = X_train.shape[1]
           3  print(input_dim)
```

# 7) MLP model

```python
# Time consumed (starts)
start_time = time.time()

#Building the Multilayer Perceptron model with back propagation.
class MLPmodel(nn.Module):
    def __init__(self):
        super(MLPmodel, self).__init__()
        self.fc1 = nn.Linear(33154,610)
        self.fc2 = nn.Linear(610,377)
        self.fc3 = nn.Linear(377,23)
        self.fc4 = nn.Linear(23,1)
        self.sigmoid = nn.Sigmoid()
        self.relu = nn.ReLU()

    def forward(self, x):
        hidden = self.relu(self.fc1(x))
        hidden = self.relu(self.fc2(hidden))
        hidden = self.relu(self.fc3(hidden))
        output = self.sigmoid(self.fc4(hidden))
        return output
```

```python
model = MLPmodel()
criterion = nn.BCELoss()
optimizer = torch.optim.Adam(model.parameters(), lr=0.00001, weight_dec

#Note: SDG was used earlier in during the experimentation, however, the
```

```python
epochs = 5000
patience = 10 #Here we define how many epochs wait until we stop
best_val_loss = float('inf') #To save the best model/early stop
patience_counter = 0 #This one starts a counter to track number of epoc

for epoch in range(epochs):
    model.train()
    optimizer.zero_grad()
  #Forward Pass
    outputs = model(X_train)
    loss = criterion(outputs.squeeze(), y_train)
  #Backward and Optimize
    loss.backward()
    optimizer.step()

  #Validation
    model.eval()
    with torch.no_grad():
        val_outputs = model(X_val)
        val_loss = criterion(val_outputs.squeeze(), y_val)

  #Early stop
    if val_loss < best_val_loss:
        best_val_loss = val_loss
        best_model_state = model.state_dict()  # Save the best model st
        patience_counter=0
    else:
        patience_counter += 1

  #Print the early stopping
    if patience_counter > patience:
        print(f'Stopping early at epoch {epoch+1}.')
        break

    if (epoch+1) % 5 == 0:
        print(f'Epoch {epoch+1}/{epochs}, Loss: {loss.item()}, Val Loss

if best_model_state:
    model.load_state_dict(best_model_state)

#Resources:
#https://pythonguides.com/pytorch-early-stopping/
#https://discuss.pytorch.org/t/can-i-deepcopy-a-model/52192
#https://machinelearningmastery.com/how-to-stop-training-deep-neural-ne
#https://stackoverflow.com/questions/71998978/early-stopping-in-pytorch
#https://github.com/Bjarten/early-stopping-pytorch
#https://debuggercafe.com/using-learning-rate-scheduler-and-early-stopp
```

```python
#Accuracy on validation set
model.eval()
with torch.no_grad():
    val_outputs = model(X_val)
    val_predicted_bin = (val_outputs.squeeze() > 0.5).int()

    #validation accuracy
    val_accuracy = accuracy_score(y_val.numpy(), val_predicted_bin.nump
    print(f'Validation set accuracy: {val_accuracy:.2f}')

# Resources (for the "bin"):
# https://scikit-learn.org/stable/modules/generated/sklearn.metrics.cla
# https://pytorch.org/docs/stable/generated/torch.Tensor.numpy.html
# https://www.youtube.com/watch?v=E35CVhVKISA&t=135s
```

```python
#Accuracy on test set
model.eval()

with torch.no_grad():
    test_outputs = model(X_test)
    predicted_bin = (test_outputs.squeeze() > 0.5).int()

#accuracy
test_accuracy = accuracy_score(y_test.numpy(), predicted_bin.numpy())
print(f'Accuracy on test set: {test_accuracy:.2f}')

#classification report
print(classification_report(y_test.numpy(), predicted_bin.numpy()))
```

```
In [ ]:    1  #Turning tensors into NumPy arrays so we can use Scikit-learn functions
           2  test_outputs_np = test_outputs.squeeze().numpy()
           3  y_test_np = y_test.numpy()
           4
           5  #Calculation of the ROC-AU
           6  fpr, tpr, thresholds = roc_curve(y_test_np, test_outputs_np)
           7  roc_auc = roc_auc_score(y_test_np, test_outputs_np)
           8
           9  #Plotting the ROC curve
          10  plt.figure(figsize=(8, 6))
          11
          12  plt.plot(fpr, tpr, color='firebrick', lw=2, label='ROC curve (area = %0
          13  plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
          14  plt.xlim([0.0, 1.0])
          15  plt.ylim([0.0, 1.0])
          16
          17  plt.xlabel('False Positive Rate')
          18  plt.ylabel('True Positive Rate')
          19  plt.title('ROC-AU Chart: Best MLP Model')
          20  plt.legend(loc="lower right")
          21
          22  plt.show()
          23
          24  #Calculation of the confusion matrix
          25  cm = confusion_matrix(y_test.numpy(), predicted_bin.numpy())
          26
          27  #Plotting the confusion matrix
          28  plt.figure(figsize=(8, 6))
          29
          30  sns.heatmap(cm, annot=True, fmt="d", cmap='seismic')
          31  plt.xlabel('Predicted labels')
          32  plt.ylabel('True labels')
          33  plt.title('C.M. Chart: Best MLP Model')
          34
          35  plt.show()
```

```
In [ ]:    1  #Total Time Consumed
           2  end_time = time.time()
           3  execution_time = end_time - start_time
           4  print(f"Total Execution Time: {execution_time} seconds")
```

# 8) SVM model

```python
# Time Consumed (starts)
start_time = time.time()

# Definition of the SVM model and hyperparameter for tuning on the trai
param_grid = {'C':  [1.25892541179416], 'kernel': ['linear'], 'tol' : [

# Hyperparameters tuning, cross-validation using training set.
grid_search = GridSearchCV(SVC(), param_grid, cv=3, scoring='accuracy',
grid_search.fit(X_train, y_train)

# Getting the best estimator
Best_SVMmodel = grid_search.best_estimator_
Best_score = grid_search.best_score_

# Print results
print(f'Best Model: {Best_SVMmodel}')
print('\n')
print(f'Best CV Score: {Best_score:.2f}')
```

```python
# Evaluation of the model on the validation set witht he best parameter
y_pred = Best_SVMmodel.predict(X_val)
val_accuracy = accuracy_score(y_val, y_pred)

# Print results
print(f'Validation set accuracy: {val_accuracy:.2f}')
```

```python
# Evaluation of the final model using the test set
y_pred = Best_SVMmodel.predict(X_test)
test_accuracy = accuracy_score(y_test, y_pred)

#Print results
print(f'Accuracy on test set: {test_accuracy:.2f}')
print(classification_report(y_test, y_pred))
```

```python
# Predict probabilities for the test set
y_pred_proba = Best_SVMmodel.decision_function(X_test)

# Calculation of the ROC-AUC
fpr, tpr, thresholds = roc_curve(y_test, y_pred_proba)
roc_auc = roc_auc_score(y_test, y_pred_proba)

# Plotting the ROC curve
plt.figure(figsize=(8, 6))
plt.plot(fpr, tpr, color='firebrick', lw=2, label='ROC curve (area = %0
plt.plot([0, 1], [0, 1], color='navy', lw=2, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.0])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('ROC-AU Chart: SVM Model')
plt.legend(loc="lower right")
plt.show()

# Calculation of the confusion matrix
cm = confusion_matrix(y_test, y_pred)

# Plotting the confusion matrix
plt.figure(figsize=(8, 6))
sns.heatmap(cm, annot=True, fmt="d", cmap='seismic')
plt.xlabel('Predicted labels')
plt.ylabel('True labels')
plt.title('C.M. Chart: SVM Model')
plt.show()
```

```python
# Total Time Consumed
end_time = time.time()
execution_time = end_time - start_time
print(f"Total Execution Time: {execution_time} seconds")
```