# *COMSATS UNIVERSITY*

**WAH CAMPUS**

## *Submitted By*

| Name: | Shayan Babar | M.Arsal | Ali Raza | Saad Rashid |
|---|---|---|---|---|
| *Registration No:* | Fa19bcs051 | Fa19-BCS-056 | Fa19-BCS-047 | Fa19-BCS-194 |

*Class/Section:*    *BSCS/6d*

## *Submitted To*

*Teacher Name:*    *Ammara Zamir*

*Date of Submission:*    *14/06/2022*

# 1. Cleaning data

```python
def clean(text):
    # Removes all special characters and numericals leaving the alphabets
    text = re.sub("[^A-Za-z]+", " ", text)
    return text
```

Because out data set was small we were able to analyze that the only irregularities it had were specials characters present in it hence using regular expressions we removed than using the above code which is called later in the actual code its self.
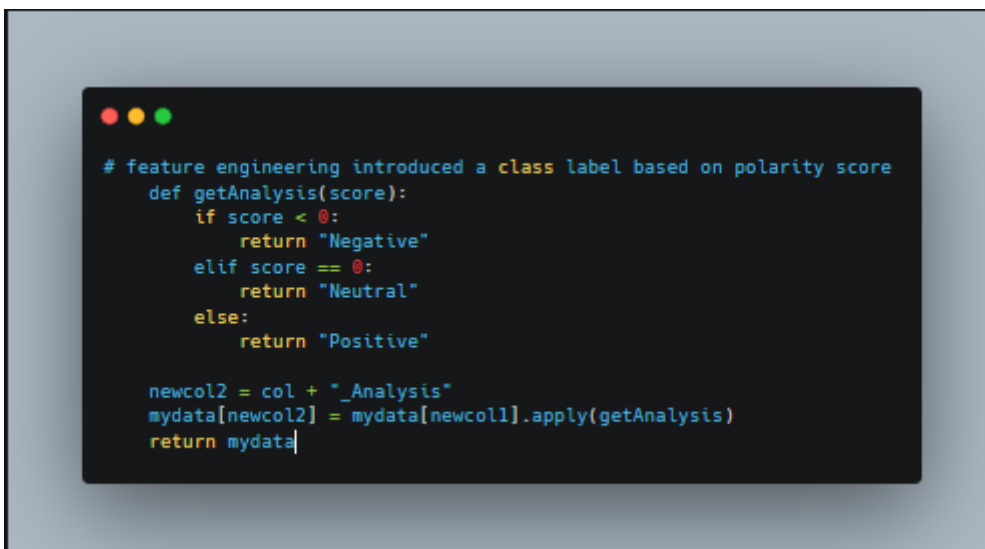
# 2. Text blob sentiment analysis and feature extraction

```python
def sentiment_analysis(mydata, col):
    def getSubjectivity(text):
        return TextBlob(text).sentiment.subjectivity

    # Create a function to get the polarity
    def getPolarity(text):
        return TextBlob(text).sentiment.polarity

    # step3 features extraction
    def features(text):
        score = SentimentIntensityAnalyzer().polarity_scores(text)
        return score["pos"], score["neg"], score["neu"], score["compound"]

    # Create two new columns 'Subjectivity' & 'Polarity'
    newcol = col + "_Subjectivity"
    mydata[newcol] = mydata[col].apply(getSubjectivity)
    newcol1 = col + "_Polarity"
    mydata[newcol1] = mydata[col].apply(getPolarity)
    newcol3 = col + "_len"
    mydata[newcol3] = mydata[col].apply(lambda x: len(x.split()))
    mydata["positive"] = mydata[col].apply(features)
    mydata["positive"], mydata["negative"], mydata["neutral"], mydata["compound"] = zip(
        *mydata[col].map(features)
    )

    # feature engineering introduced a class label based on polarity score
    def getAnalysis(score):
        if score < 0:
            return "Negative"
        elif score == 0:
            return "Neutral"
        else:
            return "Positive"

    newcol2 = col + "_Analysis"
    mydata[newcol2] = mydata[newcol1].apply(getAnalysis)
    return mydata
```

- **getSubjectivity** – finds the subjectivity of the text passed using text blob function extBlob(text).sentiment.subjectivity

- **getPolarity** --finds the Polarity of the text passed using text blob function extBlob(text).sentiment. Polarity

- **features(text)** – ectracts features( positive,negative,neutral,compound) using ntlk function SentimentIntensityAnalyzer().polarity_scores(text)

- **len**—feature is founded using  string.split()
- **data is maped using dataframe[column].apply(function_name)**
- **in feature extraction method .map() function is used to map 4 columns in data frame.**

## 3. Label class
Label class is introduced using following code (cols named is" Cleaned_essay_Analysis")

```
# feature engineering introduced a class label based on polarity score
    def getAnalysis(score):
        if score < 0:
            return "Negative"
        elif score == 0:
            return "Neutral"
        else:
            return "Positive"

    newcol2 = col + "_Analysis"
    mydata[newcol2] = mydata[newcol1].apply(getAnalysis)
    return mydata
```

Polarity is the output that lies between [-1,1], where -1 refers to negative sentiment and +1 refers to positive sentiment
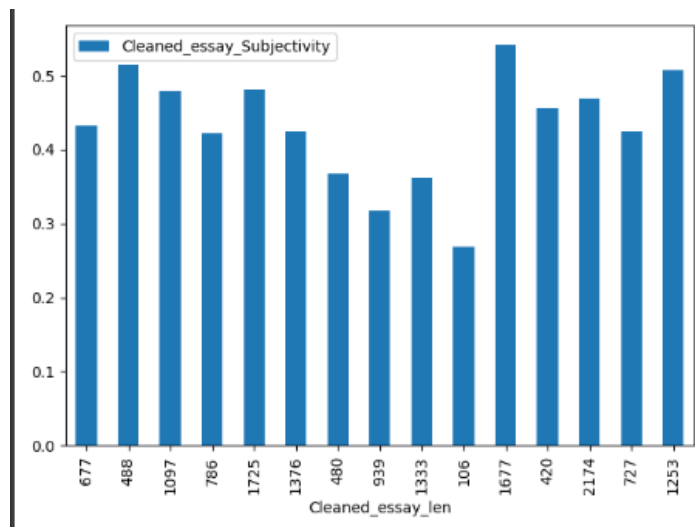
## 4. Mathplotlib



```python
def ploting(mydata):
    mydata.plot(kind="bar", x="Cleaned_essay_len", y="Cleaned_essay_Subjectivity")
    plt.show()
    mydata.plot(kind="bar", x="Cleaned_essay_len", y="Cleaned_essay_Polarity")
    plt.show()
    mydata.plot(
        kind="scatter", x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity"
    )
    plt.show()
    mydata.plot(x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity")
    plt.show()
    mydata.plot(kind="box", x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity")
    plt.show()
```
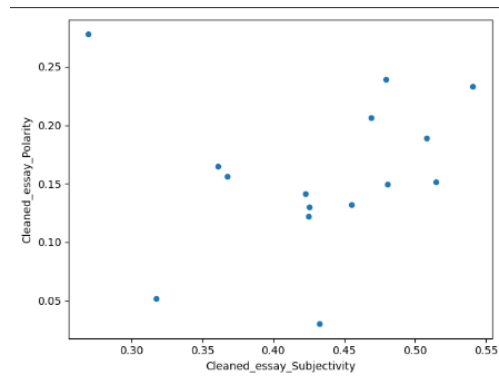




*Simple bar chart of subjectivity and length of the essay*

*Scatter plot subjectivity and polarity*

## 5. <u>Seaborn</u>

```python
def plottingsea(mydata):
    sns.countplot(x="Cleaned_essay_Analysis", data=mydata)
    plt.show()
    sns.scatterplot(
        x="Cleaned_essay_Subjectivity",
        y="Cleaned_essay_Subjectivity",
        data=mydata,
        hue="writers_names",
    )
    plt.show()
    sns.pairplot(mydata)
    plt.show()
    corr = mydata.corr()
    plt.show()
    sns.heatmap(corr)
    plt.show()
```

*Same as mathplotlip except the plot time is defined like this sns.plottype(,,,,,)*

## 6. Main code

```python
import pandas as pd
import re
from textblob import TextBlob
import matplotlib.pyplot as plt
from nltk.sentiment.vader import SentimentIntensityAnalyzer
import seaborn as sns


data = pd.read_csv("blogs.csv")
# print(data.head())

# Step 1: Cleaning the text


def clean(text):
    # Removes all special characters and numericals leaving the alphabets
    text = re.sub("[^A-Za-z]+", " ", text)
    return text


# Step 2: Text blob sentiments analysis


def sentiment_analysis(mydata, col):
    def getSubjectivity(text):
        return TextBlob(text).sentiment.subjectivity

    # Create a function to get the polarity
    def getPolarity(text):
        return TextBlob(text).sentiment.polarity

    # step3 features extraction
    def features(text):
        score = SentimentIntensityAnalyzer().polarity_scores(text)
        return score["pos"], score["neg"], score["neu"], score["compound"]

    # Create two new columns 'Subjectivity' & 'Polarity'
    newcol = col + "_Subjectivity"
    mydata[newcol] = mydata[col].apply(getSubjectivity)
    newcol1 = col + "_Polarity"
    mydata[newcol1] = mydata[col].apply(getPolarity)
    newcol3 = col + "_len"
    mydata[newcol3] = mydata[col].apply(lambda x: len(x.split()))
    mydata["positive"] = mydata[col].apply(features)
    mydata["positive"], mydata["negative"], mydata["neutral"], mydata["compound"] = zip(
        *mydata[col].map(features)
    )

    # feature engineering introduced a class label based on polarity score
    def getAnalysis(score):
        if score < 0:
            return "Negative"
        elif score == 0:
            return "Neutral"
        else:
            return "Positive"

    newcol2 = col + "_Analysis"
    mydata[newcol2] = mydata[newcol1].apply(getAnalysis)
    return mydata


# plotting the dataset using matplot
def ploting(mydata):
    mydata.plot(kind="bar", x="Cleaned_essay_len", y="Cleaned_essay_Subjectivity")
    plt.show()
    mydata.plot(kind="bar", x="Cleaned_essay_len", y="Cleaned_essay_Polarity")
    plt.show()
    mydata.plot(
        kind="scatter", x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity"
    )
    plt.show()
    mydata.plot(x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity")
    plt.show()
    mydata.plot(kind="box", x="Cleaned_essay_Subjectivity", y="Cleaned_essay_Polarity")
    plt.show()


# plotting with seaborn


def plottingsea(mydata):
    sns.countplot(x="Cleaned_essay_Analysis", data=mydata)
    plt.show()
    sns.scatterplot(
        x="Cleaned_essay_Subjectivity",
        y="Cleaned_essay_Subjectivity",
        data=mydata,
        hue="writers_names",
    )
    plt.show()
    sns.pairplot(mydata)
    plt.show()
    corr = mydata.corr()
    plt.show()
    sns.heatmap(corr)
    plt.show()


mydata = pd.DataFrame()
mydata["Cleaned_summary"] = data["article_summary"].apply(clean)
mydata["Cleaned_essay"] = data["article_essay"].apply(clean)
mydata["writers_names"] = data["article_autor_name"].apply(clean)
mydata = sentiment_analysis(mydata, "Cleaned_essay")
print(mydata.head(16))
mydata.to_csv("test.csv", index=True)
ploting(mydata)
plottingsea(mydata)
```

## 7. Csv results



Uncleaned scaped data



Cleaned data



After running the code

| | Cleaned_essay_Subjectivity | Cleaned_essay_Polarity | Cleaned_essay_len | positive | negative | neutral | compound | Cleaned_essay_Analysis |
|---|---|---|---|---|---|---|---|---|
| oll | 0.432616684 | 0.030277342 | 677 | 0.133 | 0.052 | 0.815 | 0.9956 | Positive |
| oll | 0.514882353 | 0.151347594 | 488 | 0.141 | 0.003 | 0.856 | 0.9964 | Positive |
| ooc | 0.479358885 | 0.239164007 | 1097 | 0.207 | 0.034 | 0.759 | 0.9997 | Positive |
| nn | 0.42254902 | 0.141789216 | 786 | 0.139 | 0.071 | 0.79 | 0.9966 | Positive |
| ooc | 0.480444094 | 0.149840393 | 1725 | 0.173 | 0.027 | 0.8 | 0.9998 | Positive |
| ooc | 0.425207595 | 0.129819272 | 1376 | 0.169 | 0.044 | 0.787 | 0.9996 | Positive |
| nn | 0.367592593 | 0.156216931 | 480 | 0.097 | 0.017 | 0.886 | 0.991 | Positive |
| ooc | 0.317231842 | 0.052038046 | 939 | 0.075 | 0.07 | 0.855 | 0.5558 | Positive |
| ooc | 0.360835165 | 0.164800783 | 1333 | 0.151 | 0.092 | 0.757 | 0.9986 | Positive |
| ooc | 0.269444444 | 0.277777778 | 106 | 0.145 | 0.044 | 0.811 | 0.9127 | Positive |
| ooc | 0.540861854 | 0.233568693 | 1677 | 0.209 | 0.047 | 0.743 | 0.9998 | Positive |
| ooc | 0.455294289 | 0.131969697 | 420 | 0.116 | 0.011 | 0.873 | 0.9897 | Positive |
| ot | 0.469147098 | 0.206619673 | 2174 | 0.19 | 0.048 | 0.762 | 0.9999 | Positive |
| oll | 0.425097911 | 0.121840741 | 727 | 0.099 | 0.027 | 0.874 | 0.9941 | Positive |
| ooc | 0.508204996 | 0.18908652 | 1253 | 0.171 | 0.036 | 0.793 | 0.9996 | Positive |

Values

## 8. Classification using random and accuracy with confusion matrix

```python
import pandas as pd
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.ensemble import RandomForestClassifier
import time
from sklearn.metrics import classification_report, confusion_matrix

df = pd.read_csv("test.csv")


X = df.drop(["Cleaned_essay_Analysis"], axis=1)
Y = df["Cleaned_essay_Analysis"]
X = pd.get_dummies(X)
Y = LabelEncoder().fit_transform(Y)
X = StandardScaler().fit_transform(X)


def forest_test(X, Y):
    X_Train, X_Test, Y_Train, Y_Test = train_test_split(
        X, Y, test_size=0.30, random_state=101
    )
    start = time.process_time()
    trainedforest = RandomForestClassifier(n_estimators=700).fit(X_Train, Y_Train)
    print(time.process_time() - start)
    predictionforest = trainedforest.predict(X_Test)
    print(confusion_matrix(Y_Test, predictionforest))
    print(classification_report(Y_Test, predictionforest))


forest_test(X, Y)
```

Before feeding this data into our Machine Learning models I decided to divide our data into features ($X$) and labels ($Y$)

 a function ($forest\_test$) to divide the input data into train and test sets and then train and test a Random Forest Classifier

```
"C:\Users\mians\OneDrive\Desktop\sem7\Topics in cs 1\venv\Scripts\python.exe" "C:/Users/mians/OneDrive/Desktop/sem7/Topics in cs 1/assiment#4/accuracy"
0.921875
[[2]]
              precision    recall  f1-score   support

           0       1.00      1.00      1.00         2

    accuracy                           1.00         2
   macro avg       1.00      1.00      1.00         2
weighted avg       1.00      1.00      1.00         2
```

As shown below, training a Random Forest classifier using all the features, led to 100% Accuracy in about 0.9s of training time