



Newcastle University

MACHINE LEARNING (CSC8111)

TASK – FARS, TWITTER SENTIMENT ANALYSIS

**NAME: KAILASH BALACHANDIRAN
STUDENT ID: 220243160**

FARS

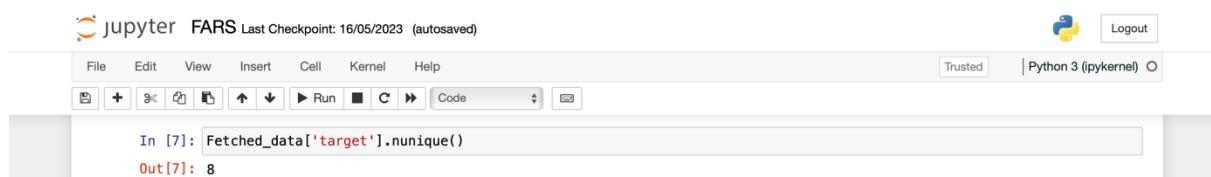
1. Introduction

The objective is to evaluate and forecast the "FARS" and "Twitter sentiment datasets" using machine learning techniques. The research is conducted using a machine learning pipeline structure that starts with data pre-processing and moves on to data cleaning, training, and prediction. A variety of machine learning models are fitted to the dataset, and the performance of each model is examined to select the most successful model with noticeably higher accuracy.

The FARS dataset is a census of fatal traffic accidents that happened inside the boundaries of all 50 states, the District of Columbia, and Puerto Rico and that resulted in at least one person's death (either a driver or a passenger) within 30 days after the accident. Several classifiers are used to predict the classes of a certain record in the FARS dataset to compare the performance of the models. The tweets in the Twitter sentiment dataset can be classified as positive, negative, or neutral depending on the content of the message. FARS dataset.

2. Exploratory Data Analysis

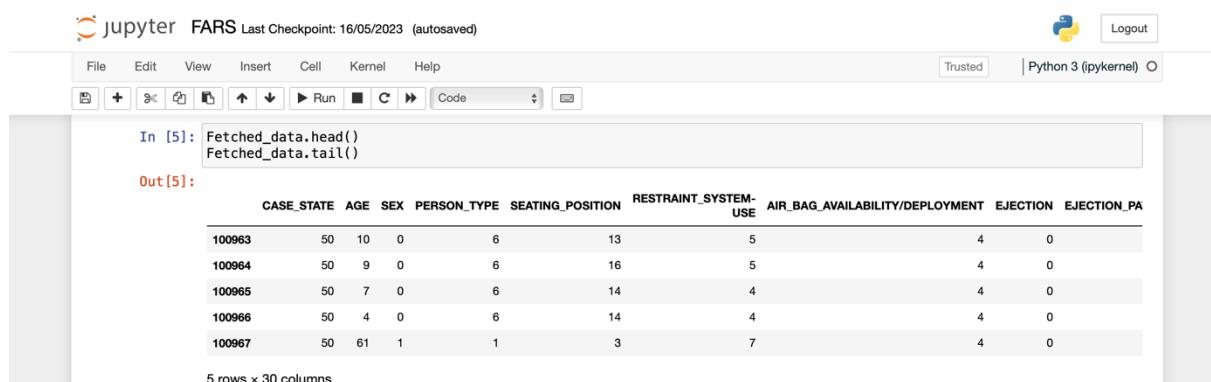
The target variable, which is likewise encoded as factors, is found in the last column of the dataset's 30 rows. There are 8 classes in the target variable.



```
jupyter FARS Last Checkpoint: 16/05/2023 (autosaved)
File Edit View Insert Cell Kernel Help
In [7]: Fetched_data['target'].nunique()
Out[7]: 8
```

Figure 1

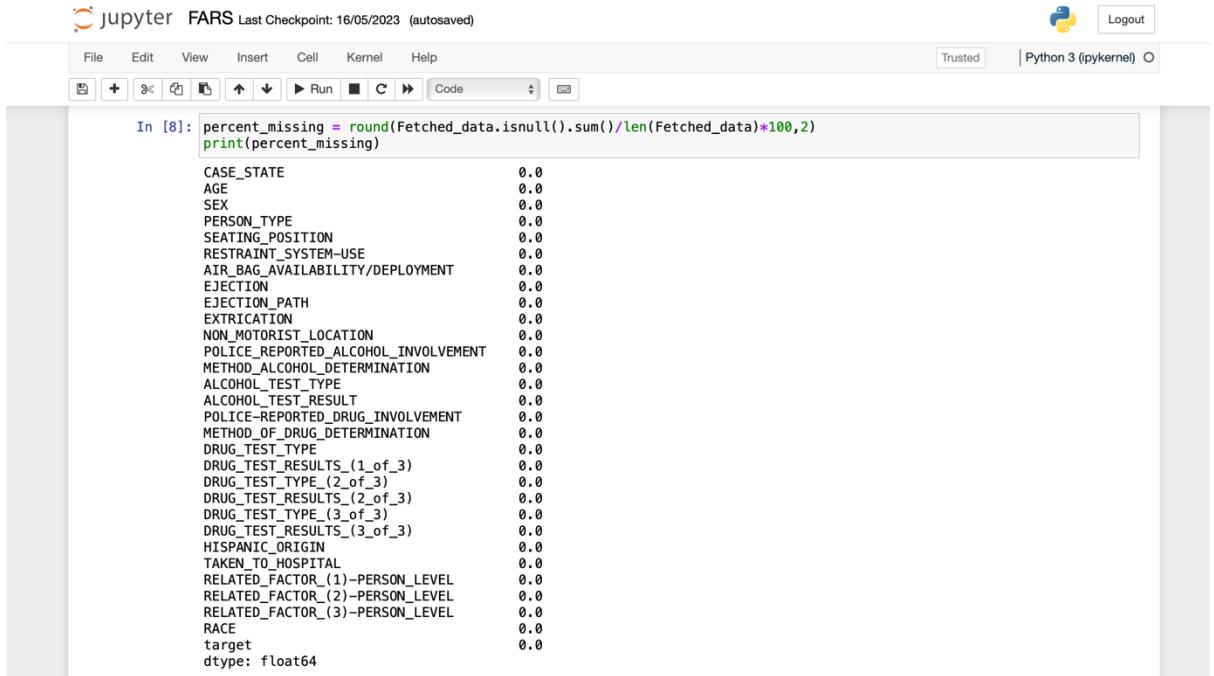
The pre-processed dataset's first five samples are listed below.



```
jupyter FARS Last Checkpoint: 16/05/2023 (autosaved)
File Edit View Insert Cell Kernel Help
In [5]: Fetched_data.head()
Fetched_data.tail()
Out[5]:
CASE_STATE AGE SEX PERSON_TYPE SEATING_POSITION RESTRAINT_SYSTEM-USE AIR_BAG_AVAILABILITY/DEPLOYMENT EJECTION EJECTION_PA
100963 50 10 0 6 13 5 4 0
100964 50 9 0 6 16 5 4 0
100965 50 7 0 6 14 4 4 0
100966 50 4 0 6 14 4 4 0
100967 50 61 1 1 3 7 4 0
5 rows x 30 columns
```

Figure 2

The data's values are as follows and don't contain any null values.



```
In [8]: percent_missing = round(Fetched_data.isnull().sum()/len(Fetched_data)*100,2)
print(percent_missing)

CASE_STATE          0.0
AGE                0.0
SEX                0.0
PERSON_TYPE         0.0
SEATING_POSITION   0.0
REstraint_System-USE 0.0
AIR_BAG_AVAILABILITY/DEPLOYMENT 0.0
EJECTION            0.0
EJECTION_PATH       0.0
EXTRICATION        0.0
NON_MOTORIST_LOCATION 0.0
POLICE_REPORTED_ALCOHOL_ININVOLVEMENT 0.0
METHOD_ALCOHOL_DETERMINATION 0.0
ALCOHOL_TEST_TYPE  0.0
ALCOHOL_TEST_RESULT 0.0
POLICE_REPORTED_DRUG_ININVOLVEMENT 0.0
METHOD_OF_DRUG_DETERMINATION 0.0
DRUG_TEST_TYPE     0.0
DRUG_TEST_RESULTS_(1_of_3) 0.0
DRUG_TEST_TYPE_(2_of_3) 0.0
DRUG_TEST_RESULTS_(2_of_3) 0.0
DRUG_TEST_TYPE_(3_of_3) 0.0
DRUG_TEST_RESULTS_(3_of_3) 0.0
HISPANIC_ORIGIN    0.0
TAKEN_TO_HOSPITAL 0.0
RELATED_FACTOR_(1)-PERSON_LEVEL 0.0
RELATED_FACTOR_(2)-PERSON_LEVEL 0.0
RELATED_FACTOR_(3)-PERSON_LEVEL 0.0
RACE               0.0
target             0.0
dtype: float64
```

Figure 3

The data collection consists of 30 features and 100968 observations. 29 independent variables and 1 dependent variable make up the 30 characteristics.



```
In [6]: Fetched_data.shape

Out[6]: (100968, 30)
```

Figure 4

3. Numerical data analysis

The screenshots below illustrate that the only discrete value is age and that there is little variation between the 75th percentile and the median. Outliers are therefore hard to come by. Most of the traits are hence symmetrically distributed with low skewness. The maximum age is 99, which is peculiar. Most of the other features in the data set are categorical, and there won't be any outliers. Each component's differences act as a reminder of the importance of homogeneity.

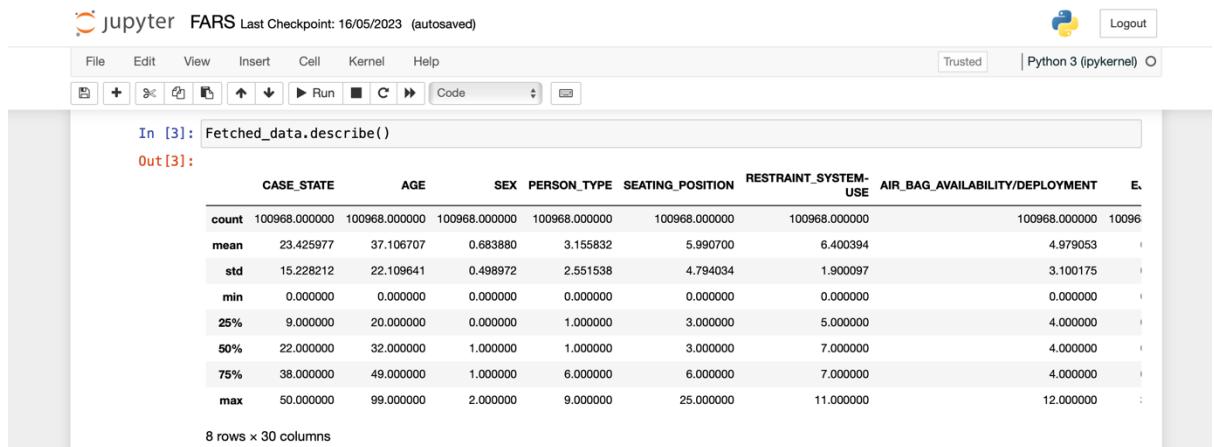


Figure 5

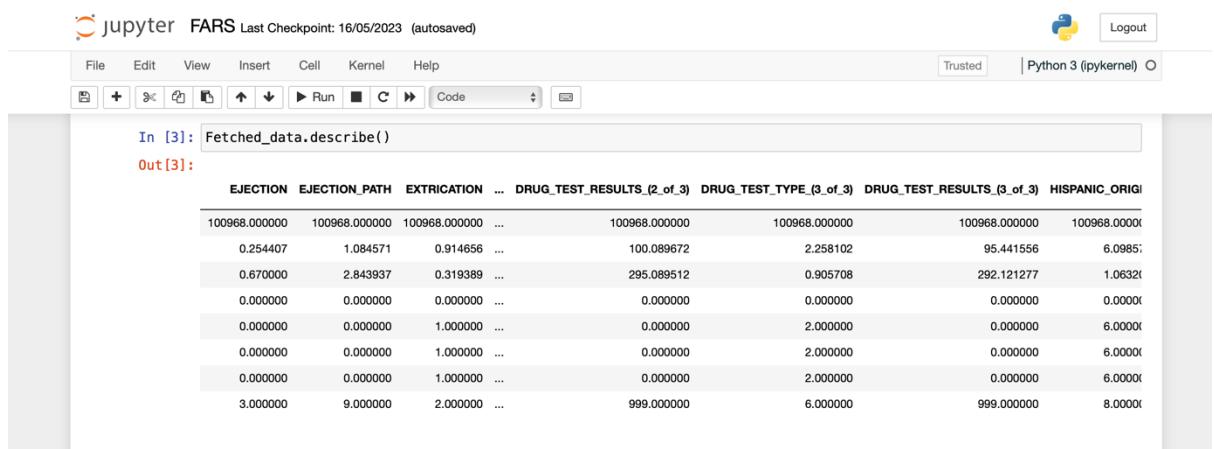


Figure 6

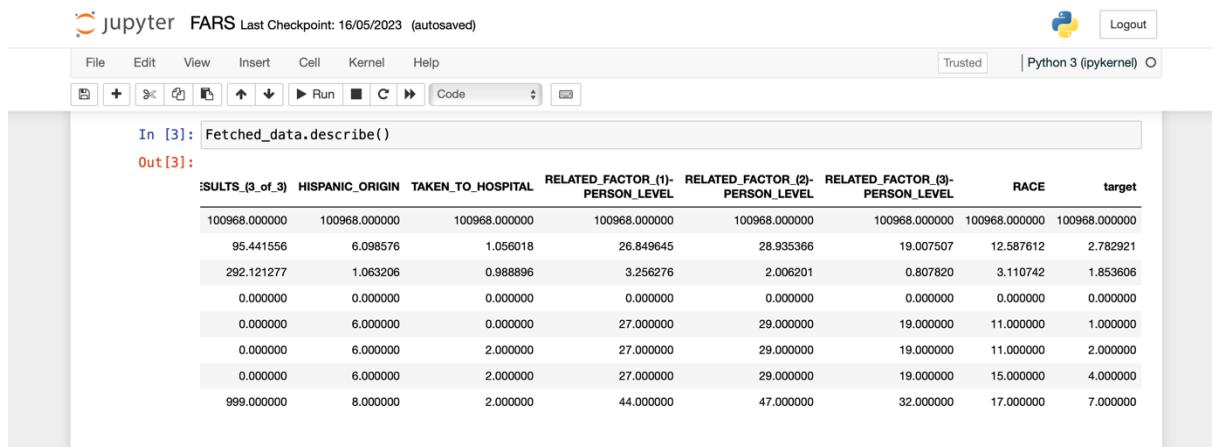


Figure 7

4. Data Pre-processing and cleaning

There are 28 discrete and category features in the data that was acquired. Pre-processing of the data has already taken place, including the transformation of categorical text data into categorical numeric data. Numerical features are used to represent category data

because machine learning models can only work with numbers to find patterns. In the dataset, there are 100968 samples.

CASE_STATE	AGE	SEX	PERSON_TYPE	SEATING_POSITION	RESTRAINT_SYSTEM-USE	AIR_BAG_AVAILABILITY/DEPLOYMENT	EJECTION	EJECTION_PA
100963	50	10	0	6	13	5	4	0
100964	50	9	0	6	16	5	4	0
100965	50	7	0	6	14	4	4	0
100966	50	4	0	6	14	4	4	0
100967	50	61	1	1	3	7	4	0

Figure 8

The data set's columns are as follows:

```
In [2]: Fetched_data.columns
Out[2]: Index(['CASE_STATE', 'AGE', 'SEX', 'PERSON_TYPE', 'SEATING_POSITION',
       'RESTRAINT_SYSTEM-USE', 'AIR_BAG_AVAILABILITY/DEPLOYMENT', 'EJECTION',
       'EJECTION_PATH', 'EXTRICATION', 'NON_MOTORIST_LOCATION',
       'POLICE_REPORTED_ALCOHOL_ININVOLVEMENT', 'METHOD_ALCOHOL_DETERMINATION',
       'ALCOHOL_TEST_TYPE', 'ALCOHOL_TEST_RESULT',
       'POLICE_REPORTED_DRUG_ININVOLVEMENT', 'METHOD_OF_DRUG_DETERMINATION',
       'DRUG_TEST_TYPE', 'DRUG_TEST_RESULTS_(1_of_3)',
       'DRUG_TEST_TYPE_(2_of_3)', 'DRUG_TEST_RESULTS_(2_of_3)',
       'DRUG_TEST_TYPE_(3_of_3)', 'DRUG_TEST_RESULTS_(3_of_3)',
       'HISPANIC_ORIGIN', 'TAKEN_TO_HOSPITAL',
       'RELATED_FACTOR_(1)-PERSON_LEVEL', 'RELATED_FACTOR_(2)-PERSON_LEVEL',
       'RELATED_FACTOR_(3)-PERSON_LEVEL', 'RACE', 'target'],
      dtype='object')
```

Figure 9

There are data samples in the dataset that are not balanced with respect to the objective variable. The model may be trained to predict data from the majority class but perform badly on data from the minority class as a result. SMOTE, a sampling method, is used to effectively address this and improve productivity. SMOTE is used to enrich the minority class samples, enhancing model recall, precision, and accuracy.

The number of samples for each class is displayed in the screenshot below. There are only 9 samples in minority class samples, such as "0." It is balanced with the majority class using the oversampling method.

```
In [4]: column = Fetched_data.groupby('target').count()
column
```

CASE_STATE	AGE	SEX	PERSON_TYPE	SEATING_POSITION	RESTRAINT_SYSTEM-USE	AIR_BAG_AVAILABILITY/DEPLOYMENT	EJECTION	EJECTION_I
target								
0	9	9	9	9	9	9	9	9
1	42116	42116	42116	42116	42116	42116	42116	42116
2	15072	15072	15072	15072	15072	15072	15072	15072
3	299	299	299	299	299	299	299	299
4	20007	20007	20007	20007	20007	20007	20007	20007
5	13890	13890	13890	13890	13890	13890	13890	13890
6	8674	8674	8674	8674	8674	8674	8674	8674
7	901	901	901	901	901	901	901	901

8 rows × 29 columns

Figure 10

```
In [10]: #SMOTE - oversampling technique to increase the train data for all the classes
from collections import Counter
from imblearn.over_sampling import SMOTE

smote = SMOTE()

# fit predictor and target variable
x_smote, y_smote = smote.fit_resample(X, y)

print('Original dataset shape', Counter(y))
print('Resample dataset shape', Counter(y_smote))
print(type(x_smote))
```

Original dataset shape Counter({1: 42116, 4: 20007, 2: 15072, 5: 13890, 6: 8674, 7: 901, 3: 299, 0: 9})
Resample dataset shape Counter({1: 42116, 2: 42116, 4: 42116, 6: 42116, 7: 42116, 5: 42116, 3: 42116, 0: 42116})<class 'numpy.ndarray'>

Figure 11

The variance of the characteristics is normalised using the usual scaler method. A standard scaler is needed since the dataset values differ on a larger scale. Age spans from 40 to 80, but most of the category entries have values of 1 or 0. The performance of the machine learning model could be harmed by this change. Following standardisation, the scaled values are as follows:

The screenshot shows a Jupyter Notebook interface. The code in cell [11] is:

```
In [11]: #normalisation technique - standard scaler
from sklearn.preprocessing import StandardScaler
scaler = StandardScaler()
from sklearn import preprocessing
scaler.fit(x_smote)
scaled_features = scaler.transform(x_smote)
df=pd.DataFrame(scaled_features)
print(df)
```

The output of the code is a DataFrame with 336923 rows and 21 columns. The first few rows of the output are:

	0	1	2	3	4	5	6	7	8	9	...	19	20	21
0	-1.626149	-0.477292	0.369947	-0.798372	-0.573967	0.108052	-1.004346	2.489578	3.180075	0.161378	...	-0.276116	-0.348155	-0.282208
1	-1.626149	-0.948436	0.369947	-0.798372	-0.573967	0.108052	1.071374	2.489578	3.180075	0.161378	...	-0.276116	-0.348155	-0.282208
2	-1.626149	-0.174414	0.369947	-0.798372	-0.573967	-0.768138	-0.411283	2.489578	3.180075	0.161378	...	-0.276116	-0.348155	-0.282208
3	-1.626149	-0.342680	-1.178985	1.265683	-0.034362	-0.768138	-0.411283	2.489578	3.180075	0.161378	...	-0.276116	-0.348155	-0.282208
4	-1.626149	0.061158	0.369947	-0.798372	-0.573967	-0.768138	1.071374	2.489578	3.180075	0.161378	...	-0.276116	-0.348155	-0.282208

Figure 12

5. Logistic regression

Logistic regression is not the ideal method when there are more than two groups to divide the data sample into. The following classification problem, which entails dividing the data into eight different groups, can be resolved using logistic regression, though. The multi class classification approach or logistic regression is used to categorise the data samples. Each class will be handled as a binary classification problem, with 1 denoting that the sample belongs to a certain class and 0 denoting that it does not. The one vs. all method is the name given to the approach detailed below. A sigmoid function is used in the logistic regression method, and it produces a result of 1 if its return value exceeds 0.5 and a result of 0 if it is less than 0.5.

$$h = \frac{1}{1+e^{-z}}$$

where z is the input features multiplied by a theta value that has been initialised randomly.

The regression model receives the standardised data as input. The train/test split technique from sklearn is then used to divide the input data into test and training sets. After that, the model is trained using training data, and test data are used to calculate its correctness. The Predict function is used to calculate accuracy.

The accuracy and confusion matrix are displayed below.

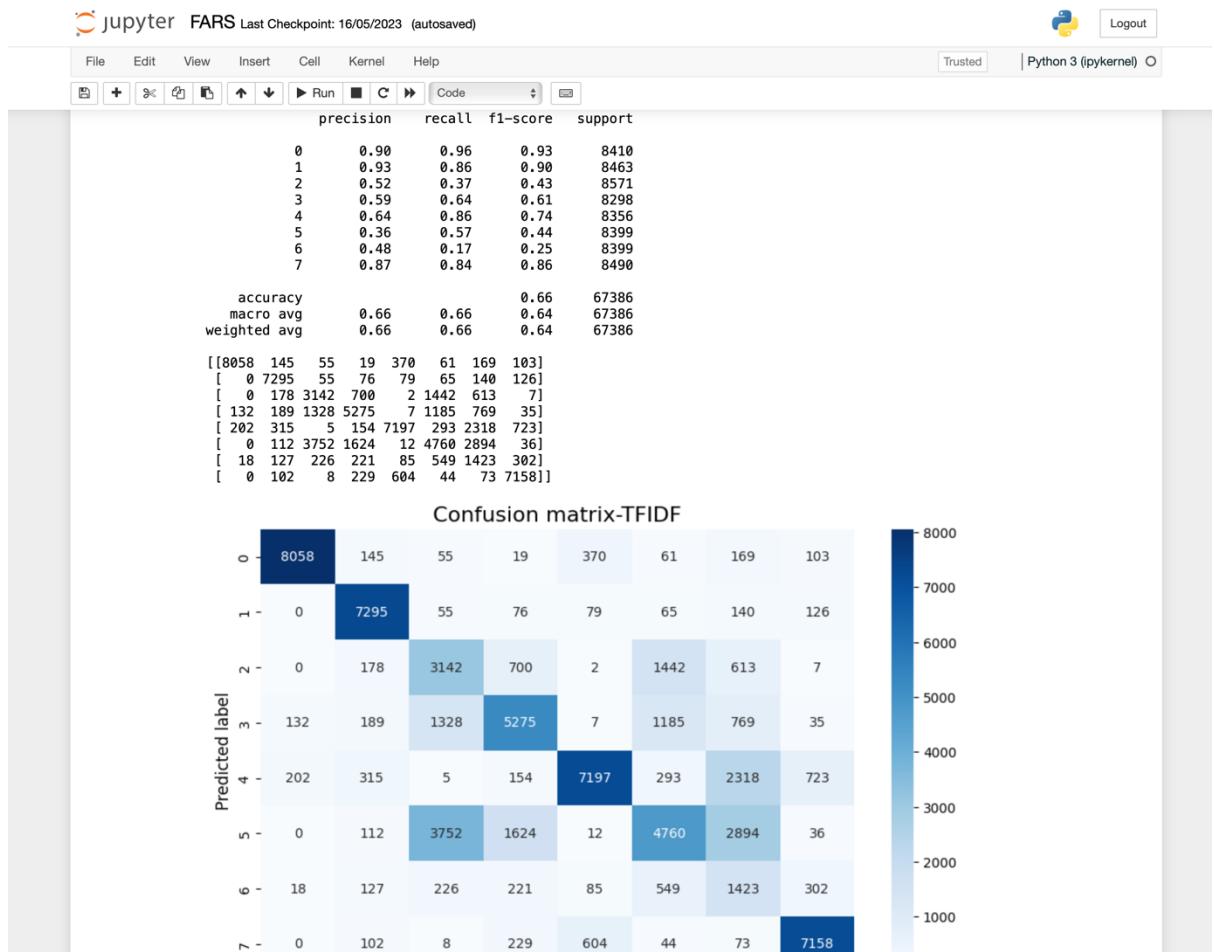


Figure 13

The accuracy is 66, which is lower than we could have anticipated, and we can see that there are a few false negative readings.

6. K Nearest Neighbour

Although KNN is not the best solution when dealing with a huge dataset, it can produce very high classification accuracy scores. KNN only functions with numerical data, and standardisation aids in enhancing accuracy. The KNN model receives the standardised value as input and trains.

The nearest annotated point is used in KNN to classify the data points. The Euclidean, Manhattan, Hamming, and Minkowski distances are used to determine how far apart two points are from one another.

KNN with SMOTE

The accuracy is 88, and the KNN prediction results following SMOTE are as follows:

```
In [13]: # testing accuracy - KNN after class imbalance correction
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report,confusion_matrix
X_train, X_test, y_train, y_test = train_test_split(df,y_smote,
                                                    test_size=0.20)
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train,y_train)
pred = knn.predict(X_test)
print(pred)
print(classification_report(y_test,pred))

[3 7 7 ... 0 5 5]
precision    recall   f1-score   support
          0       1.00     1.00      1.00     8364
          1       1.00     1.00      1.00     8379
          2       0.76     0.75      0.75     8421
          3       0.98     0.99      0.99     8394
          4       0.88     0.87      0.88     8520
          5       0.71     0.70      0.71     8493
          6       0.72     0.73      0.73     8386
          7       0.98     0.99      0.98     8429

accuracy                           0.88      67386
macro avg                           0.88      0.88      67386
weighted avg                          0.88      0.88      67386
```

Figure 14

K is initially set to 1, and KNN is applied without the use of SMOTE directly to the data. It was clear that the SMOTE technique produced high accuracy because KNN without SMOTE performed worse than KNN utilising SMOTE data.

```
In [12]: # testing accuracy -KNN before class imbalance correction
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report,confusion_matrix
X_train, X_test, y_train, y_test = train_test_split(X,y,
                                                    test_size=0.20)
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train,y_train)
pred = knn.predict(X_test)
print(pred)
print(classification_report(y_test,pred))

[2 2 2 ... 5 5 4]
precision    recall   f1-score   support
          0       0.00     0.00      0.00       2
          1       0.97     0.97      0.97     8404
          2       0.46     0.46      0.46     3094
          3       0.34     0.33      0.34       57
          4       0.77     0.76      0.76     3939
          5       0.37     0.38      0.37     2770
          6       0.24     0.24      0.24     1729
          7       0.60     0.56      0.58     199

accuracy                           0.47      20194
macro avg                           0.47      0.47      20194
weighted avg                          0.47      0.47      20194
```

Figure 15

Hyper parameter tuning

The best value for KNN is discovered by hyper parameter tuning. When K is 1, the error rate is very low. In compared to K = 1, the mistake rate is significantly higher for other K values.

```

jupyter FARS Last Checkpoint: 16/05/2023 (autosaved)
File Edit View Insert Cell Kernel Help
Trusted | Python 3 (ipykernel) Logout
In [15]: #hyperparameter tuning wit KNN
error_rate = []
# Will take some time
for i in range(1,20):
    knn = KNeighborsClassifier(n_neighbors=i)
    knn.fit(X_train,y_train)
    pred_i = knn.predict(X_test)
    error_rate.append([i, np.mean(pred_i != y_test)]) #average error rate
plt.figure(figsize=(10,6))
plt.plot(range(1,40),error_rate,color='blue', linestyle='dashed', marker='o',
markerfacecolor='red', markersize=10)
plt.title('Error Rate vs. K Value')
plt.xlabel('K')
plt.ylabel('Error Rate')
Out[15]: Text(0, 0.5, 'Error Rate')

```

Figure 16

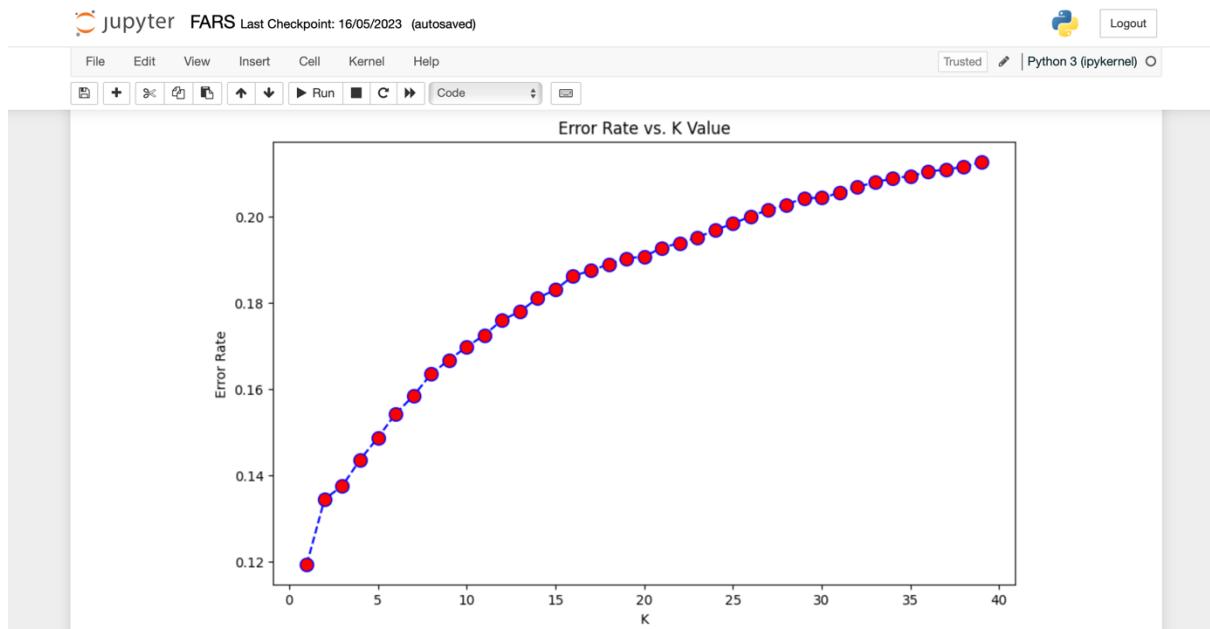


Figure 17

When K=1 accuracy is as below, 88,

```
In [12]: # testing accuracy -KNN before class imbalance correction
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report,confusion_matrix
X_train, X_test, y_train, y_test = train_test_split(X,y,
                                                    test_size=0.20)
from sklearn.neighbors import KNeighborsClassifier
knn = KNeighborsClassifier(n_neighbors=1)
knn.fit(X_train,y_train)
pred = knn.predict(X_test)
print(pred)
print(classification_report(y_test,pred))

[2 2 2 ... 5 5 4]
precision    recall   f1-score   support
      0       0.00     0.00     0.00        2
      1       0.97     0.97     0.97    8404
      2       0.46     0.46     0.46    3094
      3       0.34     0.33     0.34      57
      4       0.77     0.76     0.76    3939
      5       0.37     0.38     0.37    2770
      6       0.24     0.24     0.24    1729
      7       0.60     0.56     0.58    199
accuracy                           0.70    20194
macro avg       0.47     0.46     0.47    20194
weighted avg    0.70     0.70     0.70    20194
```

Figure 18

7. Decision tree classifier

A decision tree is just a list of choices made to accomplish a specific goal. Data samples are categorised in a decision tree from root to leaf node based on the decision rule inferred from the training data. Using a decision tree technique, one may forecast class labels by training the data and learning the fundamental decision-making principles from previous training data. The KNN model's accuracy score has already increased significantly because of splitting the oversampled SMOTE data. Using SMOTE data allows the model to be trained with a sizable sample size from each of the eight classes while maintaining balance. The accuracy rating was good, and the 88 obtained by using sample data. The decision tree classifier's results are as follows.

In [17]:

```
#decision tree classifier
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
X_train, X_test, y_train, y_test = train_test_split(x_smote, y_smote, test_size=0.20)
dtree = DecisionTreeClassifier()
dtree.fit(X_train,y_train)
predictions = dtree.predict(X_test)
from sklearn.metrics import classification_report,confusion_matrix
print(classification_report(y_test,predictions))
print(confusion_matrix(y_test,predictions))
```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	8585
1	1.00	1.00	1.00	8319
2	0.75	0.76	0.76	8564
3	0.98	1.00	0.99	8429
4	0.88	0.89	0.88	8339
5	0.71	0.70	0.70	8340
6	0.75	0.72	0.73	8420
7	0.97	0.99	0.98	8390
accuracy			0.88	67386
macro avg	0.88	0.88	0.88	67386
weighted avg	0.88	0.88	0.88	67386

```
[[8583  0  0  0  1  1  0  0]
 [ 0 8304  2  0  2  2  1  8]
 [ 0  3 6491  49  6 1454  561  0]
 [ 0  0 14 8392  4  10  6  3]
 [ 6  1  7 11 7407 109  646 152]
 [ 1  0 1446  53 149 5844  824  23]
 [ 2  0 637  48 797 836 6063  37]
 [ 0  0  3  0  57  8 16 8306]]
```

Figure 19

8. Random forest classifier

A random forest classifier is an illustration of an ensemble learning method in which a prediction is made using the results of various decision tree algorithms. An 8:2 split of the SMOTE data between the train and test datasets is fed into the model. By default, the N estimator is set to 100, meaning that 100 decision trees will be utilised for each prediction. Predictions are verified using test data once the model has been trained using simulated data. The model scored highly for accuracy, 90. As can be seen here, the confusion matrix contains a significant amount of real positive value.

The screenshot shows a Jupyter Notebook interface with the following details:

- Title Bar:** jupyter FARS Last Checkpoint: 16/05/2023 (autosaved)
- Toolbar:** File, Edit, View, Insert, Cell, Kernel, Help, Run, Stop, Kernel, Help, Code, Cell Type, Cell Kernel, Help, Trusted, Python 3 (ipykernel), Logout.
- Code Cell (In [19]):**

```
#random forest classifier
from sklearn.ensemble import RandomForestClassifier
X_train, X_test, y_train, y_test = train_test_split(x_smote, y_smote, test_size=0.20)
rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(X_train, y_train)
rfc_pred = rfc.predict(X_test)
print("CONFUSION MATRIX:")
print(confusion_matrix(y_test, rfc_pred))
print("ACCURACY REPORT:")
print(classification_report(y_test, rfc_pred))
```
- Output:**

CONFUSION MATRIX:

[8276 0 0 0 1 0 1 0]
[0 8450 0 0 0 0 0 10]
[0 0 6690 23 3 1181 430 2]
[0 0 8 8419 2 6 6 0]
[4 0 4 2 7779 51 473 116]
[0 0 1280 45 203 6155 711 14]
[0 0 542 32 759 742 6489 24]
[0 0 2 0 31 2 11 8407]

ACCURACY REPORT:

	precision	recall	f1-score	support
0	1.00	1.00	1.00	8278
1	1.00	1.00	1.00	8460
2	0.78	0.80	0.79	8329
3	0.99	1.00	0.99	8441
4	0.89	0.92	0.90	8429
5	0.76	0.73	0.74	8408
6	0.80	0.76	0.78	8588
7	0.98	0.99	0.99	8453
accuracy			0.90	67386
macro avg	0.90	0.90	0.90	67386
weighted avg	0.90	0.90	0.90	67386

Figure 20

9. Cross Validation

Cross validation is a technique for confirming that a model is a legitimate machine learning model and has a high prediction score. As testing is entirely conducted on a new data set that the model is completely unaware of during training, the idea of dividing the data into test and train is a cross validation strategy. All of the models I've used use the splitting technique, but I've also used the K-fold cross validation method, in which we alternate splitting the training data.

In my eight iterations, the model will randomly accept train and test data, and the train and test data will change each time. The results of cross validation using the K-fold technique are displayed below. The model will be trained and predicted on multiple train and test data sets concurrently.

```
jupyter FARS Last Checkpoint: 16/05/2023 (autosaved)
File Edit View Insert Cell Kernel Help
Trusted | Python 3 (ipykernel) Logout
[print('CONFUSION MATRIX:', print(classification_report(y_test, rfc_pred)))
CONFUSION MATRIX:
[[8276   0   0   0   1   0   1   0]
 [ 0 8450   0   0   0   0   0   10]
 [ 0   0 6690  23   3 1181  430   2]
 [ 0   0  8419   2   6   6   0]
 [ 4   0   4   2 7779  51 473 116]
 [ 0   0 1280  45 203 6155  711 14]
 [ 0   0  542  32 759 742 6489  24]
 [ 0   0   2   0  31   2 11 8407]]
ACCURACY REPORT:
      precision    recall   f1-score   support
          0       1.00     1.00     1.00     8278
          1       1.00     1.00     1.00     8460
          2       0.78     0.80     0.79     8329
          3       0.99     1.00     0.99     8441
          4       0.89     0.92     0.90     8429
          5       0.76     0.73     0.74     8408
          6       0.80     0.76     0.78     8588
          7       0.98     0.99     0.99     8453
   accuracy                           0.90    67386
macro avg       0.90     0.90     0.90    67386
weighted avg    0.90     0.90     0.90    67386
In [20]: #cross validation
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score,KFold
rfc = RandomForestClassifier()
kf=KFold(n_splits=8)
score=cross_val_score(rfc,x_smote,y_smote,cv=kf)
print("Cross Validation Scores are {}".format(score))
print("Average Cross Validation score :{}".format(score.mean()))
Cross Validation Scores are [0.87211511 0.87703011 0.94747839 0.76315415 0.82885364 0.74237819 0.46108367 0.63467566]
Average Cross Validation score :0.7658461154905498
```

Figure 21

The model has a high accuracy rating, demonstrating its ability to accurately anticipate any new dataset.

10. Conclusion

All the models produced very accurate results, however the best classifier for the FARS dataset was, on balance, Random Forest. The accuracy results for several classifiers are shown below.

	Logistic Regression	KNN without SMOTE	KNN with SMOTE	Decision tree	Random forest
Accuracy	66	70	88	88	90

I used all the models to forecast the 8 various class labels, and SMOTE assisted me in achieving an exceptional accuracy score. Even though Random Forest has a very high accuracy score, I choose KNN because, for a particular sample of test data, Random Forest accuracy during cross validation is just 45.

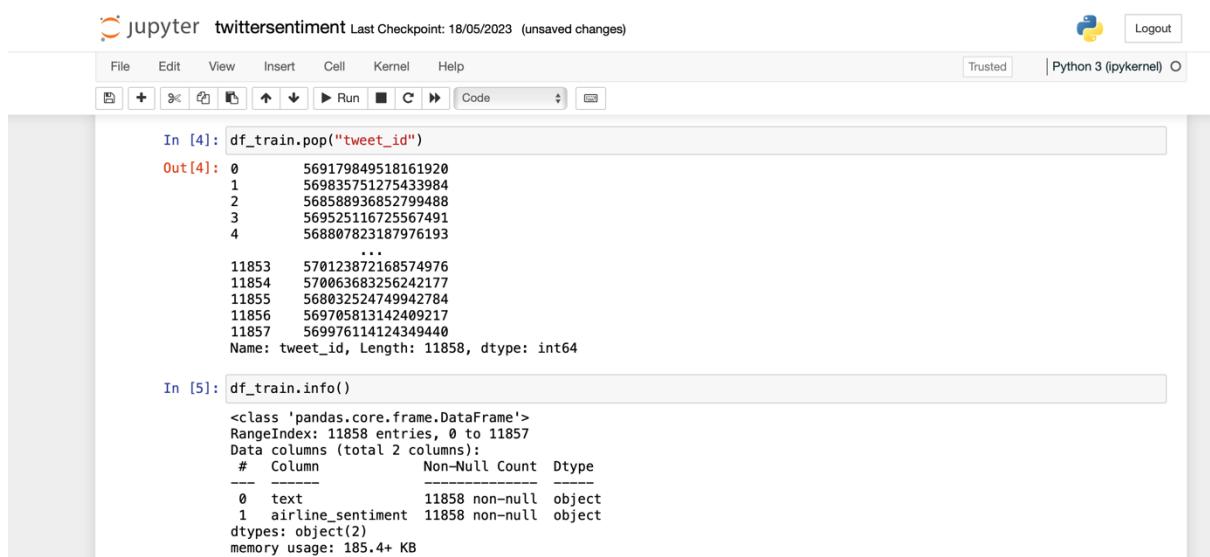
TWITTER SENTIMENT ANALYSIS

1. Introduction

The data set include both user experiences and tweets concerning the airline sector. The work entails preparing Twitter text input and converting category target labels to integers. As shallow machine learning techniques, SVM, Nave Bayes, and the Random Forest classifier were used. The convolution neural network and basic neural network, both of which use deep learning, were also used.

2. Exploratory data analysis

The airline sentiment is the only dependent variable in the dataset, with the tweets and tweet id serving as the only independent factors. The dependent target column displays labels for the positive, negative, and neutral classes. The tweet id field was removed because it included only float values and was irrelevant to my investigation.



```
In [4]: df_train.pop("tweet_id")
Out[4]: 0      569179849518161920
        1      569835751275433984
        2      568588936852799488
        3      569525116725567491
        4      568807823187976193
        ..
       11853    570123872168574976
       11854    570063683256242177
       11855    568832524749942784
       11856    569705813142409217
       11857    569976114124349440
Name: tweet_id, Length: 11858, dtype: int64

In [5]: df_train.info()
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11858 entries, 0 to 11857
Data columns (total 2 columns):
 #   Column            Non-Null Count  Dtype  
--- 
 0   text              11858 non-null   object 
 1   airline_sentiment 11858 non-null   object 
dtypes: object(2)
memory usage: 185.4+ KB
```

Figure 22

The dataset contains 11858 tweet samples with the corresponding sentiment. Additionally, the collection contains no null values. In train data, there are 7434 unfavourable tweets, compared to only 2510 neutral and 1914 positive tweets. The predictor variables are not evenly distributed among the model train samples.

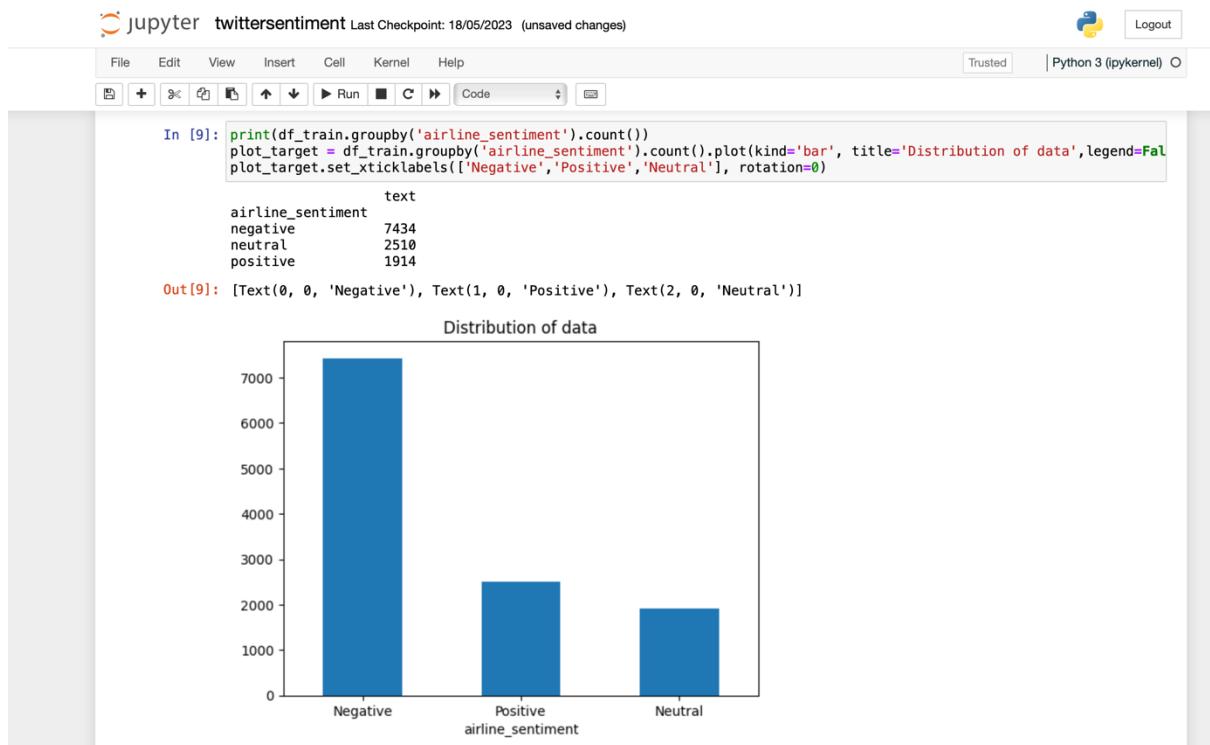


Figure 23

3. Data cleaning

The below-listed loading of the test and training data into the pandas format.

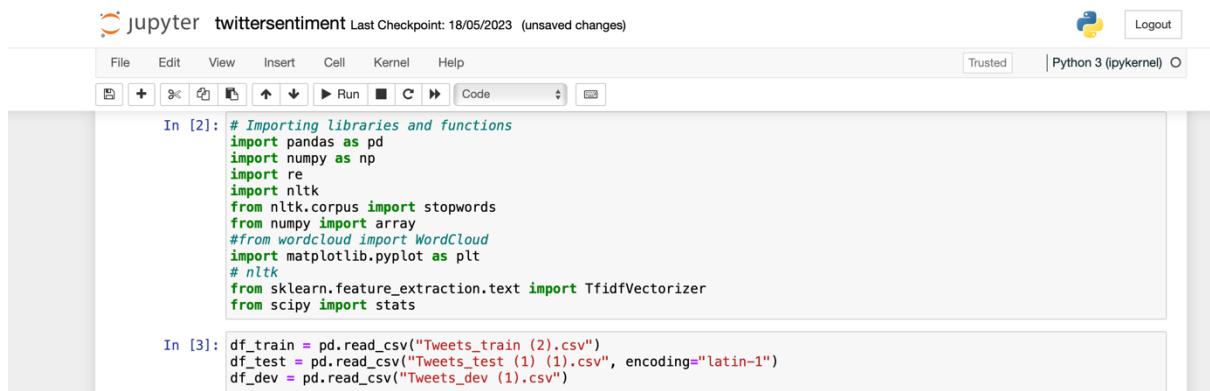
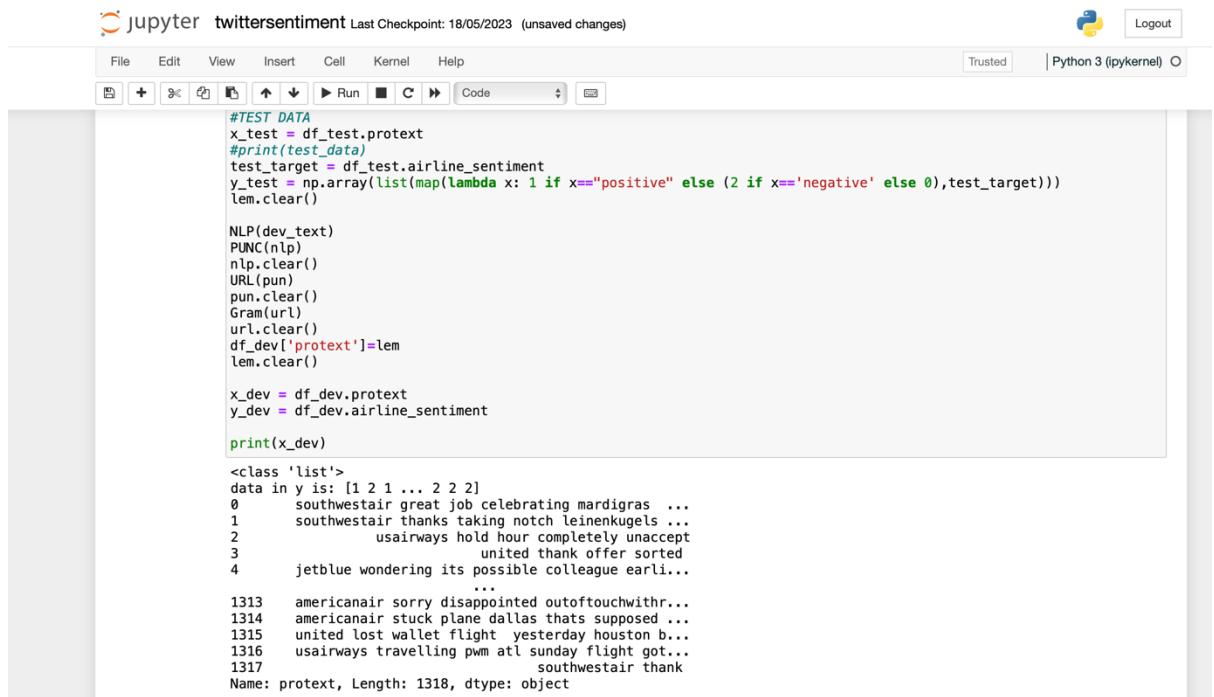


Figure 24

The data is subjected to various data cleaning processes because natural language will impair the accuracy score if supplied straight to machine learning models. Stop words are eliminated as the first step in the data cleaning process. For this, the "ENGLISH STOP WORDS" sklearn package is utilised. The term is eliminated if it appears in the library database since it is a stop word. The regex functionality is then used to remove the punctuation. URLs—which are just text that starts with "http"—can also be deleted using regex functionality. After that, the data is enhanced by the lemmatization and stemming algorithms from nltk, which are used to clean up the noisy words. The

procedure of cleansing the data yielded the following findings. Additionally, the categorical predictor variable is transformed into components as part of the cleaning process so that the machine learning model can understand it.



```
#TEST DATA
x_test = df_test.protext
#print(test_data)
test_target = df_test.airline_sentiment
y_test = np.array(list(map(lambda x: 1 if x=="positive" else (2 if x=='negative' else 0),test_target)))
lem.clear()

NLP(dev_text)
PUNC(nlp)
nlp.clear()
URL(pun)
pun.clear()
Gram(url)
url.clear()
df_dev['protext']=lem
lem.clear()

x_dev = df_dev.protext
y_dev = df_dev.airline_sentiment

print(x_dev)

<class 'list'>
data in y is: [1 2 1 ... 2 2 2]
0    southwestair great job celebrating mardigras ...
1    southwestair thanks taking notch leinenkugels ...
2    usairways hold hour completely unaccept...
3    united thank offer sorted
4    jetblue wondering its possible colleague earli...
...
1313   americanair sorry disappointed outoftouchwithr...
1314   americanair stuck plane dallas thats supposed ...
1315   united lost wallet flight yesterday houston b...
1316   usairways travelling pwm atl sunday flight got...
1317   southwestair thank
Name: protext, Length: 1318, dtype: object
```

Figure 25

Data purification frequently requires converting text data to numerical data because machine learning algorithms can only comprehend numbers. For this, the GLOVE model is employed. The text data is tokenized using the Keras Tokenizer function before the GLOVE model is applied. Thus, there are 12409 terms in the lexicon.



```
In [12]:
word_tokenizer = Tokenizer()
word_tokenizer.fit_on_texts(x_train)

X_train = word_tokenizer.texts_to_sequences(x_train)
X_test = word_tokenizer.texts_to_sequences(x_test)

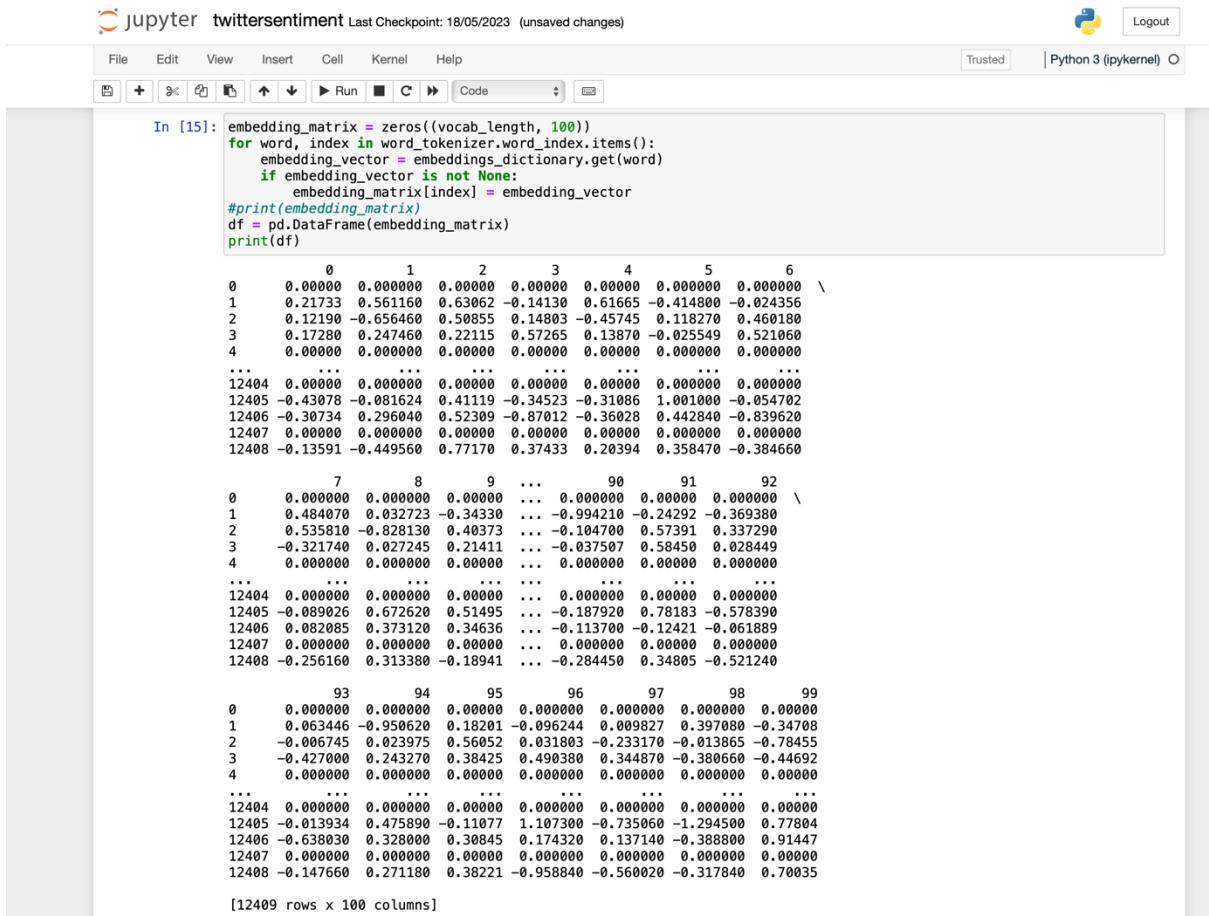
vocab_length = len(word_tokenizer.word_index) + 1
vocab_length

Out[12]: 12409
```

Figure 26

The divided words are then padded to eliminate the balancing challenge. Because of padding, each word has a vector of 100 characters. It will only include 100 '0's, which represents an empty character, if there is no character. For this, the keras library pad Sequences is employed. A list of sequences (lists of integers) with length num samples is converted into a 2D Numpy array with shape (num samples, num timesteps) by this function. The GLOVE file is then imported, which contains a vector for each word. The

dataset word is replaced with its relevant vector data if the tokenized words are available in the GLOVE file.



```
In [15]: embedding_matrix = zeros((vocab_length, 100))
for word, index in word_tokenizer.word_index.items():
    embedding_vector = embeddings_dictionary.get(word)
    if embedding_vector is not None:
        embedding_matrix[index] = embedding_vector
#print(embedding_matrix)
df = pd.DataFrame(embedding_matrix)
print(df)

      0         1         2         3         4         5         6
0  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
1  0.21733  0.56160  0.63062 -0.14130  0.61665 -0.414800 -0.024356
2  0.12190 -0.656460  0.50855  0.14803 -0.45745  0.118270  0.460180
3  0.17280  0.247460  0.22115  0.57265  0.13870 -0.025549  0.521060
4  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
...
12404 0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
12405 -0.43078 -0.081624  0.41119 -0.34523 -0.31086  1.001000 -0.054702
12406 -0.30734  0.296040  0.52309 -0.87012 -0.36028  0.442840 -0.839620
12407  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
12408 -0.13591 -0.449560  0.77170  0.37433  0.20394  0.358470 -0.384660

      7         8         9         ...         90        91        92
0  0.00000  0.00000  0.00000  ...  0.00000  0.00000  0.00000
1  0.484070  0.032723 -0.34330  ... -0.994210 -0.24292 -0.369380
2  0.535810 -0.328130  0.40373  ... -0.104700  0.57391  0.337290
3  -0.321740  0.027245  0.21411  ... -0.037507  0.58450  0.028449
4  0.00000  0.00000  0.00000  ...  0.00000  0.00000  0.00000
...
12404 0.00000  0.00000  0.00000  ...  0.00000  0.00000  0.00000
12405 -0.089026  0.672620  0.51495  ... -0.187920  0.78183 -0.578390
12406  0.082885  0.373120  0.34636  ... -0.113700 -0.12421 -0.061889
12407  0.000000  0.000000  0.00000  ...  0.000000  0.00000  0.00000
12408 -0.256160  0.313380 -0.18941  ... -0.284450  0.34805 -0.521240

      93        94        95        96        97        98        99
0  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
1  0.063446 -0.950620  0.18201 -0.096244  0.009827  0.397080 -0.34708
2  -0.006745  0.023975  0.56052  0.031803 -0.233170 -0.013865 -0.78455
3  -0.427000  0.243270  0.38425  0.490380  0.344870 -0.380660 -0.44692
4  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
...
12404 0.00000  0.00000  0.00000  0.00000  0.00000  0.00000  0.00000
12405 -0.013934  0.475890 -0.11077  1.107300 -0.735060 -1.294500  0.77800
12406 -0.638830  0.328000  0.30845  0.174320  0.137140 -0.388800  0.91447
12407  0.000000  0.000000  0.00000  0.000000  0.000000  0.000000  0.00000
12408 -0.147660  0.271180  0.38221 -0.958840 -0.560020 -0.317840  0.70035

[12409 rows x 100 columns]
```

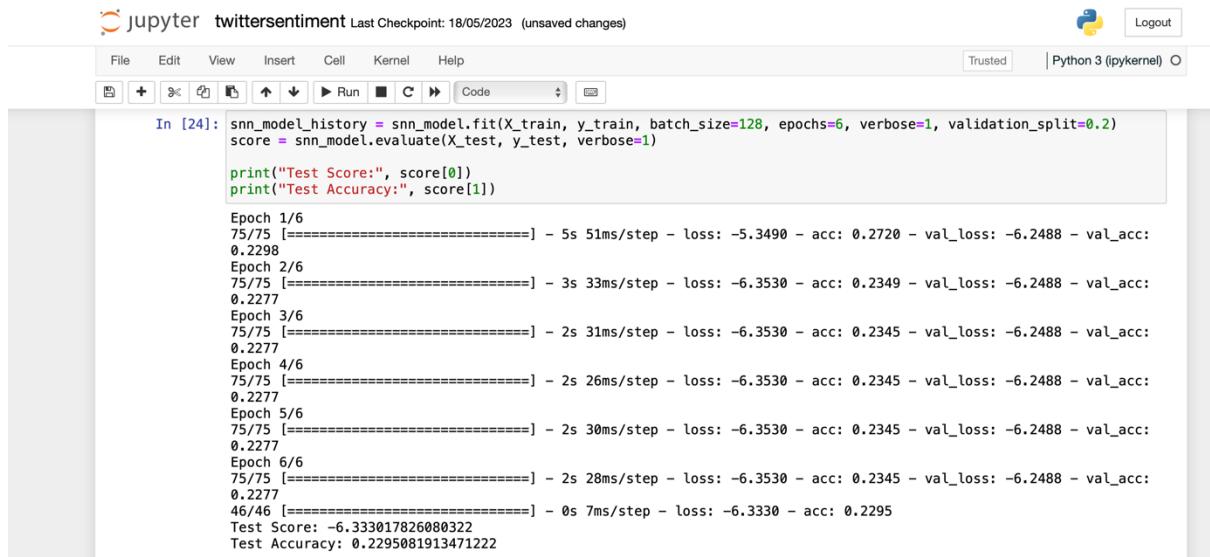
Figure 27

4. Deep Learning – Simple Neural Network

A fundamental neural network is the deep learning strategy applied in this instance. Deep learning is a technique that learns categories in a step-by-step manner using its hidden layer architecture, starting with low-level categories like letters, then progressing to slightly higher-level categories like words, and finally to higher level categories like sentences. A comprehensive representation of the image is formed by the network's neurons and nodes, each of which represents a different part of the entire network. I used deep learning because it provides highly accurate outcomes when given natural language input.

The model utilised here is a sequential model, which implies that the output of one layer is immediately fed into the next layer in all deep learning approaches. Next is the embedding layer. The Embedding layer contains the specification of the input data. Each word can be as long as 100 characters, which was established after padding. Both the entire data that is received immediately after Glove encoding and the entire vector matrix that results from the GLOVE encoding model are supplied as input. There are

three thick layers, and RELU is the activation function used on each layer. The output is altered, and categorisation is aided by the weights that are randomly created for each dense layer. The accuracy I was able to attain for the data below is 63 using 6 epochs. Entropy classification is the loss metric used, and the report includes the data at each epoch is as follows.



```
In [24]: snn_model.history = snn_model.fit(X_train, y_train, batch_size=128, epochs=6, verbose=1, validation_split=0.2)
score = snn_model.evaluate(X_test, y_test, verbose=1)

print("Test Score:", score[0])
print("Test Accuracy:", score[1])

Epoch 1/6
75/75 [=====] - 5s 51ms/step - loss: -5.3490 - acc: 0.2720 - val_loss: -6.2488 - val_acc: 0.2298
Epoch 2/6
75/75 [=====] - 3s 33ms/step - loss: -6.3530 - acc: 0.2349 - val_loss: -6.2488 - val_acc: 0.2277
Epoch 3/6
75/75 [=====] - 2s 31ms/step - loss: -6.3530 - acc: 0.2345 - val_loss: -6.2488 - val_acc: 0.2277
Epoch 4/6
75/75 [=====] - 2s 26ms/step - loss: -6.3530 - acc: 0.2345 - val_loss: -6.2488 - val_acc: 0.2277
Epoch 5/6
75/75 [=====] - 2s 30ms/step - loss: -6.3530 - acc: 0.2345 - val_loss: -6.2488 - val_acc: 0.2277
Epoch 6/6
75/75 [=====] - 2s 28ms/step - loss: -6.3530 - acc: 0.2345 - val_loss: -6.2488 - val_acc: 0.2277
46/46 [=====] - 0s 7ms/step - loss: -6.3330 - acc: 0.2295
Test Score: -6.333017826080322
Test Accuracy: 0.2295081913471222
```

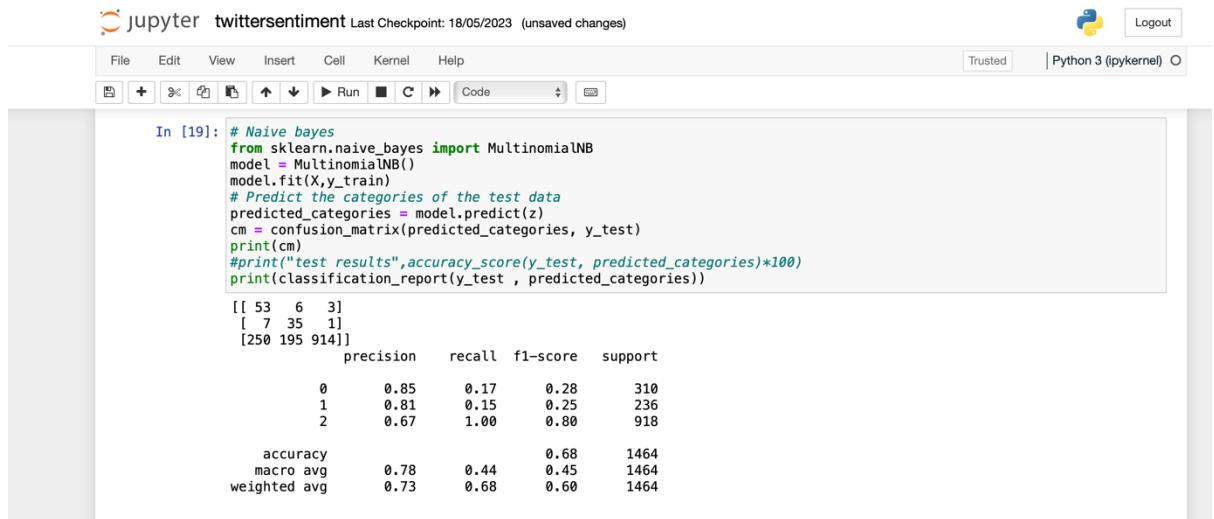
Figure 28

5. Shallow learning – Naïve bayes

A classification technique called Naive Bayes was created using the Bayes theorem. The Bayes theorem, developed by Thomas Bayes, determines the likelihood that an event will occur based on previously known conditions. Its foundation is the following equation:

$$P(A|B) = P(A) * P(B|A)/P(B)$$

The train and test parts of the model are vectorized using the TFID vectorizer. The vectorized test data is the prediction data. The accuracy of the F1 score is 68, which is less. Because Nave Bayes is a straightforward method based on probability, it has a poor accuracy.



```
In [19]: # Naive bayes
from sklearn.naive_bayes import MultinomialNB
model = MultinomialNB()
model.fit(X,y_train)
# Predict the categories of the test data
predicted_categories = model.predict(z)
cm = confusion_matrix(predicted_categories, y_test)
print(cm)
#print("test results",accuracy_score(y_test, predicted_categories)*100)
print(classification_report(y_test , predicted_categories))

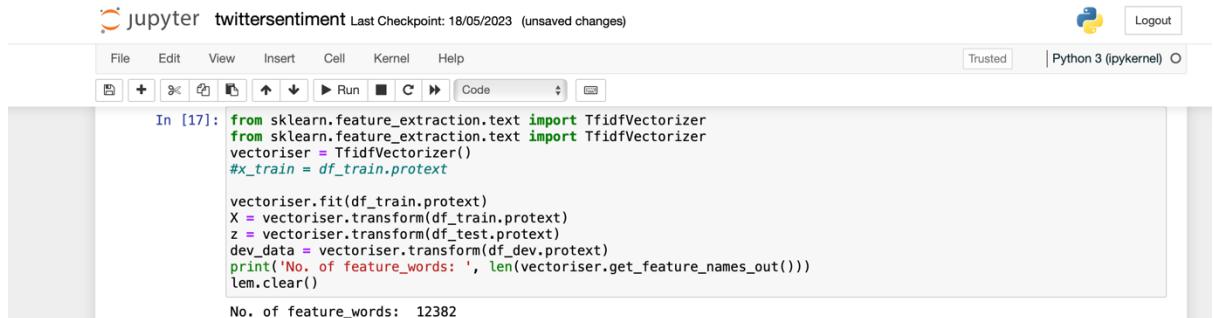
[[ 53   6   3]
 [ 7  35   1]
 [250 195 914]]
precision    recall  f1-score   support
      0       0.85     0.17     0.28      310
      1       0.81     0.15     0.25      236
      2       0.67     1.00     0.80      918
   accuracy                           0.68      1464
  macro avg       0.78     0.44     0.45      1464
weighted avg       0.73     0.68     0.60      1464
```

Figure 29

6. Shallow Learning – SVM

A supervised learning algorithm known as SVM does well with little data. It more accurately assigns each sample to a particular label. The initial training data, which was vectorized early using TFID Vectorizer as follows, is used to fit the model.

On a graph, the data is shown, and the hyper plane divides each class into groups. When there are more than two classes in SVM, this method, known as the one-to-one strategy, is utilised.

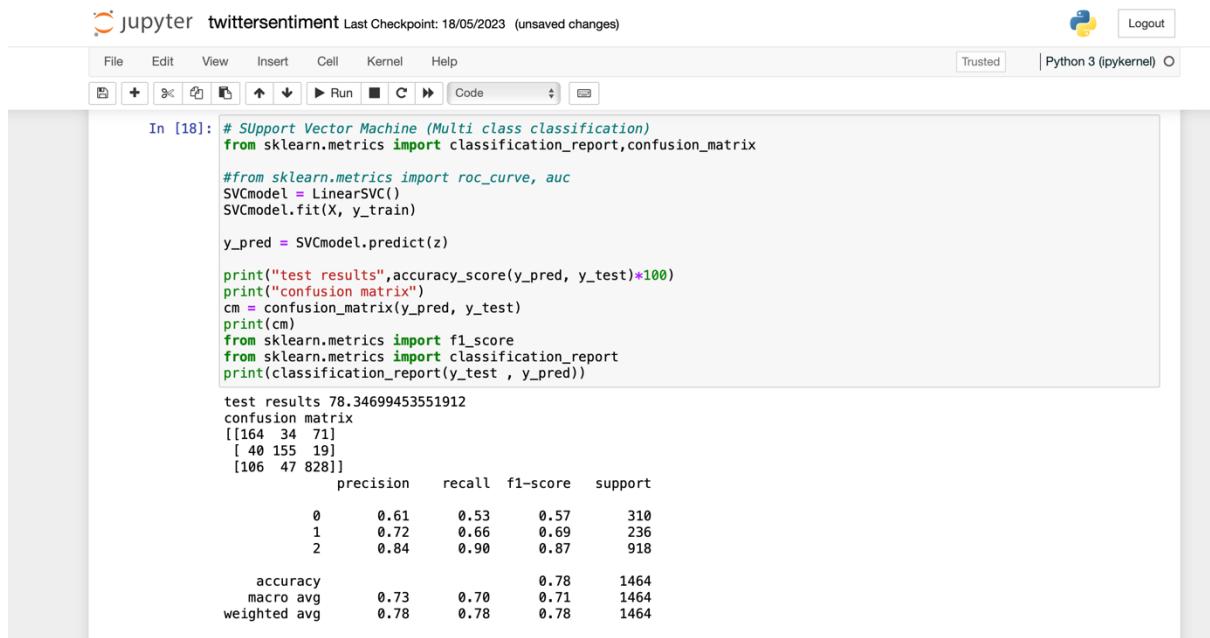


```
In [17]: from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
vectoriser = TfidfVectorizer()
#X_train = df_train.protext

vectoriser.fit(df_train.protext)
X = vectoriser.transform(df_train.protext)
z = vectoriser.transform(df_test.protext)
dev_data = vectoriser.transform(df_dev.protext)
print('No. of feature_words: ', len(vectoriser.get_feature_names_out()))
lem.clear()

No. of feature_words:  12382
```

Figure 30



```

jupyter twittersentiment Last Checkpoint: 18/05/2023 (unsaved changes)
File Edit View Insert Cell Kernel Help
In [18]: # Support Vector Machine (Multi class classification)
from sklearn.metrics import classification_report,confusion_matrix

#from sklearn.metrics import roc_curve, auc
SVCmodel = LinearSVC()
SVCmodel.fit(X, y_train)

y_pred = SVCmodel.predict(z)

print("test results",accuracy_score(y_pred, y_test)*100)
print("confusion matrix")
cm = confusion_matrix(y_pred, y_test)
print(cm)
from sklearn.metrics import f1_score
from sklearn.metrics import classification_report
print(classification_report(y_test , y_pred))

test results 78.34699453551912
confusion matrix
[[164  34  71]
 [ 40 155  19]
 [106  47 828]]
precision    recall   f1-score   support
      0       0.61      0.53      0.57      310
      1       0.72      0.66      0.69      236
      2       0.84      0.90      0.87      918
      accuracy                           0.78      1464
      macro avg       0.73      0.70      0.71      1464
      weighted avg     0.78      0.78      0.78      1464

```

Figure 31

Here, applying the deep learning technique to the vectorized data used for the GLOVE encoding method only produced 62 accuracy. After vectorizing the data, the SVM algorithms produced an accuracy of 78 that is commendable. SVM generated an extremely high accuracy score due to the 1k samples it had, even though deep learning methods are meant to produce outcomes with high accuracy.

7. Shallow classifier – Random Forest

In an ensemble learning method called random forest, the output of numerous decision trees is considered when labelling a particular sample.

```
In [20]: #random forest classifier
from sklearn.ensemble import RandomForestClassifier
#X_train, X_test, y_train, y_test = train_test_split(x_smote, y_smote, test_size=0.20)
rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(X, y)
rfc_pred = rfc.predict(z)
rfc_pred2 = rfc.predict(dev_data)

print("ACCURACY REPORT:")
print(classification_report(test_target, rfc_pred))
print("dev data:")
print(classification_report(y_dev, rfc_pred2))

ACCURACY REPORT:
              precision    recall   f1-score   support
negative        0.79     0.94     0.86      918
neutral         0.67     0.40     0.50      310
positive        0.71     0.55     0.62      236

accuracy          -         -       0.76     1464
macro avg       0.72     0.63     0.66     1464
weighted avg    0.75     0.76     0.74     1464

dev data:
              precision    recall   f1-score   support
negative        0.76     0.94     0.84      826
neutral         0.61     0.34     0.44      279
positive        0.78     0.50     0.61      213

accuracy          -         -       0.74     1318
macro avg       0.72     0.59     0.63     1318
weighted avg    0.73     0.74     0.72     1318
```

Figure 32

In the n_estimator option, which has the value 100, it is stated that 100 decision trees have been run in parallel, with the results being combined to forecast the class labels. As a part of the fit technique, the model begins to learn from the train data. With the help of TFID vectorizer, the train data is transformed into vectors that are then sent to the random forest classifier. The results of the test data are predicted by the machine learning model after it has learned from the train data. In comparison to most models, the F1 score accuracy is 77, which is greater.

As shown below, the random forest classifier's findings for development data.

```
In [20]: #random forest classifier
from sklearn.ensemble import RandomForestClassifier
#X_train, X_test, y_train, y_test = train_test_split(x_smote, y_smote, test_size=0.20)
rfc = RandomForestClassifier(n_estimators=100)
rfc.fit(X, y)
rfc_pred = rfc.predict(z)
rfc_pred2 = rfc.predict(dev_data)

print("ACCURACY REPORT:")
print(classification_report(test_target, rfc_pred))
print("dev data:")
print(classification_report(y_dev, rfc_pred2))

ACCURACY REPORT:
              precision    recall   f1-score   support
negative        0.79     0.94     0.86      918
neutral         0.67     0.40     0.50      310
positive        0.71     0.55     0.62      236

accuracy          -         -       0.76     1464
macro avg       0.72     0.63     0.66     1464
weighted avg    0.75     0.76     0.74     1464
```

Figure 33

8. Conclusion

Since SVM with TFID produced a very high accuracy score compared to other models, I prefer it.

	SNN with GLOVE encoding	SVM with TFID vectorizer	Naïve bayes with TFID vectorizer	Random forest with TFID vectorizer
Accuracy	63	78	68	77

All machine learning models' accuracy is determined by the quantity of true positive values they produce from classes with negative labels and poor scores for other classes. As a result, the model is unable to accurately forecast neutral and positive values, while having learned how to do so with negative values. This results from an imbalance in the data samples. All machine learning models' accuracy is declining because negative sample sizes are disproportionately bigger than those of the other two classes.

References

- [1] – Understanding FARS dataset.
<https://health.gov/healthypeople/objectives-and-data/data-sources-and-methods/data-sources/fatality-analysis-reporting-system-fars>
- [2] - <https://towardsdatascience.com/an-introduction-to-preprocessing-data-for-machine-learning-8325427f07ab>
- [3] - <https://scikit-learn.org/stable/>
- [4] - <https://medium.com/analytics-vidhya/understanding-embedding-layer-in-keras-bbe3ff1327ce#:~:text=Embedding%20layer%20enables%20us%20to,way%20along%20with%20reduced%20dimensions.>
- [5] - <https://towardsdatascience.com/stemming-corpus-with-nltk-7a6a6d02d3e5> stemming
- [6] - <https://towardsdatascience.com/tf-idf-simplified-aba19d5f5530> regex numpy pandas
- [7] - <https://medium.com/swlh/k-nearest-neighbor-ca2593d7a3c4>
- [8] - <https://www.baeldung.com/cs/svm-multiclass-classification>
- [9] - <https://www.tensorflow.org/>
- [11] - [https://www.analyticsvidhya.com/blog/2021/06/exploratory-data-analysis-using-data-visualization-techniques/#:~:text=A%20histogram%20is%20a%20value,at%20the%20center\(mean%\).](https://www.analyticsvidhya.com/blog/2021/06/exploratory-data-analysis-using-data-visualization-techniques/#:~:text=A%20histogram%20is%20a%20value,at%20the%20center(mean%).)