# Problem Statement 1: Tweet Classification

By,
TEAM: OkayXD
Prathmesh Deshmukh
Gavin Henry Lewis
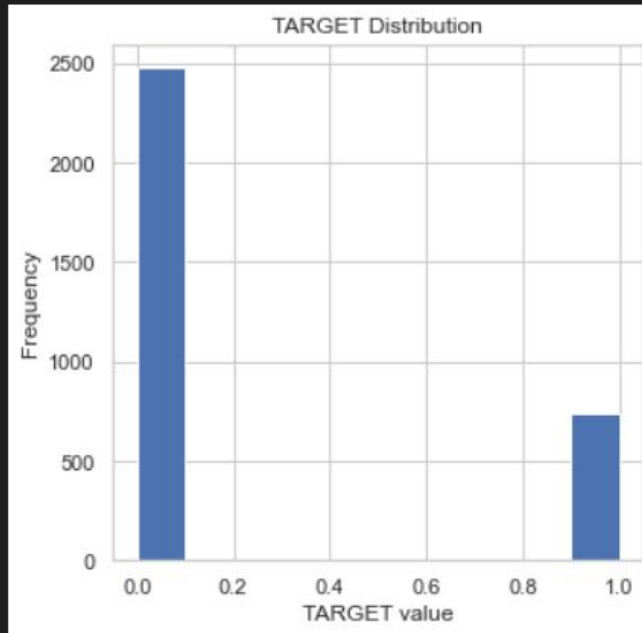
# Exploratory Data Analysis

1. Understanding the dataset.
2. Data Cleaning and removal of 'NA' values
3. Compressing the Variables and the Model:
   a. Reducing the memory impact of the variables by compressing and converting the int types in python
   b. This 'memreduce' function created reduces the memory used by the dataframe from 0.66 to 0.34 as seen in the notebook submitted.
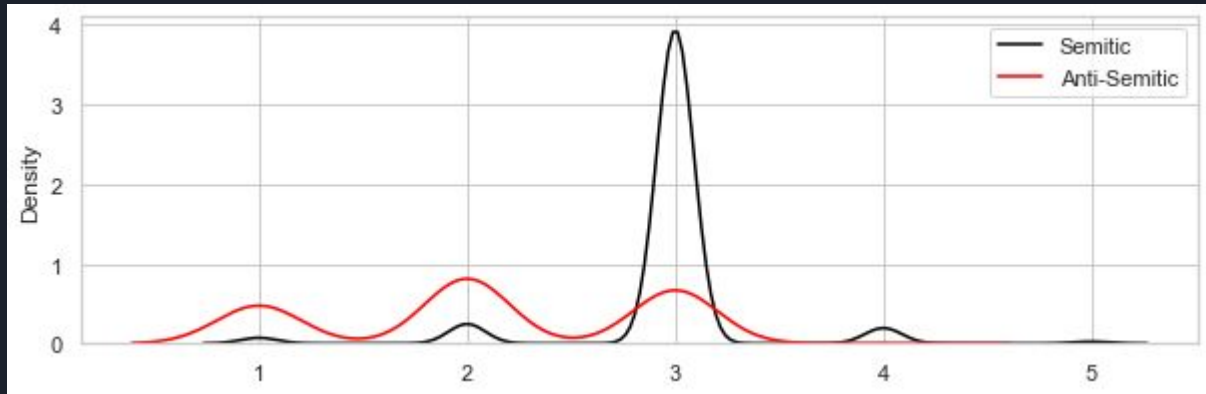
# Data Visualization

1. Understanding the skew in the Target Variable.
   a. Plotted the target variable to understand the distribution of target variable.
   b. This plot was done with respect to the 'Target' variable.

# Data Visualization

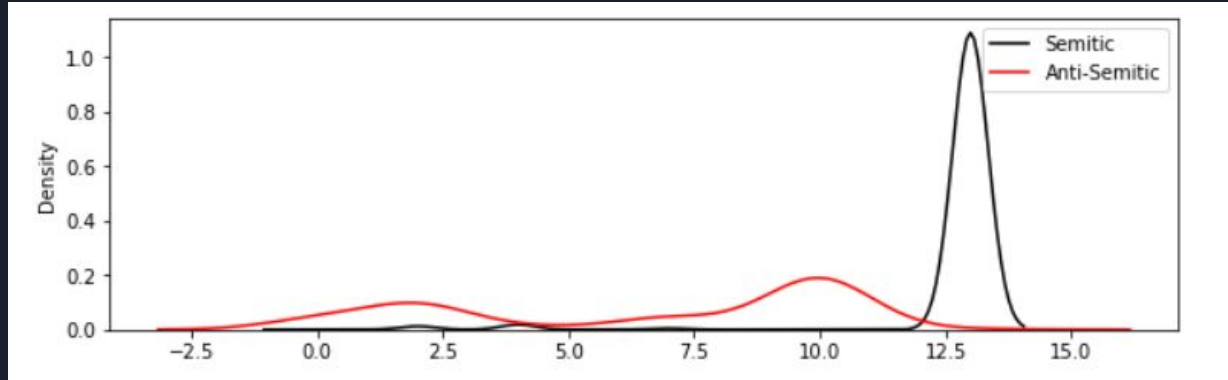2. Density Plots of the Numerical features:

    a.    Sentiment Rating:

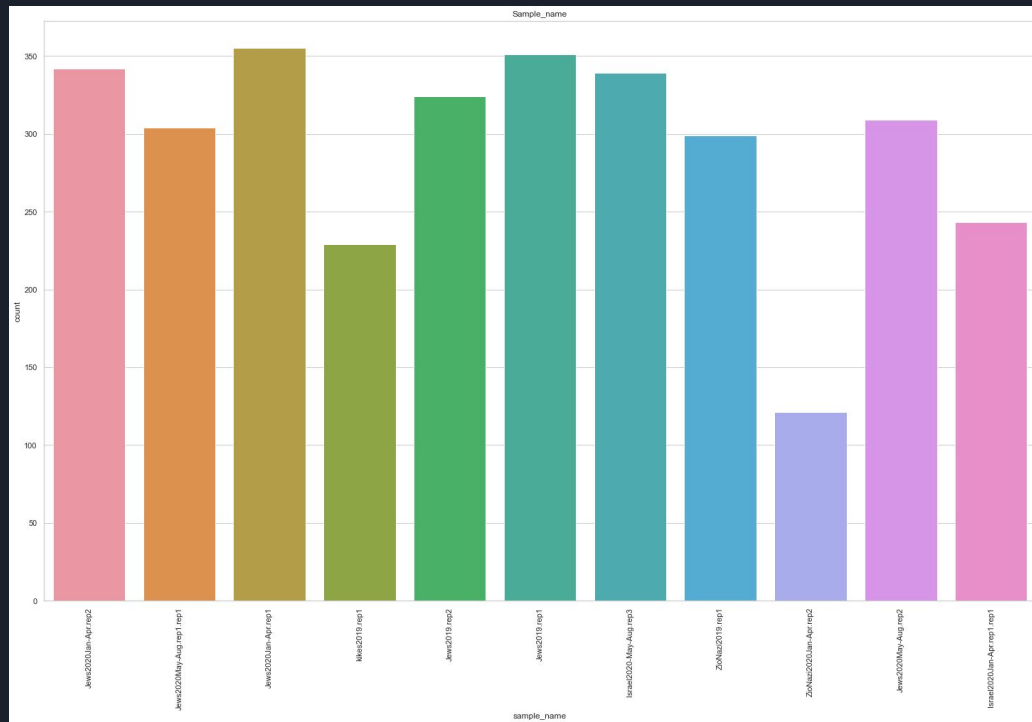# Data Visualization

2. Density Plots of the Numerical features:

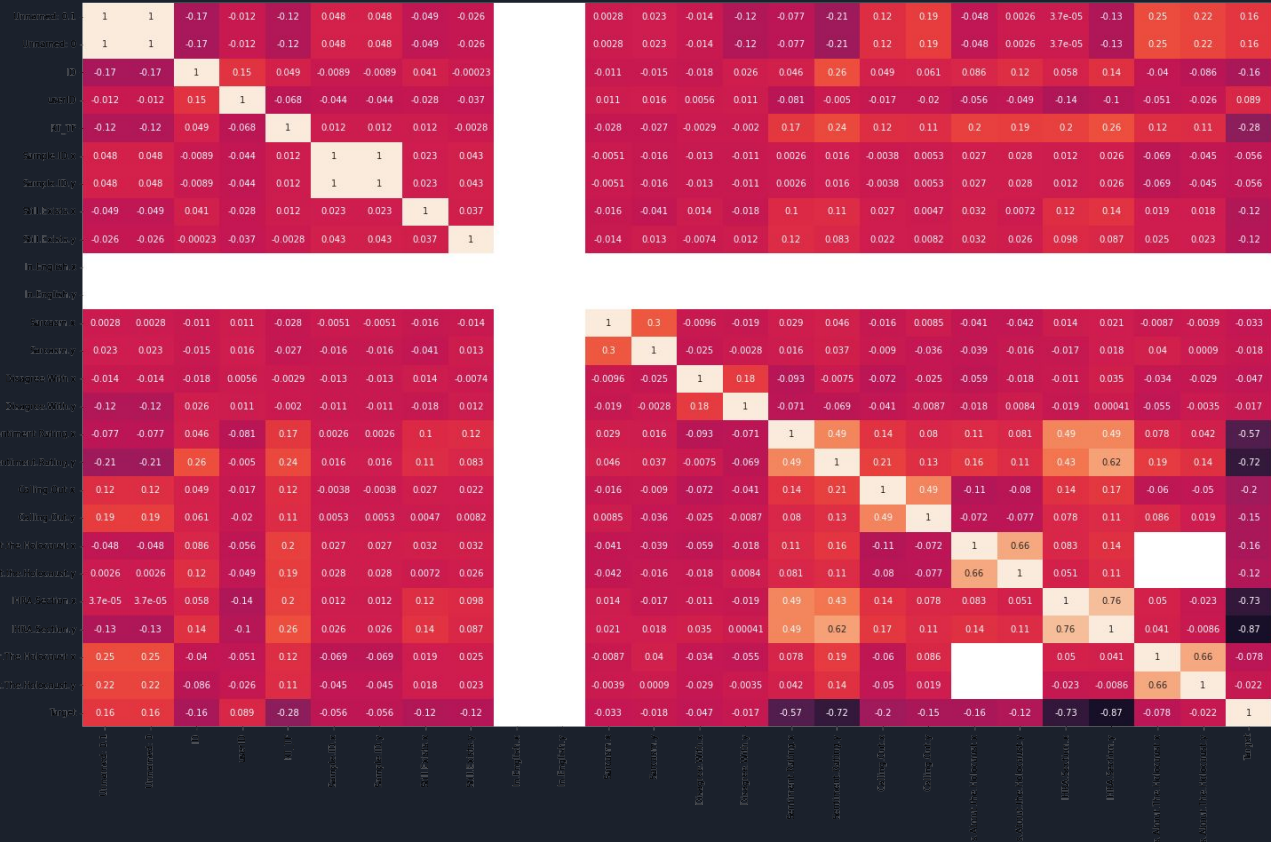      b. IHRA Section:

# Data Visualization

3. Distribution of Categorical variables:

   Distribution of Sample_name

# Correlation Matrix

# Models Used

Conducted model training with the given processed data setwith the models:

1.  XGBoost
2.  SVD
3.  Random Forest
4.  KNN

The results are stored in the expLog variable in the code.

# Experiment Log

| | exp_name | Train Acc | Test Acc | Train AUC | Test AUC | Train F1 | Test F1 |
|---|---|---|---|---|---|---|---|
| 0 | XGB | 0.9949 | 0.9938 | 0.9914 | 0.9913 | 0.9949 | 0.9938 |
| 1 | SVD_LOG | 0.9996 | 1.0000 | 0.9997 | 1.0000 | 0.9996 | 1.0000 |
| 2 | RF | 1.0000 | 0.9984 | 1.0000 | 0.9990 | 1.0000 | 0.9984 |
| 3 | KNN | 1.0000 | 0.9969 | 1.0000 | 0.9933 | 1.0000 | 0.9969 |
| 4 | XGB_b | 0.9916 | 0.9933 | 0.9905 | 0.9916 | 0.9916 | 0.9933 |
| 5 | SVD_LOG_b | 0.9983 | 1.0000 | 0.9986 | 1.0000 | 0.9983 | 1.0000 |
| 6 | SVD_LOG_b | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 7 | SVD_LOG_b | 1.0000 | 0.9966 | 1.0000 | 0.9958 | 1.0000 | 0.9966 |
| 8 | RF_b | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 9 | KNN_b | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 10 | XGB_b | 0.9958 | 0.9832 | 0.9951 | 0.9798 | 0.9958 | 0.9831 |
| 11 | SVD_LOG_b | 1.0000 | 0.9958 | 1.0000 | 0.9949 | 1.0000 | 0.9958 |
| 12 | RF_b | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 | 1.0000 |
| 13 | KNN_b | 1.0000 | 0.9958 | 1.0000 | 0.9949 | 1.0000 | 0.9958 |

# Conclusion

From the experiment log we can conclude that Random Forest model performs better than other models considered in testing.

This is because he test AUC of this model is higher than all others on which tests were conducted.