# TABULAR AND TWEET DATASET CLASSIFICATION USING MACHINE LEARNING

Akhil Raj Rajan
210409183

# Contents

# INTRODUCTION

This project discusses the development of machine learning models for two datasets namely Fars and text data using a variety of techniques. The objectives are to increase the prediction accuracy of the two data by implementing a combination of different approaches including pre-processing techniques, shallow and deep classifiers, ensemble approaches, machine learning approaches and data augmentation.

**PART A- FARS DATA SET**

A brief description of the data set is given below.

1. FARS DATASET
   This data set contains 29 features for predicting the severity of the accident. There are 8 types of severity classes.

# DATASET EXPLORATION

1. FARS DATA SET
   The distribution of each class is shown in the figure below, and it is clear that the data has to be balanced. Furthermore, no missing data were discovered in the data. Figure 2 also demonstrates several relationships between features, underlining the importance of feature selection.
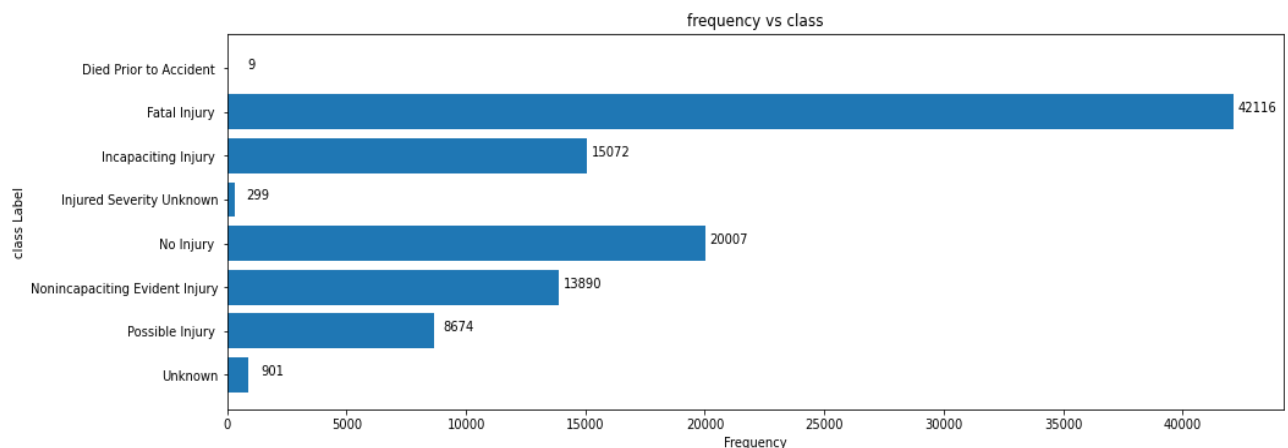


*Figure 1 Frequency of each class*
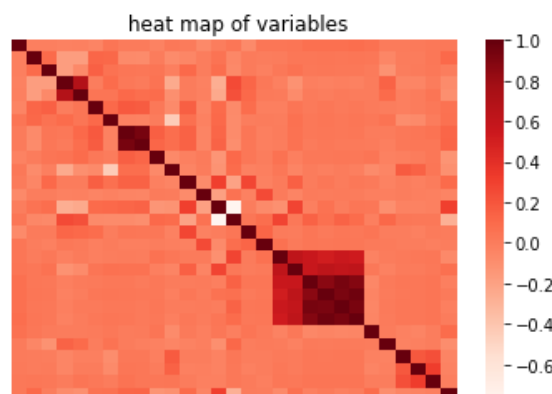


*Figure 2 Heat map of variables*

# METHODOLOGY

The Diagram given below shows a quick understanding of the processes involved in model creation.
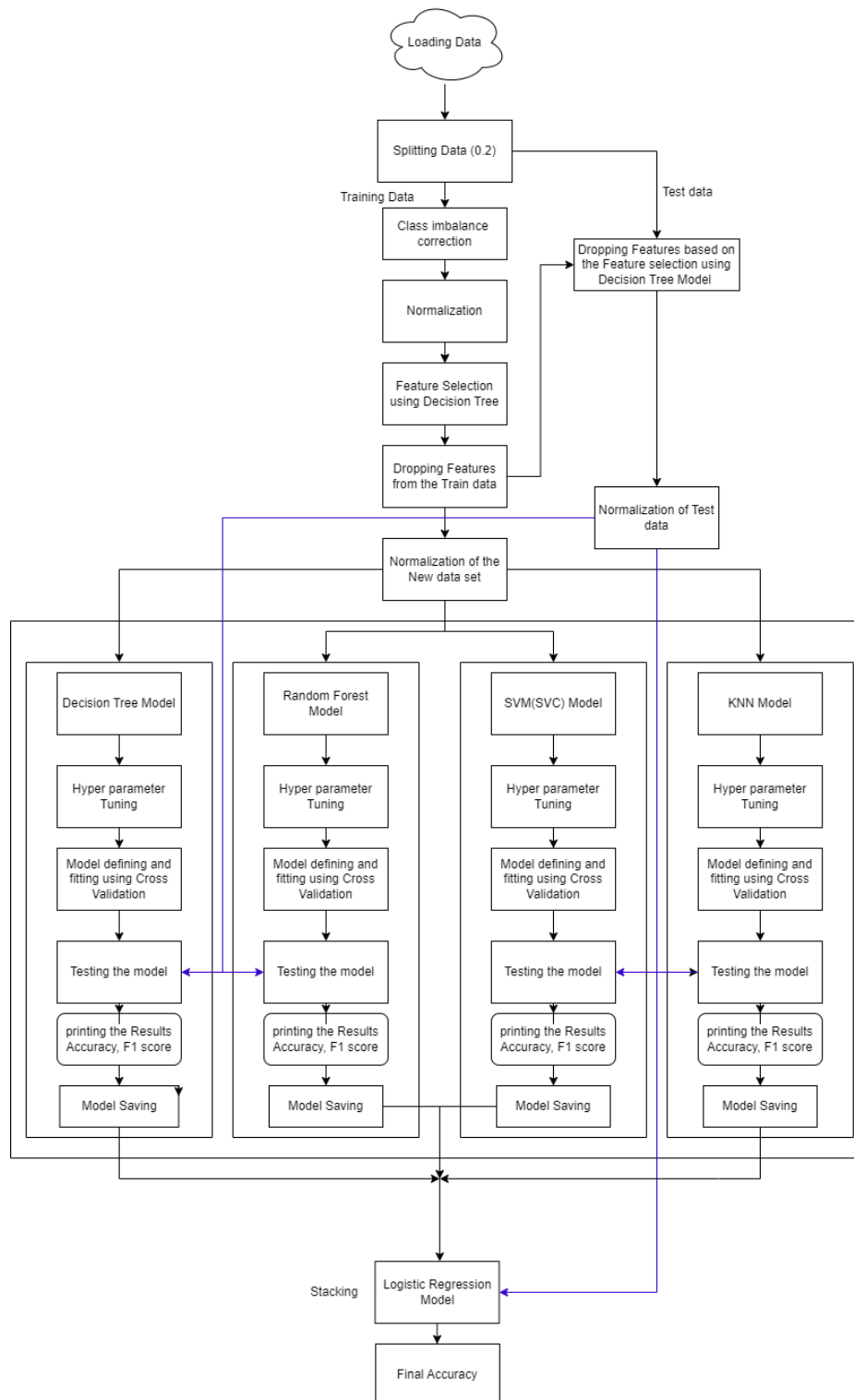


*Figure 3 Flow diagram showing the process involved in Fars data set Machine learning.*

## PRE-PROCESSING OF THE DATA

1) The main data set was split into a train data set and a test data set in the ratio 0.2. *Random_state* is set to 10 to generate the same value in each run.
2) There was an imbalance in the train data generated as the class labelled *1(Fatal injury)* contains 42116 whereas 9 samples are contained in the class labelled *0 (Died prior to accident).* The details can be found in figure 1. This imbalance is handled by SMOT algorithm (see Appendix).
3) The train data is then Normalized which means all the data is scaled between 1 and 0 without losing much information within it, by using the below formula (see Appendix).

## FEATURE SELECTION USING DECISION TREE MODEL

Decision Tree based Model is used to find the best number of features that can predict the class. The algorithm used is:

1) Select the starting number of 10, which tells the algorithm to select 10 combinations of features from the total of 29 features. So, for this case there are 29C10 combinations 20030010).
2) For each combination, a model is fitted using 5-fold cross-validation by passing pre-processed training data. Consequently, accuracy is calculated.
3) The average of the accuracies estimates the model performance which is then stored in an array corresponding to the value 10.
4) Then the number of features is incremented by 1 and iterate above steps till it reaches 25.
5) Select the number of features which showed maximum accuracy and less standard deviation.

Figure (4) shown below show the number of feature vs accuracies.



*Figure 4 Number of Features vs accuracy*

*The number of features as 18 is selected considering the maximum accuracy and minimum standard deviation*

RFE (see appendix) is a module which gives the best feature based on the provided input number. So, with the help of this module, the given features were selected with 18 as input.

The Rank represents how important that feature is in predicting the class, where 1 denotes the highest priority.

| Variables | Remarks | Ranks |
|---|---|---|
| CASE_STATE | TRUE | 1 |
| AGE | TRUE | 1 |
| SEX | TRUE | 1 |
| PERSON_TYPE | TRUE | 1 |
| SEATING_POSITION | TRUE | 1 |
| RESTRAINT_SYSTEM-USE | TRUE | 1 |
| AIR_BAG_AVAILABILITY/DEPLOYMENT | TRUE | 1 |
| EJECTION_PATH | TRUE | 1 |
| POLICE_REPORTED_ALCOHOL_INVOLVEMENT | TRUE | 1 |
| METHOD_ALCOHOL_DETERMINATION | TRUE | 1 |
| ALCOHOL_TEST_TYPE | TRUE | 1 |
| ALCOHOL_TEST_RESULT | TRUE | 1 |
| POLICE-REPORTED_DRUG_INVOLVEMENT | TRUE | 1 |
| DRUG_TEST_RESULTS_(1_of_3) | TRUE | 1 |
| DRUG_TEST_TYPE_(3_of_3) | TRUE | 1 |
| TAKEN_TO_HOSPITAL | TRUE | 1 |
| RELATED_FACTOR_(1)-PERSON_LEVEL | TRUE | 1 |
| RACE | TRUE | 1 |

*Table 1 Selected Variables*

The same feature selection and normalization are applied in the test data as well.

## MODELLING

Initially, the best parameter is decided by hyperparameter tuning using GridSearchCv and BayesSearchCV. (see appendix)

However, compared to random values, these strategies assist in identifying the superior hyperparameters. The best outcomes are produced by knowledge and comprehension of the chosen Model, the significance of the hyperparameter, their relationships, and the tuning procedure.

Four Models were chosen for this data set, namely**, Decision Tree, Random Forest, Support Vector Classifier and K- Nearest Neighbour**. Furthermore, they are all supervised learning techniques. The models were fitted using cross-validation as there are no separate validation data available.

An ensemble method is also used to combine the ability of each model for predicting a particular class.

## RESULTS OF ANALYSIS

### Decision Tree Model

The best hyperparameter combination is found using GridSearchCV and they are :

*Criterion(evaluates the split quality): entropy*

*max_depth (Depth of Tree): 11*

*min_impurity_decrease (minimum allowable value of impurity) : 0.0*

*splitter: best*

Test Result

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Died Prior to Accident | 0 | 0 | 0 | 0 | 2 |
| Fatal Injury | 1 | 1 | 1 | 1 | 8436 |
| Incapacitating Injury | 2 | 0.58 | 0.4 | 0.47 | 3029 |
| Injured Severity Unknown | 3 | 0.14 | 0.58 | 0.23 | 50 |
| No Injury | 4 | 0.85 | 0.9 | 0.87 | 3943 |
| Non incapacitating Evident Injury | 5 | 0.42 | 0.63 | 0.51 | 2848 |
| Possible Injury | 6 | 0.32 | 0.1 | 0.15 | 1711 |
| Unknown | 7 | 0.3 | 0.78 | 0.44 | 175 |
| | | | | | |
| | Accuracy | | | 0.76 | 20194 |
| | F1 macro | 0.45 | 0.55 | 0.46 | 20194 |
| | F1weighted | 0.76 | 0.76 | 0.75 | 20194 |

*Table 2 Decision tree Test Result*

## Random Forest Model

The best hyperparameter combination is found using GridSearchCV and they are:

*Criterion(evaluates the split quality): entropy*

*max_depth (Depth of Tree): 11*

*min_impurity_decrease (minimum allowable value of impurity) : 0.0*

*splitter: best*

Test Results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Died Prior to Accident | 0 | 0 | 0 | 0 | 2 |
| Fatal Injury | 1 | 1 | 1 | 1 | 8436 |
| Incapacitating Injury | 2 | 0.53 | 0.73 | 0.61 | 3029 |
| Injured Severity Unknown | 3 | 0.13 | 0.58 | 0.21 | 50 |
| No Injury | 4 | 0.84 | 0.95 | 0.89 | 3943 |
| Non incapacitating Evident Injury | 5 | 0.47 | 0.39 | 0.43 | 2848 |
| Possible Injury | 6 | 0.46 | 0.05 | 0.09 | 1711 |
| Unknown | 7 | 0.51 | 0.83 | 0.63 | 175 |
| | | | | | |
| | Accuracy | | | 0.78 | 20194 |
| | macro | 0.49 | 0.57 | 0.48 | 20194 |
| | weighted | 0.77 | 0.78 | 0.76 | 20194 |

*Table 3 Test Result of Random Forest Model*

## Support Vector Machine Model (SVC)

In this model, since hyperparameter tuning takes more than 3-4 hrs. Bayesian optimization (see appendix) technique is used for tuning.

The best parameters obtained are:

*C (Regularization parameter) =2,*

*gamma (Kernel coefficient) =0.75,*

*kernel (type of kernel used) ='poly',*

*degree (degree of polynomial kernel) =3*

Test results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Died Prior to Accident | 0 | 0 | 0 | 0 | 2 |
| Fatal Injury | 1 | 1 | 0.99 | 1 | 8436 |
| Incapacitating Injury | 2 | 0.52 | 0.58 | 0.55 | 3029 |
| Injured Severity Unknown | 3 | 0.42 | 0.22 | 0.29 | 50 |
| No Injury | 4 | 0.84 | 0.99 | 0.9 | 3943 |
| Non incapacitating Evident Injury | 5 | 0.44 | 0.54 | 0.49 | 2848 |
| Possible Injury | 6 | 0.53 | 0.04 | 0.08 | 1711 |
| Unknown | 7 | 0.7 | 0.67 | 0.69 | 175 |
| | | | | | |
| | accuracy | | | 0.78 | 20194 |
| | macro | 0.56 | 0.5 | 0.5 | 20194 |
| | weighted | 0.77 | 0.78 | 0.76 | 20194 |

*Table 4 Test Result of SVM Model*

## KNN (k nearest Neighbour) Model

Hyperparameter tuning indicated the following parameter details.

*n_neighbors (Number of neighbours used) :3*

*Metric (Distance calculator): manhattan*

Test Results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Died Prior to Accident | 0 | 0 | 0 | 0 | 2 |
| Fatal Injury | 1 | 0.98 | 0.95 | 0.97 | 8436 |
| Incapacitating Injury | 2 | 0.49 | 0.58 | 0.53 | 3029 |
| Injured Severity Unknown | 3 | 0.17 | 0.48 | 0.25 | 50 |
| No Injury | 4 | 0.84 | 0.87 | 0.85 | 3943 |
| Non incapacitating Evident Injury | 5 | 0.43 | 0.38 | 0.4 | 2848 |
| Possible Injury | 6 | 0.27 | 0.23 | 0.25 | 1711 |
| Unknown | 7 | 0.46 | 0.71 | 0.56 | 175 |
| | | | | | |
| | accuracy | | | 0.73 | 20194 |
| | F1 macro | 0.46 | 0.52 | 0.48 | 20194 |
| | F1 weighted | 0.74 | 0.73 | 0.73 | 20194 |

*Table 5 Test Result of KNN Model*

Comments: with the F1 score in the range between 0.73 and 0.76 all the above models could predict the classes except the class "Died Prior to Accident".

## Ensemble Method (Stacking) using Logistic Regression Model

A logistic regression model is trained by taking the best features of the other 3 models that is Random Forest, SVM and KNN to predict the class better than any other individual models.

Test Result

| class labels | | | precision | recall | f1-score | support |
|---|---|---|---|---|---|---|
| Died Prior to Accident | 0 | | 0 | 0 | 0 | 2 |
| Fatal Injury | 1 | | 1 | 1 | 1 | 8436 |
| Incapacitating Injury | 2 | | 0.59 | 0.61 | 0.6 | 3029 |
| Injured Severity Unknown | 3 | | 0.92 | 0.22 | 0.35 | 50 |
| No Injury | 4 | | 0.84 | 0.99 | 0.91 | 3943 |
| Non incapacitating Evident Injury | 5 | | 0.46 | 0.59 | 0.51 | 2848 |
| Possible Injury | 6 | | 0.52 | 0.06 | 0.1 | 1711 |
| Unknown | 7 | | 0.71 | 0.67 | 0.69 | 175 |
| | | | | | | |
| | accuracy | | | | 0.79 | 20194 |
| | macro | | 0.63 | 0.52 | 0.52 | 20194 |
| | weighted | | 0.79 | 0.79 | 0.77 | 20194 |

*Table 6 Ensembled model Test Result*

Cross-validation Result

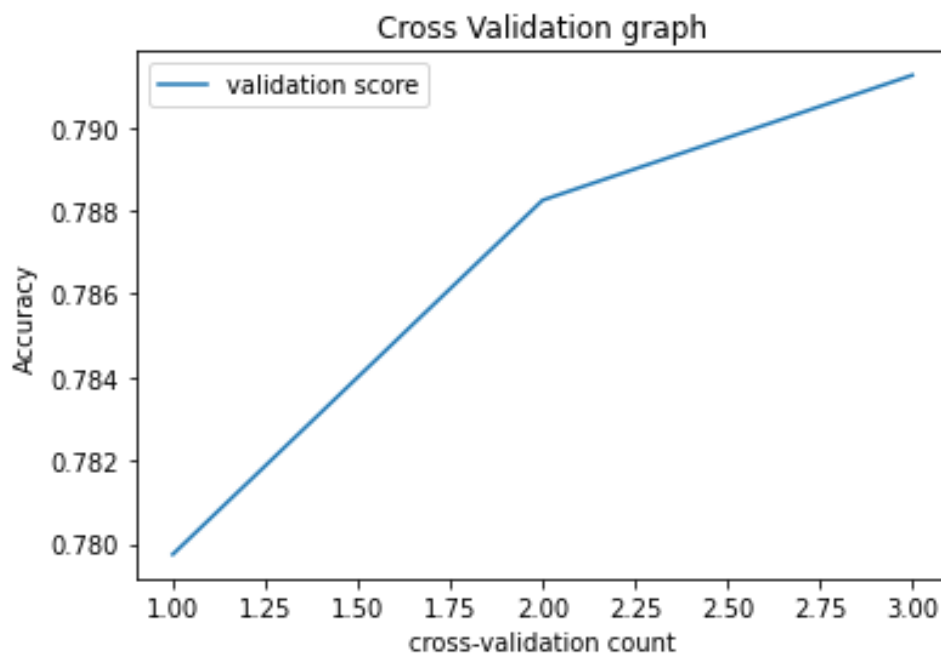The graph indicates an increase in validation accuracy (3-fold).



*Figure 5 3-fold cross validation result*

comments: significant improvement can be seen in the validation score(Accuracy). High prediction can be found in the classes 'no injury', 'Fatal injury' where as it is difficult to predict *Injured Severity Unknown* class and *Possible Injury* class.
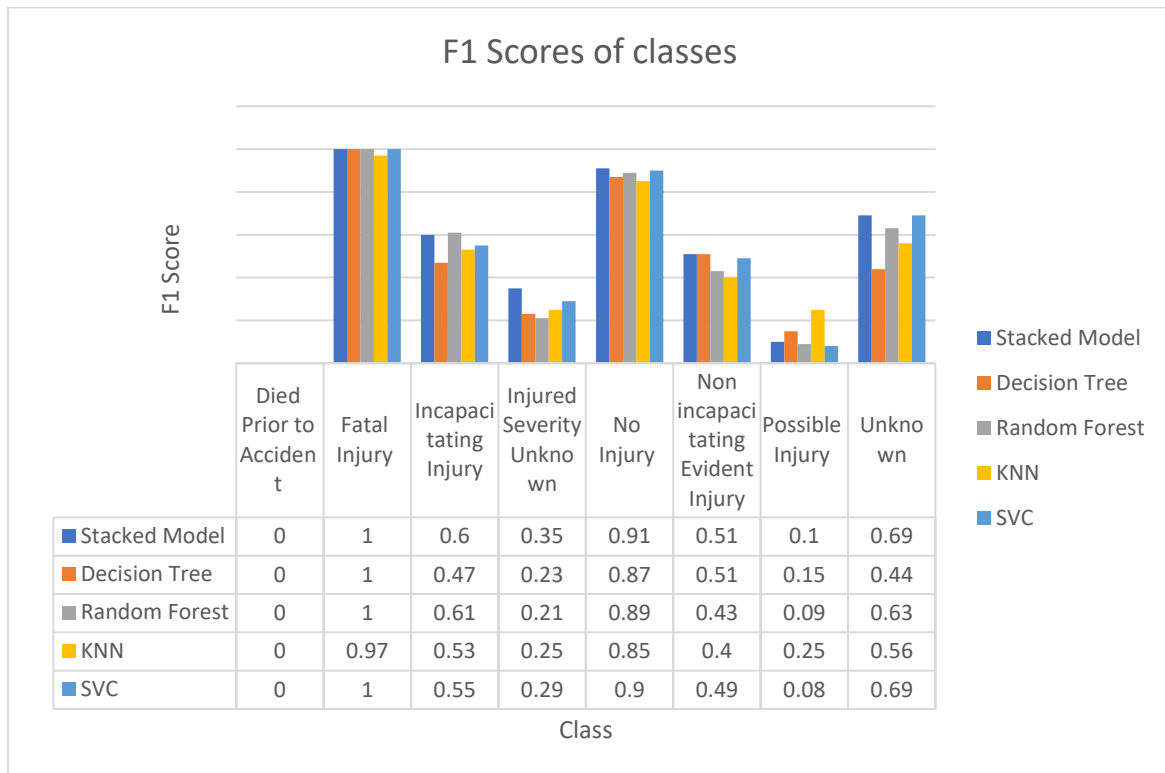
# DISCUSSION OF ANALYSIS



| | Died Prior to Accident | Fatal Injury | Incapacitating Injury | Injured Severity Unknown | No Injury | Non incapacitating Evident Injury | Possible Injury | Unknown |
|---|---|---|---|---|---|---|---|---|
| Stacked Model | 0 | 1 | 0.6 | 0.35 | 0.91 | 0.51 | 0.1 | 0.69 |
| Decision Tree | 0 | 1 | 0.47 | 0.23 | 0.87 | 0.51 | 0.15 | 0.44 |
| Random Forest | 0 | 1 | 0.61 | 0.21 | 0.89 | 0.43 | 0.09 | 0.63 |
| KNN | 0 | 0.97 | 0.53 | 0.25 | 0.85 | 0.4 | 0.25 | 0.56 |
| SVC | 0 | 1 | 0.55 | 0.29 | 0.9 | 0.49 | 0.08 | 0.69 |

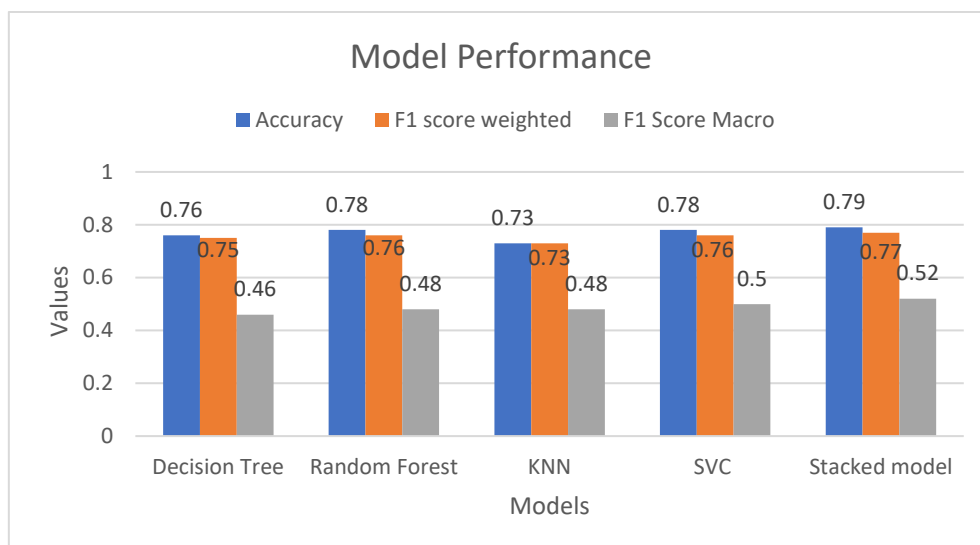*Figure 6 Summary of the Report*



*Figure 7 Performance comparison of Models*

Comments:

- The stacked model can predict accuracy better than any other model, as it could utilize the performance of all other models with a weighted F1 score of 0.77.
- Since the test data is not balanced, the F1 score should be considered for the model analysis.
- KNN produced the least accuracy whereas Random Forest showed high accuracy which is also an ensembled model on its own.
- The class "Died prior to Accident" could not be predicted by any of these models at the same time all the models can predict the class "Fatal Injury" with an accuracy of around 100%. This is due to the difference in the availability of data for these classes. Moreover, it is a good sign that these models could predict fatal injury as it is a critical factor in the real world.

- Furthermore, the classes "No Injury", "Unknown" and "incapacitating injury" can be forecasted by the model with sufficient accuracies.

**PART B**

Tweet data set.

The idea is to predict the category to which each tweet in the dataset belongs by using a Machine learning model that is using sentimental analysis.

Given below is the example data taken from the train set.

| index | tweet_id | text | airline_sentiment |
|-------|----------|------|-------------------|
| 0 | 5.6918E+17 | @united you're good. Thank you! | positive |

*Table 7 An example of training data*

# DATA EXPLORATION

## CLASS IMBALANCE AND OTHER QUALITIES CHECKING

The bar plot in figure 8 shows the distribution of 3 distinct classes (Positive, negative and neutral) in the training data.



*Figure 8 Class-wise distribution of the training data*

Comments: The bar plot shows for classes negative, neutral and positive the count of records are 7434, 2510 and 1914 respectively. Since there is an imbalance in data distribution, the data need to be balanced for getting better performance by augmentation technique.

Further exploration of the data shows no missing values or null values present in the data. However, after pre-processing there are chances of populating the columns with null values.

*Figure 9 Histogram showing the length of string*

The above histogram shows that there are strings with less number of characters. This will help in pre-processing steps later.

The below figures show how certain words contribute to each class



*Figure 10 Words Contribution to Neutral class*

It can be seen that words like need, tomorrow, know, help can be linked to neutral sentiment.

Plot representing the most occuring words negative class



*Figure 11 Words Contribution to the negative class*

Customer service, luggage, hold, time, help, delayed never etc are the most frequent words found in the negative class.

Plot representing the most occuring words positive class



*Figure 12 Word contribution to the positive class*

The words like great, thank, love, well, awesome etc are found more in the positive class

# METHODOLOGY

## PRE-PROCESSING

1) Each sentence in the data frame is processed by implementing the following steps.
   - Removing URLs and special characters and numbers.
   - Removing stop words, custom stops words and characters which depend on the nature of the data.
   - Replacing multiple spaces with a single space.
   - Changing the sentence into lowercase and lemmatization of the text data.
   - The data generated may contain some empty string, two or three characters strings in the place of a useful sentence. So it should be removed.
   - Eliminating samples with empty strings, single characters and two characters in addition to null records.

2) DATA IMBALANCE CORRECTION
   The pre-processed data were balanced using two methods, they are BERT-BASED DATA AUGMENTATION ( see appendix) and using Text Attack. Since Text attack generated more useful data, all the models are trained using this method.

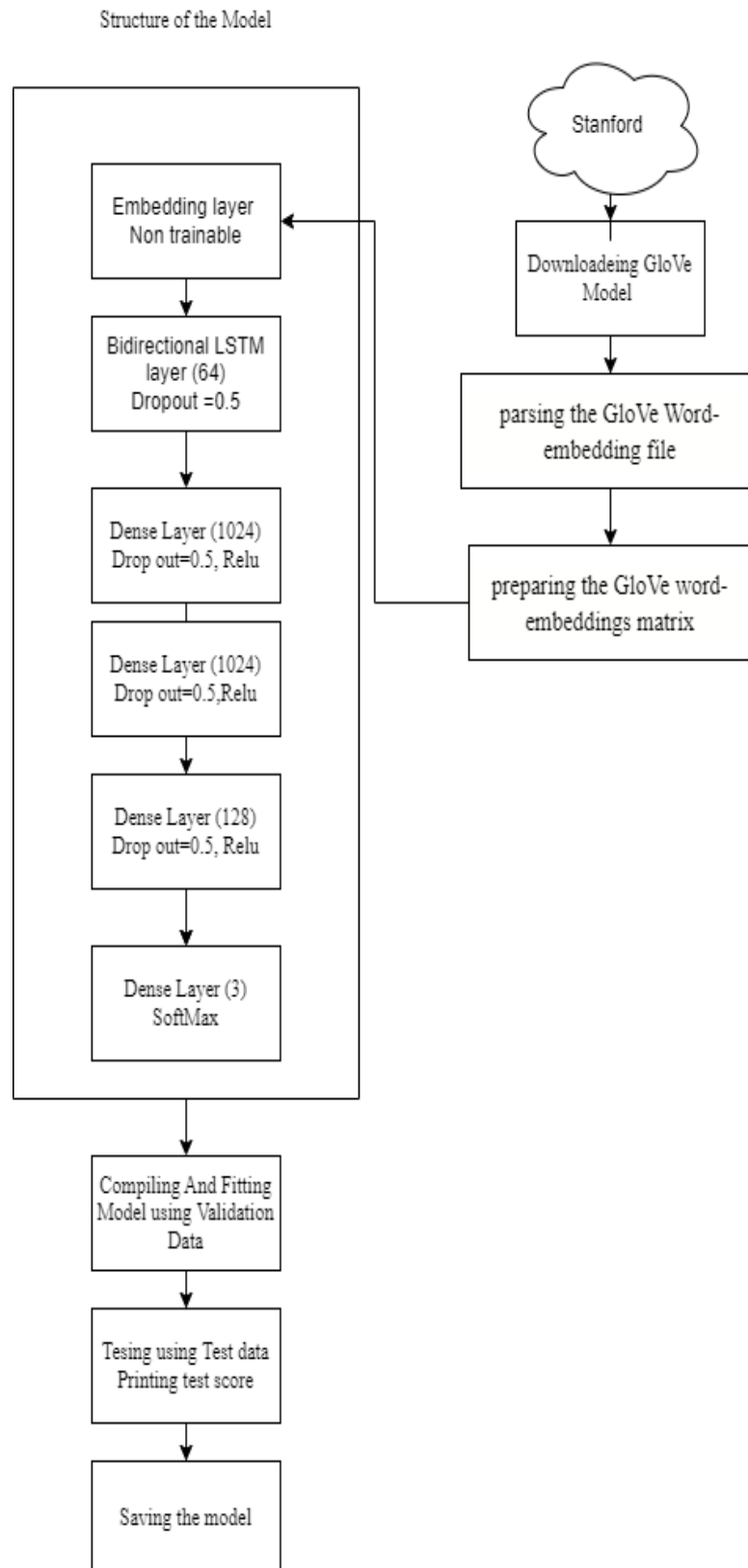   USING TEXT ATTACK AUGMENTOR

   In this method, the new data or sentence for the minority class is created by substituting a   word.

   with its synonym. Apart from this, they are other methods to generate new data using Text attacks.

3) The generated balanced data from these two methods are saved. They were again fed into, imbalance checking. Also, records with null/empty/ records with single characters were removed and saved the two data set (Bert and text attack) were separately for loading in future.

4) The data synthesized by text attack is used for training the models as it showed better accuracy.

5) The class names  'positive', 'neutral' and 'negative' are encoded to 1,0, and 2 respectively.

6) The classes are also encoded using on hot encoding method to align with the neural network model using Tensor flow API.

7) Tokenize the data by selecting the 10,000 most commonly occurring words. After that, vectorize and pad the data generated to a length of 100. These are done using the Keras package.

8) Text data are also encrypted using tf-idf and bag of words methods.

# MODELLING
## Bidirectional LSTM (RNN) Neural Network

The model structure is given below. The glove embedding is used to preserve the semantic relationship between words.



Structure of the Model

The following hyperparameters were used for training.

Optimizer- RMSprop()

Loss Function – categorical_crossentropy

Through different iterations, the best model is trained at the end of $3^{rd}$ epoch with a batch size of 32.

Apart from this model, Random Forest and Naïve bayes were implemented. Finally, an ensemble model is used for the best prediction.

## Multinomial Naïve Bayes Model
- Countvectorizer() method is used which creates a matrix of number of words in each document. Stop words were eliminated and all the words were converted to lowercase.
- The best hyperparameter is selected using GridSearchCV and trained the model.

## K-means Model (unsupervised learning)
- tf-idf method is used for vectorizing the words because this method could represent how important a particular word is in the document.

## Ensemble method (stacked) using Logistic Regression
Since SVM and naïve bayes model were selected for an ensemble model using logistic regression with a solver as 'sag'

A cross-validation fold of 5 is used for all the above models.

In addition to these models, Random forest and support vector classifiers are also implemented.

# RESULT OF ANALYSIS
## Bidirectional LSTM (RNN) Neural Network

Test Results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.33 | 0.38 | 0.35 | 277 |
| Positive | 1 | 0.2 | 0.32 | 0.25 | 213 |
| Negative | 2 | 0.72 | 0.57 | 0.63 | 825 |
| | | | | | |
| | Accuracy | | | 0.49 | 1315 |
| | F1 macro | 0.42 | 0.42 | 0.41 | 1315 |
| | F1 weighted | 0.55 | 0.49 | 0.51 | 1315 |

*Table 8 Summary of the Testing*

Comments: Different models were trained using Simple RNN, LSTM and Bidirectional LSTM methods, and with different hyperparameters like epochs, batch size, dropouts, optimizer and multiple layer combinations. However, the F1 score was not improving. The model described above produced the best results with an F1 score of 0.51.

## Multinomial Naïve Bayes Model

The learning rate of 0.002 was chosen after hyperparameter Tuning.

Test Results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.53 | 0.51 | 0.52 | 277 |
| Positive | 1 | 0.53 | 0.6 | 0.56 | 213 |
| Negative | 2 | 0.83 | 0.82 | 0.83 | 825 |
| | | | | | |
| | Accuracy | | | 0.72 | 1315 |
| | F1 macro | 0.63 | 0.64 | 0.64 | 1315 |
| | F1 weighted | 0.72 | 0.72 | 0.72 | 1315 |

*Table 9 Test summary.*

## Random Forest

The model was trained using Hyperparameters.

criterion : 'entropy'

max_depth: 34

Test Results

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.38 | 0.28 | 0.32 | 277 |
| Positive | 1 | 0.27 | 0.24 | 0.26 | 213 |
| Negative | 2 | 0.72 | 0.8 | 0.76 | 825 |
| | | | | | |
| | Accuracy | | | 0.6 | 1315 |
| | F1 macro | 0.46 | 0.44 | 0.45 | 1315 |
| | F1 weighted | 0.58 | 0.6 | 0.59 | 1315 |

*Table 10 Test summary -Random Forest*

## Support Vector Machine (SVC)

Test Result

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.5 | 0.58 | 0.54 | 277 |
| Positive | 1 | 0.67 | 0.67 | 0.67 | 213 |
| Negative | 2 | 0.86 | 0.82 | 0.84 | 825 |
| | | | | | |
| | accuracy | | | 0.74 | 1315 |
| | F1 macro | 0.68 | 0.69 | 0.68 | 1315 |
| | F1 weighted | 0.75 | 0.74 | 0.75 | 1315 |

*Table 11 SVM Result Summary*

SVM model showed better results using the *rbf kernel and C as 1.0.(gamma= scaled) vectorized using tf-idf*
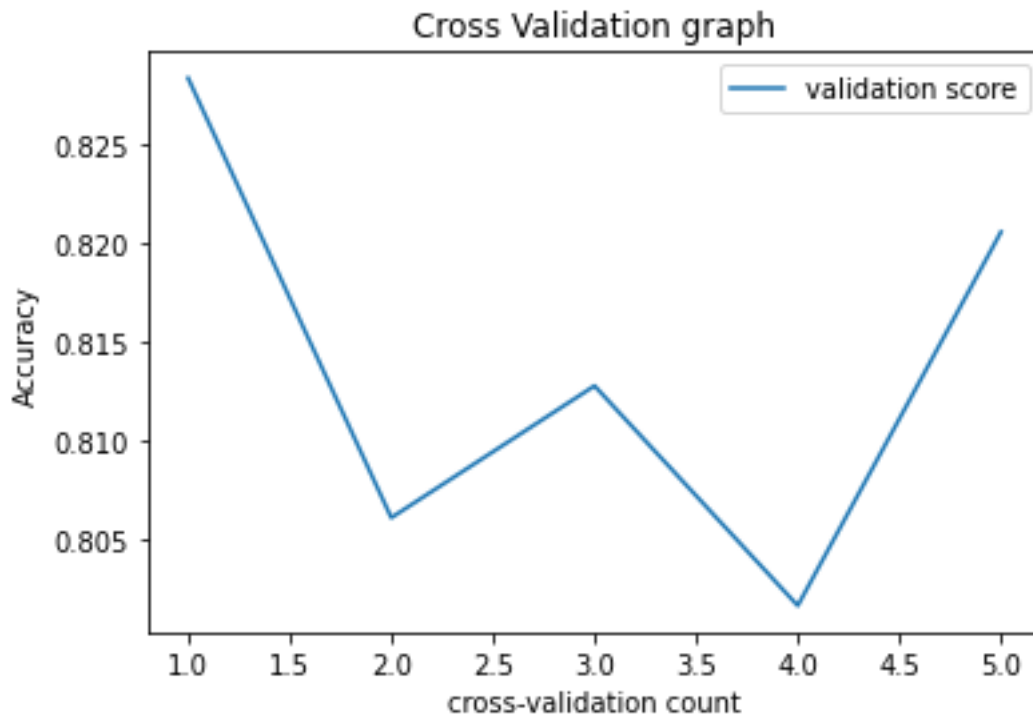
## Ensemble method (stacked) using Logistic Regression

Using the logistic regression model, SVM and multinomial Naïve Bayes Model were combined to generate the ensemble / stacked model.

The results are given below.

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.53 | 0.51 | 0.52 | 277 |
| Positive | 1 | 0.65 | 0.67 | 0.66 | 213 |
| Negative | 2 | 0.84 | 0.85 | 0.85 | 825 |
| | | | | | |
| | accuracy | | | 0.75 | 1315 |
| | F1 macro | 0.68 | 0.68 | 0.68 | 1315 |
| | F1 weighted | 0.75 | 0.75 | 0.75 | 1315 |

*Table 12 Summary of the ensemble model*

Five-fold Cross-validation graph



Comments: The validation score fluctuates between 0.8 and 0.83 while evaluating the model which can be treated as a sign of good performance.

K- means

Test Result

| class labels | | precision | recall | f1-score | support |
|---|---|---|---|---|---|
| Neutral | 0 | 0.1 | 0.03 | 0.05 | 277 |
| Positive | 1 | 0.93 | 0.07 | 0.12 | 213 |
| Negative | 2 | 0.66 | 0.97 | 0.78 | 825 |
| | | | | | |
| | Accuracy | | | 0.62 | 1315 |
| | F1 macro | 0.56 | 0.35 | 0.32 | 1315 |
| | F1 weighted | 0.59 | 0.62 | 0.52 | 1315 |

*Table 13 Test summary of K-means model*

Comments:  The model produced better prediction capability for the negative class alone when compared to other models. It is vectorized using tf-idf  and n_clusters as 3.
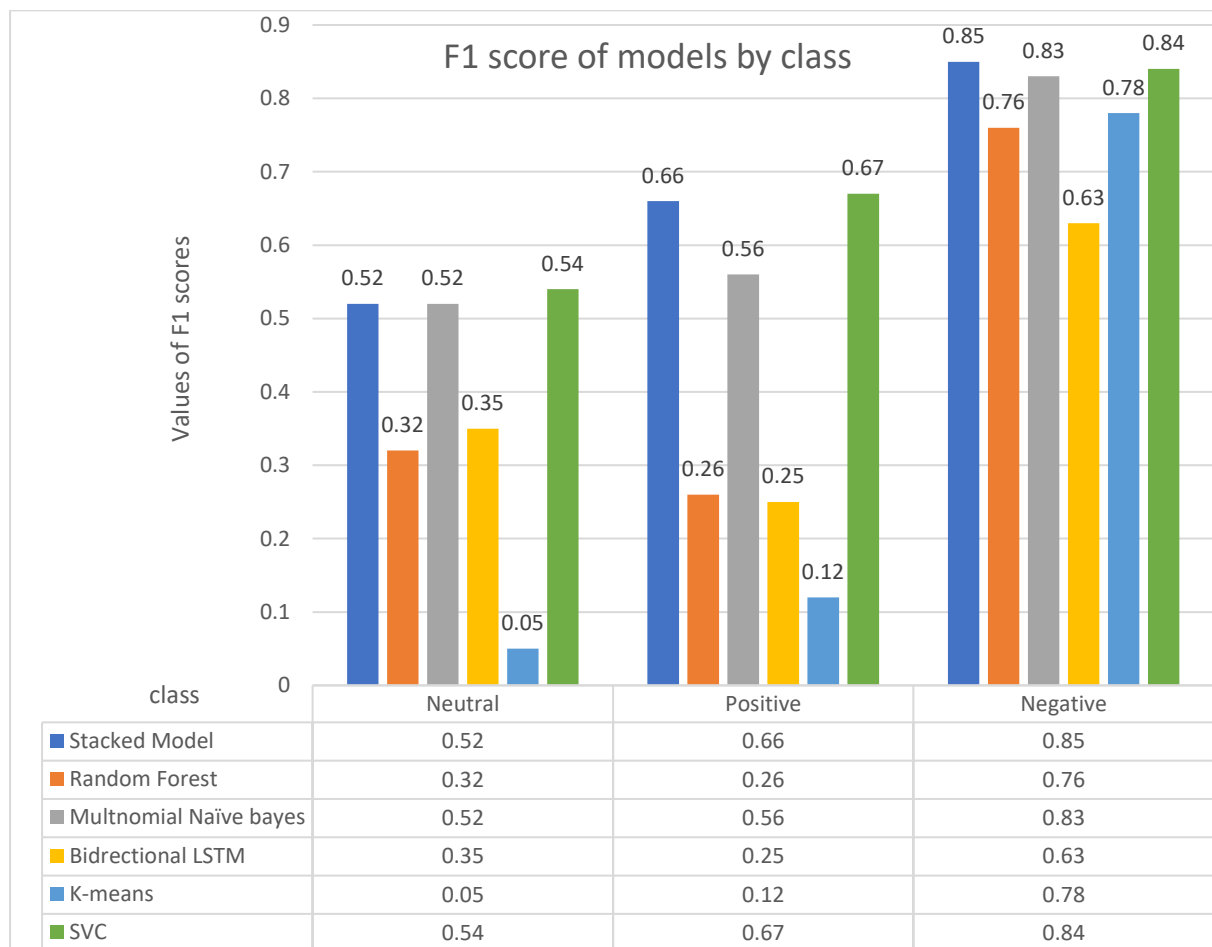
# Discussion of analysis



*Figure 13 combined F1 sores of all the models.*

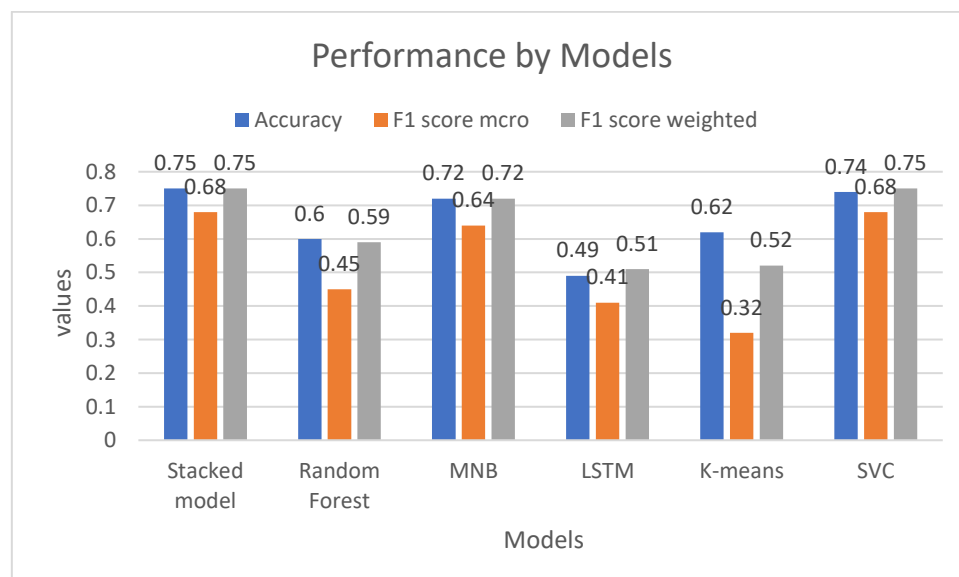| class | Neutral | Positive | Negative |
|---|---|---|---|
| Stacked Model | 0.52 | 0.66 | 0.85 |
| Random Forest | 0.32 | 0.26 | 0.76 |
| Multnomial Naïve bayes | 0.52 | 0.56 | 0.83 |
| Bidrectional LSTM | 0.35 | 0.25 | 0.63 |
| K-means | 0.05 | 0.12 | 0.78 |
| SVC | 0.54 | 0.67 | 0.84 |



*Figure 14 Performance analysis of models*

- The model which performed the best is the ensemble model (SVM, Naïve Bayes) in terms of class-wise f1 score, weighted F1 score (75%) and overall accuracy (75%).
- Since the test data is not balanced, the F1 score should be considered for the model analysis.

- All these four models could predict the negative class with high probability whereas the model struggles to predict other classes comparatively.
- K-means could predict the class with a 0.52 F1 score, which is a significant result among these models. However, the k-means could not predict positive and neutral effectively.

# Conclusion

The models generated for these two data sets (FARS, Tweet) could predict the respective classes with a sufficient F1 score. The best model fitted for the FARS data is the Stacked Model (F1 score- 0.77) combining Random forest, KNN, SVM, decision tree and Random forest. However, the model could not predict the class "Died prior to Accident". With sufficient data, this can be overcome in future.

The ensemble/stacked model is the best fit for the tweet data as it showed a weighted F1 score of 0.75. Also, the unsupervised model K-mean produced the least result (F1 score of 0.52) at the same time the model could predict the negative class with high accuracy. The model selection, training and performance significantly depend upon the quality of the data, data augmentation, pre-processing, hyperparameter combination etc. Hence, by the trial-error method, the performance can be improved in future works.

# Reference

1) Mckinney, W. (2017). Python for Data Analysis. 2nd ed. O'Reilly.
2) scikit-learn (2019). scikit-learn: machine learning in Python. [online] Scikit-learn.org. Available at: https://scikit-learn.org/stable/.
3) CholletF. (2018). Deep learning with Python. Shelter Island, NY: Manning Publications.

# Appendix

SMOTE Algorithm (Synthetic Minority Over-sampling Technique)

Based on the below process, the algorithm creates synthetic samples of minority classes.

i) Based on the below process, the algorithm creates synthetic samples of minority classes
ii) Identify a minority group to sample.
iii) Find its k-nearest minority neighbours.
iv) The nearest neighbour of k is randomly selected.
v) The difference between the feature values of the selected sample and its nearest neighbour is measured.
vi) A random number between 0 and 1 is multiplied by the difference obtained.
vii) Synthetic samples are created by adding the results to the feature values.

The data generated is similar but not identical. Hence, this helps in reducing the overfitting of the model.

Normalizing formula

$$\frac{x - \min(X)}{\min(X) - \max(X)}$$

X: sample value under a particular feature.

Min(X): Minimum value of the sample.

Max(X): Maximum value of the sample.

This could make the data contribute equally while processing it.

RFE (Recursive Feature Elimination)

RFE is a wrapper-style feature selection method. This indicates that a distinct machine learning algorithm is provided and employed in the method's core, is wrapped with RFE, and is used to aid in feature selection. In contrast, filter-based feature choices assign a score to each feature and choose the ones with the highest (or lowest) value.

## HYPERPARAMETER TUNING

Hyperparameter defines the model architecture, and they are not optimized during training instead they were added while defining the model. The performance of the model depends on the values chosen for the hyperparameter and this process of model calibration is called Hyperparameter tuning. In this report, two modules were used. They are GridSearchCV and BayesSearchCV.

### GridSearchCV

It is included in the scikit-learn class used for tuning the model hyperparameter. It takes the model and relevant values of the hyperparameter and fits the model by going through all the given combinations of parameter values. The best hyperparameter is shown based on accuracy. This method suffers from taking a long time (brute force method) and cannot consider the relationship between the hyperparameters.

### BayesSearchCV

Bayesian optimization builds a probabilistic model of the function that maps hyperparameters to model performance to find the optimal hyperparameter. It is taken from the package called sci-kit-optimize. It saves comparatively much time compared to other algorithms as it considers past knowledge for optimization.

### lemmatization

It is the process of reducing the word to its base form.

### Bert-Based Data augmentation

In this method, the minority class is balanced by generating new data for that class with the help of the Bert Model (Transformer based) model. The synthetic data is similar to the data available in the minority class. This method was implemented, but the accuracy /F1 score was not satisfactory. So, the next method given below is used.