

DEATH IN AMERICA

Erim Hayretci
Vasco de Noronha
Deniz Can Dinkci
Andrew Ize-Iyamu

Investigating Death in America

Erim Hayretci, Vasco de Noronha, Deniz Can Dinkci and Andrew Ize-Iyamu

Imperial College London, England

This report was compiled on June 24, 2024

Abstract

This project aims to investigate mortality causes in the United States through a diverse set of machine learning models applied to a data frame that includes 2.8 million recorded deaths from 2015.

Our team members analysed four main causes of death: alcohol and drug overdose, homicide, suicide, and motor accidents. For alcohol and drug overdose, a Support Vector Machine (SVM) model with an RBF kernel achieved a test accuracy of 77.8%, precision of 73.1%, and recall of 81.0%. Homicide predictions used a Logistic Regression model, yielding a validation accuracy of 82.8%, precision of 83.5%, and recall of 81.6%. Suicide risk was assessed using a Decision Tree model, resulting in a test accuracy of 75.6%, precision of 75.7%, and recall of 75.2%. Motor accidents were predicted using a Random Forest model with an accuracy of 78.4%, precision of 75.8%, and recall of 83.4%

Considering the independence and randomness of the death types examined in this study, these models prove that with correct data-driven approaches, the majority of fatalities may be precluded by implementing correct strategies for the targeted risk groups.

Keywords: Alcohol, Drug Overdose, Homicide, Suicide, Motor Accidents

Contents

1	Introduction	1
2	Methodology	2
3	Alcohol and Drug Overdose - Erim Hayretci	2
3.1	Aim	2
3.2	Introduction	2
3.3	Methodology	2
	Filtering, Balancing & Splitting Data • SVM Algorithm • Defining Hyperparameters • Optimisation	
3.4	Results	4
3.5	Discussion	4
	Validation • Suggested Improvements	
4	Homicide - Vasco de Noronha	4
4.1	Methodology	4
	Filtering the Data • Balancing the Data • Splitting the Data	
4.2	Results	4
	Feature Selection • Validation • Final Results	
4.3	Discussion	5
5	Suicide - Deniz Can Dinkci	5
5.1	Aim	5
5.2	Methodology	5
	Data Filtering and Pre-Processing • Splitting the Data • Hyperparameter Optimisation	
5.3	Results	7
	The Decision Tree • Final Results	
5.4	Discussion	7
6	Motor Accidents - Andrew Ize-Iyamu	7
6.1	Aim	7
6.2	Methodology	7
	Model Selection • Feature Engineering	
6.3	Hyperparameter Tuning	8
	Prioritising Key Performance Metric • Hyperparameters • Model 1: Benchmarking • Model 2: Recall against Maximum Depth • Model 3: Optimising the Hyperparameters	
6.4	Results	9
	Final Model • Analysis	
6.5	Discussion	9
	Feature Importance • Limitations of the Predictive Model	

7	Conclusion	10
8	Appendix	10
8.1	Preprocessing Dataset	10
	Importing PANDAS library and filtering CSV file • Converting data labels from strings to integers • Removing all rows that contain NaN values in the dataset • Binarization of dataset • Reimporting the binarized dataset • Binarizing the cause of death column for motor accidents • Combining the education columns of the data	
8.2	Alcohol and Drug Overdose - Erim Hayretci	12
8.3	Homicide - Vasco de Noronha	14
	Filtering the dataset • Balancing the dataset • Splitting the dataset • Forward Selection algorithm • Validation	
8.4	Suicide - Deniz Can Dinkci	18
	Creating the CSV as required • Splitting the variables • Accuracy vs. Maximum Depth • Precision vs. Maximum Depth • Recall vs. Maximum Depth • Accuracy vs. Minimum Impurity Decrease • Precision vs. Minimum Impurity Decrease • Recall vs. Minimum Impurity Decrease	
8.5	Motor Accidents - Andrew Ize-Iyamu	20
	Balancing and exporting the dataset • Importing binarised data • Standardising and splitting the dataset • Model 1: Benchmarking • Model 2: Plotting Recall of Model 2 against Maximum Depth • Plotting Estimator 4 Decision Tree for Model 2 • Plotting Random Forest for Model 2 • Using Random Search Cross Validation to find optimal combination of hyperparameters • Final Model • Plotting Estimator 2 Decision Tree for Final Model • Plotting Random Forest for Final Model • Feature Importance Plot	

1. Introduction

Death in the USA is a significant and complex issue, with various causes contributing to the overall mortality rate. In 2020, there were approximately **3.4 million deaths** in the USA, with heart disease, cancer, and COVID-19 being the leading causes [1]. The preventability of many of these deaths highlights the importance of early detection, effective treatment, and lifestyle modifications. However, the accessibility and cost of healthcare services remain substantial barriers. In 2019, Americans spent an average of **\$11,072 per person** on healthcare, making it the most expensive healthcare system in the world [2].

This paper aims to use different models and techniques to prevent death in America via successful diagnoses and preventative measures. By leveraging data analysis and machine learning, we can identify high-risk individuals and provide timely interventions, potentially saving countless lives and reducing healthcare costs.

2. Methodology

The dataset for this analysis was obtained from the Centres for Disease Control and Prevention (CDC), which annually releases comprehensive data detailing every death in the USA. The dataset includes 2.8 million samples and 77 attributes, providing an extensive resource for analysing mortality in America. It is important to note that this dataset is specific to the United States and cannot be generalised to other countries. The primary objective of the models developed in this report was to identify the significant variables associated with each cause of death, such as age and gender. Each of the four research questions addressed different causes of death, employing distinct modelling techniques. Consequently, much of the data filtering was conducted individually, tailored to each model's requirements.

However, certain columns were deemed irrelevant to all our research questions as a group and were consequently removed. These columns included Hispanic recode, race bridge, injury at work, autopsy, burial, and activity. This collaborative filtering ensured that our models focused on the most pertinent variables, enhancing the accuracy and relevance of our findings.

3. Alcohol and Drug Overdose - Erim Hayretci

3.1. Aim

Aim of this study is to predict the likelihood of male death caused by overdose on alcohol and psychoactive drugs, based on our dataset covering 2.8 million recorded deaths in the United States for year 2015.

3.2. Introduction

Global trends indicate a concerning rise in alcohol and drug overconsumption across various demographics. According to NIH, combined deaths of overdose on alcohol and drugs have escalated from 62.30 to 97.72 per 100,000 from 2000 to 2009 [3]. The gender-specific effects are also notable. The WHO estimates that of the 67,000 deaths from alcohol poisoning worldwide, men outnumbered women by a ratio of 4 to 1.[4] Regarding drug-related cases, data suggests that in 2015, men died of overdose at twice the rate of women, which increased to 3.2 by 2021[5]. These trends that are influenced by sociocultural, behavioural, and biological factors that differentiate how genders interact with substances, indicate that focusing a machine learning algorithm on particularly male individuals may help uncover subtle patterns and risk factors, ensuring the identification of at-risk individuals and providing targeted precautions and support.

3.3. Methodology

3.3.1. Filtering, Balancing & Splitting Data

The dataset was initially filtered, revealing an all-male subset that contained 1,308,641 data points. Subsequently the set was then split into two groups based on whether individuals died from substance overconsumption, or other reasons. As expected, the dataset was extremely imbalanced which was confirmed via a confusion matrix, revealing 7,868 positive cases against 1,300,773 negative cases. To balance the data for training our algorithm, undersampling has been performed using the stratified sampling method. The major dataset (negative class) was divided into 18 age-differentiated strata. Random sampling was then conducted within each stratum, proportional to their representation in the minority group. This approach ensures the exclusion of underage individuals and prevents age-based sampling bias, leading to a more representative dataset.

Given the abundance of data points for training the model effectively, the final data is split into 60% for training and 40% for evaluation. While 9,441 data points are used to train the model by fitting its parameters, the remaining is divided into equal parts for validation and testing. The validation set is used to fine-tune the model by predicting responses for its observations, while the test set provides

an unbiased evaluation of the model's final performance based its optimised parameters.

3.3.2. SVM Algorithm

To achieve our aim, a Support Vector Machine (SVM) algorithm is preferred, that uses our final data to find a hyperplane in an N-dimensional space that classifies tested data points. A kernelised model is employed to handle non-linear relationships within the data, enabling more precise differentiation between unbiased positive and negative cases, which, in our scenario, represent deaths caused by overdoses and other reasons.

3.3.3. Defining Hyperparameters

Essential hyperparameters for the kernelised SVM model include:

- **Kernel:** Transforms the input data into a higher-dimensional space with a set of reduction rules to make it easier to find a hyperplane that bounds the data points effectively. Kernel settings are expressed as mathematical function types.
- **Gamma γ :** Determines how far the influence of a single data point reaches. It affects the shape of the decision boundary where a high gamma value means that each point has a narrow area of influence, making the boundary tightly surround the points which can cause overfitting. Usually set between $1/n$ to $6/n$ where n represents the number of input fields.
- **C (Control):** Controls the trade-off between achieving a low error on the training data and controlling regularization to avoid overfitting. Covers all positive values starting with 0.

3.3.4. Optimisation

Kernel Selection



Figure 1. Accuracy of the Training and Validation Sets

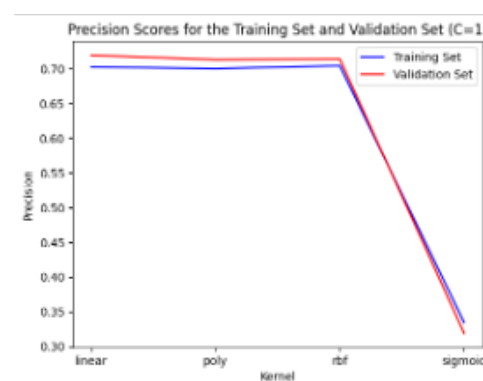


Figure 2. Precision of the Training and Validation Sets

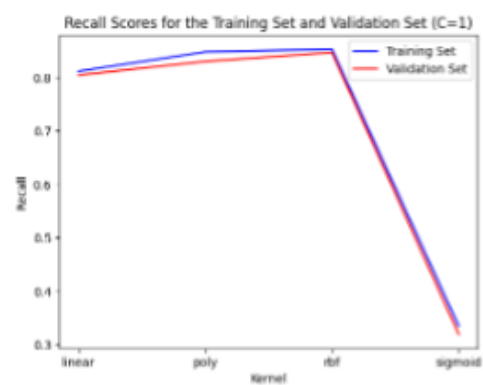


Figure 3. Recall of the Training and Validation Sets

Apart from the sigmoid function, all three kernel types provide similar decision boundaries regarding the accuracy, precision, and recall scores of the model. However, through further investigation, RBF (Radial Basis Function) is found to offer the best performance with 75.6% accuracy and 71.3% precision in detecting correct individuals dying from substance overdose with a recall score of 85.3%. This is noted to be the default setting of the SVM algorithm.

C Value Selection

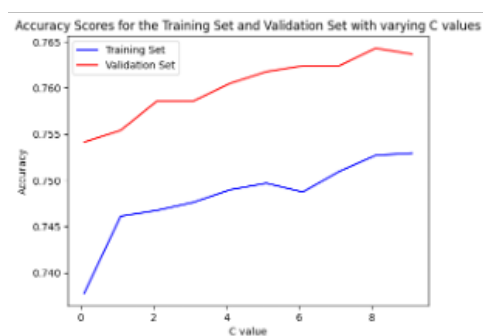


Figure 4. Accuracy of the Training and Validation Sets with varying C values



Figure 5. Precision of the Training and Validation Sets with varying C values



Figure 6. Recall of the Training and Validation Sets with varying C values

Due to the computational constraints of the device used for modelling, the continuous control variable, C, is tested at 10 evenly spaced points between 0 and 10. It is generally observed that increasing values of C lead to lower error margins. However, as large values of control can cause overfitting by reducing regularization, a C value of 3.0 is selected for the testing model.

Gamma Value Selection

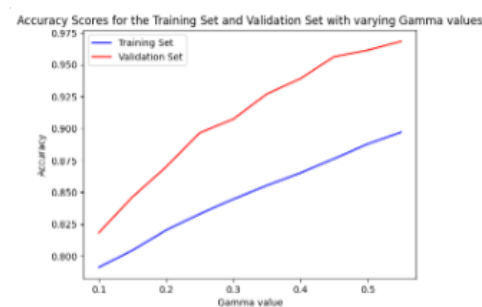


Figure 7. Accuracy of the Training and Validation Sets with varying Gamma values

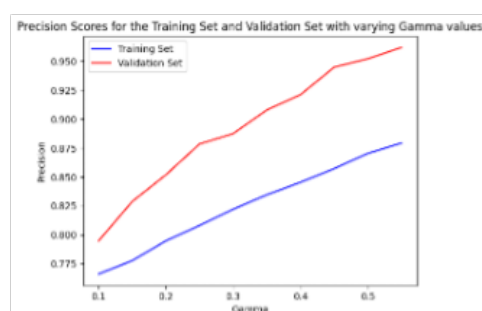


Figure 8. Precision of the Training and Validation Sets with varying Gamma values



Figure 9. Recall of the Training and Validation Sets with varying Gamma values

The default gamma value for the RBF kernel function is set to 'scaled,' which is 0.1 in our case. Testing gamma values between $1/n$ and $6/n$ showed a positive correlation with performance metrics. However, despite the positive trend, the high number of data points justifies using a simpler model with a value of 0.1 to avoid overfitting, where each point would excessively influence the model.

3.4. Results

Table 1. SVM algorithm results with selected hyperparameters

	Training Set	Validation Set	Testing Set
Accuracy	0.746	0.756	0.778
Precision	0.704	0.713	0.731
Recall	0.853	0.847	0.810

3.5. Discussion

3.5.1. Validation

The optimisation of the hyperparameters in our SVM model is observed to significantly improve its performance. By tuning the C and gamma variables and selecting the RBF kernel structure, we achieved a testing accuracy of 77.8% and precision of 73.1% which are both 2-3% higher than the scores of training and validation sets. This indicates that the model has effectively learned the underlying patterns in the data, despite the inherent randomness of predicting deaths due to overdosing.

Precision stands out as the lowest score amongst other, which refers to the risk of approaching males that do not carry the likelihood of overdosing, as at-risk individuals. The high recall score (81%), on the other hand, ensures that majority of risky cases could be accurately detected. In our scenario this trade-off is favourable, since offering a more widespread set of precautions is better than overlooking certain parts of the population which is likely to have fatal results.

3.5.2. Suggested Improvements

Future work on this model must implement a more extensive parameter optimisation process with higher computational power, testing a wider range of C and gamma values to further optimize the model's performance. Additionally, the population sample's coverage must be extended, including cases from other countries, for a more comprehensive dataset. This expansion would account for different policies and cultural attitudes towards alcohol and drug consumption, potentially improving model's generalizability across different demographics.

4. Homicide - Vasco de Noronha

4.1. Methodology

In 2015, **18,882 people** in America died from homicide-related causes. To address this, a Logistic Regression Model was developed, aiming

to find associations between independent variables and the binary dependent variable of homicide-related death. This model's goal is to help mitigate future homicides in the USA by identifying the most significant predictors and informing laws and policies. Forward selection was employed to create a model with high sensitivity, identifying the features that most strongly correlate with homicide-related deaths.

4.1.1. Filtering the Data

America has one of the **highest homicide rates among affluent nations**, with approximately 5 homicides per 100,000 people, a rate stable since 2000, according to a paper from the NCJ [6]. To analyse the factors affecting this, we processed the data by sorting the manner of death column, coding homicides as 1 and other deaths as 0. We further filtered the data by removing unnecessary columns, such as autopsy status and burial information. Education index values from pre- and post-1983 were combined, and rows with any missing values or unknown variables were removed. This resulted in a dataset of 2,147,246 data points for the forward selection algorithm.

4.1.2. Balancing the Data

Due to the computational power of my device, we decided to randomly sample 20,000 deaths from the 2.15 million in our filtered dataset. We performed a confusion matrix on the dataset to test for imbalance. When examining the manner of death, specifically homicides, the dataset was highly imbalanced, with significantly more non-homicide cases (19,828) than homicide cases (172). This imbalance is likely due to the high proportion of other related deaths compared to homicides. As a result, the dataset would almost always predict a non-homicide death. To better understand the variables that influence homicide-related deaths, we needed to balance the data. Therefore, a new dataset was created, containing 10,000 homicide deaths and 10,000 non-homicide-related deaths. Despite the large size of the dataset, the effect of overfitting is likely to be minimal even though we are under-sampling.

4.1.3. Splitting the Data

The filtered and balanced data was randomly split into 50% for training, 25% for validation, and 25% for testing. Additionally, the data was split in 80/10/10, 70/15/15, and 60/20/20 proportions, as the split ratio affects the model's final accuracy. These various splits have been further analysed in the validation section 4.2.2.

4.2. Results

4.2.1. Feature Selection

A forward selection algorithm was implemented to identify the features with the highest predictive power. Initially, a model with no features was created. The feature with the largest correlation to homicide-related deaths was then added to the model. This process was repeated until no more variables remained or until the next variable changed the model's validation accuracy by less than 0.005. The final features were then saved and used to calculate their accuracy on the test set to prevent overfitting the model.

4.2.2. Validation

Many different data splits were used to find the most appropriate one. Figures 10–13 illustrate the relationship between average model accuracy and model complexity (number of variables). A common trend across all figures is that the validation data has a slightly higher mean accuracy than the training data. While this could be due to random variations, it is unlikely given the significant amount of data points, which results in a very accurate model. Figure 11, representing the 60/20/20 split, fits the model best up to a model complexity of 5, where the validation and training data accuracies begin to diverge slightly. Therefore, the 60/20/20 split was used during the logistic regression forward selection to select key features and create a model.

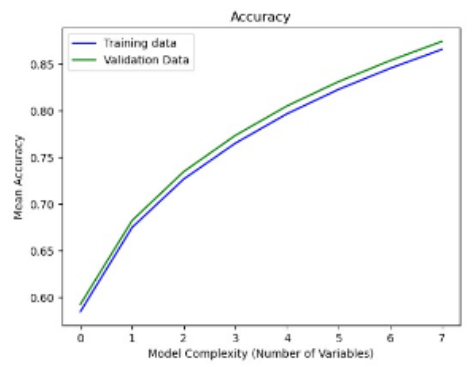


Figure 10. 50/25/25 split Model Complexity vs Model Accuracy

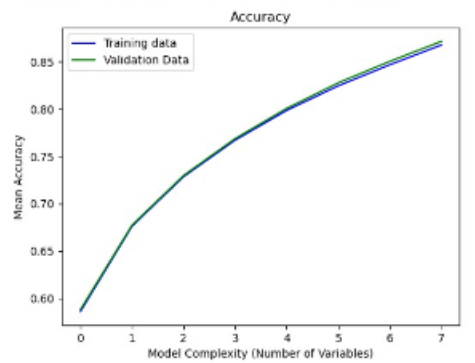


Figure 11. 60/20/20 split Model Complexity vs Model Accuracy

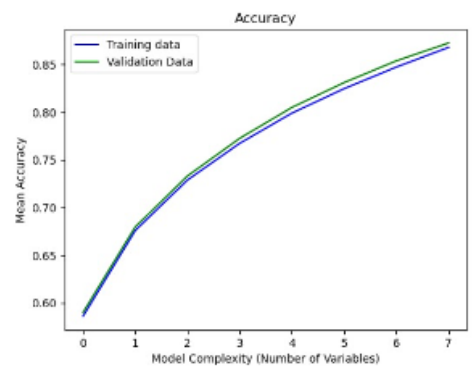


Figure 12. 70/15/15 split Model Complexity vs Model Accuracy

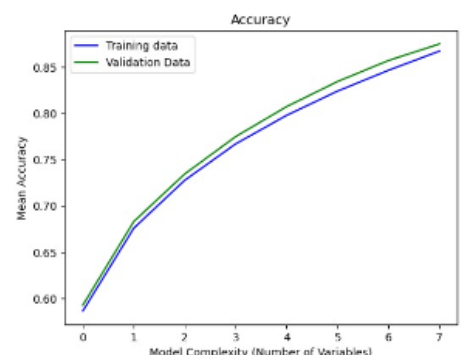


Figure 13. 80/10/10 split Model Complexity vs Model Accuracy

4.2.3. Final Results

Table 2. Results of Logistic Regression Forward Selection Model

	Accuracy	Precision	Recall
4 Training	0.798	0.807	0.788
4 Validation	0.801	0.812	0.785
5 Train	0.824	0.829	0.817
5 Validation	0.828	0.835	0.816

4.3. Discussion

The model identifies the features that most affected homicide-related deaths in the USA in 2015. These features, in order, are: Age, Race, Resident Status, Education, Marital Status, Day of the Week of Death, and Month of Death. Interestingly, the final model included all these features except the month of death, which was surprising, as it was assumed that the season would significantly affect homicide rates.

As the selected features were added, the validation accuracy increased from 0.869 to 0.878. The model's accuracy peaked at a model complexity of 4 with an accuracy of 0.867. The accuracies don't diverge greatly at complexities greater than 4, suggesting minimal overfitting in this model. This aligns with the graph in Figure 11.

The model should have high sensitivity, as it is better to overpredict the likelihood of homicide-related deaths. In the USA, a homicide has an estimated cost of **\$9 million** per case [7], making it more cost-effective to overpredict and target individuals at risk of homicide rather than underpredict and miss opportunities for prevention. As mentioned, the model strongly fits the validation data up to a complexity of 4 or 5, so the results for accuracy, precision, and recall were plotted in Table 2.

For both the validation and training sets, the model values for accuracy, precision, and recall are high. However, at a complexity of 5, the model's accuracy starts to diverge slightly from the validation dataset. Therefore, a model complexity of 4 was selected to minimize error. The model achieved a relatively high recall of 0.788, resulting in high sensitivity.

To conclude, the model produced is satisfactory. However, logistic regression can still suffer from issues such as linearity and other limitations. Therefore, it is suggested to use a different model for further research to better understand and mitigate homicides in the USA.

5. Suicide - Deniz Can Dinkci

5.1. Aim

The goal of this study is to predict the probability of an individual committing suicide by investigated their background and demographics, based on an inclusive data having all deaths in the USA in 2015.

5.2. Methodology

Decision trees are non-parametric, tree-like flowchart models that are utilized for classification and regression. Each internal node represents a feature or attribute, each branch denotes a decision rule related on that attribute, and each leaf node indicates an outcome or prediction. In the case of predicting suicides, it will be used as a classification. The structurization of the nodes is done to maximize information gain. Furthermore, as the tree splits down further, the Gini impurity decreases which ensures that the classification become progressively more reliable.

5.2.1. Data Filtering and Pre-Processing

To use a decision tree, all individual parameters are first binarized. The "manner of death" column is also binarized, assigning a value of 1 for suicides and 0 for other types of deaths. Next, unnecessary columns, such as the "method of disposition," are removed from the

dataset. Due to the significant imbalance in the data, with only 43,029 suicides compared to over 2.7 million other deaths, the dataset is undersampled to create a more balanced dataset. The dataset comprises exactly 20000 records, which 10000 of people who committed suicide and 10000 people who died from other reasons is selected to prevent overfitting.

5.2.2. Splitting the Data

The filtered data was first divided into a training set and a secondary set, with 70% allocated to the training set and 30% to the secondary set. The secondary set was then split evenly, creating a test set and a validation set, each containing half of the data from the secondary set.

This 70/30 split was chosen over the 80/20 split because of the considerably large dataset, allowing a larger portion to be allocated to validation and testing. A larger validation and testing set reduces overfitting and results in more reliable metrics.

5.2.3. Hyperparameter Optimisation

To ensure the model did not overfit, graphs displaying accuracy, precision, and recall were plotted against maximum depth and minimum impurity decrease.

Maximum Depth

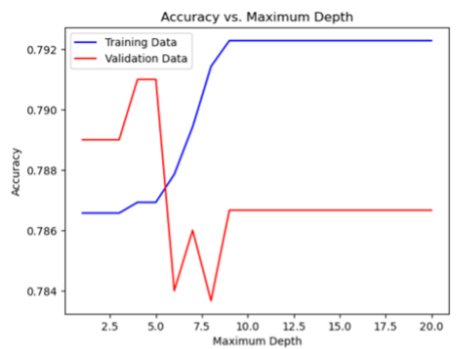


Figure 14

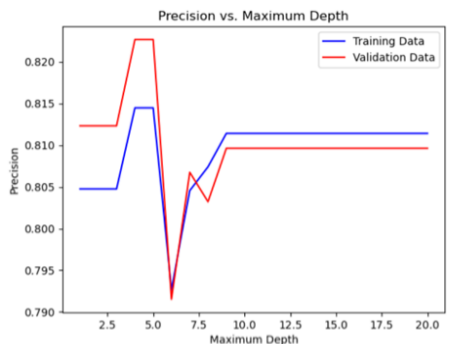


Figure 15

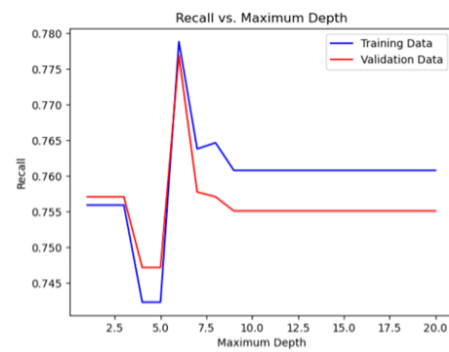


Figure 16

Figure 14 shows that the validation data peaks between maximum depths of 3 and 4, with a decline afterwards and overfitting increases after this point. Figure 15 indicates a spike in precision around the same depth values, making a depth of 4 a logical choice to pick.

Minimum Impurity Decrease

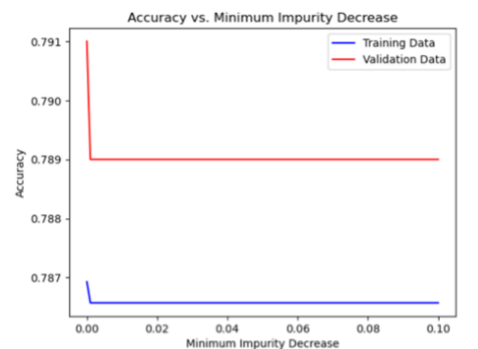


Figure 17

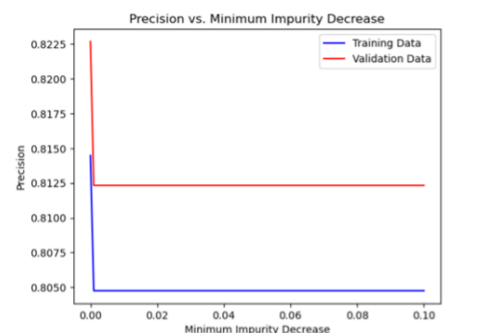


Figure 18

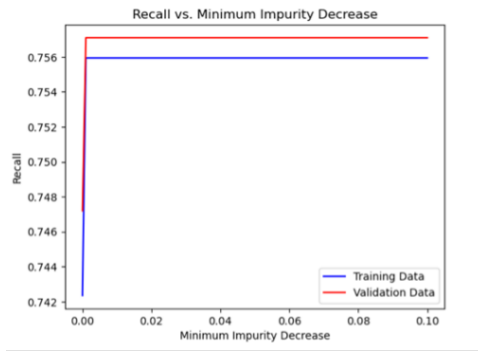


Figure 19

Minimum impurity decrease regulates the creation of new leaf nodes in a decision tree. The accuracy is constant across different values, peaking between 0 and 0.001. Although precision is lowest from 0.001 to 0.1, this comes at the cost of a high recall value. So, choosing the value as 0.001 is logical.

5.3. Results

5.3.1. The Decision Tree

The decision tree was constructed using the chosen hyperparameters: a minimum impurity decrease of 0.001 and a maximum depth of 4. The Gini coefficient on the leaf nodes ranged from 0 to 0.5, but none of the end leaves had a value of 0, indicating that there is a chance of misclassification for new data at every branch.

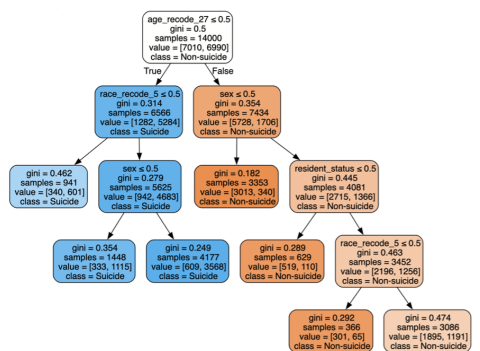


Figure 20. Decision Tree

Each of the condition statements in the nodes of the decision tree is based on whether a condition is less than or equal to 0.5. This is because the data has been binarized. If the data point is less than 0.5, the person does not exhibit that characteristic. If it is greater, they do.

5.3.2. Final Results

Table 3. Metrics of the Performance

	Training Set	Validation Set	Testing Set
Accuracy	0.787	0.805	0.756
Precision	0.789	0.812	0.757
Recall	0.789	0.812	0.752

Table 3 summarizes the performance metrics since all values in the table exceeds 70%, it indicates a great and realistic performance. The small differences between the testing dataset and validation and training datasets suggest that the model is representative of multiple datasets. Although the testing dataset has slightly lower scores, it indicates that the dataset does not overfit.

5.4. Discussion

The decision tree analysis reveals several factors that significantly influence the likelihood of suicide. Key factors identified include race recode, sex, and resident status. Each of these variables has been binarized, with the decision points in the tree indicating whether a condition is less than or equal to 0.5.

The model's decision process starts with the race recode variable. If the race recode is less than or equal to 0.5, the model proceeds to evaluate sex. If sex is less than or equal to 0.5, it further evaluates based on resident status. This hierarchical decision-making indicates the relative importance of these factors in predicting suicide.

The Gini impurity values give important insights into the model's predictions capability. The lower Gini values are higher purity and less chance of misclassification at that node occurs. For instance, a Gini impurity of 0.249 for $\text{sex} \leq 0.5$ and $\text{samples} = 4177$ shows a relatively high certainty in predicting suicide for this subgroup. Meaning sex is the leading candidate of suicide and males as can be seen from the decision tree are suiciding more relative to women. Conversely, higher Gini values, such as 0.474 for $\text{race_recode} \leq 0.5$ and $\text{samples} = 3068$, indicate a higher chance of misclassification, suggesting less reliability in these predictions.

Though the model performed well, there are limitations to using a decision tree. Firstly, there is no guarantee that the sequence of locally optimal decisions will lead to a global optimum. Secondly, decision trees are unstable; small changes in the data split can significantly impact the resulting model. This instability can be mitigated by using the random forest method. While random forests have a higher computational cost, this would not be an issue for this relatively small dataset.

To improve the dataset, a bigger data set can be used, even though we had the possibility to use 80000 persons, we decided to go with 20000 because of the computational power of the computer we use. Secondly, this data set is taken from 2015 in USA meaning it might not be correct to implement the model to other countries with different demographics.

6. Motor Accidents - Andrew Ize-Iyamu

6.1. Aim

The model aims to determine the likelihood of motor accidents being the cause of death for individuals in the US. With crash deaths rising by 20% from 2011 to 2021 and costing the US \$340 billion in 2022 [8], this research identifies key factors influencing motor accidents. By informing government bodies like the IIHS with data-driven evidence, we aim to enhance road safety regulations, reduce annual fatalities, and alleviate the socioeconomic burden of road-related accidents.

6.2. Methodology

6.2.1. Model Selection

Given the non-linear nature of predicting motor accident fatalities, a classification ML algorithm seemed the most appropriate. The available options included Decision Tree, Random Forests, and SVM. The Random Forests method was selected due to its ensemble approach, which reduces overfitting — a common issue with single decision trees. Despite SVMs' sensitivity to errors and outliers, Random Forests offer easier hyperparameter tuning and computational efficiency, making it a more cost-effective and reliable method for government bodies to generate predictions.

6.2.2. Feature Engineering

The "manner of death" column was filtered out from the updated Kaggle dataset due to its generality, while the "cause of death" column provided more detailed information. Two recoded education columns were combined to address NaN values, as Random Forests cannot handle them. Despite Random Forests' ability to manage discrete variables, remaining features were binarized to simplify decision

boundaries, enhance accuracy, and reduce overfitting and computational costs. Post-binarization, there were 36,233 motor accident fatalities and 2,481,907 other fatalities, indicating a significant class imbalance. To counteract this, we under-sampled the majority class, thus 36,233 samples were randomly selected from the majority. Then the processed dataset was standardized then split into training (80%), validation (10%), and test sets (10%) to avoid overfitting and ensure optimal model performance.

6.3. Hyperparameter Tuning

6.3.1. Prioritising Key Performance Metric

Keeping track of the ML model's performance is important, as we need to ensure that the model can provide stable and reliable predictions on whether the cause of death is related to a motor accident. Our performance metrics include:

- **Accuracy:** measures the overall correctness of the model.
- **Precision:** evaluates the correctness of positive predictions.
- **Recall:** measures the model's ability to identify all relevant instances.

In our context, the model must highlight key features impacting the likelihood of death via motoring incidents. Minimizing false negatives is crucial, as failing to identify high-risk individuals may lead to inadequate resource allocation by the US government. Consequently, recall is our primary performance metric for benchmarking and comparing the model during hyperparameter tuning. While false positives are less critical because they may still enhance safety, precision remains important. Excessive false positives could strain the government's budget by spreading resources too thin, thus reducing policy effectiveness. Accuracy is unsuitable for this classification model as it doesn't differentiate between type I and type II errors, failing to account for the skewed costs heavily favouring type I errors.

6.3.2. Hyperparameters

For the Random Forest ML method, the following parameters were available to modify:

- ***N Estimators:*** Number of decision trees in the random forest.
- ***Maximum Depth:*** Maximum depth of each decision tree.
- ***Minimum Samples Leaf:*** Minimum number of samples required to be a leaf node
- ***Minimum Samples Split:*** Minimum number of samples required to split an internal node.
- ***Maximum Leaf Nodes:*** Maximum number of leaf nodes in each tree.

6.3.3. Model 1: Benchmarking

Table 4. Model 1 Performance Metrics against the Training and Validation Set

Performance Metric (%)	Training Set	Validation Set
Accuracy	78.5	77.5
Precision	77.4	75.9
Recall	81.0	79.6

- ***Max Depth:*** Number of decision trees in the random forest.
- ***Minimum Samples Split:*** 0.1 (i.e. 10% of the samples).

In the table above, we can see that for the validation set, the model's recall (79.6%) is relatively higher than the model's precision (75.9%), and only marginally higher than the model's mean accuracy (77.5%). The purpose of the first model is to act as a benchmark for any future iterations of the ensemble ML method, so that we can clearly see any improvements in the model's performance post-revision of hyperparameters.

6.3.4. Model 2: Recall against Maximum Depth

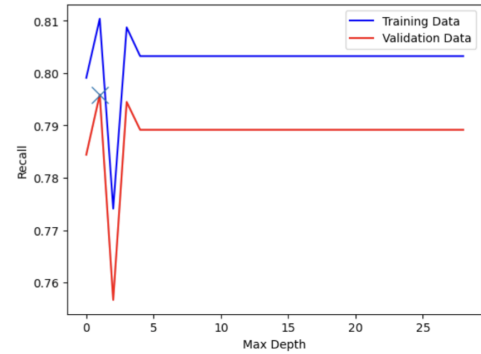


Figure 21. Recall of Model 2 against Maximum Depth

As seen in the figure above, the maximum recall for the model occurred when the maximum depth was 1. However, this would lead to underfitting within our dataset, because only variable is used by the model to make a prediction on whether the cause of death of an individual a result of a motor accident is. Thus, in exchange for marginally lower recall of 75.5% (as seen in Table 5), our model's maximum depth was increased to 4, which borders the point at which the data's recall remains unchanged at maximum depths greater than 4. However, the model's precision for the validation set has marginally improved by 1.8%. This will ensure that the model avoids underfitting the test data.

Table 5. Model 2 Performance Metrics against the Training and Validation set

Performance Metric (%)	Training Set	Validation Set
Accuracy	77.4	75.5
Precision	79.0	77.7
Recall	78.4	77.3

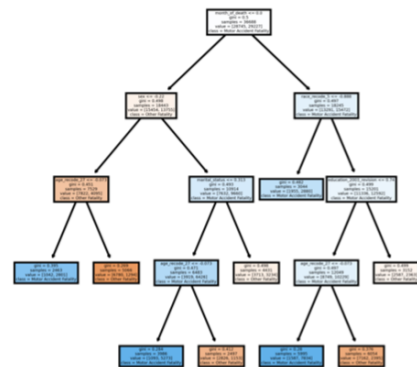


Figure 22. Estimator 4 Decision tree for the Random Forest in Model 2

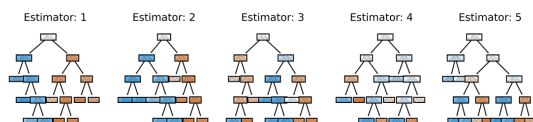


Figure 23. Random Forest in Model 2

6.3.5. Model 3: Optimising the Hyperparameters

To ensure the final model generated has the highest recall possible, the use of a very powerful computational function in the sklearn library called the *RandomizedSearchCV* allows the computer to run multiple random forests simultaneously, testing various combinations of the remaining hyperparameters until the most optimal is attained. So, after running 1000 different fits for the model with randomly assigned hyperparameters, the following optimal hyperparameters were selected:

- *N Estimators*: 39
- *Maximum Depth*: 4
- *Minimum Samples Leaf*: 33
- *Minimum Samples Split*: 0.2
- *Maximum Leaf Nodes*: 20

Some of the hyperparameters seemed nonsensical due to the randomised nature of the method to select the most optimal outcome, such as the maximum leaf nodes variable (820). The variable was reduced to a more appropriate value, since the maximum leaf nodes possible for a max depth of 4 is $2^4 = 16$. Thus, maximum number of leaf nodes was changed to 20.

6.4. Results

6.4.1. Final Model

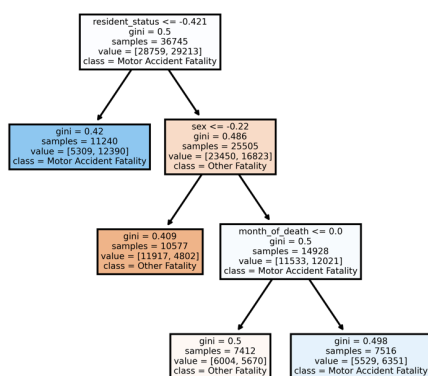


Figure 24. Estimator 2 Decision tree for the Random Forest in Model 3



Figure 25. Random Forest in Model 3

6.4.2. Analysis

Compared to the previous models, there's been a considerable improvement in the recall for the validation set by 6.5%. However, the

precision had fallen marginally by 2.7%, while the accuracy remained unchanged. This suggested that the model could perform well against unseen data, as seen with how the high recall and accuracy of the model on the test set is close to the performance on the training set.

Table 6. Model 3 Performance Metrics against the Training, Validation and Test set

Performance Metric (%)	Training	Validation	Testing
Accuracy	78.5	77.3	78.4
Precision	76.5	75.0	75.8
Recall	83.4	82	83.4

6.5. Discussion

6.5.1. Feature Importance

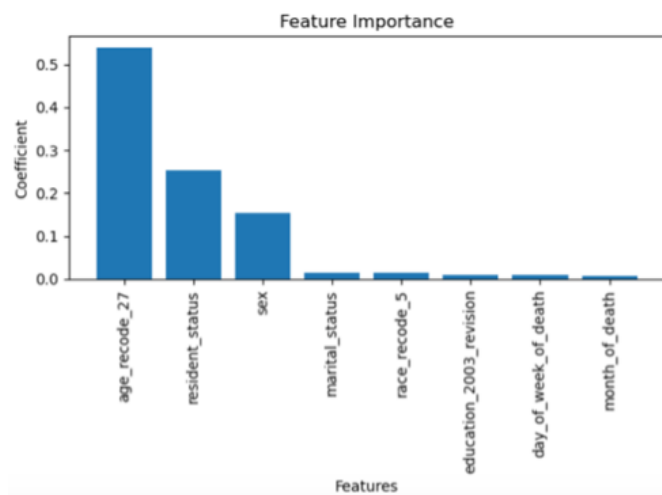


Figure 26. Feature Importance Chart for Final Model

As seen in Figure 26, the 3 most significant features for our Random Forest classifier are, in order of most significant:

1. Age
2. Resident Status
3. Sex

According to Statista [9], the largest group of licensed drivers in the US in 2021 was aged 30-34, underscoring the importance of age in our model. Younger individuals are more likely to die in motor accidents, suggesting higher recklessness among this demographic. The US government should raise road safety awareness among the younger generation to deter reckless driving and encourage adherence to safety rules. Non-residents in the US are less likely to die in motor accidents, likely due to high barriers for obtaining insurance and licenses. These barriers may result in a higher proportion of safer drivers among non-residents. Implementing similar barriers for US residents could reduce motor accident fatalities. Men are more likely than women to die in motor accidents, as they drive more miles annually and engage in riskier behaviours, such as not wearing seat belts or ignoring traffic signals. Consequently, men face a higher likelihood of encountering accidents. Road safety campaigns specifically targeting young men could help minimize and deter unlawful behaviours on roads and highways, further reducing fatalities. By addressing these factors, the US government can make informed decisions to improve road safety and reduce the socioeconomic burden of motor accidents.

6.5.2. Limitations of the Predictive Model

Our data and predictive model have limitations. The dataset lacks variables like weather, state, city, and years as a licensed driver, which

could provide more context and help US road safety bodies target critical regions in the country. Additionally, our Random Forest model could benefit from more rigorous hyperparameter tuning. Increasing iterations from 1,000 to 1,000,000 and narrowing the range of hyperparameter values could improve optimisation and ensure more sensible values for the optimised model.

7. Conclusion

The team created a series of models to predict deaths caused by suicide, homicide, motor accidents, and alcohol and drug abuse. Since each model addresses different questions, it's not feasible to determine which model is most relevant to this dataset. All models could benefit from larger datasets. The models' specificity to subsets of the data and the fact that they were trained on data from a specific country and demographic group means they may not be as relevant to the wider population.

The highest accuracy was achieved with the logistic regression forward selection model, which had a test accuracy of 0.788 and it follows by the random forest method with an test accuracy of 0.784. Through all methods, age and sex was the most significant factors in predicting these types of deaths.

This paper aims to use different models and techniques to prevent death by applying data analysis and machine learning techniques. Implementing those techniques can make us quickly identify high-risk individuals and provide timely interventions, potentially saving countless lives and reducing healthcare costs. While the current models show promise, further improvements with larger and more diverse datasets will improve the methods in addressing the complex issue of mortality in the USA.

References

- [1] J. Xu, S. L. Murphy, K. D. Kochanek, and E. Arias, "Deaths: Final Data for 2019," *National Vital Statistics Reports*, vol. 70, no. 8, Jul. 2021. [Online]. Available: <https://www.cdc.gov/nchs/data/nvsr/nvsr70/nvsr70-08-508.pdf>.
- [2] R. Tikkanen and M. K. Abrams, "U.S. health care from a global perspective, 2019: Higher spending, worse outcomes?," Jan. 2020. [Online]. Available: <https://www.commonwealthfund.org/publications/issue-briefs/2020/jan/us-health-care-global-perspective-2019>.
- [3] A. M. White, R. W. Hingson, I.-j. Pan, and H. yi, "Hospitalizations for Alcohol and Drug Overdoses in Young Adults Ages 18–24 in the United States, 1999–2008: Results From the Nationwide Inpatient Sample," *Journal of Studies on Alcohol and Drugs*, vol. 72, no. 5, pp. 774–786, 2011. [Online]. Available: [https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3357438/#:~:text=Rates%20of%20drug%20overdoses%20were,\(Accessed: Jun. 19, 2024\)](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3357438/#:~:text=Rates%20of%20drug%20overdoses%20were,(Accessed: Jun. 19, 2024)).
- [4] W. H. Organization, "Global Status Report on Alcohol and Health 2018," vol. 70, no. 8, Jul. 2018. [Online]. Available: https://books.google.co.uk/books?hl=en&lr=&id=qnOyDwAAQBAJ&oi=fnd&pg=PR7&ots=a2luNGwiho&sig=_JI5uwqwc75IyRgBVXnWyBPG0C8&redir_esc=y#v=onepage&q&f=false, (Accessed: Jun. 19, 2024).
- [5] N. I. of Health (NIH), "Men died of overdose at 2-3 times greater a rate than women in the U.S. in 2020-2021," Jun. 2023. [Online]. Available: <https://www.nih.gov/news-events/news-releases/men-died-overdose-2-3-times-greater-rate-women-us-2020-2021>.
- [6] R. Roth, *Why Is the United States the Most Homicidal Nation in the Affluent World?* | Office of Justice Programs, Dec. 2013. [Online]. Available: <https://www.ojp.gov/nij/virtual-library/abstracts/why-united-states-most-homicidal-nation-affluent-world>.
- [7] K. E. McCollister, M. T. French, and H. Fang, "The cost of crime to society: New crime-specific estimates for policy and program evaluation," *Drug and Alcohol Dependence*, vol. 108, no. 1–2, pp. 98–109, Apr. 2010. DOI: <https://doi.org/10.1016/j.drugalcdep.2009.12.002>.
- [8] I.-H. crash testing and highway safety, "Fatality facts 2022: Yearly snapshot," Jun. 2022. [Online]. Available: [https://www.iihs.org/topics/fatality-statistics/detail/yearly-snapshot#:~:text=A%20total%20of%2042%2C514%20people,\(Accessed: Jun. 19, 2024\)](https://www.iihs.org/topics/fatality-statistics/detail/yearly-snapshot#:~:text=A%20total%20of%2042%2C514%20people,(Accessed: Jun. 19, 2024)).
- [9] Statista, *U.S.: licensed drivers by age group*. [Online]. Available: [https://www.statista.com/statistics/206311/total-number-of-us-licensed-drivers-in-2010-by-age/#:~:text=In%202021%2C%20about%2017.44%20million,\(accessed: Jun.22.2024\)](https://www.statista.com/statistics/206311/total-number-of-us-licensed-drivers-in-2010-by-age/#:~:text=In%202021%2C%20about%2017.44%20million,(accessed: Jun.22.2024)).

8. Appendix

8.1. Preprocessing Dataset

8.1.1. Importing PANDAS library and filtering CSV file

```
1 import pandas as pd # Import the pandas library
   for data manipulation
2
3 # Read the entire CSV file into a DataFrame
4 data = pd.read_csv('2015_data.csv')
5
6 # Print the entire DataFrame to the console
7 print(data)
8
9 # Read the CSV file again, but this time only
   load specific columns into a new DataFrame
10 filtered_dataset = pd.read_csv('2015_data.csv',
   usecols=[
11     'resident_status', 'education_1989_revision',
12     'education_2003_revision',
13     'education_reporting_flag', 'month_of_death',
14     'sex', 'age_recoded_27',
15     'marital_status', 'day_of_week_of_death',
16     'manner_of_death', '358_cause_recoded',
17     'race_recoded_5'
18 ])
19
20 # Create a new DataFrame from the filtered
   dataset
21 df27 = pd.DataFrame(filtered_dataset)
22
23 # Save the new DataFrame to a new CSV file
24 df27.to_csv('new_data27.csv')
```

8.1.2. Converting data labels from strings to integers

```
1 # Create a copy of the df27 DataFrame
2 df27_1 = df27.copy()
3
4 # Replace 'M' with 1 and 'F' with 0 in the 'sex'
   column
5 df27_1['sex'].replace('M', 1, inplace=True)
6 df27_1['sex'].replace('F', 0, inplace=True)
7
8 # Replace marital status codes with numeric
   values
9 # 'W' (Widowed) -> 0
10 df27_1['marital_status'].replace('W', 0, inplace
   =True)
11 # 'S' (Single) -> 1
12 df27_1['marital_status'].replace('S', 1, inplace
   =True)
13 # 'M' (Married) -> 2
14 df27_1['marital_status'].replace('M', 2, inplace
   =True)
15 # 'D' (Divorced) -> 3
16 df27_1['marital_status'].replace('D', 3, inplace
   =True)
```

```

17 # 'U' (Unknown) -> 4
18 df27_1['marital_status'].replace('U', 4, inplace
19     =True)
20 # Save the modified DataFrame to a new CSV file
21 df27_1.to_csv('new_data27_1.csv', index=False)
22
23 # Create another copy of the modified DataFrame
24 df27_2 = df27_1.copy()

```

```

42 df27_8 = df27_7[df27_7.manner_of_death != 'Blank
43     ']
44 # Print rows where 'manner_of_death' still
45 # equals 'Blank' (should be empty)
46 print(df27_8.loc[df27_8['manner_of_death'] == '
47     Blank'])
48
49 # Save the final DataFrame to a new CSV file
50 df27_8.to_csv('new_data27_2.csv', index=False)

```

8.1.3. Removing all rows that contain NaN values in the dataset

```

1 # Drop the 'education_reporting_flag' column
2 # from df27_2
3 df27_2.drop('education_reporting_flag', axis=1,
4     inplace=True)
5
6 # Filter out rows where 'education_1989_revision'
7 # equals 99
8 df27_3 = df27_2[df27_2.education_1989_revision
9     != 99]
10
11 # Print rows where 'education_1989_revision'
12 # still equals 99 (should be empty)
13 print(df27_3.loc[df27_3['education_1989_revision
14     '] == 99])
15
16 # Filter out rows where 'education_2003_revision'
17 # equals 9
18 df27_4 = df27_3[df27_3.education_2003_revision
19     != 9]
20
21 # Print rows where 'education_2003_revision'
22 # still equals 9 (should be empty)
23 print(df27_4.loc[df27_4['education_2003_revision
24     '] == 9])
25
26 # (This line is redundant and can be removed
27 # since df27_4 is already created above)
28 # Filter out rows where 'education_2003_revision'
29 # equals 9
30 df27_4 = df27_3[df27_3.education_2003_revision
31     != 9]
32
33 # Print rows where 'education_2003_revision'
34 # still equals 9 (should be empty)
35 print(df27_4.loc[df27_4['education_2003_revision
36     '] == 9])
37
38 # Filter out rows where 'age_recode_27' equals
39 # 27
40 df27_5 = df27_4[df27_4.age_recode_27 != 27]
41
42 # Print rows where 'age_recode_27' still equals
43 # 27 (should be empty)
44 print(df27_5.loc[df27_5['age_recode_27'] == 27])
45
46 # Filter out rows where 'marital_status' equals
47 # 4
48 df27_6 = df27_5[df27_5.marital_status != 4]
49
50 # Print rows where 'marital_status' still equals
51 # 4 (should be empty)
52 print(df27_6.loc[df27_6['marital_status'] == 4])
53
54 # Filter out rows where 'day_of_week_of_death'
55 # equals 9
56 df27_7 = df27_6[df27_6.day_of_week_of_death !=
57     9]
58
59 # Print rows where 'day_of_week_of_death' still
60 # equals 9 (should be empty)
61 print(df27_7.loc[df27_7['day_of_week_of_death']
62     == 9])
63
64 # Filter out rows where 'manner_of_death' equals
65 # 'Blank'

```

8.1.4. Binarization of dataset

```

1 # Create a copy of df27_8
2 df27_9 = df27_8.copy()
3
4 # Replace 'resident_status' values 2, 3, and 4
5 # with 0
6 df27_9['resident_status'].replace([2,3,4], 0,
7     inplace=True)
8
9 # Replace 'education_1989_revision' values 0-12
10 # with 0 and values 13-17 with 1
11 df27_9['education_1989_revision'].replace
12 ([0,1,2,3,4,5,6,7,8,9,10,11,12], 0, inplace=
13     True)
14 df27_9['education_1989_revision'].replace
15 ([13,14,15,16,17], 1, inplace=True)
16
17 # Replace 'education_2003_revision' values 0-4
18 # with 0 and values 5-8 with 1
19 df27_9['education_2003_revision'].replace
20 ([0,1,2,3,4], 0, inplace=True)
21 df27_9['education_2003_revision'].replace
22 ([5,6,7,8], 1, inplace=True)
23
24 # Replace 'month_of_death' values 10-3 with 0
25 # and values 4-9 with 1
26 df27_9['month_of_death'].replace
27 ([10,11,12,0,1,2,3], 0, inplace=True)
28 df27_9['month_of_death'].replace([4,5,6,7,8,9],
29     1, inplace=True)
30
31 # Replace 'age_recode_27' values 1-17 with 0 and
32 # values 18-26 with 1
33 df27_9['age_recode_27'].replace
34 ([1,2,3,4,5,6,7,8,9,10,11,12,13,14,15,16,17],
35     0, inplace=True)
36 df27_9['age_recode_27'].replace
37 ([18,19,20,21,22,23,24,25,26], 1, inplace=
38     True)
39
40 # Replace 'marital_status' values 0, 1, and 3
41 # with 0, and value 2 with 1
42 df27_9['marital_status'].replace([0,1,3], 0,
43     inplace=True)
44 df27_9['marital_status'].replace(2, 1, inplace=
45     True)
46
47 # Replace 'day_of_week_of_death' values 1 and 7
48 # with 0, and values 2-6 with 1
49 df27_9['day_of_week_of_death'].replace([1,7], 0,
50     inplace=True)
51 df27_9['day_of_week_of_death'].replace
52 ([2,3,4,5,6], 1, inplace=True)
53
54 # Replace 'race_recode_5' values 0, 2, 3, and 4
55 # with 0
56 df27_9['race_recode_5'].replace([0,2,3,4], 0,
57     inplace=True)
58
59 # Print the final DataFrame to the console
60 print(df27_9)
61
62 # Save the final binarised DataFrame to a new
63 # CSV file
64 df27_9.to_csv('binarised_data.csv', index=False)

```


8.1.5. Reimporting the binarized dataset

```

1 import pandas as pd # Import the pandas library
  for data manipulation
2 import numpy as np # Import the numpy library
  for numerical operations
3
4 # Read the CSV file into a DataFrame
5 data = pd.read_csv('binarised_data.csv')
6
7 # Create a new DataFrame from the loaded data
8 df1 = pd.DataFrame(data)
9
10 # Drop the 'Unnamed: 0' column from the
   DataFrame
11 df1.drop('Unnamed: 0', axis=1, inplace=True)
12
13 # Print the DataFrame to the console
14 print(df1)

```

8.1.6. Binarizing the cause of death column for motor accidents

```

1 # Define a function to replace specific values
2 def replace_values(x):
3     # If x equals 384, 389, 390, or 391, return
   0
4     if x == 384 or x == 389 or x == 390 or x ==
   391:
5         return 0
6     # If x is less than 384 or greater than 399,
   return 0
7     elif x < 384 or x > 399:
8         return 0
9     # Otherwise, return 1
10    else:
11        return 1
12
13 # Apply the replace_values function to the '358
   _cause_recode' column
14 df1['358_cause_recode'] = df1['358_cause_recode']
   .apply(replace_values)
15
16 # Print the DataFrame after replacing values
17 print("After Replacing:\n", df1)

```

8.1.7. Combining the education columns of the data

```

1 # Set reference and target column names
2 reference_column = 'education_1989_revision'
3 target_column = 'education_2003_revision'
4
5 # Define a function to apply the mapping
6 def apply_mapping(row):
7     # Check if the target column value is NaN
8     if pd.isna(row[target_column]):
9         # Convert reference column value to
   string
10        left_value = str(row[reference_column])
11        # If reference column value is "0" or
   "1", return it
12        if left_value == "0" or left_value == "1"
   ":
13            return left_value
14        # Otherwise, return the original target
   column value
15        return row[target_column]
16
17 # Apply the apply_mapping function to the target
   column
18 df1[target_column] = df1.apply(apply_mapping,
   axis=1)
19
20 # Drop the 'manner_of_death' column from the
   DataFrame
21 df1.drop('manner_of_death', axis=1, inplace=True)

```

```

22
23 # Drop the 'education_1989_revision' column from
   the DataFrame
24 df1.drop('education_1989_revision', axis=1,
   inplace=True)
25
26 # Drop rows with any NaN values
27 df2 = df1.dropna()

```

8.2. Alcohol and Drug Overdose - Erim Hayretci

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3 from scipy import stats
4 import numpy as np
5 from sklearn import svm
6 from sklearn.metrics import confusion_matrix,
   accuracy_score
7 from sklearn.linear_model import
   LogisticRegression
8
9 from google.colab import drive
10 drive.mount('/content/drive')
11
12 file_path = '/content/drive/MyDrive/new.csv'
13
14 df = pd.read_csv(file_path)
15
16 #Filtering the dataset to take all male
   individuals
17
18 del df['Unnamed: 0']
19 del df['education_1989_revision']
20 filt1 = df['sex'] == 1
21 all_male = df[filt1]
22
23 print(all_male)
24
25
26 #Filtering out deaths from alcohol or substance
   overdose
27 filt2 = all_male['358_cause_recode'] == 178
28 filt3 = all_male['358_cause_recode'] == 179
29
30 positive = pd.concat([all_male[filt2], all_male[
   filt3]])
31
32 #Undersampling male dataframe
33
34 age_proportions = positive['age_recode_27'].
   value_counts(normalize=True)
35
36 total = []
37
38 #Stratified sampling based on age proportions of
   overdose deaths
39 for i in age_proportions.keys():
40     filt5 = all_male['age_recode_27'] == i
41     age_group = all_male[filt5]
42     specific_age = age_group.groupby('
   age_recode_27', group_keys=False).
43     apply(lambda x, frac=age_proportions[i]: x.
   sample(frac=frac))
44     total.append(specific_age)
45
46 final = pd.concat(total)
47
48 print(final)
49
50 negative = final.groupby('age_recode_27',
   group_keys=False).
51 apply(lambda x: x.sample(frac=0.09640))
52
53 print(negative)
54
55
56 dataset = pd.concat([positive, negative])
57
58 #Binarising cause of death

```

```

59 for index, row in dataset.iterrows():
60     if row['358_cause_recode'] == 178:
61         dataset.loc[index, '358_cause_recode'] = 1
62     elif row['358_cause_recode'] == 179:
63         dataset.loc[index, '358_cause_recode'] = 1
64     else:
65         dataset.loc[index, '358_cause_recode'] = 0
66
67 print(dataset)
68
69
70 from sklearn.model_selection import
71     train_test_split
72
73 x_pre = dataset[['resident_status',
74                 'education_2003_revision',
75                 'month_of_death',
76                 'age_recode_27',
77                 'marital_status',
78                 'day_of_week_of_death',
79                 'manner_of_death',
80                 'race_recode_5']]
81 x = x_pre.fillna(0)
82 y = dataset['358_cause_recode']
83
84 #Splitting a training set
85 x_train, x_test_val, y_train, y_test_val =
86     train_test_split(x, y, test_size=0.6, shuffle =
87         True, random_state = 0)
88
89 #Splitting validation and testing sets
90 x_val, x_test, y_val, y_test =
91     train_test_split(x_test_val, y_test_val,
92         test_size=0.5,
93         shuffle = True, random_state =
94             0)
95
96 print(x_train.shape)
97 print(x_val.shape)
98 print(x_test.shape)
99
100 print(y_train.shape)
101 print(y_val.shape)
102 print(y_test.shape)
103
104 from sklearn.metrics import accuracy_score,
105     precision_score, recall_score
106
107 #Accuracy, preicision and recall scores with
108     different kernel types
109 kernel_list = ['linear', 'poly', 'rbf', 'sigmoid
110     ']
111
112 accuracy_training = []
113 accuracy_validation = []
114 precision_training = []
115 precision_validation = []
116 recall_training = []
117 recall_validation = []
118
119 for i in kernel_list:
120     model = svm.SVC(kernel=i)
121     model.fit(x_train, y_train)
122     y_pred = model.predict(x_train)
123     accuracy_training.append(accuracy_score(
124         y_train, y_pred))
125     precision_training.append(precision_score(
126         y_train, y_pred))
127     recall_training.append(recall_score(y_train,
128         y_pred))
129
130 for i in kernel_list:
131     model = svm.SVC(kernel=i)
132     model.fit(x_val, y_val)
133     y_pred = model.predict(x_val)
134     accuracy_validation.append(accuracy_score(
135         y_val, y_pred))
136     precision_validation.append(precision_score(
137         y_val, y_pred))
138     recall_validation.append(recall_score(y_val,
139         y_pred))
140
141 print(accuracy_training)
142 print(accuracy_validation)
143 print(precision_training)
144 print(precision_validation)
145 print(recall_training)
146 print(recall_validation)
147
148 #Kernel - Accuracy
149 plt.plot(kernel_list, accuracy_training, 'b',
150     label = 'Training Set')
151 plt.plot(kernel_list, accuracy_validation, 'r',
152     label = 'Validation Set')
153 plt.ylabel('Accuracy')
154 plt.xlabel('Kernel')
155 plt.title('Accuracy Scores for the Training Set
156     and Validation Set (C=1)')
157 plt.legend()
158
159 #Kernel - Precision
160 plt.plot(kernel_list, precision_training, 'b',
161     label = 'Training Set')
162 plt.plot(kernel_list, precision_validation, 'r',
163     label = 'Validation Set')
164 plt.ylabel('Precision')
165 plt.xlabel('Kernel')
166 plt.title('Precision Scores for the Training Set
167     and Validation Set (C=1)')
168 plt.legend()
169
170 #Kernel - Recall
171 plt.plot(kernel_list, recall_training, 'b', label
172     = 'Training Set')
173 plt.plot(kernel_list, recall_validation, 'r',
174     label = 'Validation Set')
175 plt.ylabel('Recall')
176 plt.xlabel('Kernel')
177 plt.title('Recall Scores for the Training Set
178     and Validation Set (C=1)')
179 plt.legend()
180
181 #Accuracy, precision and recall scores with
182     different C values
183 accuracy_training_c = []
184 accuracy_validation_c = []
185 precision_training_c = []
186 precision_validation_c = []
187 recall_training_c = []
188 recall_validation_c = []
189 y_data_c = []
190
191 for i in np.arange(0.1,10,1):
192     y_data_c.append(i)
193
194 for i in y_data_c:
195     model = svm.SVC(C=i)
196     model.fit(x_train, y_train)
197     y_pred = model.predict(x_train)
198     accuracy_training_c.append(accuracy_score(
199         y_train, y_pred))
200     precision_training_c.append(precision_score(
201         y_train, y_pred))
202     recall_training_c.append(recall_score(y_train,
203         y_pred))
204
205 for i in y_data_c:
206     model = svm.SVC(C=i)
207     model.fit(x_val, y_val)
208     y_pred = model.predict(x_val)
209     accuracy_validation_c.append(accuracy_score(
210         y_val, y_pred))
211     precision_validation_c.append(precision_score(
212         y_val, y_pred))
213     recall_validation_c.append(recall_score(y_val,
214         y_pred))
215
216 print(accuracy_training_c)
217 print(accuracy_validation_c)
218 print(precision_training_c)
219 print(precision_validation_c)
220 print(recall_training_c)

```

```

189 print(recall_validation_c)
190
191 #C - Accuracy
192 plt.plot(y_data_c, accuracy_training_c, 'b',
193         label = 'Training Set')
194 plt.plot(y_data_c, accuracy_validation_c, 'r',
195         label = 'Validation Set')
196 plt.ylabel('Accuracy')
197 plt.xlabel('C value')
198 plt.title('Accuracy Scores for the Training Set
199         and Validation Set
200         'with varying C values')
201 plt.legend()
202
203 #C - Precision
204 plt.plot(y_data_c, precision_training_c, 'b',
205         label = 'Training Set')
206 plt.plot(y_data_c, precision_validation_c, 'r',
207         label = 'Validation Set')
208 plt.ylabel('Precision')
209 plt.xlabel('C value')
210 plt.title('Precision Scores for the Training Set
211         and Validation Set
212         'with varying C values')
213 plt.legend()
214
215 #C - Recall
216 plt.plot(y_data_c, recall_training_c, 'b', label
217         = 'Training Set')
218 plt.plot(y_data_c, recall_validation_c, 'r',
219         label = 'Validation Set')
220 plt.ylabel('Recall')
221 plt.xlabel('C value')
222 plt.title('Recall Scores for the Training Set
223         and Validation Set
224         'with varying C values')
225 plt.legend()
226
227 #Accuracy, precision and recall scores with
228         different gamma values
229 accuracy_training_g = []
230 accuracy_validation_g = []
231 precision_training_g = []
232 precision_validation_g = []
233 recall_training_g = []
234 recall_validation_g = []
235 y_data_g = []
236
237 for i in np.arange(0.1,0.6,0.05):
238     y_data_g.append(i)
239
240 for i in y_data_g:
241     model = svm.SVC(gamma=i)
242     model.fit(x_train, y_train)
243     y_pred = model.predict(x_train)
244     accuracy_training_g.append(accuracy_score(
245         y_train, y_pred))
246     precision_training_g.append(precision_score(
247         y_train, y_pred))
248     recall_training_g.append(recall_score(y_train,
249         y_pred))
250
251 for i in y_data_g:
252     model = svm.SVC(gamma=i)
253     model.fit(x_val, y_val)
254     y_pred = model.predict(x_val)
255     accuracy_validation_g.append(accuracy_score(
256         y_val, y_pred))
257     precision_validation_g.append(precision_score(
258         y_val, y_pred))
259     recall_validation_g.append(recall_score(y_val,
260         y_pred))
261
262 print(accuracy_training_g)
263 print(accuracy_validation_g)
264 print(precision_training_g)
265 print(precision_validation_g)
266 print(recall_training_g)
267 print(recall_validation_g)
268
269 #Gamma - Accuracy
270 plt.plot(y_data_g, accuracy_training_g, 'b',
271         label = 'Training Set')
272 plt.plot(y_data_g, accuracy_validation_g, 'r',
273         label = 'Validation Set')
274 plt.ylabel('Accuracy')
275 plt.xlabel('Gamma value')
276 plt.title('Accuracy Scores for the Training Set
277         and Validation Set
278         'with varying Gamma values')
279 plt.legend()
280
281 #Gamma - Precision
282 plt.plot(y_data_g, precision_training_g, 'b',
283         label = 'Training Set')
284 plt.plot(y_data_g, precision_validation_g, 'r',
285         label = 'Validation Set')
286 plt.ylabel('Precision')
287 plt.xlabel('Gamma')
288 plt.title('Precision Scores for the Training Set
289         and Validation Set
290         'with varying Gamma values')
291 plt.legend()
292
293 #Gamma - Recall
294 plt.plot(y_data_g, recall_training_g, 'b', label
295         = 'Training Set')
296 plt.plot(y_data_g, recall_validation_g, 'r',
297         label = 'Validation Set')
298 plt.ylabel('Recall')
299 plt.xlabel('Gamma')
300 plt.title('Recall Scores for the Training Set
301         and Validation Set
302         'with varying Gamma values')
303 plt.legend()
304
305 #Testing Set
306 model = svm.SVC(kernel='rbf',C=3,gamma=0.1)
307 model.fit(x_train,y_train)
308 y_pred = model.predict(x_test)
309 print('test')
310 print(accuracy_score(y_test, y_pred))
311 print(precision_score(y_test, y_pred))
312 print(recall_score(y_test, y_pred))

```

8.3. Homicide - Vasco de Noronha

8.3.1. Filtering the dataset

```

1 # Combining to make 1 education column
2
3 import pandas as pd
4
5 # json code from data set was used to create
6         this map and to make education one column
7         which would have worked better for my model,
8         it does this by finding nan values in the
9         education column that we kept and looking at
10        the older education index to fill in the
11        gaps and preserve as many data points as
12        possible.
13
14 mapping = {
15     '10': '3',
16     '11': '3',
17     '12': '3',
18     '13': '4',
19     '14': '5',
20     '15': '6',
21     '16': '7',
22     '17': '8',
23     '99': '9',
24     '00': '9',
25     '01': '1', '02': '1', '03': '1', '04': '1',
26     '05': '1', '06': '1', '07': '1', '08': '1',
27     '09': '3'
28 }
29
30 }
31

```

```

22 df = pd.read_csv('2015_filtered.csv', low_memory
23                 =False)
24
25 reference_column = 'education_1989_revision'
26 target_column = 'education_2003_revision'
27
28 def apply_mapping(row):
29     if pd.isna(row[target_column]):
30         left_value = str(row[reference_column])
31         if left_value in mapping:
32             return mapping[left_value]
33         return row[target_column]
34
35 # Apply the function to the target column
36 df[target_column] = df.apply(apply_mapping, axis
37                             =1)
38
39 df = df.drop('education_1989_revision', axis=1)
40
41 df.to_csv('2015_filtered_2.csv', index=False)
42
43 print(df.tail())
44
45 # Checking for Nan Values
46
47 import pandas as pd
48
49 df = pd.read_csv('2015_filtered_2.csv',
50                 low_memory=False)
51
52 has_nan = df.isnull().values.any()
53 nan_counts = df.isnull().sum()
54 rows_with_nan = df[df.isnull().any(axis=1)]
55
56 print(f"DataFrame contains NaN values: {has_nan}
57       ")
58 print("Count of NaN values in each column:")
59 print(nan_counts)
60 print("Rows with NaN values:")
61 print(rows_with_nan)
62
63 if has_nan:
64     print("The dataset has NaN values.")
65 else:
66     print("The dataset has zero NaN values.")
67
68 # Removing all columns with Nan Values
69
70 import pandas as pd
71
72 df = pd.read_csv('2015_filtered_2.csv',
73                 low_memory=False)
74
75 # Define the columns to check for NaN values
76 columns_to_check = ['education_2003_revision', '
77                     manner_of_death'] # replace with your
78                     actual column names
79
80 # Remove rows with NaN values in the specified
81 columns
82 df_cleaned = df.dropna(subset=columns_to_check)
83
84 df_cleaned.to_csv('2015_filtered_3.csv', index=
85                  False)
86
87 df = pd.read_csv('2015_filtered_3.csv',
88                 low_memory=False)
89
90 has_nan = df.isnull().values.any()
91 nan_counts = df.isnull().sum()
92 rows_with_nan = df[df.isnull().any(axis=1)]
93
94 print(f"DataFrame contains NaN values: {has_nan}
95       ")
96 print("Count of NaN values in each column:")
97 print(nan_counts)
98 print("Rows with NaN values:")
99 print(rows_with_nan)
100
101 if has_nan:

```

```

92     print("The dataset has NaN values.")
93 else:
94     print("The dataset has zero NaN values.")

```

8.3.2. Balancing the dataset

```

1 import pandas as pd
2
3 input_file_path = '2015_filtered_3.csv'
4 df = pd.read_csv(input_file_path, low_memory=
5                 False)
6
7 # low memory set to false due to large dataset
8
9 # Separate the data into two groups based on
10 manner of death
11 homicides = df[df['manner_of_death'] == 3]
12 non_homicides = df[df['manner_of_death'] != 3]
13
14 homicides_count = homicides.shape[0]
15 non_homicides_count = non_homicides.shape[0]
16
17 homicides_sample = homicides.sample(n=10000,
18                                     random_state=42)
19 non_homicides_sample = non_homicides.sample(n
20                                             =10000, random_state=42)
21
22 combined_sample = pd.concat([homicides_sample,
23                               non_homicides_sample])
24
25 output_file_path = '2015_20000_filtered_even10.
26                   csv'
27
28 combined_sample.to_csv(output_file_path, index=
29                       False)

```

8.3.3. Splitting the dataset

```

1 import pandas as pd
2
3 input_file_path = '2015_20000_filtered_even10.
4                   csv'
5 df = pd.read_csv(input_file_path, low_memory=
6                 False)
7
8 # Separate the data into two groups based on
9 manner of death
10 # 3 was a homicide
11 homicides = df[df['manner_of_death'] == 3]
12 non_homicides = df[df['manner_of_death'] != 3]
13
14 # Determine the sample sizes
15 total_samples = 20000
16 train_size = int(total_samples * 0.6)
17 val_test_size = (total_samples - train_size) //
18                 2
19
20 #splitting into different dfs
21 homicides_train = homicides.sample(n=train_size,
22                                   random_state=42)
23 non_homicides_train = non_homicides.sample(n=
24 train_size, random_state=42)
25 train_sample = pd.concat([homicides_train,
26                           non_homicides_train])
27
28 remaining_homicides = homicides.drop(
29     homicides_train.index)
30 remaining_non_homicides = non_homicides.drop(
31     non_homicides_train.index)
32
33 homicides_val = remaining_homicides.sample(n=
34 val_test_size, random_state=42)
35 non_homicides_val = remaining_non_homicides.
36 sample(n=val_test_size, random_state=42)
37 val_sample = pd.concat([homicides_val,
38                         non_homicides_val])
39
40 homicides_test = remaining_homic

```



```

28 # Save all dfs to different files
29 train_sample.to_csv('60.20.20_train', index=
30 False)
31 homicides_test.to_csv('60.20.20_test', index=
32 False)
33 val_sample.to_csv('60.20.20_val', index=False)

```

8.3.4. Forward Selection algorithm

```

1 import numpy as np
2 import pandas as pd
3 from sklearn.metrics import accuracy_score
4 from sklearn.linear_model import
5     LogisticRegression
6 import warnings
7 warnings.filterwarnings('ignore')
8
9 def select_column_to_add(X_train, y_train, X_val
10 , y_val, columns_in_model, columns_to_test):
11     acc_best = 0
12     column_best = None
13
14     if len(columns_in_model) == 0:
15         acc_best = 0
16     elif len(columns_in_model) == 1:
17         mod = LogisticRegression(C=1e9).fit(
18             X_train[columns_in_model].values.reshape(-1,
19             1), y_train)
20         acc_best = accuracy_score(y_val, mod.
21             predict(X_val[columns_in_model].values.
22             reshape(-1, 1)))
23     else:
24         mod = LogisticRegression(C=1e9).fit(
25             X_train[columns_in_model], y_train)
26         acc_best = accuracy_score(y_val, mod.
27             predict(X_val[columns_in_model]))
28
29     for column in columns_to_test:
30         mod = LogisticRegression(C=1e9).fit(
31             X_train[columns_in_model + [column]],
32             y_train)
33         y_pred = mod.predict(X_val[
34             columns_in_model + [column]])
35         acc = accuracy_score(y_val, y_pred)
36
37         if acc > acc_best:
38             acc_best = acc
39             column_best = column
40
41     if column_best is not None:
42         print('Adding {} to the model'.format(
43             column_best))
44         print('The new best validation accuracy
45         is {}'.format(acc_best))
46         columns_in_model_updated =
47         columns_in_model + [column_best]
48     else:
49         print('Did not add anything to the model
50         ')
51         columns_in_model_updated =
52         columns_in_model
53
54     return columns_in_model_updated, acc_best
55
56 def auto_forward_selection(X_train, y_train,
57 X_val, y_val, max_num_of_features):
58     columns_to_test = list(X_train.columns)
59     columns_in_model = []
60     current_acc = 0
61
62     for i in range(0, max_num_of_features):
63         columns_in_model, acc_best =
64         select_column_to_add(X_train, y_train, X_val
65         , y_val, columns_in_model, columns_to_test)
66
67         if acc_best == current_acc:
68             break

```

```

51     else:
52         current_acc = acc_best
53         for feature in columns_in_model:
54             if feature in columns_to_test:
55                 columns_to_test.remove(
56                     feature)
57
58         return columns_in_model, acc_best
59
60 df_train = pd.read_csv('60.20.20_train.csv',
61 low_memory=False)
62 df_val = pd.read_csv('60.20.20_val.csv',
63 low_memory=False)
64
65 df_train['manner_of_death'] = df_train['
66 manner_of_death'].apply(lambda x: 1 if x ==
67 3 else 0)
68 df_val['manner_of_death'] = df_val['
69 manner_of_death'].apply(lambda x: 1 if x ==
70 3 else 0)
71
72 X_train = df_train.drop(columns=['
73 manner_of_death', 'Unnamed: 0', '358
74 _cause_recode'])
75 y_train = df_train['manner_of_death']
76 X_val = df_val.drop(columns=['manner_of_death',
77 'Unnamed: 0', '358_cause_recode'])
78 y_val = df_val['manner_of_death']
79
80 max_num_of_features = 9
81
82 columns_selected, best_accuracy =
83 auto_forward_selection(X_train, y_train,
84 X_val, y_val, max_num_of_features)
85
86 print("Final selected columns:",
87 columns_selected)
88 print("Best validation accuracy:", best_accuracy
89 )

```

8.3.5. Validation

```

1 import itertools
2 import numpy as np
3 import pandas as pd
4 import matplotlib.pyplot as plt
5 from sklearn.linear_model import
6     LogisticRegression as logreg
7 from sklearn.metrics import accuracy_score,
8     precision_score, recall_score
9
10 df_train = pd.read_csv('50.25.25_train.csv',
11 low_memory=False)
12 df_val = pd.read_csv('50.25.25_val.csv',
13 low_memory=False)
14 df_test = pd.read_csv('50.25.25_test.csv',
15 low_memory=False)
16
17 # binarizing manner of death, homicide =1, non-
18 homicide death = 0
19 df_train['manner_of_death'] = df_train['
20 manner_of_death'].apply(lambda x: 1 if x ==
21 3 else 0)
22 df_val['manner_of_death'] = df_val['
23 manner_of_death'].apply(lambda x: 1 if x ==
24 3 else 0)
25 df_test['manner_of_death'] = df_test['
26 manner_of_death'].apply(lambda x: 1 if x ==
27 3 else 0)
28
29 # Split features and target variable
30 # removed 358_cause_recode due to the cause
31 being dependent on manner of death
32 X_train = df_train.drop(columns=['
33 manner_of_death', 'Unnamed: 0', '358
34 _cause_recode'])
35 y_train = df_train['manner_of_death']
36 X_val = df_val.drop(columns=['manner_of_death',
37 'Unnamed: 0', '358_cause_recode'])

```

```

22 y_val = df_val['manner_of_death']
23 X_test = df_test.drop(columns=['manner_of_death',
24                               , 'Unnamed: 0', '358_cause_recode'])
25 y_test = df_test['manner_of_death']
26
27 def combinations(attributes, group_size):
28     return list(itertools.combinations(
29         attributes, group_size))
30
31 # Initialize plot data dictionary
32 plot_data = {
33     'Acc_train': [], 'Prec_train': [], '
34     Rec_train': [],
35     'Acc_val': [], 'Prec_val': [], 'Rec_val': []
36 }
37
38 # List of columns to test
39 # All columns left after filtering
40 columns_to_test = [
41     'age_recode_27', 'race_recode_5', '
42     resident_status', 'education_2003_revision',
43     'month_of_death', 'sex',
44     'marital_status', 'day_of_week_of_death' ]
45
46 # Function gets plot data
47 def get_plot_data(column_to_test):
48     for i in range(len(column_to_test)):
49         trial_predictors = combinations(
50             column_to_test, i + 1)
51         acc_list = []
52         prec_list = []
53         rec_list = []
54         acc_val_list = []
55         prec_val_list = []
56         rec_val_list = []
57
58         for col_comb in trial_predictors:
59             X = X_train[list(col_comb)]
60             X_Test = X_test[list(col_comb)]
61             y = y_train
62             mylr = logreg(C=1e9)
63             mylr.fit(X, y)
64
65             predicted_y = mylr.predict(X)
66             acc = accuracy_score(y, predicted_y)
67             prec = precision_score(y,
68                 predicted_y)
69             rec = recall_score(y, predicted_y)
70
71             predicted_y_test = mylr.predict(
72                 X_Test)
73             acc_val = accuracy_score(y_test,
74                 predicted_y_test)
75             prec_val = precision_score(y_test,
76                 predicted_y_test)
77             rec_val = recall_score(y_test,
78                 predicted_y_test)
79
80             acc_list.append(acc)
81             prec_list.append(prec)
82             rec_list.append(rec)
83             acc_val_list.append(acc_val)
84             prec_val_list.append(prec_val)
85             rec_val_list.append(rec_val)
86
87             if acc_list and acc_val_list:
88                 plot_data['Acc_train'].append(np.
89                     mean(acc_list))
90                 plot_data['Prec_train'].append(np.
91                     mean(prec_list))
92                 plot_data['Rec_train'].append(np.
93                     mean(rec_list))
94                 plot_data['Acc_val'].append(np.mean(
95                     acc_val_list))
96                 plot_data['Prec_val'].append(np.mean(
97                     prec_val_list))
98                 plot_data['Rec_val'].append(np.mean(
99                     rec_val_list))
100
101 get_plot_data(columns_to_test)

```

```

86 df = pd.DataFrame(plot_data)
87 print(df)
88
89 # Plot the results on a graph
90 plt.figure()
91 plt.plot(df['Acc_train'], 'b')
92 plt.plot(df['Acc_val'], 'g')
93 plt.title('Accuracy')
94 plt.xlabel('Model Complexity (Number of
95     Variables)')
96 plt.ylabel('Mean Accuracy')
97 plt.legend(['Training data', 'Validation Data'])
98 plt.show()
99
100 # Second model
101
102 import numpy as np
103 import pandas as pd
104 from sklearn.metrics import accuracy_score,
105     precision_score, recall_score
106 from sklearn.linear_model import
107     LogisticRegression
108 import warnings
109 import itertools
110 import matplotlib.pyplot as plt
111
112 warnings.filterwarnings('ignore')
113
114 def select_column_to_add(X_train, y_train, X_val,
115     y_val, columns_in_model, columns_to_test):
116     acc_best = 0
117     column_best = None
118
119     if len(columns_in_model) == 0:
120         acc_best = 0
121     elif len(columns_in_model) == 1:
122         mod = LogisticRegression(C=1e9).fit(
123             X_train[columns_in_model].values.reshape(-1,
124                 1), y_train)
125         acc_best = accuracy_score(y_val, mod.
126             predict(X_val[columns_in_model].values.
127                 reshape(-1, 1)))
128     else:
129         mod = LogisticRegression(C=1e9).fit(
130             X_train[columns_in_model], y_train)
131         acc_best = accuracy_score(y_val, mod.
132             predict(X_val[columns_in_model]))
133
134     for column in columns_to_test:
135         mod = LogisticRegression(C=1e9).fit(
136             X_train[columns_in_model + [column]],
137             y_train)
138         y_pred = mod.predict(X_val[
139             columns_in_model + [column]])
140         acc = accuracy_score(y_val, y_pred)
141
142         if acc > acc_best:
143             acc_best = acc
144             column_best = column
145
146     if column_best is not None:
147         print('Adding {} to the model.
148             Validation accuracy: {}'.format(column_best,
149                 acc_best))
150         columns_in_model_updated =
151             columns_in_model + [column_best]
152     else:
153         print('Did not add anything to the model
154             ')
155         columns_in_model_updated =
156             columns_in_model
157
158     return columns_in_model_updated, acc_best
159
160 def auto_forward_selection(X_train, y_train,
161     X_val, y_val, max_num_of_features):
162     columns_to_test = list(X_train.columns)
163     columns_in_model = []
164     current_acc = 0
165     plot_data = {

```

```

147     "Acc_train": [], "Prec_train": [], "
148     Rec_train": [],
149     "Acc_val": [], "Prec_val": [], "Rec_val"
150     : []
151     }
152
153     for i in range(0, max_num_of_features):
154         columns_in_model, acc_best =
155         select_column_to_add(X_train, y_train, X_val
156         , y_val, columns_in_model, columns_to_test)
157
158         if acc_best == current_acc:
159             break
160         else:
161             current_acc = acc_best
162             for feature in columns_in_model:
163                 if feature in columns_to_test:
164                     columns_to_test.remove(
165                     feature)
166
167             # Train the model with the selected
168             features and record metrics
169             mod = LogisticRegression(C=1e9).fit(
170             X_train[columns_in_model], y_train)
171             y_train_pred = mod.predict(X_train[
172             columns_in_model])
173             y_val_pred = mod.predict(X_val[
174             columns_in_model])
175
176             plot_data["Acc_train"].append(
177             accuracy_score(y_train, y_train_pred))
178             plot_data["Prec_train"].append(
179             precision_score(y_train, y_train_pred))
180             plot_data["Rec_train"].append(
181             recall_score(y_train, y_train_pred))
182             plot_data["Acc_val"].append(
183             accuracy_score(y_val, y_val_pred))
184             plot_data["Prec_val"].append(
185             precision_score(y_val, y_val_pred))
186             plot_data["Rec_val"].append(recall_score
187             (y_val, y_val_pred))
188
189             return columns_in_model, plot_data
190
191 df_train = pd.read_csv('60.20.20_train.csv',
192 low_memory=False)
193 df_val = pd.read_csv('60.20.20_val.csv',
194 low_memory=False)
195 df_test = pd.read_csv('60.20.20_test.csv',
196 low_memory=False)
197
198 df_train['manner_of_death'] = df_train['
199 manner_of_death'].apply(lambda x: 1 if x ==
200 3 else 0)
201 df_val['manner_of_death'] = df_val['
202 manner_of_death'].apply(lambda x: 1 if x ==
203 3 else 0)
204 df_test['manner_of_death'] = df_test['
205 manner_of_death'].apply(lambda x: 1 if x ==
206 3 else 0)
207
208 X_train = df_train.drop(columns=['
209 manner_of_death', 'Unnamed: 0', '358
210 _cause_recode'])
211 y_train = df_train['manner_of_death']
212 X_val = df_val.drop(columns=['manner_of_death',
213 'Unnamed: 0', '358_cause_recode'])
214 y_val = df_val['manner_of_death']
215 X_test = df_test.drop(columns=['manner_of_death',
216 'Unnamed: 0', '358_cause_recode'])
217 y_test = df_test['manner_of_death']
218
219 # number of columns left essentially
220 max_num_of_features = 8
221
222 columns_selected, plot_data =
223 auto_forward_selection(X_train, y_train,
224 X_val, y_val, max_num_of_features)
225
226 print("Final selected columns:",
227 columns_selected)
228 print("Validation accuracies:", plot_data["
229 Acc_val"])
230
231 # Create DataFrame from plot data
232 df_plot = pd.DataFrame(plot_data)
233 print(df_plot)
234
235 # Plot the results
236 plt.figure()
237 plt.plot(df_plot["Acc_train"], "b")
238 plt.plot(df_plot["Acc_val"], "g")
239 plt.title("Accuracy")
240 plt.xlabel("Model Complexity (Number of
241 Variables)")
242 plt.ylabel("Mean Accuracy")
243 plt.legend(["Training data", "Validation Data"])
244 plt.show()
245
246 # Test the final model on the test set
247 means = X_train.mean(axis=0)
248 stds = X_train.std(axis=0)
249 X_train_standardised = (X_train - means) / stds
250 X_test_standardised = (X_test - means) / stds
251
252 X_final_train = X_train_standardised[
253 columns_selected]
254 X_final_test = X_test_standardised[
255 columns_selected]
256
257 final_model = LogisticRegression(C=1e9)
258 final_model.fit(X_final_train, y_train)
259 predicted_y_test = final_model.predict(
260 X_final_test)
261
262 print("Test set accuracy:", accuracy_score(
263 y_test, predicted_y_test))
264 print("Test set precision:", precision_score(
265 y_test, predicted_y_test))
266 print("Test set recall:", recall_score(y_test,
267 predicted_y_test))

```

8.4. Suicide - Deniz Can Dinkci

8.4.1. Creating the CSV as required

```

1 import pandas as pd
2 from sklearn.model_selection import
3 train_test_split
4 from sklearn.impute import SimpleImputer
5 from sklearn.preprocessing import StandardScaler
6 import numpy as np
7 from sklearn.model_selection import
8 train_test_split
9 from sklearn.tree import DecisionTreeClassifier
10 from sklearn.metrics import accuracy_score
11 import matplotlib.pyplot as plt
12
13 # Load the original dataset
14 df = pd.read_csv('binarised_data.csv')
15
16 df = df.drop(columns=['Unnamed: 0', '358
17 _cause_recode'])
18
19 mapping = {
20     '10': '3', '11': '3', '12': '3', '13': '4',
21     '14': '5', '15': '6', '16': '7', '17': '8',
22     '99': '9', '00': '9', '01': '1', '02': '1',
23     '03': '1', '04': '1', '05': '1', '06': '1',
24     '07': '1', '08': '1', '09': '3'
25 }
26
27 # merging two education columns
28 reference_column = 'education_1989_revision'
29 target_column = 'education_2003_revision'

```

```

28
29 def apply_mapping(row):
30     if pd.isna(row[target_column]):
31         left_value = str(row[reference_column])
32         if left_value in mapping:
33             return mapping[left_value]
34     return row[target_column]
35
36 df['education'] = df.apply(apply_mapping, axis
37                             =1)
38
39 df = df.drop(columns=[reference_column,
40                       target_column])
41
42 # Binarize the 'manner_of_death' for my case
43 df['manner_of_death'] = df['manner_of_death'].
44     apply(lambda x: 1 if x == 2.0 else 0)
45
46 non_manner_2_df = df[df['manner_of_death'] == 0]
47 manner_2_df = df[df['manner_of_death'] == 1]
48
49 # Randomly select 10,000 examples from each both
50 # suicide and other deaths
51 random_non_manner_2_sample = non_manner_2_df.
52     sample(n=10000, random_state=42)
53 random_manner_2_sample = manner_2_df.sample(n
54     =10000, random_state=42)
55
56 final_dataset = pd.concat([
57     random_non_manner_2_sample,
58     random_manner_2_sample])
59
60 # Saving the final csv
61 final_dataset.to_csv('final_dataset_20000.csv',
62                     index=False)
63
64 print(final_dataset.head())

```

8.4.2. Splitting the variables

```

1 df = pd.read_csv('final_dataset_20000.csv')
2
3
4 X = df.drop(columns=['manner_of_death'])
5 y = df['manner_of_death']
6
7
8 imputer = SimpleImputer(strategy='mean')
9 X_imputed = imputer.fit_transform(X)
10
11
12 scaler = StandardScaler()
13 X_scaled = scaler.fit_transform(X_imputed)
14
15 # Split the data into training (70%) and
16 # temporary (30%) sets
17 X_train, X_temp, y_train, y_temp =
18     train_test_split(X_scaled, y, test_size=0.3,
19                     shuffle=True, random_state=0)
20
21 # Further split the temporary set into
22 # validation 15% and 15%
23 X_val, X_test, y_val, y_test = train_test_split(
24     X_temp, y_temp, test_size=0.5, random_state
25     =0)
26
27 print(f'Training set: {X_train.shape}, {y_train.
28     shape}')
29 print(f'Validation set: {X_val.shape}, {y_val.
30     shape}')
31 print(f'Test set: {X_test.shape}, {y_test.shape}
32     ')

```

8.4.3. Accuracy vs. Maximum Depth

```

1 max_depths = range(1, 21) # From 1 to 20
2 for max_depth in max_depths:
3     dt = DecisionTreeClassifier(max_depth=
4     max_depth, random_state=0)
5     dt.fit(X_train, y_train)
6     y_pred_train = dt.predict(X_train)
7     y_pred_val = dt.predict(X_val)
8
9     acc_train.append(accuracy_score(y_train,
10     y_pred_train))
11     acc_val.append(accuracy_score(y_val,
12     y_pred_val))
13
14 # Plot the results
15 plt.figure()
16 plt.plot(max_depths, acc_train, 'b', label='
17     Training Data')
18 plt.plot(max_depths, acc_val, 'r', label='
19     Validation Data')
20 plt.xlabel('Maximum Depth')
21 plt.ylabel('Accuracy')
22 plt.legend(['Training Data', 'Validation Data'])
23 plt.title('Accuracy vs. Maximum Depth')
24 plt.show()

```

8.4.4. Precision vs. Maximum Depth

```

1 max_depths = range(1, 21) # From 1 to 20
2 for max_depth in max_depths:
3     dt = DecisionTreeClassifier(max_depth=
4     max_depth, random_state=0)
5     dt.fit(X_train, y_train)
6     y_pred_train = dt.predict(X_train)
7     y_pred_val = dt.predict(X_val)
8
9     prec_train.append(precision_score(y_train,
10     y_pred_train))
11     prec_val.append(precision_score(y_val,
12     y_pred_val))
13
14 # Plot the results
15 plt.figure()
16 plt.plot(max_depths, prec_train, 'b', label='
17     Training Data')
18 plt.plot(max_depths, prec_val, 'r', label='
19     Validation Data')
20 plt.xlabel('Maximum Depth')
21 plt.ylabel('Precision')
22 plt.legend(['Training Data', 'Validation Data'])
23 plt.title('Precision vs. Maximum Depth')
24 plt.show()

```

8.4.5. Recall vs. Maximum Depth

```

1 max_depths = range(1, 21) # From 1 to 20
2 for max_depth in max_depths:
3     dt = DecisionTreeClassifier(max_depth=
4     max_depth, random_state=0)
5     dt.fit(X_train, y_train)
6     y_pred_train = dt.predict(X_train)
7     y_pred_val = dt.predict(X_val)
8
9     rec_train.append(recall_score(y_train,
10     y_pred_train))
11     rec_val.append(recall_score(y_val,
12     y_pred_val))
13
14 # Plot the results
15 plt.figure()
16 plt.plot(max_depths, rec_train, 'b', label='
17     Training Data')
18 plt.plot(max_depths, rec_val, 'r', label='
19     Validation Data')
20 plt.xlabel('Maximum Depth')

```



```

16 plt.ylabel('Recall')
17 plt.legend(['Training Data', 'Validation Data'])
18 plt.title('Recall vs. Maximum Depth')
19 plt.show()

```

8.4.6. Accuracy vs. Minimum Impurity Decrease

```

1 min_impurity_decreases = [i / 1000 for i in
   range(0, 101)] # From 0 to 0.1 in
   increments of 0.001
2 for min_impurity_decrease in
   min_impurity_decreases:
3     dt = DecisionTreeClassifier(max_depth=5,
   min_impurity_decrease=min_impurity_decrease,
   random_state=0)
4     dt.fit(X_train, y_train)
5     y_pred_train = dt.predict(X_train)
6     y_pred_val = dt.predict(X_val)
7
8     acc_train.append(accuracy_score(y_train,
   y_pred_train))
9     acc_val.append(accuracy_score(y_val,
   y_pred_val))
10
11 # Plot the results
12 plt.figure()
13 plt.plot(min_impurity_decreases, acc_train, 'b',
   label='Training Data')
14 plt.plot(min_impurity_decreases, acc_val, 'r',
   label='Validation Data')
15 plt.xlabel('Minimum Impurity Decrease')
16 plt.ylabel('Accuracy')
17 plt.legend(['Training Data', 'Validation Data'])
18 plt.title('Accuracy vs. Minimum Impurity
   Decrease')
19 plt.show()

```

8.4.7. Precision vs. Minimum Impurity Decrease

```

1 min_impurity_decreases = [i / 1000 for i in
   range(0, 101)] # From 0 to 0.1 in
   increments of 0.001
2 for min_impurity_decrease in
   min_impurity_decreases:
3     dt = DecisionTreeClassifier(max_depth=5,
   min_impurity_decrease=min_impurity_decrease,
   random_state=0)
4     dt.fit(X_train, y_train)
5     y_pred_train = dt.predict(X_train)
6     y_pred_val = dt.predict(X_val)
7
8     prec_train.append(precision_score(y_train,
   y_pred_train))
9     prec_val.append(precision_score(y_val,
   y_pred_val))
10
11 # Plot the results
12 plt.figure()
13 plt.plot(min_impurity_decreases, prec_train, 'b',
   label='Training Data')
14 plt.plot(min_impurity_decreases, prec_val, 'r',
   label='Validation Data')
15 plt.xlabel('Minimum Impurity Decrease')
16 plt.ylabel('Precision')
17 plt.legend(['Training Data', 'Validation Data'])
18 plt.title('Precision vs. Minimum Impurity
   Decrease')
19 plt.show()

```

8.4.8. Recall vs. Minimum Impurity Decrease

```

1 min_impurity_decreases = [i / 1000 for i in
   range(0, 101)] # From 0 to 0.1 in
   increments of 0.001

```

```

2 for min_impurity_decrease in
   min_impurity_decreases:
3     dt = DecisionTreeClassifier(max_depth=5,
   min_impurity_decrease=min_impurity_decrease,
   random_state=0)
4     dt.fit(X_train, y_train)
5     y_pred_train = dt.predict(X_train)
6     y_pred_val = dt.predict(X_val)
7
8     rec_train.append(recall_score(y_train,
   y_pred_train))
9     rec_val.append(recall_score(y_val,
   y_pred_val))
10
11 # Plot the results
12 plt.figure()
13 plt.plot(min_impurity_decreases, rec_train, 'b',
   label='Training Data')
14 plt.plot(min_impurity_decreases, rec_val, 'r',
   label='Validation Data')
15 plt.xlabel('Minimum Impurity Decrease')
16 plt.ylabel('Recall')
17 plt.legend(['Training Data', 'Validation Data'])
18 plt.title('Recall vs. Minimum Impurity Decrease')
19 plt.show()

```

8.5. Motor Accidents - Andrew Ize-Iyamu

8.5.1. Balancing and exporting the dataset

```

1 # Filter df2 to create df1_1 containing rows
   where '358_cause_recode' equals 1
2 df1_1 = df2.loc[df2['358_cause_recode'] == 1]
3
4 # Filter df2 to create df1_0 containing rows
   where '358_cause_recode' equals 0
5 df1_0 = df2.loc[df2['358_cause_recode'] == 0]
6
7 # Randomly sample 36233 rows from df1_0 to
   create reduced_df1_0
8 reduced_df1_0 = df1_0.sample(n=36233)
9
10 # Concatenate df1_1 and reduced_df1_0 to create
   the final DataFrame
11 final_df = pd.concat([df1_1, reduced_df1_0])
12
13 # Print the final DataFrame to the console
14 print(final_df)
15
16 # Save the final DataFrame to a CSV file
17 final_df.to_csv('true_final_data.csv', index=
   False)

```

8.5.2. Importing binarised data

```

1 import pandas as pd # Import the pandas library
   for data manipulation
2 import numpy as np # Import the numpy library
   for numerical operations
3 import matplotlib.pyplot as plt # Import the
   matplotlib library for plotting
4
5 # Read the CSV file into a DataFrame
6 data = pd.read_csv('true_final_data.csv')
7
8 # Create a new DataFrame from the loaded data
9 df = pd.DataFrame(data)
10
11 # Drop the 'Unnamed: 0' column from the
   DataFrame
12 df.drop('Unnamed: 0', axis=1, inplace=True)
13
14 # Print the DataFrame to the console
15 print(df)

```

8.5.3. Standardising and splitting the dataset

```

1 from sklearn.model_selection import
  train_test_split
2 import pandas as pd
3
4 # Assuming df is already defined and contains
  your data
5
6 # Define features (X) and target (y)
7 X = df[['resident_status', '
  education_2003_revision',
8         'month_of_death', 'sex', 'age_recode_27',
9         'marital_status', 'day_of_week_of_death',
10        'race_recode_5']]
11 y = df['358_cause_recode']
12
13 # Split the data into train, validation, and
  test sets
14 # First split: train and combined validation/
  test (X_test_val, y_test_val)
15 X_train, X_test_val, y_train, y_test_val =
  train_test_split(X, y, test_size=0.2,
16                  shuffle=True, random_state=0)
17
18 # Second split: split X_test_val and y_test_val
  into validation and test sets
19 X_val, X_test, y_val, y_test = train_test_split(
  X_test_val, y_test_val, test_size=0.5,
20                  random_state=0)
21
22 # Calculate means and standard deviations of
  features from training set
23 X_means = X_train.mean(axis=0)
24 X_stds = X_train.std(axis=0)
25
26 # Standardize the data using the mean and
  standard deviation from the training set
27 X_train = (X_train - X_means) / X_stds
28 X_val = (X_val - X_means) / X_stds
29 X_test = (X_test - X_means) / X_stds
30
31 # Print the shapes of the datasets to verify the
  splits and standardization
32 print('X shape: {}'.format(X.shape))
33 print('X_train shape: {}'.format(X_train.shape))
34 print('X_test shape: {}'.format(X_test.shape))
35 print('X_val shape: {}'.format(X_val.shape))

```

8.5.4. Model 1: Benchmarking

```

1 from sklearn.ensemble import
  RandomForestClassifier
2 from sklearn.metrics import accuracy_score,
  precision_score, recall_score
3
4 # Define the Random Forest classifier with
  specified parameters
5 model = RandomForestClassifier(max_depth=2,
  min_samples_split=0.1, random_state=0)
6
7 # Train the model on the training data
8 model = model.fit(X_train, y_train)
9
10 # Predictions on the training, validation, and
  test sets
11 ypred_train = model.predict(X_train)
12 ypred_val = model.predict(X_val)
13 ypred_test = model.predict(X_test)
14
15 # Calculate accuracy, precision, and recall
  scores for the training set
16 acc_train = accuracy_score(y_train, ypred_train)
17 prec_train = precision_score(y_train,
  ypred_train)
18 rec_train = recall_score(y_train, ypred_train)

```

```

20 # Calculate accuracy, precision, and recall
  scores for the validation set
21 acc_val = accuracy_score(y_val, ypred_val)
22 prec_val = precision_score(y_val, ypred_val)
23 rec_val = recall_score(y_val, ypred_val)
24
25 # Calculate accuracy, precision, and recall
  scores for the test set
26 acc_test = accuracy_score(y_test, ypred_test)
27 prec_test = precision_score(y_test, ypred_test)
28 rec_test = recall_score(y_test, ypred_test)
29
30 # Print the results for the training set
31 print('Performance Metrics for Model 1 against
  Training set:')
32 print('Accuracy Training: {}'.format(acc_train))
33 print('Precision Training: {}'.format(prec_train))
34 print('Recall Training: {}'.format(rec_train))
35 print('')
36
37 # Print the results for the validation set
38 print('Performance Metrics for Model 1 against
  Validation set:')
39 print('Accuracy Validation: {}'.format(acc_val))
40 print('Precision Validation: {}'.format(prec_val))
41 print('Recall Validation: {}'.format(rec_val))
42 print('')

```

8.5.5. Model 2: Plotting Recall of Model 2 against Maximum Depth

```

1 import matplotlib.pyplot as plt
2 import numpy as np
3 from sklearn.ensemble import
  RandomForestClassifier
4 from sklearn.metrics import accuracy_score,
  precision_score, recall_score
5
6 # Lists to store recall scores for training and
  validation sets
7 rec_train = []
8 rec_val = []
9
10 # Loop through different max_depth values from 1
  to 14
11 for i in range(1, 15):
12     # Initialize RandomForestClassifier with
  current max_depth and other parameters
13     dt = RandomForestClassifier(max_depth=i,
  min_samples_split=0.1, random_state=0)
14
15     # Train the model on the training data
16     dt.fit(X_train, y_train)
17
18     # Predictions on training and validation
  sets
19     ypred_train = dt.predict(X_train)
20     ypred_val = dt.predict(X_val)
21
22     # Calculate recall scores and append to
  rec_train and rec_val lists
23     rec_train.append(recall_score(y_train,
  ypred_train))
24     rec_val.append(recall_score(y_val, ypred_val))
25
26 # Plot recall scores for training and validation
  sets across different max_depth values
27 plt.figure()
28 plt.plot(rec_train, 'b', label='Training Data')
29 plt.plot(rec_val, 'r', label='Validation Data')
30 plt.xlabel('Max Depth')
31 plt.ylabel('Recall')
32 plt.legend()
33 plt.xticks(np.arange(0, 14, step=1))
34 plt.grid(True)
35 max_recall = np.max(rec_val)

```

```

36 plt.plot(np.argmax(rec_val), max_recall, marker=
    'x', markersize=15) # Highlight maximum
    recall point
37 plt.title('Recall vs Max Depth')
38 plt.show()
39
40 # Choose the best max_depth based on validation
    performance
41 best_max_depth = np.argmax(rec_val) + 1 # +1
    because max_depth starts from 1
42
43 # Train the final model using the best max_depth
    on the entire training set
44 model = RandomForestClassifier(max_depth=
    best_max_depth, min_samples_split=0.1,
    random_state=0)
45 model.fit(X_train, y_train)
46
47 # Predictions on training, validation, and test
    sets using the final model
48 ypred_train = model.predict(X_train)
49 ypred_val = model.predict(X_val)
50
51 # Calculate accuracy, precision, and recall
    scores for the training set
52 acc_train = accuracy_score(y_train, ypred_train)
53 prec_train = precision_score(y_train,
    ypred_train)
54 rec_train = recall_score(y_train, ypred_train)
55
56 # Calculate accuracy, precision, and recall
    scores for the validation set
57 acc_val = accuracy_score(y_val, ypred_val)
58 prec_val = precision_score(y_val, ypred_val)
59 rec_val = recall_score(y_val, ypred_val)
60
61 # Print the results for the training set
62 print('Performance Metrics for Model 2 against
    Training set:')
63 print('Accuracy Training: {}'.format(acc_train))
64 print('Precision Training: {}'.format(prec_train
    ))
65 print('Recall Training: {}'.format(rec_train))
66 print('')
67
68 # Print the results for the validation set
69 print('Performance Metrics for Model 2 against
    Validation set:')
70 print('Accuracy Validation: {}'.format(acc_val))
71 print('Precision Validation: {}'.format(prec_val
    ))
72 print('Recall Validation: {}'.format(rec_val))
73 print('')

```

8.5.6. Plotting Estimator 4 Decision Tree for Model 2

```

1 from sklearn import tree
2 import matplotlib.pyplot as plt
3
4 # Assuming model is your trained
    RandomForestClassifier
5 # and you want to plot the 4th estimator (index
    3) as an example
6 # Change the index (3) to plot a different tree
    from the Random Forest
7
8 plt.figure(figsize=(4, 4), dpi=800) # Set the
    figure size and DPI
9
10 # Plot the individual tree
11 tree.plot_tree(model.estimators_[3],
12                 feature_names=X_train.columns,
13                 class_names=['Other Fatality', '
    Motor Accident Fatality'],
14                 filled=True);
15
16 plt.savefig('rf_individualtree1.png') # Save
    the figure
17 plt.show() # Display the plot

```

8.5.7. Plotting Random Forest for Model 2

```

1 import matplotlib.pyplot as plt
2 from sklearn import tree
3
4 # Assuming model is your trained
    RandomForestClassifier
5 # and you want to plot the first 5 estimators
6 # Adjust the range (0, 5) to plot different sets
    of trees
7
8 fig, axes = plt.subplots(nrows=1, ncols=5,
9                           figsize=(10, 2), dpi=900) # Adjust figsize
    and dpi as needed
9
10 for index in range(0, 5):
11     tree.plot_tree(model.estimators_[index],
12                    feature_names=X_train.columns
13                    ,
14                    class_names=['Other Fatality',
15                                'Motor Accident Fatality'],
16                    filled=True,
17                    ax=axes[index])
18
19 axes[index].set_title('Estimator: ' + str(
    index + 1), fontsize=11)
20
21 fig.savefig('rf_5trees1.png') # Save the figure
    with all trees
22 plt.show() # Display the plot

```

8.5.8. Using Random Search Cross Validation to find optimal combination of hyperparameters

```

1 from sklearn.ensemble import
    RandomForestClassifier
2 from sklearn.model_selection import
    RandomizedSearchCV
3 from scipy.stats import randint, uniform
4 from sklearn.metrics import recall_score
5
6 # Initialize the Random Forest model
7 rf = RandomForestClassifier(random_state=42)
8
9 # Define the parameter distribution for
    Randomized Search
10 param_dist = {
11     'n_estimators': randint(1, 500), #
    Number of trees in the forest
12     'max_depth': randint(1, 10), #
    Maximum depth of the trees
13     'min_samples_leaf': randint(1, 100), #
    Minimum number of samples required to be at
    a leaf node
14     'min_samples_split': uniform(0.01, 0.49), #
    Minimum number of samples required to split
    an internal node
15     'max_leaf_nodes': randint(1, 300) #
    Maximum number of leaf nodes
16 }
17
18 # Perform Randomized Search with Cross-
    Validation
19 random_search = RandomizedSearchCV(
20     estimator=rf,
21     param_distributions=param_dist,
22     n_iter=200, # Number of parameter settings
    that are sampled
23     scoring='recall', # Use recall as the
    scoring metric
24     cv=5, # Number of cross-validation
    folds
25     verbose=2,
26     random_state=42,
27     n_jobs=-1 # Use all available cores
28 )
29

```

```

30 # Fit the Randomized Search to find the best
    parameters
31 random_search.fit(X_train, y_train)
32
33 # Get the best parameters and best score
34 best_params = random_search.best_params_
35 best_recall = random_search.best_score_
36
37 print(f"Best Parameters: {best_params}")
38 print(f"Best Recall: {best_recall}")

```

```

53 # Print the results for the test set
54 print('Performance Metrics for Model 3 against
    Test set:')
55 print('Accuracy Test: {}'.format(acc_test))
56 print('Precision Test: {}'.format(prec_test))
57 print('Recall Test: {}'.format(rec_test))
58 print('')

```

8.5.10. Plotting Estimator 2 Decision Tree for Final Model

8.5.9. Final Model

```

1 from sklearn.ensemble import
    RandomForestClassifier
2 from sklearn.metrics import accuracy_score,
    precision_score, recall_score
3
4 # Define and initialize the
    RandomForestClassifier with specified
    hyperparameters
5 model = RandomForestClassifier(max_depth=4,
6                               max_leaf_nodes
7                               =36,
8                               min_samples_leaf
9                               =33,
10                              min_samples_split
11                              =0.2,
12                              n_estimators=39,
13                              random_state=42)
14 # Set random_state for reproducibility
15
16 # Fit the model on the training data
17 model.fit(X_train, y_train)
18
19 # Predictions on training set
20 ypred_train = model.predict(X_train)
21
22 # Predictions on validation set
23 ypred_val = model.predict(X_val)
24
25 # Predictions on test set
26 ypred_test = model.predict(X_test)
27
28 # Calculate evaluation metrics for training set
29 acc_train = accuracy_score(y_train, ypred_train)
30 prec_train = precision_score(y_train,
31                              ypred_train)
32 rec_train = recall_score(y_train, ypred_train)
33
34 # Calculate evaluation metrics for validation
    set
35 acc_val = accuracy_score(y_val, ypred_val)
36 prec_val = precision_score(y_val, ypred_val)
37 rec_val = recall_score(y_val, ypred_val)
38
39 # Calculate evaluation metrics for test set
40 acc_test = accuracy_score(y_test, ypred_test)
41 prec_test = precision_score(y_test, ypred_test)
42 rec_test = recall_score(y_test, ypred_test)
43
44 # Print the results for the training set
45 print('Performance Metrics for Model 3 against
    Training set:')
46 print('Accuracy Training: {}'.format(acc_train))
47 print('Precision Training: {}'.format(prec_train
48 ))
49 print('Recall Training: {}'.format(rec_train))
50 print('')
51
52 # Print the results for the validation set
53 print('Performance Metrics for Model 3 against
    Validation set:')
54 print('Accuracy Validation: {}'.format(acc_val))
55 print('Precision Validation: {}'.format(prec_val
56 ))
57 print('Recall Validation: {}'.format(rec_val))
58 print('')

```

```

1 from sklearn import tree
2 import matplotlib.pyplot as plt
3
4 # Assuming 'model' is your trained
    RandomForestClassifier and you want to plot
    the second estimator (index 1)
5 fig, axes = plt.subplots(nrows=1, ncols=1,
6                           figsize=(4, 4), dpi=800)
7
8 # Plot the tree
9 tree.plot_tree(model.estimators_[1], # Change
10                index to plot different trees
11                feature_names=X_train.columns,
12                class_names=('Other Fatality', '
13                Motor Accident Fatality'),
14                filled=True,
15                ax=axes)
16
17 axes.set_title('Estimator 2', fontsize=11) #
18 Set title for the plot
19
20 # Save the figure
21 fig.savefig('rf_individualtree2.png')
22
23 # Display the plot (optional)
24 plt.show()

```

8.5.11. Plotting Random Forest for Final Model

```

1 import matplotlib.pyplot as plt
2 from sklearn import tree
3
4 # Assuming 'model' is your trained
    RandomForestClassifier
5 fig, axes = plt.subplots(nrows=1, ncols=5,
6                           figsize=(20, 4), dpi=900) # Adjust figsize
7 as needed
8
9 # Plot each tree in a separate subplot
10 for index in range(5): # Plotting 5 trees as
11 per your code
12     tree.plot_tree(model.estimators_[index],
13                    feature_names=X_train.columns
14                    ,
15                    class_names=('Other Fatality',
16                    , 'Motor Accident Fatality'),
17                    filled=True,
18                    ax=axes[index])
19
20 axes[index].set_title('Estimator: ' + str(
21 index + 1), fontsize=11)
22
23 # Adjust layout to prevent overlapping
24 plt.tight_layout()
25
26 # Save the figure
27 fig.savefig('rf_5trees2.png')
28
29 # Display the plot (optional)
30 plt.show()

```

8.5.12. Feature Importance Plot

```

1 import numpy as np
2 import matplotlib.pyplot as plt

```



```
3
4 # Assuming 'model' is your trained
   RandomForestClassifier
5 importances = model.feature_importances_
6
7 # Sort the feature importance indices in
   descending order
8 sorted_indices = np.argsort(importances)[::-1]
9
10 # Plotting the feature importances
11 plt.figure(figsize=(10, 6)) # Adjust figsize as
   needed
12 plt.bar(range(X_train.shape[1]), importances[
   sorted_indices], align='center', alpha=0.8)
13 plt.xticks(range(X_train.shape[1]), X_train.
   columns[sorted_indices], rotation=90)
14 plt.xlabel('Features')
15 plt.ylabel('Importance Score')
16 plt.title('Feature Importances')
17
18 plt.tight_layout()
19 plt.show()
```