

Using Machine Learning to Predict Traffic Accidents

Lewis Quick – 22016949, William Forber – 22015706, Yamin Shwe Yi Htay - 23019880

Abstract

In this study, we attempted to predict the severity of traffic accidents in Bristol, UK using machine learning classification. We were aiming to provide a classification between slight and severe accidents, based on the number of vehicles involved in the incident. After selecting Support Vector Machines and Random Forests as our algorithms to use, we tuned the parameters of each to obtain the strongest classification possible. The Random Forest model provided the best performance, with an average test accuracy of 0.609. These results were not as conclusive as we had hoped, but factors of the dataset contributed towards this.

Introduction

As recently as 2022, road traffic accidents attributed to 135,480 casualties in the UK (GOV.UK, 2023). This remains a significant number of people, despite the trend being a general decrease over recent decades. Machine Learning technologies enable versatile classification of many different problems, including humanitarian issues such as this. In this study, we aim to identify factors of traffic accidents and their severity in Bristol, UK using machine learning methods. We aimed to predict the severity of a car accident depending on the number/type of vehicles involved.

Related Work

Support Vector Machines (SVM)

Breast Cancer Prediction (Min-Wei, H., et al., 2017)

This study used SVMs to predict breast cancer susceptibility, focusing on different configurations and kernels, including linear, polynomial, and RBF kernels. The study found linear SVM with bagging worked best for smaller datasets. This aligns with our project, as it also uses a relatively small dataset.

Mountain Freeways in China (Li, J., et al., 2023)

This study applied SVMs to predict the accident severity on mountain freeways in Yunnan Province, China. It also used Random Forest feature selection to boost accuracy, with Ada_SVM achieving high precision and accuracy.

Artificial Neural Networks (ANN)

Road Traffic Accident Prediction (Gatarić, D., et al., 2023)

This study predicted road traffic accidents in Eastern Europe, comparing ANN to logistic regression. The ANN showed positive results, with road length being a significant factor in predicting accidents. In contrast, our project does not use ANN, focusing on SVM for smaller datasets.

Deep Neural Network (Formosa, N., et al., 2020)

This study employed a Deep Neural Network (DNN) to predict accidents based on vehicle telemetry data, using complex inputs such as speed and distance from other vehicles. This study's use of DNN differs from our project, which primarily uses text-based data and focuses on SVM.

Random Forest and Other Algorithms

UK Traffic Accident Prediction (Obasi, I.C. and Benson, C., 2023)

This study analysed 10 years of UK traffic accident data, using Random Forest, Naive Bayes, Logistic Regression, and ANN. Random Forest and Logistic Regression outperformed the other models, achieving high accuracy rates. This study also performed feature importance analysis using Random Forest, identifying key factors like engine capacity and vehicle age. Compared to our project, this study has a broader dataset and uses various machine learning algorithms.

75 **Indian Highways Accident Prediction**
76 **(Khanum, H., Garg, A. and Faheem, M.I.,**
77 **2023)**

78 This study used Random Forest to predict accident
79 severity on Indian highways, achieving moderate
80 accuracy on the training set but lower accuracy on the
81 test set due to possible data imbalances. Our project
82 focuses on SVM with a smaller dataset and k-fold
83 cross-validation to prevent overfitting, unlike the
84 Indian study.

85 **Traffic Accident Severity in China (Yang, J.,**
86 **Han, S. and Chen, Y., 2023)**

87 This study used a large Chinese dataset to predict
88 traffic accident severity, achieving high accuracy with
89 the Random Forest algorithm. It highlighted collision
90 patterns and car structure as key factors,
91 recommending improvements in road infrastructure
92 and driver training. This study's broader scope and
93 higher accuracy with larger datasets differ from our
94 project's more localized focus with SVM on a smaller
95 dataset.

96 ---

97 Compared to these studies, the major differences
98 between all these cases and our project include:
99 **Dataset Size:** Many studies use large datasets covering
100 broader regions (like China, India, and the UK). In
101 contrast, our project uses a smaller dataset from
102 Bristol, UK.

103 **Machine Learning Algorithms:** While other studies
104 use various algorithms (such as Random Forest, ANN,
105 Logistic Regression, and DNN), our project focuses
106 primarily on SVM and Random Forest emphasizing
107 their effectiveness with smaller datasets.

108 **Scope and Focus:** Other studies often have a broader
109 scope, exploring various factors and environmental
110 conditions. Our project is localized, focusing on a
111 specific city with a smaller dataset.

112 **Evaluation Techniques:** Our project uses k-fold cross-
113 validation to prevent overfitting, while other studies
114 might use different evaluation techniques, such as
115 feature importance analysis or SHAP analysis for
116 model interpretation.

117 **Data**

118 The dataset we used for this project was available from
119 the Bristol Council website (Open Data Bristol, 2017).
120 It is a CSV file containing over 4000 records of Bristol
121 traffic accidents between 2017-2021. Each record
122 represents an incident, with several fields to represent
123 details:

Field	Description
Date	Date of accident.
Time	Time of accident.
Severity	Integer value of severity. (3, 2, 1).

Severity Description	Label of severity. (Slight, severe, fatal).
Accident Type	Accident code. (LC, A, HO, ...)
Accident Description	Description of accident code. (Loss of Control, Adult Pedestrian, Head On, ...).
Vehicles	Number of vehicles involved.
Casualties	The number of casualties involved.
Pedestrian	The number of pedestrians involved.
Cycles	Number of cycles involved.
MCycles	The number of motorcycles involved.
Children	Number of children involved.
OAPs	Number of OAPs involved.
X	X coordinates of location.
Y	Y coordinates of location.
Render	The main cause of the accident. (Cars, Cyc, A, ...) {Cars, Cycles, Adult pedestrian, ...} e.g. Many different accident types can be attributed to 'Cars'.

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

Pre-processing the data was required since we were
predicting the severity depending on the number and
types of vehicles involved. To adjust the data to the
classification problem we were solving, we removed
all columns other than:

- Severity.
- Number of vehicles.
- Number of pedestrians involved.
- Number of cycles involved.
- Number of motorcycles involved.

As is accustomed to pre-processing data, we also
checked for any null values in the dataset to ensure that
the models were not trained on erroneous data.

To create a classification on the severity of incidents,
we found it best to convert the 3 categories of severity
into a binary classification. Seeing as there were very
few records classed as 'fatal', we combined these with
the much more abundant 'serious' records to reduce
the number of severity categories to 2. This could in
turn be made into a binary classification – an accident
would be classed as slight or 'not slight' (serious/fatal).
This is optimal for training machine learning
algorithms as it's a very simple classification.

Furthermore, we identified that since there was a large
imbalance in the dataset, with 3861 slight accidents
and only 405 other examples. This would have caused
misclassifications of severe accidents. Therefore, we
under-sampled the dataset removing slight cases until
there was an even split between severe and slight
accidents.

158 **Methods**

159 In consideration of the related works, we found that
160 Support Vector Machines, Random Forests and
161 Artificial Neural Networks should be considered for
162 our research.

163 **Support Vector Machine**

164 A Support Vector Machine is “a supervised machine
165 learning algorithm that classifies data by finding an
166 optimal line or hyperplane that maximizes the distance
167 between each class in an N-dimensional space” (IBM,
168 2023). SVMs are essentially an extension of logistic
169 regression, which is one of the earliest ideas of
170 classification algorithms. They are known for their
171 effective classification across small to medium-sized
172 datasets, via a relatively simple implementation. Large
173 datasets are to be avoided, as the algorithm can be
174 computationally expensive. This is particularly true
175 when the classification is completed on a multi-
176 dimensional dataset. To improve this performance,
177 most SVMs utilise kernels. Kernels supplement SVMs
178 with mathematical functions in the algorithm which
179 simulate distances between data points, rather than an
180 actual calculation of said distances. Different kernels
181 are found to produce different results, which is why
182 kernel selection is a key part of SVM experimentation.

183 **Random Forest**

184 An ensemble is a collection of classifiers each trained
185 using different parts of a dataset, resulting in an
186 aggregation of the output of each of the classifiers. A
187 common type of ensemble method is a Random Forest
188 Classifier. This method utilises an ensemble of
189 **Decision Trees** and combines their results with the
190 ambition of producing a classification (IBM, 2023).
191 The nuance in the method is in the behaviour of the
192 decision trees, and how the data is handled between the
193 trees. Random Forests are known for their low
194 computational cost and effectiveness on relatively
195 simple classifications. This makes them a good fit for
196 quick and easy classification of small to medium
197 datasets.

198 **Artificial Neural Network**

199 ANNs are a highly sophisticated classification
200 methodology. They consist of multiple ‘layers’ which
201 represent the data flow. Between an input and output
202 layer, there are many ‘hidden layers’. The specification
203 of these hidden layers dictates the model’s output. For
204 example, certain hidden layers can be set to analyse
205 certain portions of the input. Traits like this make the
206 methodology useful for the classification of highly
207 complex input, such as images. ANNs are best used on
208 larger datasets, as they’ve been seen to overfit with too
209 little data (BuiltIn, 2023).

210 ---

211 In consideration of the previous 3 methods, we chose
212 to use SVMs and Random Forest classifiers on our
213 dataset. The ANNs would have proved to be a lengthy

214 and difficult implementation and would have yielded
215 sub-par results on our small dataset.

216 **Grid Search**

217 To optimize the performance of both SVM and
218 Random Forest models we used grid search, which is a
219 common technique for identifying the best
220 hyperparameters. Grid search involves defining a
221 range of potential values for various hyperparameters
222 and systematically exploring these combinations to
223 identify the optimal set of variables. Considering it is
224 computationally expensive, grid search is only viable
225 for a relatively small number of possible variables.
226 However, it must be said that it is far more effective
227 than a manual ‘trial and error’ test of the variables.

228
229 For each SVM kernel variant (linear, polynomial, and
230 Radial Basis Function), we defined parameter grids for
231 the key hyperparameters. These typically included the
232 regularization (C), the degree of polynomial kernel,
233 and the gamma parameter for non-linear kernels. By
234 conducting a comprehensive search across different
235 values, we aimed to identify the combination that
236 yielded the best performance.

237
238 The grid search process was paired with cross-
239 validation to ensure that the tuning was robust and not
240 overly reliant on specific data splits. This approach
241 helped us identify hyperparameters that consistently
242 performed well across different training and testing
243 datasets, leading to models with greater accuracy and
244 reduced risk of overfitting.

245 **Performance Evaluation**

246 **K-fold Cross Validation**

247 K-fold cross-validation is where the dataset is split k-
248 number of times and the model is trained on each split
249 dataset (Anguita *et al.*, 2012), the results are then
250 averaged giving a more stable and reliable estimate of
251 model accuracy and precision. Using k-fold cross-
252 validation helped us ensure that our evaluation process
253 was robust and minimized the risk of overfitting,
254 which is when a model performs well on training data
255 but poorly on unseen data and allowed us to compare
256 how each model performs with different data and
257 consequently how the model would perform in a real-
258 world scenario.

259 **Confusion Matrix**

260 We used a confusion matrix coupled with k-fold cross-
261 validation to evaluate different models. A confusion
262 matrix is used to measure the performance of a
263 classifier numbering the occurrences of true positives
264 against false positive cases and true negatives against
265 false negative cases. This allowed us to easily compare
266 the classification performance of each model.

267 To enhance our evaluation, we calculated additional
268 metrics like precision, recall, F1-score, and accuracy.
269 Precision shows the proportion of correct positive

270 predictions, recall reflects how many actual positives
271 were identified, and the F1-score balances precision
272 and recall. Accuracy provides an overall measure of
273 correct predictions.

274 Languages & Libraries

275 For software tools and libraries, we used Python as the
276 primary language, with a range of supporting libraries
277 for data analysis, machine learning, and data
278 visualization. Python's '**Scikit-learn**' served as a
279 highly useful module for all aspects of machine
280 learning. It provides easy application of classifiers like
281 '**SVC**', '**RandomForestClassifier**', and
282 '**AdaBoostClassifier**', offers utilities for data splitting
283 ('**train_test_split**'), hyperparameter
284 tuning ('**GridSearchCV**'), and model evaluation
285 ('**confusion_matrix**', '**accuracy_score**', and
286 '**classification_report**'). For data manipulation, we
287 used '**pandas**', which allowed us to pre-process, filter,
288 and manage the dataset effectively. '**Numpy**' was
289 employed for numerical operations, while '**seaborn**'
290 and '**matplotlib**' were used for data visualization and
291 plotting results. Together, these tools created a
292 comprehensive and efficient environment for our
293 project, enabling us to build, evaluate and optimise our
294 models with flexibility and ease.

295 Ethical Considerations

296 When working with data related to traffic accidents,
297 ensuring ethical standards is critical. Our project used
298 a dataset from the Bristol Council, which included
299 details about traffic accidents, such as date, time,
300 severity, vehicle types, and casualties. Given the
301 potential sensitivity of this information, we took
302 several steps to ensure the data was anonymized and
303 that no personally identifiable information (PII) could
304 be used to discriminate against any individual or
305 group.

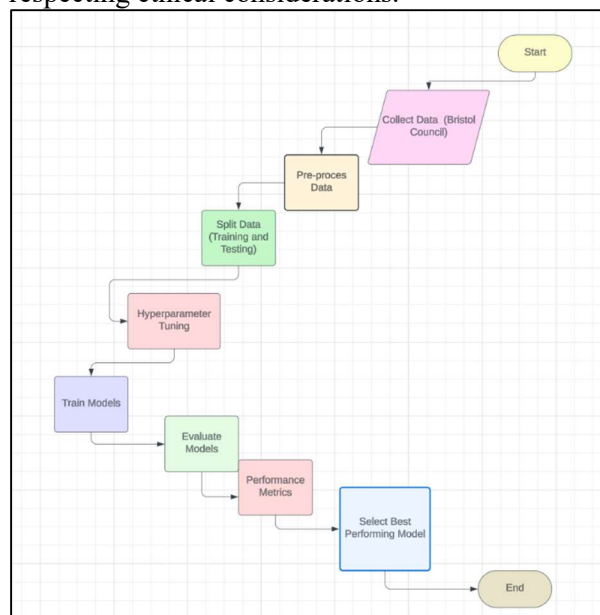
306 Anonymization of Data (Regulation (EU) 307 2016/679 of the European Parliament and of 308 the Council, 2016)

309 Any personal information, such as names, addresses,
310 or vehicle registration numbers, was not included in
311 the dataset. While location data was present, it was
312 generalized to avoid pinpointing specific addresses or
313 exact locations where accidents occurred. Any
314 demographic data that could identify individuals, such
315 as age or gender, was not described to prevent
316 discrimination against specific groups.

317 Avoiding Discrimination

318 The data was processed and analysed without
319 reference to protected characteristics, focusing solely
320 on factors related to traffic accidents. We avoided any
321 analysis that could lead to biased conclusions or
322 discriminatory practices (e.g. We did not use
323 demographic data to conclude accident causes.). The
324 project complied with ethical guidelines and

325 regulations, ensuring that all data handling, processing,
326 and analysis were conducted concerning privacy and
327 non-discrimination.
328 By focusing on anonymization and avoiding
329 discrimination, we aimed to contribute to a broader
330 understanding of traffic accident patterns while
331 respecting ethical considerations.



332 Fig 1.0 – The machine learning workflow which we
333 planned to follow.
334

335 Experiments

336 We experimented using different kernels for the
337 Support Vector Machine and comparing them to a
338 Random Forest classifier.

339 We started by training each model using
340 *train_test_split* from *Sklearns*, which splits the dataset
341 into training and test data, the test data to be used to
342 evaluate each model. Furthermore, the random state
343 parameter was set to 0 so we could tune each model
344 using the same data to get repeatable results.
345 Afterwards, we used the *make_pipeline* feature from
346 *Sklearns* and scaled the data using the standardization
347 technique. To get a control baseline each model was
348 evaluated with no hyperparameters specified using the
349 test data mentioned previously. We then started tuning
350 the hyperparameters for each model.
351

352 Hyperparameter Tuning

353 Linear Kernel SVM

354 The linear kernel Support Vector Machine only
355 hyperparameter that can be tuned is the **regularization**
356 parameter (how much error is allowed). We started by
357 running a grid search which incremented the
358 regularization parameter in factors of 10, which
359 allowed us to narrow down on a range of values.

360 `parameter_grid = {'C': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}`
361 Here, the best value for the hyperparameter was 1.
362 With this information, we used the *linspace* feature

from *Numpy* to search through 100 values in the range 0.1 to 10.

```
{'C': np.linspace(0.1, 10, num=100, dtype=float)}
```

This search found that the best value was:

- regularization, C = 3.**

Polynomial SVM

This kernel is affected by more than one hyperparameter, being the **regularization, degree** (the complexity of the model), and **gamma** (kernel coefficient). Since there were more parameters to tune than the previous method, to get the best possible hyperparameters using *linspace* will simply take too long to complete. Therefore, we instead created a grid search with the **regularization** and **gamma** hyperparameters incrementing in factors of 10 and the degree hyperparameter incrementing in steps of 1.

```
{'C': [0.001, 0.01, 0.1, 1, 10, 100],  
'degree': [1, 2, 3, 4],  
'gamma': [0.01, 0.1, 1, 10]}
```

This grid search found that the best hyperparameters for the polynomial kernel were:

- regularization, C = 0.01**
- degree = 2**
- gamma = 1**

Radial Basis Function SVM

The RBF SVM has two hyperparameters which can be adjusted, being the **regularization** and **gamma**. This means that we can follow the same method as for the linear kernel for tuning the hyperparameters, starting with narrowing down the range of possible best values and finishing using *linspace* to find the best possible combination. So, we started by creating a grid search to search the **regularization** and **gamma** parameters incrementing in factors of 10.

```
{'C': [0.01, 0.1, 1, 10, 100, 1000],  
'gamma': [0.001, 0.01, 0.1, 1, 10, 100, 1000]}
```

This search found the best combination to be **C = 10**, and **gamma = 0.1**. We then used *linspace* feature to generate 100 values between 1 and 100 for **C**, and between 0.01 and 1 for the **gamma**.

```
{'C': np.linspace(1, 100, num=100, dtype=float)  
'gamma': np.linspace(0.01, 1, num=100, dtype=float)}
```

This grid search found that the best hyperparameters for the RBF SVM were:

- regularization, C = 18.**
- gamma = 0.09.**

Random Forest

Random Forests have 4 different hyperparameters to tune, being **n_estimators** (number of trees in the forest), **max_features** (the number of features to consider when looking for the best split), **max_depth** (the maximum depth of the tree) and **max_leaf_nodes** (grows the tree by the number of max_leaf_nodes).

Because there are many hyperparameters in Random Forests, it was infeasible to use *linspace* as it would have taken too much time. We used a grid search with **n_estimators** in increments of 25, with **max_depth** and **max_leaf_nodes** searched through in increments of 3.

```
'n_estimators': [25, 50, 100, 150],  
'max_features': ['sqrt', 'log2', None],  
'max_depth': [3, 6, 9],  
'max_leaf_nodes': [3, 6, 9],
```

The results from this search concluded that the best parameters for the Random Forest classifier were:

- max_depth = 3.**
- max_features = 'sqrt'.**
- max_leaf_nodes = 9.**
- n_estimators = 25.**

Performance Evaluation

To evaluate the models, we used k-fold cross-validation coupled with a confusion matrix (see more detail in [Methods](#)). We focused the most on the f1-score, which is the combined average of precision and recall (GeeksForGeeks, 2023), and average test accuracy.

Results

Model	F1 score positive	F1 score negative	Avg. Test Accuracy
Linear SVM	0.59	0.64	0.584
Polynomial SVM	0.59	0.62	0.584
RBF SVM	0.65	0.62	0.602
Random Forest	0.64	0.65	0.609

The best model was the Random Forest classifier (ensemble model), producing the highest f1-score of 0.65 for negative classifications and 0.64 for positive classifications. The Random Forest classifier also had the highest average test accuracy of 0.609.

The runner-up model was the Radial Basis Function Support Vector Machine, which had a lower f1-score of 0.62 for negative classifications and 0.65 for positive classifications. The model had a lower test average accuracy of 0.602.

The worst model was the polynomial SVM, producing the lowest f1-score of 0.62 for negative classifications and 0.59 for positive classifications. This was coupled with the tied lowest test average accuracy of 0.584 with the linear kernel SVM, which had a higher f1-score for negative classifications of 0.64 but had the same f1-score for positive classifications of 0.59.

Conclusion

Through the use of both SVMs and Random Forest classifiers, we were unable to make as clear of a classification as we'd have liked. This is even though

both models were optimised effectively using grid search. However, the models did produce correct classifications for the majority of cases.

A lesson we have learned is to pay closer attention to the dataset, and the range of data which it holds. The results of the models may have been negatively affected by the under-sampling used on the dataset. While this was necessary for the training on this dataset, ideally it would have been avoided. Clearly, it is best to have the highest number of cases possible when training a model, as small datasets are always prone to producing inaccurate or invalid results.

It is unlikely that our choice of algorithm affected the results. An ANN would have likely suffered more than our SVM and RF because, as mentioned in Methods, they are prone to overfitting on small datasets.

References

Anguita, D., Ghelardoni, L., Ghio, A., Oneta, L. and Ridella, S. (2012) The 'K' in K-fold Cross Validation. *Esann 2012 Proceedings, European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning* [online]. [Accessed 15 April 2024].

BuiltIn (2023) 4 Disadvantages of Neural Networks. Available from: <https://builtin.com/data-science/disadvantages-neural-networks> [Accessed 15/03/2024].

GeeksforGeeks (2023) F1 Score in Machine Learning. Available from: <https://www.geeksforgeeks.org/f1-score-in-machine-learning/>.

Formosa, N., Quddus, M., Ison, S., Abdel-Aty, M. and Yuan, J. (2020) Predicting Real-time Traffic Conflicts Using Deep Learning. *Accident Analysis & Prevention* [online]. 136 [Accessed 23 April 2024].

Gatarić, D., Ruškić, N., Aleksić, B., Đurić, T., Pezo, L., Lončar, B., Pezo, M. (2023) Predicting Road Traffic Accidents - Artificial Neural Network Approach. *Algorithms* [online]. pp. 257. 16 (5). [Accessed 17 April 2024].

GOV.UK (2023) Reported road casualties Great Britain, annual report: 2022. Available from: <https://www.gov.uk/government/statistics/reported-road-casualties-great-britain-annual-report-2022/reported-road-casualties-great-britain-annual-report-2022#:~:text=In%20reported%20road%20collisions%20in,of%2012%25%20compared%20to%202019> [Accessed 23/04/2024].

IBM (2023) What are Support Vector Machines (SVMs)? Available from:

<https://www.ibm.com/topics/support-vector-machine> [Accessed 30/04/2024].

IBM (2023) What is Random Forest? Available from: <https://www.ibm.com/topics/random-forest> [Accessed 30/04/2024].

Khanum, H., Garg, A. and Faheem, M.I. (2023) Accident Severity Prediction Modelling For Road Safety Using Random Forest Algorithm: *An Analysis of Indian Highways*. F1000research [online]. 2 [Accessed 23 April 2024].

Li, J., Guo, F., Zhou, Y., Yang, W. and Ni, D. (2023) Predicting the Severity of Traffic Accidents on Mountain Freeways with Dynamic Traffic and Weather Data. *Transportation Safety and Environment* [online]. 5 (4) [Accessed 23 April 2024].

Min-Wei, H., Chen, C., Wei-Chao, L., Shih-Wen, K. and Chih-Fong, T. (2017) SVM and SVM Ensembles in Breast Cancer Prediction. *PLoS One* [online]. 12 (1) [Accessed 16 April 2024].

Obasi, I.C. and Benson, C. (2023) Evaluating the Effectiveness of Machine Learning Techniques in Forecasting the Severity of Traffic Accidents. *Heliyon* [online]. 9 (8) [Accessed 23 April 2024].

Open Data Bristol (2017) Traffic Accidents. Available from: <https://opendata.bristol.gov.uk/datasets/bcc::traffic-accidents-1/about> [Accessed 05/03/2024].

Regulation (EU) 2016/679 of the European Parliament and of the Council [online]. Chapter 1. (2016) legislation.gov.uk. Available from: <https://www.legislation.gov.uk/eur/2016/679/contents#/> [Accessed 04 April 2024].

Yang, J., Han, S. and Chen, Y. (2023) Prediction of Traffic Accident Severity Based on Random Forest. *Journal of Advanced Transportation* [online]. [Accessed 23 April 2024].