



THE CHANCE OF A TRAFFIC COLLISION IS PREDICTED USING MACHINE LEARNING

J C Achutha^{1*}, Raveena Y², Rekha H³, Rekha S⁴, Premalatha V⁵, Pavithra D K⁶

Article History: Received: 15.08.2022

Revised: 16.10.2022

Accepted: 22.12.2022

Abstract

Vehicle collisions are still as among the largest sources of fatalities, disabling injuries, and being hospitalised in the nation. Due to this, it is crucial to foresee the possibility that collisions while driving with the aim to prevent them from happening and protect people's lives. To achieve a single objective, a number of models—from conventional statistical techniques to contemporary modeled inspired by the development of automated learning—have ever been suggested. Plenty of those models are compared in this work, we want to evaluate and draw an effective technique for circulation incident risk predictions. The research investigation intends to anticipate the likelihood of traffic accidents. for the reason passengers control the path they take. to drivers by examining information they would be aware of beforehand, such as vehicle type age, gender, time of day, weather, and so forth. The use of Ideal The categorization Plants is a framework that could end in findings that will make obvious sense to the chauffeur in addition to using Logistic Regression and Random Forest. Additionally, geospatial analysis of data utilising the K-means clustering technique can offer information about high-accident locations. The road has gotten more difficult in the design and management sectors as the number of cars upon the highway has grown. Given that Automobile collisions are a major global source of worry since they have a considerable impact on people's safety, health, and well-being. According to the World Health Organisation (WHO), 1.35 million people die in automobile accidents each period. As a consequence, they represent a substantial field of research for the use of cutting-edge methodologies as well as prediction and assessment algorithm design. Despite how common automobile crashes are caused by external factors, some are caused by the driver. Unfavourable weather conditions, such as Operating a vehicle on these types of roads is challenging and frequently fatal because of factors like poor visibility caused by rainfall skies, and fog. The system of model for forecasting used today evaluated only many probable causal factors.

Keywords: Road Accident, Traffic Accident, Machine Learning, K-Means, Geo-Location.

^{1*}Asst. Professor, Department of MCA, the Oxford College of Engineering, Bengaluru, Karnataka, India – 560068
^{2,3,4,5,6}MCA Final Year, ⁶Department of MCA, The Oxford College of Engineering, Bengaluru, Karnataka, India – 560068

Corresponding Author: J C Achutha^{1}

Email: ^{1*}achutha.sir@gmail.com

^{1*}Asst. Professor, Department of MCA, the Oxford College of Engineering, Bengaluru, Karnataka, India – 560068

DOI: 10.31838/ecb/2022.11.12.48

1. Introduction

The endeavour for predicting the risk of traffic accidents is designed to create a system that uses Machine Learning algorithms to forecast Potential for collisions between vehicles. There are many collisions while driving, public safety problem, with millions of incidents happening each year throughout the world, resulting in thousands of fatalities and injuries. Existing traffic accident prediction systems frequently depend on statistical models and heuristic-based techniques that may be incapable of capturing the underlying data's complexity. Machine Learning algorithms, contrasted with, provide a more data-driven approach that can analyse vast datasets and uncover patterns and correlations that are not immediately obvious. The suggested system would analyse historic traffic collision Through the use of a random forest classifier on this information algorithm, which will include elements such as location, time, weather conditions, road ailments, among others important aspects. The algorithm will be optimised to deliver reliable forecasts of the likelihood of a crash with traffic will feature a user-friendly interface that will allow users to enter pertinent information and obtain a risk score indicating the chance of a traffic accident occurring. The technology has the ability to improve road safety and minimise the frequency of collection mischance by giving useful insights on the risk variables associated with traffic accidents to drivers, transportation authorities, and other stakeholders. Furthermore, the system can aid in the creation of targeted initiatives and policies aimed at alleviating traffic congestion.

Literature Survey

Vehicle crashes [1] are A major field of study is studying the usage of new algorithms and methodologies to evaluate and anticipate vehicle crashes, and also pinpoint the major contributing elements, because these incidents are a single leading causes of injury and death globally. road accidents. The goal of road accident prediction research is to answer to the problem of creating a more safe transportation environment and, eventually, saving lives. The purpose of this paper is to give a description of the contemporary in vehicle collision detection using deep learning frameworks, training computational methods, as sophisticated data assessment methods, including convolutional neural networks and long short-term memory networking. Also, this page compiles and studies the greatest popular information places for road accident forecasting. Also proposed is a classification based on its origin and features, such as open data, measuring technology, onboard equipment, and social media data. For information analysis, the many methods used to forecast give a description on the contemporary in vehicle collision detection

leveraging deep learning structures, training computational methods, as well as sophisticated data assessment methods, including convolutional neural networks and long short-term memory networking. Also, In this study, a systematic technique for identifying No connection between the prevalence of strong stop events [2] as well as where and when they happened is given. The technique being suggested combines previous and current information in order to predict the timing and locations of incidents involving disastrous interruption. It depends on computations for batch grouping and real-time grouping. Subtractive gathering and fuzzy c-means clustering are used in conjunction with sequential sorting to create groupings that mirror the preliminary association signals. No connection between the prevalence of strong stop events [2] as well as where and when they happened is given. The technique being suggested combines previous and current information to anticipate when and where disastrous interruption incidents will occur. It relies on computations for batches grouping and real-time grouping. Subtractive grouping and fuzzy c-means clustering are used in conjunction in sequential sorting to create groupings that mirror the preliminary association signals. Road and traffic accidents [3] are a major worry all over the world. Road accidents not only endanger public health by causing varying degrees of harm, but they also cause property loss. Data analysis can determine the many causes of road accidents, such as traffic factors, weather characteristics, road features, and so on. Numerous researches on road accident data analysis have previously demonstrated its significance. Some research concentrated on identifying factors related with accident severity, whereas others concentrated on identifying elements connected to accident occurrence. Traditional statistical approaches as symptomless as Data analysis strategies were utilised in these study investigations. In current study, data mining is a popular tool for analysing road accident data. Another major study field in the world of road accidents is trend analysis.

Data mining [4] showed his competence a viable approach for analysing road accidents and producing useful findings. The majority of road accident data analysis use data mining techniques, with the goal of discovering characteristics that influence the severity of an accident. Any harm caused by a car accident, however, is always undesirable at relation to health, property loss, and other economic issues. It has been shown that traffic accidents occur more often in some areas. The examination of these locations might assist in locating important crashes causes that eventually lead to them frequency of road accidents in these areas. Association rule mining is a prominent data mining approach for identifying correlations in various aspects of a traffic

accident. We originally used in this paper Road traffic accidents[5] are one among the main causes of death and injury worldwide. Abu Dhabi, United Arab Emirates. in 2014, 971 road incidents occurred, resulting in 121 fatalities and 135 serious injuries. Several variables, including driver-related factors, road-related factors, and accident-related factors, all contribute to the degree of seriousness for a wound. In the present investigation, systems (classifiers)

were built based on 5,973 traffic accident records in Abu Dhabi between 2008 and 2013, in order to estimate the degree of harm of each individual lesion. accident with fair accuracy. Furthermore, the study attempted to develop A compilation of guidelines the United Arab Emirates (UAE) Traffic Agencies may use to identify the key factors affecting crash damage. Known as WEKA (Waikato Environment for Knowledge Acquisition).

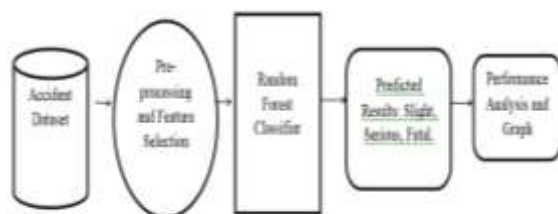


Fig. 1. Proposed Architecture

Existing Model

The current system study analysed and studied several approaches for applying it and found that collection misadventure hazard projections generated from previously known to drivers elements such as personal descriptors, vehicle descriptors, and location made a lot of intuitive sense. Once individuals are aware of such high risk variables, they have a certain degree of power to lessen the danger. Because Those who are that were driving charge on the road, having this information can help them make better decisions about their trips, reducing the likelihood potential preventing deaths from car crashes. The current system in use The K-means clustering technique is used. The present method analyses accident geolocation data using clustering algorithms such as K-means to categorise them into high risk hotspots in a specific region. Once gathered, clusters can be subjected to a classification algorithm to discover which characteristics are responsible for raising the risk. These factors might include time of day, road conditions, weather conditions, and so on.

Sensitivity to initialization: K-means clustering is extremely sensitive to centroids' initial location. If the starting centroids are not correctly positioned, the method may converge to a suboptimal solution that does not adequately represent the underlying data distribution.

K-means clustering has a limited application because it is only applicable to datasets with a spherical or circular form. It might not be appropriate for datasets with uneven forms or clusters with changing densities.

Scalability: K-means clustering is not scalable to huge datasets since processing massive amounts of

data demands a substantial amount of memory and computer resources.

Cluster size and form bias: K-means clustering assumes isotropic clusters with equal variances. This may not be the case for all datasets, resulting in a skewed cluster size.

Outlier sensitivity: K-means clustering is sensitive to outliers and noisy data points. Outliers might cause misleading clusters to develop or the algorithm to converge to poor solutions.

Difficulty in establishing the appropriate number of clusters: The optimal number of clusters is not always known a priori and may need trial and error. This might take time and may need specialist knowledge or subject competence.

2. Proposed Methodology

The suggested system is This Random Forest Classifier-based method for predicting the chances of crash events to occur algorithm. The system tries to increase the precision that traffic accident risk prediction compared to previous systems, and it also includes a geo-location component, which is significant. The suggested system would make use of a collection of the past collision data, acquired through the Kaggle archive, with location services, date, time, weather, roadway conditions, and other relevant details. The data will be cleaned and transformed into features that capture the data's underlying patterns and trends. Upon python extracted statistics, the method for the randomly generated forest Classifiers will be generated. And will generate predictions using an ensemble of decision trees. The algorithm will be tuned to maximise accuracy while minimising mistake rate. It uses a forest-based Compressor. Algorithm's performance will be tested using several

measurements like F1-score, recall, reliability, and exactitude. The system will also have a user interface that will allow users to enter important information such as location, time, the environment, and road conditions to obtain a risk score indicating the possibility of a traffic accident occurring.

Improved accuracy: The Random Forest Classifier algorithm is very accurate and can accurately forecast the risk of traffic accidents. This can help minimise the amount of traffic accidents a car conditioner depending in forests is used. And save lives.

Robustness: Because the Random Forest Classifier method is resistant to noise and outliers in the data, it is less prone to mistakes and can handle a larger variety of data.

Non-linear connections: The method can handle non-linear relationships that may exist in the underlying data between characteristics and the output variable. As a result, the system can capture more complicated patterns and enhance its accuracy.

Interpretability: The algorithm gives insights into the most essential factors that contribute to traffic accident risk. This information is available.

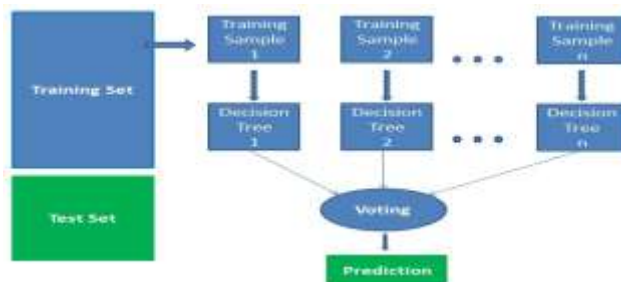
Scalability: Because the Random Forest Classifier method is very scalable, it can handle big datasets with millions of data points. This enables more precise and efficient prediction of traffic accident risk.

The system has a user-friendly interface that allows users to enter important information and obtain a risk score indicating the possibility of a traffic accident occurring.

Implementation

Data Collection:

Collecting data is the first real step towards the actual construction of a machine learning model. This is a vital phase that will have a knock-on effect on how successful the model is; the more and better



Feature

A random forest is also an effective feature selection indication. With the model, Scikit-learn includes an additional variable that reflects the relative value or contribution of each attribute to the prediction. During the training phase, it automatically computes the relevance score of each feature. The relevance is then scaled down so that the total of all ratings is 1.

data we have, the better our model will perform. There are numerous methods for gathering data, including online scraping and manual interventions. The dataset may be found in the model folder. The dataset is from the well-known dataset repository kaggle.

Preprocessing of Data:

Gather and arrange data for training. Clean up everything that needs it (remove duplicates, rectify mistakes, deal with missing numbers, normalisation, data type conversions, and so on).

Randomise data to remove the impacts of the sequence in which we acquired and/or otherwise prepared our data.

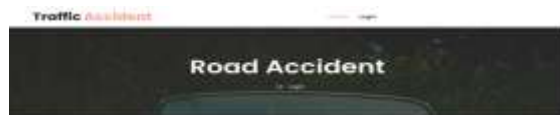
Visualise data to aid in the detection of meaningful correlations between variables or class imbalances (bias alert!), or do other exploratory analysis.

Sets are divided into training and assessment sets.

Model Selection:

- The Random Nature Classifier machine learning algorithm was utilised. We applied this method after achieving an accuracy of 99.0% on the training set.
- The Algorithm of Random Forests
- Let's go out the algorithm in layman's words. Assume you want to go on a trip and you want to go somewhere you would love.
- So, how do you go about finding a place you'll like? You may conduct an internet search, read reviews on travel blogs and websites, or ask your friends.
- Assume you chose to question your pals and asked them about their previous trip experiences to various locations. Every buddy will give you some recommendations. You must now construct a list of the recommended locations.

This score will assist you in selecting the most vital characteristics and eliminating the least important ones for model development. Random forest calculates the relevance of each feature using Gini importance or mean reduction in impurity (MDI). Gini significance is often referred to as the complete decrease in node impurity.



Road Accident Prediction

Input fields for prediction, including a red 'Predict' button.



Road Accident Prediction

Input fields for prediction, including a red 'Predict' button.

Road Accident Prediction

id	latitude	longitude	age_group_of_driver	sex_of_driver	vehicle_type	road_type	junction_conditions	weather_conditions	light_conditions	road_surface_conditions	urban_or_rural_area	accident_type	severity
1	40.7128	-87.6298	21-30	Male	Car	Highway	Clear	Sunny	Daylight	Asphalt	Urban	Minor	Low
2	40.7128	-87.6298	31-40	Female	Truck	City Street	Clear	Sunny	Daylight	Asphalt	Urban	Minor	Low
3	40.7128	-87.6298	41-50	Male	Car	Highway	Clear	Sunny	Daylight	Asphalt	Urban	Minor	Low
4	40.7128	-87.6298	51-60	Female	Truck	City Street	Clear	Sunny	Daylight	Asphalt	Urban	Minor	Low

Road Accident Prediction

Form with input fields for: Latitude, Longitude, Age_group_of_driver, Sex_of_driver, Vehicle_Type, Road_Type, Junction_Conditions, Weather_Conditions, Light_Conditions, Road_Surface_Conditions, and Urban_or_Rural_Area.



Road Accident Prediction

Form with input fields for: Latitude, Longitude, Age_group_of_driver, Sex_of_driver, Vehicle_Type, Road_Type, Junction_Conditions, Weather_Conditions, Light_Conditions, Road_Surface_Conditions, and Urban_or_Rural_Area.

Slight: Slightly Injured Will Occur So Drive Carefully



Road Accident Prediction



3. Conclusions

Finally, the Traffic Accident Risk Prediction using Machine Learning research has demonstrated considerable promise for enhancing road safety and lowering the frequency of traffic accidents. The suggested method can reliably forecast the likelihood of traffic accidents and identify the most relevant risk variables by analysing historical traffic accident data and applying powerful Machine Learning techniques such as Random Forest Classifier. Future work might concentrate on increasing system performance by combining other data sources, such as traffic flow data and driver behaviour data, as well as integrating real-time data to offer up-to-date risk rankings. Furthermore, the system might be improved to deliver more personalised risk rankings based on individual driver profiles, perhaps encouraging safer driving behaviour. Overall, the Traffic Accident Risk Prediction project has shown the anticipated benefits of neural networks for boosting crash prevention and lowering collisions, plus relevance of adopting data-driven techniques to address complex transportation concerns.

4. References

1. Camilo Gutierrez-Osorio and César Pedraza, "Modern data sources and techniques for analysis and forecasting of road accidents: A review." 7.4 (2020): 432-446 in the Journal of Traffic & Transportation Engineering (English Edition).
2. G. Cao, J. Michelini, K. Grigoriadis, B. Ebrahimi, and M. A. Franchek, "Cluster-based correlation of severe braking events with time and location," 2015, pp. 187-192, doi: 10.1109/SYSOSE.2015.7151986.
3. Kumar, S., and D. Toshniwal (2016). Hierarchical clustering and the cophenetic correlation coefficient (CPCC) were used to analyse hourly traffic accident numbers. 1-11 in Journal of Big Data, 3(1).
4. Kumar, S., and D. Toshniwal (2016). A data mining strategy to characterising the sites of traffic accidents. 62-72 in Journal of Modern Transportation.
5. Taamneh, M., S. Alkheder, and S. Taamneh (2017). In the United Arab Emirates, data mining techniques are being used to model and forecast traffic accidents. Transportation Safety & Security, 9(2), pp. 146-166.
6. Tiwari, P., Dao, H., and G. N. Nguyen (2017). On traffic accident analysis, the performance of slow, decision tree classifier, and multilayer perceptron is evaluated. 41(1), Informatica.
7. Ait-Mlouk, A., and T. Agouti (2019). A case study on a road accident using DM-MCDA, a web-based tool for data mining and multiple criteria decision analysis. SoftwareX, vol. 10, no. 100323.