# Road Accident Analysis and Prediction using Machine Learning Algorithmic Approaches

## Venkata Koteswara Rao Ballamudi

Solution Architect, Software Engineering Architecture Department, Cyber Global, Inc., **USA**

**\*Email for correspondence:**
koteswarabv@gmail.com

## ABSTRACT

Ongoing studies have anticipated that in 2030, car crashes will be the fifth driving reason for death around the world. The main cause of car crashes is difficult to decide these days because of a complex mix of qualities like the mental condition of the driver, road conditions, climate conditions, traffic, and infringement of traffic rules to give some examples. The expenses of fatalities and driver wounds because of car crashes incredibly influence the general public. The use of machine learning methods in the field of road accidents is picking up speed nowadays. The organization of machine learning classifiers has swapped conventional data mining methods for creating higher outcomes and exactness. This work presents a review of different existing businesses related to accident prediction utilizing the machine learning area. Wounds because of road accidents are one of the most pervasive reasons for death separated from health-related issues. The investigation of road accident seriousness was finished by running an accident dataset through a few machine learning arrangement calculations to see which model played out the best in characterizing the accidents into severity classes, for example, slight, extreme, and fatal. It was seen that calculated relapse to perform multilevel order gave the most noteworthy exactness score. It was additionally seen that variables, for example, the number of vehicles, lighting conditions, and road highlights assumed a part in deciding the seriousness of the accident. Engineers and analysts in the car business have attempted to plan and manufacture more secure vehicles, yet auto collisions are unavoidable. Examples associated with hazardous accidents could be identified by building up a prediction model that naturally orders the sort of injury severity of different traffic accidents. These social and roadway designs are valuable in the improvement of traffic security control strategies. Significantly, estimates be founded on logical and target reviews of the reasons for accidents and the seriousness of injuries. This paper presents a few models to predict the seriousness of the injury that happened during traffic accidents utilizing machine-learning paradigms. We considered networks prepared to utilize machine learning methods. Analysis results uncover that among the machine learning ideal models considered different standards paradigm approaches.

**Keywords:** Road accident, machine learning algorithmic, paradigm approaches, seriousness of the injury

## INTRODUCTION

Road Accident is the most unwanted and unforeseen thing to happen to a road user, however, they happen regularly. Tragically, we can see a minatory ascent of road accidents everywhere in the world, prominently highroad accidents over the past barely any years. It massively affects society just as in the economy of our nation as there is a huge expense of fatalities and injuries. Lately, traffic accident analysis caused extensive to notice the scientists decide the elements that fundamentally influence traffic accidents. Yet, tragically, the greatest exploration techniques depend on measurable records or by doing some basic overview dependent on meetings or polls. Yet, it is unimaginable to expect to improve and certain arrangement by utilizing these sorts of basic methodologies. The primary riddle is that social highlights in traffic accidents are very hard to concentrate on these sorts of conventional research techniques. Since accidents are generally eccentric and impromptu, so immediate perception is very troublesome. Therefore, getting 100% exact information is next to impossible. Execution of a serious technique that can give better examination results is a crying need here. Machine learning is one of the most exceptional logical fields of AI that can be applied here to improve results. The prime objective of this paper is to analyze road accidents and decide the seriousness of an accident by applying progressed machine learning strategies (Bulbul & Unsal, 2011). There exist so many created strategies in machine learning to look at this area. In this paper, we perform traffic accident analysis, by applying four progressed and most famous managed learning strategies of machine learning because of their demonstrated precision in this area. Those methodologies are-Decision Tree, K-Nearest Neighbors (KNN), Naïve Bayes, and Adaptive Boosting (AdaBoost). Accident seriousness is controlled by the harm coming about because of the accident as far as real harm (fatal accidents being the most extreme) and arranges the accident. Accident severity is simply identified with the speed of the vehicle at the hour of the impact. Albeit monetary misfortune and property harm are too plentiful with the more extreme accidents, actual damage is centered around and considered for consistency in the accident expectation.

The goal of this paper is to introduce a survey of the cutting edge in the expectation of road accidents through calculations and progressed strategies for dissecting data and the joining of new information sources, which were absent in distributed audits on road accident prediction and analysis. In this report, attention was set on archives distributed from the year 2015 onwards and the survey of novel techniques to dissect and figure traffic accidents. It must be viewed as that techniques and algorithms, for example, natural language processing (NLP) and deep learning are generally utilized in different fields of data, sentiment analysis using NLP, or advanced image processing utilizing deep learning methods, however, their application in traffic accident analysis is the most recent. In this paper, a part is devoted to auditing the most pertinent information sources utilized by analysts, including new data sources on traffic accidents, for example, social media and open information given by governments.

## ROAD ACCIDENT DATA SOURCES

Numerous authors refer to the data given by public or private organizations that satisfy the elements of operational traffic control in every nation, regardless of whether they are road security offices or roadway police, while different people utilize freely available datasets on the internet. These informational indexes contain segment data about those engaged with the accident, a variable degree of insight about road conditions and

ecological settings, specialized insights concerning the vehicles in question and their geographical position, the level of the severity of the wounds on drivers, travelers, and people on foot, among other pertinent factors. Then again, a few authors assemble their information by introducing gear on vehicles, for instance, satellite positional frameworks (Donepudi, 2019), cameras and sensors, to collect information like increasing speed, startling slowing down occasions, unexpected path changes, and data about the driver conduct and status like languor and level of pressure. Another arising information source reasonable for proposing models of road accident expectation is social media. The omnipresence and accessibility of social media make it achievable to get progressively data revealed by road users that can't be found in other information sources, for example, road framework disintegration, parked vehicles, minor traffic occurrences, and episodes around the road.
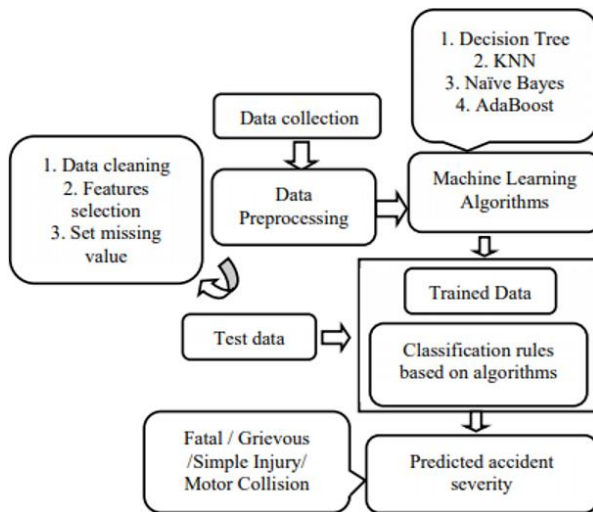


Figure 1: Working Mechanism of Proposed Approaches

The collaboration of new sources of data, for example, mobile applications, internet of things (IoT) devices and more canny instrumentation and hard-product accessible in new vehicles, for example, on-board PCs, GPS, and sensors, it is normal that there will be broader accessibility of information that might be defenseless to broader and inside and out dissects, that includes extra data about drivers, traveler and walkers, and their propensities and everyday movement, and more definite information about ecological, climatic and lighting conditions, as well as data about events and episodes around the road that could influence the security of people on foot, drivers and travelers (Ballamudi, 2016).

**Government & Open Data Sources**

Government information refers to those informational indexes that are produced, gathered, protected, stored, and made accessible to people in general by government substances or those that are designated to practice elements of control, execution, or revealing of data concerning road accidents. Among these offices can be incorporated police bodies, traffic police, and road concessionaires. Government information can be presented as verifiable, since it contains information crossing a very long while, and can be considered as solid since it is upheld by the care cycle of the elements answerable for the

information. Government information is generally the specialized help for the age of public strategy in every nation with respect to angles like road framework plan and road security plans. One perspective to consider is that not all the factors of the informational collection might be accessible for free, for example, specific segment data, sex, and identity, in compliance with the data protection laws currently in every nation. Open Data can be characterized, as per past analysis as the information that is created and financed with public cash, that is made accessible and available without limitation to the public taking into thought security and privacy matters. The most utilized technique to make accessible open data substance to the overall population is by methods for committed sites that empower search and uncover interfaces, including web services, to allow automated utilization of the information and incorporate it to other sites or applications. Road traffic data is normally one of the most accessible information, among different points, for example, population, economic measurements, and geographic data. Relevant instances of open data gateways are the United States government open information index for traffic accidents, that incorporates data from all the nation, United Kingdom open information inventory (Habibullah et al., 2019) and Australia open information inventory, that incorporates constant travel, traffic accidents, and public vehicle time tables.

## Social Media

Social media can be considered the freshest created information source in traffic and road accident-related examinations, and at present, the most utilized information source comes from Waze, Inrix, Google Maps, and Twitter streams as revealed by Sinnott. A social media information can be named as untrustworthy, one-sided, and hard to decipher. Social media data is temperamental because isn't anything but difficult to survey the reliability of its root or distributer; social media is hard to decipher because the users use local language to post their content and the content may contain spelling and linguistic blunders; lastly, it tends to be expressed that social media information can be one-sided since, on account of road accidents, not all the relevant data about accidents are reported by the road users (Donepudi, 2017). The untrustworthiness of social media can be located by utilizing strategies that remove the time, captivity, and subject of the report to connect the report to a genuine event; to manage the arrangement with the syntactic complexities of social media, a technique was proposed dependent on a deep learning engineering and a strategy was depicted dependent on a convolutional recurrent network; the two papers expect to overpass the restriction of the sack of-words and predefined set of keywords techniques that are generally utilized to handle the content of the tweets. To handle the predisposition of social media, it was suggested that the consideration of human movement data revealed in Twitter in the spatial examination of traffic accidents can improve the adequacy and execution of a road accident model.

## Onboard Equipment

Onboard equipment alludes to all gadgets introduced on a vehicle that can store or send information concerning the vehicle factors and driver conditions. The onboard equipment may incorporate global positioning units, cameras set up to record road conditions or driver circumstances, for example, sluggishness or ready status, accelerometers, vehicle condition recorders, for example, change on vehicle speed, abrupt slowing down occasion or path changes, lastly, braking and increasing speed in the event of an effect or crash. One creative methodology was recommended that utilized a development framework called chain road traffic incident to mimic vehicle impact settings dependent on the PreScan stage. Regarding the matter of the information arrangement and quality cycles needed to

collect an appropriate dataset utilizing information from onboard equipment, the most common exercises, are I) eliminate exceptions and equipment blunders announced by GPS hardware; ii) eliminate inconsequential information outside the region of study; iii) coordinate caught GPS traces or information gathered with the road portions characterized in the investigated region; and iv) channel information dependent on vehicle course or other condition required.

## THE METHODOLOGY OF PROPOSED APPROACHES

### Decision Tree

For classification issues, the decision tree is broadly utilized the supervised algorithm. The essential point of view of this algorithm is predicting the estimation of the ideal variable by learning decision standards derived from the highlights of the information and make a model of that.

Most importantly, a root node is assigned for the development of this model depends on the best feature picked by the addition approach, and the sub-nodes are then produced based on the decision taken corresponding to the status of value chose at every node. At the point when every node is decreased to a single quality status, the class is resolved toward the finish of the node; it is known as a leaf. These blueprints proceed recursively until a class is characterized toward the finish of every node (Bulbul & Unsal, 2011).

### AdaBoost

AdaBoost is primarily a boosting algorithm that is utilized with short decision trees. Each example is weighted in the preparation dataset. From the outset, the weight is set to,

$$\text{Weight } w_i = \frac{1}{n}$$

Where $w_i$ is the i'th training example weight and n is the number of training instances. Further, the primary tree is made, the performance of the tree on each training case is utilized. From that point forward, it assesses general mistakes. Next cycle loads are determined by the mistakes. More weight is given where difficult to anticipate, while less weight is given were simple to predict (Wang, 2012).

### KNN

KNN is a classification algorithm that depends on feature comparability. It examines the information and measure the distance and similitudes among information and cluster them based on K values. Distance is determined from multiple points of view, for this exploration, we utilized Euclidean distance estimation. The class of new input data is characterized by computing the distance between the clusters and allocated it to the nearest one (Zhang, 2016).

### Naïve Byes

Naïve Byes is another classification method dependent on the Bayes hypothesis. It predicts the likelihood of various class dependent on a few attributes and allocates the new class to the highest likelihood.

## ROAD ACCIDENT ANALYSIS METHODS

By utilizing analytic techniques, experts look to present the data and factors of the road accident, to find concealed examples, profile practices, produce rules, and derivations. These examples are helpful to profile drivers or drivers' conduct on the road, to delimitate hazardous territories for driving, to produce arrangement rules identified with road accident data, to perform determination of factors to be gotten continuously model of accidents and to choose relevant factors to be utilized to prepare different strategies, for example, fake neural networks and deep learning algorithms.

### Clustering Algorithms

Clustering is a technique for partitioning and gathering objects into clusters (groups) so that objects grouped in each group share normal qualities while searching for them to be unmistakably not quite the same as different objects assembled in different groups. Regular qualities can be understood as the degree of relationship of objects as indicated by the attributes on which clustering strategies are applied. In contrast to classification strategies, clustering doesn't need that the information is recently set apart with a specific class to recognize various groups inside the data. The nonappearance of these past classifications or classes demonstrates that the target of clustering is to locate a basic structure in the data and accomplish a more minimal portrayal of it as opposed to separating future information into classes (Donepudi, 2014a). The principle favorable circumstances of clustering algorithms are that they don't need earlier data training, function admirably with huge informational indexes, and their outcomes are interpretable graphically. Then again, clustering algorithms are sensitive to the chance of finding a nearby greatest rather than a worldwide most extreme on their enhancement capacities.

Clustering algorithms can be ordered by the portrayal of their outcomes and how they play out the gathering and dividing of the data collection. Clustering algorithms utilize a distance capacity to compute the closeness in attributes when they work with persistent components and a proportion of similarity for information with subjective components. Among the procedures dependent on similarity capacities, we can incorporate K-nearest neighbor and K-means clustering. On account of group methods whose similarity work depends on distribution probabilities, their activity depends on the reason that each cluster has a fundamental likelihood of circulation from which the data components are created. An illustration of this kind of calculation is Latent Class Clustering (LCC). For data collections with attributes both subjective and quantitative, clustering procedures, for example, two-step clustering might be utilized, in which a pre-portion of groups is performed utilizing a component of logarithmic distance and afterward said pre-allocation is approved by comparing their distances with a given limit value, at that point the clusters are joined if the distance esteem is more noteworthy than the defined threshold value.

### Classification Algorithms

A decision tree builds classification models as trees or dendrogram, every node speaks to one of the info factors, and every node has a few branches equivalent to the number of potential estimations of the said input variable. Similarly, each leaf node is an estimation of the objective property and speaks to the choice made dependent on the estimation of the information factors in its way from the root node to the leaf. Decision trees are valuable tools in example characterization applications. Its most prominent utility is that its area information isn't needed for its development, its strategy for investigation is

exploratory and not inferential. They can be utilized with deeply dimensional information and with informational collections with fragmented data. Rule learners and classifiers don't need earlier information handling, they function admirably with enormous data indexes and rule learners and classifiers are interpretable graphically; nonetheless, their outcomes are not as precise contrasted with different techniques.

Experts utilize a decision tree classifier, rules induction PART, multilayer perceptron, and Naive Bayes to determine the main factors reasonable for the forecast of the seriousness of a traffic accident. By looking at the changed ruled based models got, the experts reasoned that the decision tree classifier and rules acceptance had better precision, with an estimation of 0.08218. The factors that have more weight in accident casualty were age, sexual orientation, ethnicity, year of the accident, and sort of accident.

Performance of the classifier algorithms was assessed, for example, decision tree, lazy classifier, and multilayer perceptron, analyzing a dataset containing traffic accidents. It was noted that the best precision was acquired by the lazy classifier utilizing clustered data, with a precision of 0.8235. The most significant end got by the experts was that the treatment of the dataset by utilizing clustering algorithms, for this situation various leveled clustering, can prompt a superior performance of classifiers, comparing the exhibition of a similar set of classifier algorithms on a non-clustered dataset (Donepudi, 2014b).

### Natural Language Processing Algorithms

Natural language processing (NLP) algorithms are for the most part used to deal with road accidents revealed by social media and infer data, for example, geolocation, road accident highlights, and pertinent factors. NLP algorithms are utilized to perform the order of social media content, as per predefined target classes, for the most part as a double clustering, being the social media content related or not to a traffic accident. NLP specialists concentrated on how to build the viability of road traffic recognition, utilizing Twitter as a data source, using natural language processing. Their outcomes indicated that in the wake of training the obtained information, just 5% of the information was helpful, under the presumption that the tweets were traffic accident-related and having the option to geocode the information on a guide. The outcomes were approved against direct traffic accident information sources, for example, road condition report framework and police traffic report. The experts asserted that there was a strong pattern in the happening of the data posting, having a top on ends of the week. The analysts announced a precision estimation of 0.9500 for the general characterization of the dataset as traffic accident-related. Then again, the precision value for the way toward getting geocoding data from the tweets dataset was 0.5200. Another NLP expert introduced a methodology for recovering, training, and arranging data from Twitter identified with traffic accidents, by joining natural language processing and support vector machine algorithm to perform text classification. The experts will remember for their future work different procedures to improve their precision and to explore the sentiment analysis inside the content of the tweets.

## CONCLUSION

Losses in road accidents are unendurable, to the general public as well as a non-industrial nation like many. Along these lines, it has become a fundamental prerequisite to control and arrange traffic with a serious framework to diminish the number of road accidents. By avoiding potential risk, because of expectations or warnings of a complex framework may avoid traffic accidents. Also, it's an essential requirement for various countries currently,

to handle this circumstance where consistently endless individuals were killed in a traffic accident, and step by step this rate is getting expanded. The usage of machine learning is a utilitarian and an extraordinary way to deal with taking a precise choice with the experience to deal with the current circumstance and the discoveries of the investigation part can be recommended to traffic experts for reducing the number of accidents. We can utilize proposed ways to deal with actualize machine learning here because of their demonstrated and higher precision to expect traffic accident severity.

The past examination focused essentially on recognizing no-injury and injury (including casualty) classes. We stretched out the exploration to possible injury, non-incapacitating injury, incapacitating injury, and fatal injury classes. Our tests demonstrated that the model for deadly and non-fatal injury performed in a way that is better than different classes. The capacity of predicting deadly and non-deadly injury is significant since drivers' casualty has the greatest expense to society economically and socially. Notably, one of the significant variables causing different injury levels is the genuine speed that the vehicle was going when the accident occurred. Tragically, our dataset doesn't give enough data on the genuine speed since the speed for 67.68% of the data records' was obscure. If the speed was accessible, all things considered, it might have assisted with improving the presentation of models concentrated in this paper.

## REFERENCES

Ballamudi, K. R. (2016). Blockchain as a Type of Distributed Ledger Technology. *Asian Journal of Humanity, Art and Literature*, *3*(2), 127-136. https://doi.org/10.18034/ajhal.v3i2.528

Bulbul, H. I., & Unsal, Ö. (2011). Comparison of Classification Techniques used in Machine Learning as Applied on Vocational Guidance Data. 2011 10th International Conference on Machine Learning and Applications and Workshops, Honolulu, HI, pp. 298-301.

Donepudi, P. K. (2014a). Technology Growth in Shipping Industry: An Overview. *American Journal of Trade and Policy*, *1*(3), 137-142. https://doi.org/10.18034/ajtp.v1i3.503

Donepudi, P. K. (2014b). Voice Search Technology: An Overview. *Engineering International*, *2*(2), 91-102. https://doi.org/10.18034/ei.v2i2.502

Donepudi, P. K. (2017). Machine Learning and Artificial Intelligence in Banking. *Engineering International*, *5*(2), 83-86. https://doi.org/10.18034/ei.v5i2.490

Donepudi, P. K. (2019). Automation and Machine Learning in Transforming the Financial Industry. *Asian Business Review*, *9*(3), 129-138. https://doi.org/10.18034/abr.v9i3.494

Habibullah, K. M., Alam, A., Saha, S., Amin, A., & Das A. K. (12019). A Driver-Centric Carpooling: Optimal Route-Finding Model using Heuristic Multi-Objective Search. 2019 4th International Conference on Computer and Communication Systems (ICCCS), Singapore.

Wang, R. (2012). AdaBoost for Feature Selection, Classification and Its Relation with SVM. A Review. Physics Procedia (2012), vol. 25, pp. 800-807.

Zhang, Z. (2016). Introduction to machine learning: k-nearest neighbors. *Annals of Translational Medicine,* *4*(11), 218-218.

--0--