

# A Review on Road Accident in Traffic System using Data Mining Techniques

Maninder Singh<sup>1</sup>, Amrit Kaur<sup>2</sup>

<sup>1</sup>Student M. Tech, Computer Engineering Department, Punjabi University, Patiala, India

<sup>2</sup>Assistant Professor M. Tech, Computer Engineering Department, Punjabi University, Patiala, India

**Abstract:** Road traffic accidents are a major public health concern, resulting in an estimated 1.2 million deaths and 50 million injuries worldwide each year. In the developing world, road traffic accidents are among the leading cause of death and injury. The objective of this study is to evaluate a set of variable that contribute to the degree of accident severity in traffic crashes. The issue of traffic safety has raised great concerns across the sustainable development of modern traffic and transportation. The study on road traffic accident cause can identify the key factors rapidly, efficiently and provide instructional methods to the traffic accident prevention and road traffic accident reduction, which could greatly reduce personal casualty by road traffic accidents. Using the methods of traffic data analysis, can improve the road traffic safety management level effectively.

**Keywords:** Road Traffic Accidents, Data Mining, Influential Factor, WEKA, Data Mining Techniques

## 1. Introduction

In recently years, because of too much travel speed of road traffic, the accidents have been increasing on yearly basis. So the traffic safety has raised great concerns across the globe. It has become one of the key for challenging the modern traffic and transportation so that traffic administrations can be more accurately informed and better policies can be introduced. The encyclopaedia defines road traffic accident as any vehicle accident occurring on a public highway (i.e. originating on, terminating on, or involving a vehicle partially on the highway). These accidents therefore include collisions between vehicles and animals, vehicles and pedestrians, or vehicles and fixed obstacles. Single vehicle accidents, in which one vehicle alone and no other road user involve, are included. A report by [17] estimated that the cost of fatalities and injuries due to road traffic accidents have tremendous impact on societal well-being and socio-economic development. Road traffic accidents are among the leading cause of death and injury worldwide, causing an estimated 1.2 million deaths and 50 million injuries each year. Moreover, by the year 2020 road accident will be the third leading cause of death. This puts road safety well ahead of wars, HIV/AIDS, malaria and (other) acts of violence as world health problem. Among children aged 5-14years, and young people aged 15-29 years, road traffic injuries are the second-leading cause of death worldwide. In low income countries, the majority of road deaths are among pedestrians, passengers, cyclists, users of motorized two wheelers, and occupants of buses and minibuses. Globally, the risk of dying in a road crash is far higher for vulnerable road users like pedestrians, cyclists and motorcyclists than for car occupants. This paper structure in four sections; section 1 covers the basic introduction whereas section 2 presents the review of literature. Section 3 gives methods and various

data mining techniques are used. The conclusion is drawn in section 4.

## 2. Literature Survey

The cost of deaths and injuries due to traffic accidents has a great impact on society. In recent years, researchers have paid a great attention at determining the factor that significantly affects accident severity in traffic system. The author in [5] presents a random forest & rough set theory to identify the factors significantly influencing single vehicle crash severity. The author in [5] presents a decision tree which predicts cause of accident and accident prone locations. Papers [10] predict traffic accident duration of incident and driver information system. The author in [13] used various data mining techniques and tells the random forest outperforms than other classification algorithms. In paper [14], author talks about the significance of data mining classification algorithm in predicting the factor which influences road traffic accident. The author in [15] used to explore the possible application of data mining technology for developing a classification model and the result shows that developed model could classify accidents within a reasonable accuracy. It is important to analyse these datasets to extracts useful knowledge. Data mining is an effective tool for analysing data to extract useful knowledge. Table 1 shows a sample of different data mining technique and their influential factor used in traffic accident severity. The severity of injuries measured for crash records has both continuous and categorical characteristics. Hence many previous studies have used models with ordered structure to analyse risk factor and their effect on severity of injuries sustained in traffic crashes.

**Table 1:** Summary of the Pertinent Literature

<i>Author</i>	<i>Objective</i>	<i>Data Mining Techniques</i>	<i>Influential factors</i>
Ali et.al (2010)	To identify Most important factors which affect injury severity	Classification & Regression tree	Injury Severity, Gender, Age, Seat Belt, Cause Of crash, Collision type, Vehicle Type, Location type, Lighting conditions, Weather conditions, Road surface condition, Occurrence, Shoulder type, Shoulder Width
Chaozhong et.al (2009)	To identify the factors significantly influencing single vehicle crash severity.	Random Forest Rough set theory	Weather, Speed limit, Lightingconditions, collision factors, gender, Age, Experience, Safety belt, Vehicle type, Severity of svc
DipoT.Akomolafe, Akinbola Olutayo (2012)	To predict causes of accidents and accident prone locations	Decision tree: Id3, Functional tree	Vehicle Type, Time of the day, Season, Causes
Liping et.al (2010)	To predict Traffic accident duration of incident and driver information system	Artificial neural Networks	No. of trucks involved, Rollover vehicle, Facility damage, Degree of traffic jam, No. of fatalities, No. of severe personal injuries, Road pollution involvement, Hazard material, Fire Involvement, Police involvement, Patrol vehicle involvement, Fire engine involvement
S.Krishnaveni, Dr. M.Hemalatha (2011)	To predict severity of injury using data mining techniques & compare algorithm performance.	Naïve Bayes AdaBoost M1 Meta Classifier Part J48, Random Forest	Casualty, Fatal accident, Slight accident, Killed causality, Serious injury, Slight injury, Road users, Vehicles involved
S.Shanthi, R.Geetha Ramani (2012)	Significance of data mining classification algorithms in predicting the factors which influence road traffic accident.	Classification techniques: C4.5, ID3, CS-CRT, CR-T, CS MC4, Naive Bayes, Random forest	Key Value, state, County, Month, date, Time, Day, Harmful event, Manner of collision, Person type, Seating position, age, Gender, Injury Severity, Air bag, Protection System, Ejection, Ejection path, Year_ of_ death, Month_ of_ death, Alcohol Test, Drug test, Drug involvement, Accident Location, Related Factors
Tibebe et.al (2013)	To Explore the possible application of data mining technology for developing a classification model	Classification & Regression tree	Accident_Id, Driver_Age, Driv_Exp, Vehic_Age, Vehic_Type, Road_Surf_Type, Road_Cond, Weather_Con, Light_Cond, Acci_Type

Different data mining technique have been used to help traffic accident severity such as Decision tree, Naïve bayes, Artificial Neural networks, classification and regression tree, J48, PART classifier, random forest. Those most frequently used focus on: Classification, Decision tree and Artificial Neural networks. Driver related factor that affect severe injury crashes have been recognized to have a great influence in the occurrence of crashes. There are various levels of injury severity like No injury means not bodily harm from the crash but only property damage, possible injury means no visible signs of injury but complaint of pain, fatal injury

means an injury sustained in a crash that results in a death within 30 days of the crash. Each road traffic accidents records contain multiple data attributes. Each attribute value reflects a characteristic in a traffic accident. And the more data attribute data constitute the multiple dimensions of traffic accidents. In addition, data attributes are set from the basic accident information, personnel information, vehicle information, road information and environmental information.

**Table 2** Summary of different data mining tools on different datasets

<i>Author</i>	<i>Data Mining Tool</i>	<i>Work Done</i>
Ali et.al (2010)	Variable importance measure	Cause factors: seat belt, improper overtaking & speed.
Chaozhong et.al (2009)	Cross Validation Method	Cause factors: Lighting Conditions, vehicle type, driving experience, wearing belt or not. The efficiency of attribute reduction is not high.
DipoT.Akomolafe, Akinbola Olutayo (2012)	WEKA	Decision tree predict causes of accidents and accident prone locations accurately.
Liping et.al (2010)	MATLAB	The incident duration is predicted in practice, as time goes & incident information gradually increases.
S.Krishnaveni, Dr. M.Hemalatha (2011)	WEKA	Random Forest Outperforms than classification algorithms.
S.Shanthi, R.Geetha Ramani (2012)	TANAGRA	Random tree classifier using Arc-x4 Meta Classifier outperforms & also improves accuracy.
Tibebe et.al (2013)	WEKA	Results shows that developed model could classify accidents with in a reasonable accuracy.

Table 2.2 illustrates a sample of different data mining tools used in the road traffic system over different datasets. WEKA is a Waikato Environment for Knowledge Analysis. WEKA is a collection of machine learning algorithms for

data mining tasks and well suited for developing new machine learning schemes. WEKA is java based software capable of working under various operating systems. These algorithms can either be applied directly to a dataset or can

be called from your own java code. WEKA is probably the most successful open source data mining software which has inspired by the development of other programs with more sophisticated graphical user interface and better visualization methods [2][8]. In WEKA datasets should be formatted to the ARFF format. The WEKA explorer will use these automatically if it does not recognize a given file as an ARFF file the pre-process panel has facilities for importing data from a database, a CSV file, etc., and for pre-processing this data using a filtering algorithm. These filters can be used to transform the data and make it possible to delete instances and attributes according to specific criteria. TANAGRA is free data mining software for academic and research purposes. It offers several data mining methods like exploratory data analysis, statistical learning and machine learning. The first purpose of the Tanagra project is to give researchers and student easy to use data mining software. The second purpose of TANAGRA is to propose to researchers an architecture allowing them to easily add their own data mining methods, to compare their performances. The third and last purpose is that novice developers should take advantage of free access to source code, to look how this sort of software is built, the problems to avoid, the main steps of the project, and which tools and code libraries to use for. In this way, TANAGRA can be considered as a pedagogical tool for learning programming techniques [16]. Revolution is a free software programming language and software environment for statistical computing and graphics. R provides a wide variety of graphical and statistical technique such as linear and non-linear modelling, classical statistical tests, time-series analysis, classification clustering and is highly extensible. Extensibility and superb data visualization are the two main reasons for the success of R [12][9].

**Table 3:** Summary of Accident Severity accuracy crash involvements

Author	Sample Size	Accuracy
Ali et.al (2010)	169648	72.49%
Chaozhong et.al(2009)	59	0.73% 0.54%
DipoT.Akomolafe, Akinbola Olutayo (2012)	148	77.70% 70.27%
Liping et.al(2010)	170	85.35%
S.Krishnaveni,Dr. M.Hemalatha(2011)	34575	84.66% 84.64% 84.64% 85.18% 88.25%
S.Shanthi, R.Geetha Ramani (2012)	457549	99.73%
Tibebe et.al (2013))	5207	87.47%

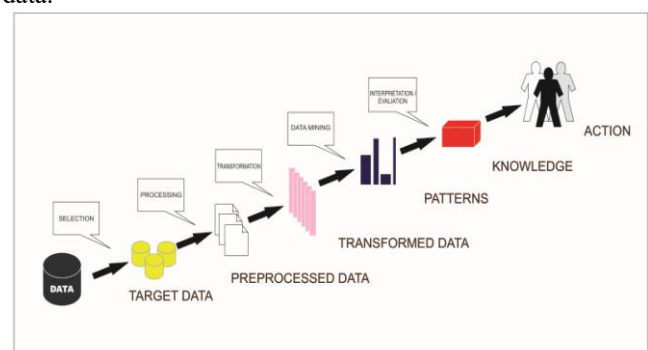
A summary of past studies where crash involvements were used in the accident severity analysis presented in table 3 above. The amount of past research demonstrates the capabilities of different methods to model crash accident severity versus a set of continuous and discrete independent variable. Interactions were found significant in several studies such as: light-weather, alcohol, seat-belt, among others. Although there is a substantial amount of literature demonstrating different uses of severity analysis, only a few studies dealt with a sample as large as the one undertaken in

this research. The crash data [1] from the records of the information and technology department of the Iranian traffic police from 2006 to 2008 was used to study hundreds of drivers who were involved in traffic crashes on the main two-lane two-way rural roads of Iran. The result indicated that seat belt is the most important factor associated with accident severity of traffic crashes and not using it significantly increase the probability of being injury or killed. The crash data [13] from the records of the transport department of Government of Hon Kong from 2008 was used. The total record set used was 14,476 investigates the performance of Naïve Bayes, j48, AdaBoostM1, PART, and Random Forest classifiers for predicting classification accuracy the attributes involved in this case are severity, District council, Hit and Run, Weather, Rain, Natural Light, Junction control, Road classification, Vehicle Movements, Types of collision, Number of vehicles involved and number of casualties injured, the result indicated that Random Forest outperforms than other classification algorithm instead of selecting all the attributes for classification.

### 3. Methods

#### 3.1 Data Mining

This fast growth and tremendous amount of data, collected and stored in large and numerous databases need a powerful tool to elicit useful information. The tool helps to get benefit from the collected data, by identifying relevant and useful information. Data mining is one of the solutions to analyse huge amount of data and turn such data into useful information and knowledge. Data mining [6] refers to extracting or mining knowledge from large amounts of data. There are some other terms which carry a similar or slightly different meaning to data mining, such as knowledge mining from data, knowledge extraction, data or pattern analysis, and data archaeology. Data mining functionalities are used to specify the kind of patterns to be found in data mining tasks. In general, it has been classified into two categories: descriptive and predictive. Descriptive mining tasks characterize the general properties of the data in the database where as predictive mining tasks perform inference on the current data in order to make predictions. Data mining, also popularly known as Knowledge Discovery in Database refers to extracting or “mining” knowledge from large amount of data.



**Figure 1:** An overview of the steps that compose the KDD process

Data mining techniques are used to operate on large volumes of data to discover hidden patterns and relationships helpful in decision making. While data mining and knowledge discovery in database are frequently treated as synonyms, data mining is actually part of the knowledge discovery process [3]. The sequences of steps identified in extracting knowledge from data are shown in Figure 1.

### 3.2 Data Mining Tasks

The cycle of data and knowledge mining comprises various analysis steps, each step focusing on a different aspect or task. Hand et. al (2001) propose the following categorization of data mining tasks.

#### 3.2.1 Description and Summarization

At the beginning of each data analysis are the wish and the need to get an overview on the data to see general trends as well as extreme values quickly

#### 3.2.2 Descriptive Modeling

Descriptive modeling tries to find models for the data. The aim of this model is to describe not to predict models. Even these models are used in the setting of unsupervised learning. Various methods of descriptive modeling are density estimation, smoothing, data segmentation, and clustering. The most widely used method of clustering is k-means clustering. The reasoning behind cluster analysis is the assumption that the data set contains natural clusters can be characterized and labeled. While for some cases it might be difficult to decide to which group they belong, we assume that the resulting group are clear-cut and carry an intrinsic meaning. In segmentation analysis, the user typically sets the number of group in advance and tries to partition all cases in homogeneous subgroup.

#### 3.2.3 Predictive Modeling

The aim of this task is to build a model that will permit the value of one variable to be predicted from the known values of other variables. Predictive modeling falls into the category of supervised learning; hence, one variable is clearly labeled as target variable and will be explained as a function of the other variables. The nature of the target variable determines the type of model: classification model, if it is a discrete variable) or regression model (if it is a continuous one. Many models are typically built to predict the behavior of new cases and to extend the knowledge to objects that are new or not yet as widely understood.

#### 3.2.4 Discovering Patterns and Rules

The area of the previous tasks has been much within the statistical tradition in describing functional relationships between explanatory variables and target variables. There are situations where such a functional relationship is either not appropriate or too hard to achieve in a meaningful way. So association rules a method originating from market basket analysis to elicit patterns of common behaviour.

#### 3.2.5 Retrieving Similar Objects

The World Wide Web contains an enormous amount of information in electronic journal articles, electronic catalogs and private and commercial homepages. Having found an

interesting article or picture, it is a common desire to find similar objects quickly. Based on key word and indexed meta-information search engines are providing us with this desired information. They can not only work on text documents, but to a certain extent also on images.

### 3.3 Data Mining Techniques

#### 3.3.1 Association

Association is the discovery of togetherness or connection of objects. Such kind of togetherness or connection is tends as association rule. An association rule reveals the associative relationship between objects, i.e. the appearance of set of objects. The association rule can be useful for marketing, commodity management, advertising, etc. for example, a retail store may discover that people tend to buy soft drinks together with potato chips and then put the potato chips on sale to promote the sale of soft drinks.

#### 3.3.2 Clustering

Clustering is the identification of class, also called clusters or groups, for a set of objects whose classes are unknown. The objected are so clustered that the interclass similarities are minimized based on some criteria defined on attributes of the objects. Once the cluster is decided, the objects are labelled with their corresponding clusters and common feature of the object in a cluster are summarized to form the class description. For example, a bank may cluster its customers into several group based on the similarities of their age, income, residence, etc. and common characteristics of the customers in a group can be used to describe that group of customers. The clusters will help the bank to understand its customer better and thus provide more suitable products and customized services.

#### 3.3.3 J48

J48 is a version of an earlier algorithm developed by J. Ross Quinlan, C4.5. Decision trees are a classic way to represent information from a machine learning algorithm, and offer a fast and powerful way to express structures in data. It is important to understand the variety of options available when using this algorithm, as they can make a significant difference in the quality of results. J48 in WEKA3.6.0 employs two pruning methods. The first is known as sub tree replacement. This means that nodes in a decision tree may be replaced with a leaf basically reducing the number of tests along a certain path. This process starts from the leaves of the fully formed tree, and works backwards toward the root. The Second type of pruning used in J48 is termed sub tree raising. In this case, a node may be moved upwards towards the root of the tree, replacing other nodes along the way. Sub tree raising often has a negligible effect on decision tree models. There is often no clear way to predict the utility of the option, though it may be advisable to try turning it off if the induction process is taking a long time. This is due to the fact that sub tree raising can be somewhat computationally complex [2].

#### 3.3.4 Classification

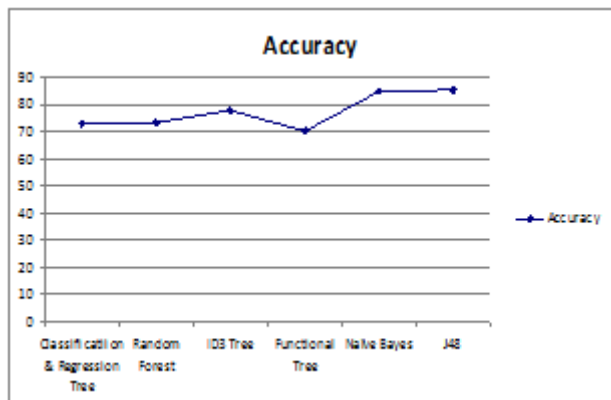
Classification tree are used to predict the classes of a categorical dependent variable. From their measurements on one or more predictor or independent variables. Decision



trees have emerged as a powerful technique for modelling general input output relationship. They are tree shaped structure that represents a series of roles that lead to sets of decision. They generate rules for the classification of a datasets and a logical model represented as a binary tree that show how the value of a target variable can be predicted by using the values of a set predictor variables. Decision tree, which are considered in a regression analysis problem are called regression tree. Thus the decision tree represents a logic model of regularities of the researched phenomenon.

#### 4. Analysis

To predict accident severity, various classification models were built using Decision tree, Random forest, ID3, Functional Tree, J48 and Naïve Bayes. Decision trees are easy to build and understand can manage both continuous and categorical variables and can perform classification as regression.



**Figure 2:** Graph of Accuracy against various Techniques

Figure 2 shows the graph of accuracy against various techniques used. The statistics shows that having a means of predicting likely accuracy of different technique base on some input value. It is evident from the line graph that value of classification and regression tree is slightly less than the random forest and ID3 tree.

Classification models are generated on the basis of the training data whose independent variables and target variable are known, to be applied for the new dataset whose objective is the prediction of the target variable. The principle of CART method in developing the classification tree is described in the following: at first all data are concentrated at the root node, situated at the top of the tree. Further, it will be divided into two child node, on the basis of an independent variable (splitter), which creates the best purity. A decision tree is a flowchart like tree structure, where each internal node denotes a test on an attribute, each branch represents an outcome of the test, and each leaf node holds a class label. It further reveals that value of Naive Bayes is slightly less than J48 technique.

**Table 4:** Detailed Accuracy By Different Techniques

S. No.	Techniques	Accuracy
1	Classification and Regression Tree	72.49
2	Random Forest	73
3	ID3	77.7
4	Functional Tree	70.27
5	Naive Bayes	84.66
6	J48	85.18

Table 4 illustrates a detailed accuracy of different data mining techniques used in the traffic accident severity over different datasets. It further reveals that value of Naive Bayes and J48 techniques are approximately same accuracy. The values are nearby 84.66% and 85.18%. In the end, I have to conclude that J48 technique has higher accuracy than other techniques.

#### 5. Conclusion

Data mining in recent year with the database and artificial intelligence developed a new technology, its aim the large amount of data from the excavated useful knowledge, to achieve the effective consumption of data resources. A through literature review revealed a gap in published studies on the relationship between road characteristics and road traffic accident severity. The study on road traffic accident cause can identify the key factor rapidly and efficiently and provide instructional methods to the traffic accident prevention and road traffic accident reduction, which could greatly reduce personal casualty and property loss by road traffic accidents. Meanwhile, it would be helpful for improving the efficiency and security service level of the road transportation system.

#### References

- [1] Ali et al., "A Data Mining Approach to identify key factors of traffic injury severity" Traffic & Transportation, Vol. 23, 2011, No. 1, 11-17.
- [2] Bouckaert Remco, Eibe Frank, Mark Hall, Richard Kirkby, Peter Reutemann, and AlexSeewald, 2008.WEKA Manual for Version 3-6-0. University of Waikato, New Zealand.
- [3] Brijesh Kumar Baradwaj, Saurabh Pal," Mining Educational Data to Analyze Students Performance" (IJACSA) International Journal of Advanced Computer Science and Applications, Vol. 2, No. 6, 2011.
- [4] Chaozhong et.al., "Severity Analyses of Single-Vehicle Crashes Based on Rough Set theory" 2009 International Conference on Computational Intelligence and Natural Computing.
- [5] DipoT.Akomolafe, Akinbola Olutayo," Using Data Mining Technique to Predict Cause ofAccident and Accident Prone Locations on Highways", American Journal of Database Theory and Application 2012, 1(3): 26-38.
- [6] Han, Jiawei and Kamber, Micheline. (2006)," Data Mining: concepts and Techniques. San Fransisco", Morgan kufman Publishers.
- [7] Hand, D.J., Mannila, H., and Smyth, P. (2001) ,"Principles of Data Mining" ,MIT Press.

- [8] <http://en.wikipedia.org/wiki/weka> (machine learning)/ accessed on May 2014
- [9] <http://rapid-i.com/content/view/181/190> accessed on May 2014
- [10] Liping et al., "Traffic Incident Duration Prediction Based on Artificial Neural Network" 2010 International Conference on Intelligent Computation Technology and Automation.
- [11] Mehmed Kantardzic (2003), "Data Mining: Concepts, Models, Methods, and Algorithms" ISBN13: 9780471228523, John Wiley & Sons Publisher.
- [12] Pasko Konjevoda and Nikola Stambuk, "Open-Source Tools for Data Mining in Social Science," Theoretical and Methodological Approaches to Social Sciences and Knowledge Management, pp.163-176.
- [13] S.Krishnaveni, Dr. M.Hemalatha, "A Perspective Analysis of Traffic Accident using Data Mining Techniques", International Journal of Computer Applications (0975 – 8887) Volume 23– No.7, June 2011.
- [14] S.Shanthi, R.Geetha Ramani "Feature Relevance Analysis and Classification of Road Traffic Accident Data through Data Mining Techniques" Proceedings of the World Congress on Engineering and Computer Science 2012 Vol. I WCECS 2012, October 24-26, 2012, San Francisco, USA.
- [15] Tanagra – a Free Data Mining Software for Teaching and Research, Available at: <http://eric.univ-lyon2.fr/~ricco/tanagra/en/tanagra.html>, (Accessed 20 June 2014).
- [16] Tibebe Shah, Shawndra Hill (2013) "Mining Road Traffic Accident Data to Improve Safety: Role of Road-related Factors on Accident Severity in Ethiopia".
- [17] WHO (2004). World report on road traffic injury prevention, Switzerland, Geneva.

## Author Profile

**Maninder Singh** completed B-Tech degree in Computer Science Engineering from Punjab Technical University, Jalandhar, Punjab and Pursuing M-Tech in Computer Engineering from Punjabi University Patiala, Punjab.

**Amrit Kaur** received her B-Tech, M-Tech degree in Computer Science Engineering respectively from Punjab Technical University Jalandhar, Punjab and Thaper University Patiala, Punjab. Her interest is primarily in the area of Data Mining.