

Data Mining and Business Intelligence (ETH)

Course Code: ITA5007

Slot: B1 + TB1

Registration No.: 19MCA0070

Name: Debayan Das

**Predictive Analysis on Road Accident Risks Based on
Heterogeneous Sparse Data**

Madhurima Chakraborty
19MCA0087

Pragya Yadav
19MCA0125

Debayan Das
19MCA0070

Literature review:

In the paper [1], the authors have used Hadoop to analyze the large data based on various criteria to predict road accidents, in order to raise precaution alarms for the same. On comparing Hadoop with other methods, it has proved to be the most efficient method for big data analysis. The algorithms used in this paper are CCMF and TCAMP which analyze the dataset effectually to predict the risks of road accidents. The proposed algorithm tries to predict the road accident risks with the help of the enormous data about vehicle movement and circumstances that favor the accidents.

This paper [2] focuses on finding and predicting road accident patterns based on the severity, road type, accident type, climate, hour of accident etc. The method of finding interesting and useful patterns from spatial database is termed as spatial data mining. Spatiotemporal algorithms are able to locate hidden patterns more easily than traditional data mining techniques. Spatial data is the data which contains location on the earth surface. Two spatiotemporal clustering algorithms are used for finding patterns, which are DBSCAN or Density Based Spatial Clustering of Applications with Noise and Grid Based algorithms. While on one hand, the data for DBSCAN must have data points and limited threshold, the grid based clustering algorithms require a multi-dimensional data structure. It has been found that Grid based algorithms like STING and CLIQUE produce more accurate results along with faster processing time.

In [3] traffic accident is inferred from heterogeneous data. Huge amount of heterogeneous data containing accident data and GPS records have been collected, to check how vehicle mobility affect the road accidents. On analyzing this data a Stack de-noise auto encoder model is prepared

which studies the features in human mobility to predict accident risks. The model on being prepared simulates real time accident risks, which can be used to warn people of the possible accident, for a safer route and journey.

Anupama Makkar et. Al has analysed the accident dataset of recent years to forecast road accidents in [4]. The proposed approach in this paper incorporates amalgamation of machine learning algorithms like Bayes Net, j48 graft and j48 decision tree in the data mining process to examine the performances of the algorithms in prediction of accidents. It has been thus noticed that the combination of such algorithms render better results than a single algorithm used. The results obtained would support in forecasting road traffic accidents and hence prevention and control can be provided.

The crucial issue of predicting road traffic and the accidents caused due to it has been addressed in this paper [5]. The sparse and heterogeneous data have been worked upon by various algorithms to find interesting patterns. Through a case study, this paper explores various algorithms to predict the amount of accidents occurring every hour. The problem is framed as a binary classification problem. Big data inclusive of features like weather, accidents, road networks etc. have been map-matched. Algorithms like support vector machine, deep neural networks, random forest and decision tree have been evaluated along with Eigen analysis.

In this paper,[6] the features causing road traffic accidents were detected using three techniques of classifications; which are: Decision trees, SVM and ANN. Based on these algorithms a prediction model is built, testing all these algorithms on real time data set. The results determine that Random Forest algorithm gave the most accurate predictions, followed by ANN and then SVM. Hence a data mining model using decision trees is made for forecasting the accidents.

Accident Prediction Model is proposed by the authors in the paper [7], where accidents on horizontal curves are anticipated based on the patterns examined in the dataset. This prediction model takes measures to lessen the accidents at some extent. The factors which cause road accidents are analyzed to make predictions and take safety measures to decrease the accident rates. In this paper, the geometric features of the road are collected at different level length, height, stretch, etc. Regression modelling is used for the data mining procedure.

This paper [8] has used two predictive models for the analysis of previous accident data and current accident data to predict the number of accidents occurred in that year. Multiple Linear Regression as well as Artificial Neural Networks has been used for the predictive analysis. After conducting the analysis it is concluded that the predicted values from regression model had greater errors. While the predictions made from Artificial Neural Network analysis was more accurate having less errors. Hence ANN was proved to be a better methodology to make predictions for accidents.

The problems addresses in the paper [9] are, predicting the number of accidents occurring on intersection of roads or any road and finding roads prone to accident risks. Algorithms have been used on heterogeneous data to mine traffic risk. The algorithm incorporated for the framework is an advanced feature based non-negative matrix factorization (FNMF). This framework is successful in predicting the traffic risks at any road or intersection more accurately than the existing algorithms or methods. Two clusters were defined that segregate the risk locations, in

which one cluster had larger roads with accident risks, and other cluster having higher risk of vehicle collision. Risky locations were ranked based on the results of the clusters.

In this paper [10], it states that reducing traffic accidents is a necessary concern for safety. Small datasets, depending on large data and not valid for real-time data are basic demerits of previous studies. The deep neural network model is used to provide a solution for real-time problems. It collects the data, integrates and obtains different attributes like traffic events, weather data and time. Using US-Accident data they perform a comparison between DAP (Deep Accident Prediction) and traditional model and resulting from that extensive model are more appropriate than the traditional approach.

In this paper [11], they surveyed on semi-automated vehicles factors that are affecting the crashes and accident. Using negative binomial regression model to decrease number of accident. As the ADT dataset are expensive and not available in whole road instead they taken sub sample dataset of same period of time. Then applying spatial; regression on the traffic accident dataset. So the result is that driver should drive more carefully to reduce the accident.

In this paper [12], stated that the spatial data is deal with the patchy data. For example, in population data the destruction is having hard edges which is patchy and have gaps in it. The comparative studies say that the approach used in spatially referenced data has ecologically utilize these attribute.

In this paper [13], it is a comparative study by seeing the condition of Bangladesh. So many accidents caused in Bangladesh due to heavy traffic applying following supervised machine learning technique they are decision tree, KNN, naive Bayes and ad boost. So it is observed that most of accident occur due to surface like if the surface is dry no accident occur and if the surface is wet they chance are more for the accident.

In this paper [14], stated that road transportation is the major form of transportation from one country to another there are few factors involve that are responsible for these accident. Using data mining approaches find the relation between these factor and perform which algorithm perform well. They basically use Decision tree, KNN, and Naïve Bayes classifier. It resulting as the KNN is the best among the other two algorithms using the confusion matrix they conclude the results.

In this paper [15], stated that the rate of accident occur in 1 hrs. in the major cities can be predicted using some machine learning algorithms. Every paper represent the factor involve in the accident and area affected but they time zoned it. They use balanced random forest and random forest algorithm. By comparing these two using various factors they concluded that random forest better than balanced random forest.

In this paper [16] the high-density areas of accident has been identified. For this GIS and KDE methodologies has been used to study the spatial patterns of injury related road accidents. And K-means clustering methodology is used for creating a classification of road accident hotspots in London and UK. Five groups and 15 clusters were created based on collision and attribute data. These clusters are discussed and evaluated according to their robustness and potential uses in road safety campaigning.

In this paper[17] they have surveyed the spanish road accident data and found the injury severity involving novice drivers in urban areas. The information root node variation (IRNV) method (based on decision trees) was used to get a rule set that provides useful information about the most probable causes of fatalities in accidents involving inexperienced drivers in urban areas. This method is based on the decision tree classifier. These rules provide useful knowledge in order to prevent these kinds of accidents.

In this article[18] the zones, roads and specific time in the CDMX in which the largest number of road traffic accidents are concentrated during 2016 has been identified. A database compiling information obtained from the social network known as Waze is built. The methodology Discovery of knowledge in the database (KDD) for the discovery of patterns in the accidents reports was used. The Maximization of Expectations (EM) algorithm was used to obtain the number ideal of clusters for the data and k-means was used for grouping method.

In this study[19], the classification of road accident on the basis of road user category has been performed. Self Organizing Map (SOM), K-modes clustering technique has been used to group the data into homogeneous segments and then Support vector machine (SVM), Naive Bayes (NB) and Decision tree models has been applied to classify the data. The classification has been performed on data with and without clustering.

In this paper[20] a learning model has been proposed to overcome the challenges that the decision makers face in order to encounter a huge number of resulting association rules that can make them unable to choose and decide rationally between these different extracted rules. The learning model is based on based on FP-growth algorithm using Apache Spark framework, in order to analyze data and extract interesting association rules by taking into account some quality measures.

Summary:

S.No	Title	Models Used	Summary	Author
1	A Data Mining Framework to Analyze Road Accident Data using Map Reduce Methods CCMF and TCAMP Algorithms	CCMF and TCAMP	The authors have used Hadoop to analyze the large data based on various criteria to predict road accidents, in order to raise precaution alarms for the same. On comparing Hadoop with other methods, it has proved to be the most efficient method for big data analysis	S. Nagendra Babu, J. Jebamalar Tamilselvi
2	ACCIDENT PREDICTION BASED ON ACCIDENT TYPES USING SPATIOTEMPORAL	DBSCAN, STING and CLIQUE	This paper focuses on finding and predicting road accident patterns based on the severity, road type, accident type, climate, hour of accident etc. The method of finding interesting and useful patterns from spatial database is termed as spatial data mining. Spatiotemporal algorithms are able to	Dara Anitha Kumari, Dr. A. Govardhan

	CLUSTERING ALGORTIHMS		locate hidden patterns more easily than traditional data mining techniques	
3	Learning Deep Representation from Big and Heterogeneous Data for Traffic Accident Inference	Stack de-noise auto encoder model	Huge amount of heterogeneous data containing accident data and GPS records have been collected, to check how vehicle mobility affect the road accidents. On analyzing this data a Stack de-noise auto encoder model is prepared which studies the features in human mobility to predict accident risks	Quanjun Chen, Xuan Song, Harutoshi Yamada, Ryosuke Shibasaki
4	A Radical Approach to Forecast the Road Accident Using Data Mining Technique	Bayes Net, j48 graft and j48 decision tree	It has been thus noticed that the combination of such algorithms render better results than a single algorithm used. The results obtained would support in forecasting road traffic accidents and hence prevention and control can be provided.	Anupama Makkar, Harpreet Singh Gill
5	Predicting Traffic Accidents Through Heterogeneous Urban Data: A Case Study	support vector machine, deep neural networks, random forest and decision tree and Eigen Analysis	The crucial issue of predicting road traffic and the accidents caused due to it has been addressed in this paper [5]. The sparse and heterogeneous data have been worked upon by various algorithms to find interesting patterns. Through a case study, this paper explores various algorithms to predict the amount of accidents occurring every hour	Xun Zhou, Tianbao Yang, Zhuoning Yuan, James Tamerius, Ricardo Mantilla
6	Data Mining Methods for Traffic Accident Severity Prediction	Decision trees, SVM and ANN	In this paper, the features causing road traffic accidents were detected using three techniques of classifications; which are: Decision trees, SVM and ANN. Based on these algorithms a prediction model is built, testing all these algorithms on real time data set.	Qasem A. Al-Radaideh and Esraa J. Daoud
7	Development of Accident Prediction Model on Horizontal Curves	Regression analysis	Accident Prediction Model is proposed by the authors in the paper [7], where accidents on horizontal curves are anticipated based on the patterns examined in the dataset. This prediction model takes measures to lessen the accidents at some extent.	Jerry Soman, Jisha Akkara
8	Study of Road Accident Prediction Model at Accident	Multiple Linear Regression, Artificial	After conducting the analysis it is concluded that the predicted values from regression model had greater errors. While the predictions made	Haikal Aiman Hartika, Mohd Zakwan Ramli, Muhamad

	Blackspot Area: A Case Study at Selangor	Neural Networks	from Artificial Neural Network analysis was more accurate having less errors. Hence ANN was proved to be a better methodology to make predictions for accidents.	Zaihafiz Zainal Abidin, Mohd Hafiz Zawawi
9	Traffic Risk Mining From Heterogeneous Road Statistics	advanced feature based non-negative matrix factorization (FNMF).	This framework is successful in predicting the traffic risks at any road or intersection more accurately than the existing algorithms or methods. Two clusters were defined that segregate the risk locations, in which one cluster had larger roads with accident risks, and other cluster having higher risk of vehicle collision. Risky locations were ranked based on the results of the clusters.	Koichi Moriya, Shin Matsushima, Kenji Yamanishi
10	Accident Risk Prediction based on Heterogeneous Sparse Data: New Dataset and Insights (2019)	Feature Vector Representation, Deep Accident Prediction (DAP) Model	it states that reducing traffic accidents is a necessary concern for safety. Small datasets, depending on large data and not valid for real-time data are basic demerits of previous studies. The deep neural network model is used to provide a solution for real-time problems. It collects the data, integrates and obtains different attributes like traffic events, weather data and time. Using US-Accident data they perform a comparison between DAP (Deep Accident Prediction) and traditional model and resulting from that extensive model are more appropriate than the traditional approach.	Sobhan Moosavi, Mohammad Hossein, Srinivasan Parthasarathy, Radu Teodorescu, Rajiv Ramnath
11	Spatial prediction of traffic accidents with critical driving events – Insights from a nationwide field study	Spatial regression	they surveyed on semi-automated vehicles factors that are affecting the crashes and accident. Using negative binomial regression model to decrease number of accident. As the ADT dataset are expensive and not available in whole road instead they taken sub sample dataset of same period of time. Then applying spatial; regression on the traffic accident dataset. So the result is that driver should drive more carefully to reduce the accident.	Benjamin Rydera,*, Andre Dahlingerb, Bernhard Gahr, Peter Zundritscha, Felix Wortmannb, Elgar Fleischa

12	Twenty years and counting with ADIE: Spatial Analysis by Distance Indices software and review of its adoption and use	Ia, index of aggregation, Patch and gap cluster indices and Red blue Plot	It stated that the spatial data is deal with the patchy data. For example, in population data the destruction is having hard edges which is patchy and have gaps in it. The comparative studies say that the approach used in spatially referenced data has ecologically utilize these attribute.	Linton Winder ¹ , Colin Alexander ² , Georgianne Griffiths ³ , John Holland ⁴ , Chris Woolley ⁵ , Joe Perry ⁶
13	Road Accident Analysis and Prediction of Accident Severity by Using Machine Learning in Bangladesh	Decision Tree, KNN, Naïve Bayes and Ada-Boost	it is a comparative study by seeing the condition of Bangladesh. So many accidents caused in Bangladesh due to heavy traffic applying following supervised machine learning technique they are decision tree, KNN, naive Bayes and ad boost. So it is observed that most of accident occur due to surface like if the surface is dry no accident occur and if the surface is wet they chance are more for the accident.	Md. Farhan Labib, Ahmed Sady Rifat, Md. Mosabbir Hossain, Amit Kumar Das, Faria Nawrine
14	Analyzing Road Accident Criticality using Data mining	KNN, Decision Tree and Naïve Bayes	It stated that road transportation is the major form of transportation from one country to another there are few factors involve that are responsible for these accident. Using data mining approaches find the relation between these factor and perform which algorithm perform well. They basically use Decision tree, KNN, and Naïve Bayes classifier. It resulting as the KNN is the best among the other two algorithms using the confusion matrix they conclude the results.	Shahsitha Siddique V*, Nithin Ramakrishnan
15	High-Resolution Road Vehicle Collision Prediction for the City of Montreal	Balanced random forest, and random forest	It stated that the rate of accident occur in 1 hrs. in the major cities can be predicted using some machine learning algorithms. Every paper represent the factor involve in the accident and area affected but they time zoned it. They use balanced random forest and random	Antoine H'ebert_, Timoth'ee Gu'edon_, Tristan Glatard, Brigitte Jaumard

			forest algorithm. By comparing these two using various factors they concluded that random forest better than balanced random forest.	
16	Kernel density estimation and K-means clustering to profile road accident hotspots	Geographical Information System, Kernel Density Estimation, K-means Clustering	In this paper the high-density areas of accident has been identified. For this GIS and KDE methodologies has been used to study the spatial patterns of injury related road accidents. And K-means clustering methodology is used for creating a classification of road accident hotspots in London and UK. Five groups and 15 clusters were created based on collision and attribute data. These clusters are discussed and evaluated according to their robustness and potential uses in road safety campaigning.	Tessa K. Anderson
17	Decision Tree Ensemble Method for Analyzing Traffic Accidents of Novice Drivers in Urban Areas	Decision Tree	In this paper they have surveyed the spanish road accident data and found the injury severity involving novice drivers in urban areas. The information root node variation (IRNV) method (based on decision trees) was used to get a rule set that provides useful information about the most probable causes of fatalities in accidents involving inexperienced drivers in urban areas. This method is based on the decision tree classifier. These rules provide useful knowledge in order to prevent these kinds of accidents.	Serafín Moral-García , Javier G. Castellano , Carlos J. Mantas , Alfonso Montella and Joaquín Abellán
18	Road Traffic Accidents Analysis in Mexico City through Crowdsourcing Data and Data Mining Techniques	KDD, EM, K-Means	In this article the zones, roads and specific time in the CDMX in which the largest number of road traffic accidents are concentrated during 2016 has been identified. A database compiling information obtained from the social network known as Waze is built. The methodology Discovery of knowledge in the database (KDD) for the discovery of patterns in the	Gabriela V. Angeles Perez, Jose Castillejos Lopez, Araceli L. Reyes Cabello, Emilio Bravo Grajales, Adriana Perez

			accidents reports was used. The Maximization of Expectations (EM) algorithm was used to obtain the number ideal of clusters for the data and k-means was used for grouping method.	Espinosa, Jose L. Quiroz Fabian
19	Road-User Specific Analysis of Traffic Accident Using Data Mining Techniques	SOM, K-Mode, SVM, Naive Bayes, Decision Tree	In this study, the classification of road accident on the basis of road user category has been performed. Self Organizing Map (SOM), K-modes clustering technique has been used to group the data into homogeneous segments and then Support vector machine (SVM), Naive Bayes (NB) and Decision tree models has been applied to classify the data. The classification has been performed on data with and without clustering.	Prayag Tiwari, Sachin Kumar, and Denis Kalitin
20	Data mining for road accident analysis in a big data context	FP-growth algorithm, Apache Spark framework	In this paper a learning model has been proposed to overcome the challenges that the decision makers face in order to encounter a huge number of resulting association rules that can make them unable to choose and decide rationally between these different extracted rules. The learning model is based on based on FP-growth algorithm using Apache Spark framework, in order to analyze data and extract interesting association rules by taking into account some quality measures.	Fatima Zahra El Mazouri, Mohammed Chaouki Abounaima, Said Najah, Khalid Zenkouar

References:

- [1] S. Nagendra Babu, J. Jebamalar Tamilselvi,” A Data Mining Framework to Analyze Road Accident Data using Map Reduce Methods CCMF and TCAMP Algorithms” , IJSSST, 2013
- [2] Dara Anitha Kumari, Dr. A. Govardhan,” ACCIDENT PREDICTION BASED ON ACCIDENT TYPES USING SPATIOTEMPORAL CLUSTERING ALGORTIHMS”, International Journal of Pure and Applied Mathematics Volume 120 No. 6 , 2018
- [3] Quanjun Chen, Xuan Song, Harutoshi Yamada, Ryosuke Shibasaki,” Learning Deep Representation from Big and Heterogeneous Data for Traffic Accident Inference”, Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence

- [4] Anupama Makkar, Harpreet Singh Gill, "A Radical Approach to Forecast the Road Accident Using Data Mining Technique", International Journal of Innovative Science and Research Technology ISSN No: - 2456 – 2165 Volume 2, Issue 8, 2017
- [5] Xun Zhou, Tianbao Yang, Zhuoning Yuan, James Tamerius, Ricardo Mantilla, "Predicting Traffic Accidents Through Heterogeneous Urban Data: A Case Study", Proceedings of 6th International Workshop on Urban Computing, Halifax, Nova Scotia, Canada, 2017
- [6] Qasem A. Al-Radaideh and Esraa J. Daoud, "Data Mining Methods for Traffic Accident Severity Prediction", INTERNATIONAL JOURNAL OF NEURAL NETWORKS and ADVANCED APPLICATIONS, 2018
- [7] Jerry Soman, Jisha Akkara, "Development of Accident Prediction Model on Horizontal Curves", International Research Journal of Engineering and Technology (IRJET) Volume: 06 Issue: 03, 2019
- [8] Haikal Aiman Hartika, Mohd Zakwan Ramli, Muhamad Zaihafiz Zainal Abidin, Mohd Hafiz Zawawi, "Study of Road Accident Prediction Model at Accident Blackspot Area: A Case Study at Selangor", International Journal of Scientific Research in Science, Engineering and Technology, 2017
- [9] Koichi Moriya, Shin Matsushima, Kenji Yamanishi, "Traffic Risk Mining From Heterogeneous Road Statistics", IEEE TRANSACTIONS ON INTELLIGENT TRANSPORTATION SYSTEMS, VOL. 19, NO. 2018
- [10] "Sobhan Moosavi, Mohammad Hossein, Srinivasan Parthasarathy, Radu Teodorescu, Rajiv Ramnath" Accident Risk Prediction based on Heterogeneous Sparse Data: New Dataset and Insights (2019)
- [11] "Benjamin Rydera*, Andre Dahlingerb, Bernhard Gahr, Peter Zundritsch, Felix Wortmannb, Elgar Fleischa" Spatial prediction of traffic accidents with critical driving events – Insights from a nationwide field study, 2018
- [12] "Linton Winder¹, Colin Alexander², Georgianne Griffiths³, John Holland⁴, Chris Woolley⁵, Joe Perry⁶" Twenty years and counting with ADIE: Spatial Analysis by Distance Indices software and review of its adoption and use, 2019
- [13] "Md. Farhan Labib, Ahmed Sady Rifat, Md. Mosabbir Hossain, Amit Kumar Das, Faria Nawrine" Road Accident Analysis and Prediction of Accident Severity by Using Machine Learning in Bangladesh, 2019
- [14] "Shahsitha Siddique V*, Nithin Ramakrishnan", Analyzing Road Accident Criticality using Data mining, 2019
- [15] "Antoine H'ebert_, Timoth'ee Gu'edon_, Tristan Glatard, Brigitte Jaumard" High-Resolution Road Vehicle Collision Prediction for the City of Montreal, 2019

[16] Tessa K. Anderson, “Kernel density estimation and K-means clustering to profile road accident hotspots”, 2008

[17] Serafín Moral-García , Javier G. Castellano , Carlos J. Mantas , Alfonso Montella and Joaquín Abellán , “Decision Tree Ensemble Method for Analyzing Traffic Accidents of Novice Drivers in Urban Areas”, 2019

[18] Gabriela V. Angeles Perez, Jose Castillejos Lopez, Araceli L. Reyes Cabello, Emilio Bravo Grajales, Adriana Perez Espinosa, Jose L. Quiroz Fabian, “Road Traffic Accidents Analysis in Mexico City through Crowdsourcing Data and Data Mining Techniques”, 2018

[19] Prayag Tiwari, Sachin Kumar, and Denis Kalitin, “Road-User Specific Analysis of Traffic Accident Using Data Mining Techniques”, 2017

[20] Fatima Zahra El Mazouri, Mohammed Chaouki Abounaima, Said Najah, Khalid Zenkouar, “Data mining for road accident analysis in a big data context”, 2019