

# ACCIDENT DATA ANALYSIS USING HADOOP HIERARCHICAL CLUSTERING

Mr.K. Venkatesh

Asst. Professor, Department of computer science  
MallaReddyEngineeringCollegeForWomen, Hyderabadh

Mr. M.Chandra Rao

Asst.professor,Department of computer science  
Mallareddy Engineering college for women,Hyderabad

K.shiva Bhavani

Asst.professor,Department of computer science  
MallaReddyEngineeringCollegeForWomen, Hyderabadh

**Keywords:** Bigdata, MapReduce, Hadoop, Trafficdataanalysis, HierarchicalClustering

## Abstract

Road accidents are more frequent today than ever before for a variety of reasons, including poor road conditions, poor weather, driver inexperience, etc. Due to these reasons, a sizable number of individuals in India have died or been hurt. Traffic accidents need to be analyzed to find issues and their influence on death rates and injury severity levels in order to increase traffic safety. In order to investigate the primary causes of accidents in all of India's states, we did cluster analysis using data sets from road accidents that were published by the Indian government. On the Hadoop platform, hierarchical clustering is used to evaluate the traffic data set for every state in India. The main agenda of Hierarchical cluster analysis is to find datasets that belong together and separating them from the other data resulting a cluster of variables.

## 1 Introduction

Each and everyday the traffic in India is increasing gradually. It is shown that there is a positive correlation among growth of automobiles, population of growth and change in income. [1] The main traffic causes two wheelers, four wheelers and three wheelers, it creates air pollution, and noise pollution. In addition to the pollution the increase in traffic is also leading to traffic accidents which are a major cause of death of people. Approximately 1.25 million people die every year and 20 to 50 million people suffer with injuries. If the preventive actions are not taken it is predicted that road traffic crashes may rise to become 7th place in the world. [2]

The number of traffic accidents is rising for a number of reasons, including driver error, two-wheeler traffic, environmental factors, driver drug usage, etc. Traffic accidents must be examined to discover potential risk variables and their influence on injury severity levels in order to increase traffic safety. Planning for the improvement of road conditions and implementing the necessary corrective measures to stop traffic accidents requires conducting a traffic study. [3]

Government of India publishing a lot of data related to various departments like health and family welfare,

Data over the last ten years on information and communications, agriculture, finance, rural areas, and traffic. With the research and analysis of this data, new business initiatives may be developed, precautionary measures may be taken, or some decision support systems may be supported.

## 2 Literature Review

After reading and examining 66 case studies about the effect of drugs on car accidents, Rune Elvik carried out a meta-analysis of them. His findings indicate that 264 estimations of the likelihood of accidents are associated with drug usage while driving. [4] nations. Based on local data and studies into the causes of accidents, safety rules should be put into practice. [6]

Ross After doing logistic regression analysis on various survey responses, Owen Phillips and Fridulv Sagberg found that factors like driving off the road, poor road conditions, long commutes, and inexperienced drivers are linked to incidents involving tired drivers. [5]

Only 0.7% of studies published on road traffic injuries are from India, according to N.N. Borse and A.A. Hyder's 2009 study of 826 PubMed articles about injuries caused by traffic accidents. More research should be supported in developing nations in order to increase traffic safety and decrease injuries from accidents. Based on local data and studies into the causes of accidents, safety rules should be put into practice. [6]

## 3 Implementation

Hadoop is a free, open-source set of tools for scalable, distributed computing. For the analysis of huge datasets, MapReduce is used as the implementation. Hadoop employs the Hadoop Distributed File System, or HDFS, to manage storage resources throughout the cluster. A primary distributed storage file system called HDFS was created to allow dispersed jobs to communicate data across various hardware and software platforms. [7][16][17]

MapReduce: Based on Java, Map Reduce is a framework and programming model for handling huge data collections. Large datasets are processed using Map Reduce on a node cluster. Map Reduce leverages the proximity of data processing principle to minimize data transmission. The map function transforms a set of data into another set of data, where each element is divided into tuples of key/value pairs. Two crucial jobs, namely Map and Reduce, are part of the Map Reduce algorithm. using Map Reduce

theworkflowofInput→Map()→Copy()/Sort()→Reduce()→Output.Mapfunctionexecutesonasetofinputandoutputsasetofrecordsintheformofkey-valuepairs,andpassthemtotheReducefunction.TheReducefunction,acceptstheintermediatekeypairsandmergestogetherthesamekeyvaluestoformasmallerasetofvaluesasfinaloutput.[8][9][3][10][15]

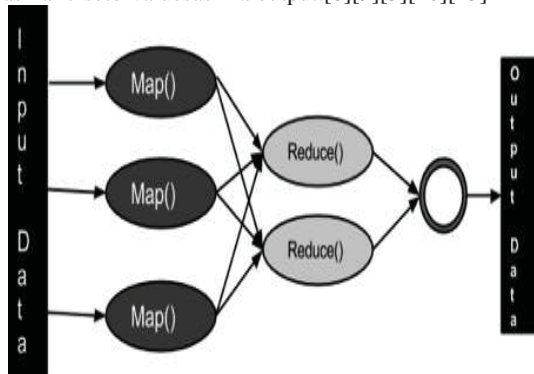


Fig.1:MapReduceframework

## 4 Methodology

A cluster analysis technique called hierarchical clustering aims to create a hierarchy of groups. It is a commonly used tool for data analysis. The comparable data are grouped together and separated from the other data using hierarchical cluster analysis. The clusters will be made up of uniform variables. Each data point is treated as its own cluster when using hierarchical clustering, and the distance between clusters is determined. The clusters that are closest to one another are combined. Based on the correlation coefficients between the variables, clusters will be created. Depending on the Euclidian distance between the variables, clusters can form. [11][12]

The road accident datasets are collected from the data.gov.in. and performed the vertical processing of data and clustering analysis is done using hierarchical algorithm.

### MapReduce

Stage of the map: The job of the map or mappers is to process the supplied data. The input data is often stored in the Hadoop file system as text files. The mapper function receives the input file line by line. The mapper breaks up the data into multiple little pieces after processing it.

Reduce stage: This step combines the shuffle stage with the previous stage. Processing the data that arrives from the mapper is the Reducer's responsibility. It produces a fresh set of output after processing.[13]

[12]

### Steps

1. The input dataset is loaded and apply preprocessing on that data
2. Partition the data by mapping
3. Taking the mapper output as an input, reduce the data

4. After reducing, clustering the data by reading jar files.

### The input dataset is loaded and apply preprocessing on that data.

The initial action taken when data enters the process state is parsing. Parsing is the process of separating or identifying data from a text file. The Indian government publishes information about road accidents that is gathered. Accidents involving two-wheelers, those caused by driver error, and those caused by vehicle defects are all included in the data collection. Preprocessing is carried out based on the data that are available. In the absence of data, zero is used in its place.

### 2. Partition the data by mapping

The input dataset is distributed to three mappers. Each mapper will separately give the result as the partitions of data. The inbuilt class for partition is MyPartitioner which is implemented from the class Partitioner. The inbuilt method is getPartition by taking the parameters as text key, text value.

Partition number

public int getPartition(Text key, Text value, int partitionerNo).[13]

Taking the fields such as cause of accidents, years and total number of accidents, data is partitioned into three sets and distributed to three mappers. Each contains the data related to number of accidents took place in a particular year in all the states.

### 3. Taking the mapper output as an input, reduce the data.

The mapper output is taken as an input to reducer and collect data from each mapper. The inbuilt class is MyReducer which is implemented from Reducer and inbuilt method is reduce. As per the year, the total number of accidents will be displayed.

The partitioned data is joined into a single file by using reduce function.

### 4. Clustering and visualizing the clusters

The data is grouped according to the type of accident that occurred in a particular year using hierarchical clustering. Bar charts are used to display the output data. The data in that graph shows the annual total number of accidents caused by three factors: two-wheeler use, driver error, and vehicle problem.

## 5 Results And Discussions

The data published by government of India is collected for the analysis.[14] The data is collected for the years from 2006 to 2011 for all the states in India. The sampled data used for the analysis is shown below in Fig.No.2.

State	Year	Reason	Count	...
Bihar	2015	2015	632	...
Bihar	2016	2016	725	...
Bihar	2017	2017	637	...
Bihar	2018	2018	722	...
Bihar	2019	2019	806	...
Bihar	2020	2020	584	...
Bihar	2021	2021	318	...
Bihar	2022	2022	267	...
Bihar	2023	2023	181	...
Bihar	2024	2024	114	...
Bihar	2025	2025	383	...
Bihar	2026	2026	382	...
Bihar	2027	2027	217	...
Bihar	2028	2028	974	...
Bihar	2029	2029	1228	...
Bihar	2030	2030	2817	...
Bihar	2031	2031	835	...

Fig.2:Sampleddataset

We considered three dimensioned data with year, state and reason for accident. After the preprocessing step the data is partitioned based on the year. and the sample partitioned data is shown in the below fig. no.3

**Partitioned data**

Year	State	Reason	Count
2015	Bihar	2015	632
2016	Bihar	2016	725
2017	Bihar	2017	637
2018	Bihar	2018	722
2019	Bihar	2019	806
2020	Bihar	2020	584
2021	Bihar	2021	318
2022	Bihar	2022	267
2023	Bihar	2023	181
2024	Bihar	2024	114
2025	Bihar	2025	383
2026	Bihar	2026	382
2027	Bihar	2027	217
2028	Bihar	2028	974
2029	Bihar	2029	1228
2030	Bihar	2030	2817
2031	Bihar	2031	835

Fig.3:Partitioneddata

In hierarchical clustering the data is clustered based on the proximity from one cluster to another. The proximity considered here is the type of reason for the accident. The resultant clusters are visualized using bar charts in Fig.4.

**Clustered output**

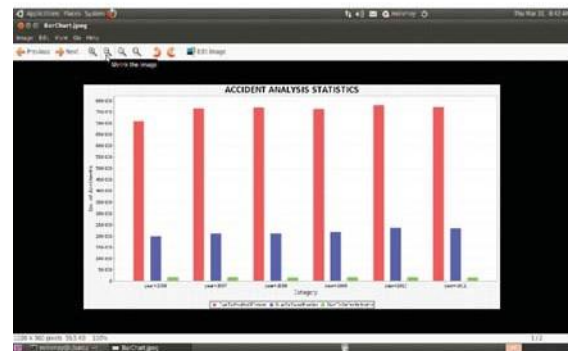


Fig.4:Bar chart showing the clusters

**6 Conclusion And Future Scope:**

As the population and earning capability of the people is increasing the rate of purchase of motor vehicles is also increasing. But, the infrastructure facilities and safety measures are not provided on par with the available vehicles. Every year a huge number of people are killed or injured in road accidents. The reasons for the road accidents are many like due to fault of driver, due to bad road conditions, due to consumption of drugs by drivers, due to two wheelers, due to defect in motor etc.. The need for analyzing road accidents due to heavy traffic is most important. The government can suggest appropriate safety measures once the cause of the event has been determined. In order to determine the primary reason for accidents, we examined the traffic statistics in this study. Hierarchical clustering is used to examine the data using the MapReduce architecture. We can infer from the facts that the majority of accidents occur as a result of driver error. The government can lower the number of fatalities and injuries caused by traffic accidents by taking the required action. In the future, we can examine which states have more accidents based on the weather and the state of the roads.

**References**

- [1] Rameshwar DASHARMA & Sandeep JAIN & Kewal SINGH "Growth rate of Motor Vehicles in India - Impact of Demographic and Economic Development" Journal of Economic and Social Studies, Volume 1, Number 2, July 2011, P137-p153
- [2] <http://www.who.int/mediacentre/factsheets/fs358/en/> (accessed in January, 2016)
- [3] Vadivel.Mand Raghunath.V "Enhancing Map-Reduce Framework for Big data with Hierarchical Clustering", International Journal of Innovative Research in Computer and Communication Engineering, March 2014.
- [4] Rune Elvik "Risk of road accident associated with the use of drugs: A systematic review and meta-analysis of evidence from epidemiological studies" Accident Analysis &

- Prevention ,Volume 60, November 2013, Pages 254–267doi:10.1016/j.aap.2012.06.017
- [5] RossOwenPhillips , FridulvSagberg,“Roadaccidentscaused by sleepy drivers: Update of a Norwegian survey”Accident Analysis & Prevention Volume 50, January 2013,Pages138–146.
- [6] N.N. Borse & A.A. Hyder, “Call for more research on injuryfromthedevelopingworld:Resultsofabibliometricanalysis”Indian JMedRes129, March2009, pp321-326.
- [7] JeffreyShafer;ScottRixner.”The[::Hadoop::]distributedfilesystem: Balancing portability and performance”, ;AlanL.CoxPublicationYear:2010,Page(s):122-133.
- [8] Ronald Taylor,“AnoverviewoftheHadoop/MapReduce/HBaseframeworkanditscurrentapplicationsin bioinformatics”, BMCBioinformatics201011(Suppl12): S1DOI: 10.1186/1471-2105-11-S12-S1
- [9] Jeffrey Dean and Sanjay Ghemawat “MapReduce: SimplifiedDataProcessingonLarge Clusters”October2013.
- [10] Dr. Siddaraju and Sowmya “Efficient Analysis of Big DataUsingMapReduceFramework” International JournalofRecentDevelopmentinEngineeringandTechnology.
- [11] ShenWangandHaimontiDutta,PARABLE:APArallelRAndom - partitionBasedHierarchicalLClusteringAlgorithmfortheMap ReduceFramework”https://www.researchgate.net/publication/266422725
- [12] Georg Foerster; “Traffic State estimation using hierarchicalclustering:apracticalapproach”YoungResearchers Seminar2007.CD-ROM:27to30May2007,Brno,CzechRepublicBrno: CDV,2007,12pp.
- [13] PrasadkumarKaleI andArtiMohanpurkar2“BigDataAnalysis usingPartitionTechnique”
- [14] www.data.gov.in
- [15] HadoopMapReducehttp://hadooptutorial.wikispaces.com/MapReduce
- [16] [www.youtube.com/hadoopinstallation](http://www.youtube.com/hadoopinstallation)
- [17] <http://mvnrepository.com/artifact/org.apache.hadoop/hadoop-common/2.2.0>

