# Big Data Analysis Application for Road Accident with Decision Tree

**R. Sravan Kumar*[1]**          **Dr.R.Radha*[2]**

[1]* Research Scholar KL University ,
[2]*Prfessor KL University.

*Abstract*— **In transportation field, an enormous measure of information gathered by IOTframeworks, remote detecting and other information accumulation apparatuses brings new difficulties; the span of this information turns out to be amazingly huge and progressively complex for customary methods of information mining. To manage this test, Apache Spark remain as an amazing huge scale circulated processing stage that can be utilized effectively for machine learning against exceptionally expansive databases. This work utilized substantial scale machine learning strategies particularly Decision Tree with Apache Spark system for enormous information investigation to assemble a model that can foresee the elements lead to street mishaps dependent on a few info factors identified with auto collisions. In view of this, the anticipating model first preprocesses the huge mishap information and investigates it to make information for a learning framework. Observational outcomes demonstrate that the proposed model could give new data that can help the leaders to break down and enhance street wellbeing.**

*Keywords*— **Deep learning, deep neural network, Internet of Things, on-device intelligence, IOT big data, fast data analytics, cloud-based analytics.**

## I.  INTRODUCTION

Information mining procedures have been intended to dispread helpful learning and justifiable examples from databases1,2. To be sure, with the unstable increment of data advances, the expansive measure of information is generated, which includes various issues, the major of which is information preparing to make a preparation dataset that required equipment assets and tedious for the examination. To manage these issues, conveyed registering is generally utilized, Hadoop and MapReduce3,4 comprise the incredible answer for one-pass computations, yet not extremely effective for use cases that require multi-pass calculations. Each progression in the information procedure work process has one Map stage and one Reduce stage and you'll have to change over any utilization case into Map Reduce designs. Thus, this methodology will in general be moderate because of the colossal space utilization by each activity. As of late, there has been significant research in structuring huge information models In 2009, AMP Lab created Apache Spark8 as an open source huge information preparing system worked around speed, convenience, and modern explanatory, publicly released in 2010 as an Apache venture. Flash takes Map Reduce to the following dimension with more affordable rearranges in the information procedure. With capacities like in-memory information stockpiling and close ongoing preparing, the execution can be a few times quicker than other enormous information advances.

The framework of this paper is sorted out as pursues: Section 2 portrays an outline of related work about mishap investigation, and Section 3 clarifies Spark outline work and choice tree strategy. In Section 4, clarifies the proposed methodology and the experimental examination, at last this paper finished by an outcome exchange and closing area.

## II.  RELATED WORK

Street mishaps have developed as an essential general medical issue on the planet, as per World Health organization 1, 24 million individuals kick the bucket in street crashes every year and upwards of 50 million are harmed. In the writing survey, information mining procedures are broadly used to break down street mishap.

utilized CART and MARS to break down of an epidemiological case-control investigation of wounds coming about because of engine vehicle mishaps and they recognized potential regions of hazard generally brought about by the driver circumstance. utilized calculated relapse models to dissect the accident elements, and they found that the shopping locales are more risky than town destinations. Creator in12 utilized three strategies of information mining, for example, choice tree, neural systems, and calculated relapse for finding significant factors for Korea Road traffic seriousness. In this way, utilized choice tree to investigate the seriousness of auto collision, and they found that lethal damage brought about by numerous components among them safety belts, liquor, and light conditions. built up a CART model to analyze the connection between drivers, damage seriousness and parkway condition variable. utilized Binary Logistic Regression, Logistic Regression Diagnostics to controlling the impacts of statistic and street condition. Likewise, utilized bunching, order trees to cover intuitive investigations dependent on brushing and connecting strategies to identify and recognize intriguing examples. examined the spatial examples of street mishap damage and results from the patterns so as to make a grouping of street mishap hotspots. Additionally, creator in18 utilized distinctive methodologies to find mishap seriousness factors, they found that an unsafe mishap brought about by a blend of contrasting factors. Considered the driver duty by utilizing ID3, J48, and MLP calculations to find the related elements and they found that numerous variables directly affect seriousness mishap, for example, permit grades, driver age and experience. utilized CART and Multinomial Logistic Regression (MLR) to ponder the pretended by drivers' attributes in the subsequent accident seriousness, and they found that the CART technique gave more outcomes that are exact.

In a similar rationale, utilized remote detecting for territorial scale examination and viable administration of the ecological, this innovation can be valuable for aiding in the counteractive action of some sort of mishaps. Utilized the Global Positioning System (GPS) in the prevention of the impact mishaps.

Likewise, Author in23 reasoned that the non-utilization of safety belts and insufficient preparing were likewise two imperative elements. Analyzes the primary driver of those mishaps by utilizing Bayesian classifiers and choice tree.

As of late, with the quick advancement of data advances, information investigation turns out to be increasingly more unpredictable since the information are incredibly substantial. To handle this issue Google Company25 proposed Map Reduce as a programming model and a related implementation for preparing and creating substantial datasets with a parallel and circulated calculation. In a similar rationale, proposed a foreseeing model dependent on C5.0 to gain huge data successfully between the foot issue gatherings and biomechanical parameters identified with manifestations. Likewise, proposed a methodology of grouping to fabricate an expectation display that can resolve the issue of huge information by utilizing Hadoop system and mahout to process and break down car crash.

### III.      KNOWLEDGE DISCOVERY IN BIG DATA

#### A. DECISION TREE (DT)

As of late, with the quick advancement of data advances, information investigation turns out to be increasingly more unpredictable since the information are incredibly substantial. To handle this issue Google Company25 proposed Map Reduce as a programming model and a related implementation for preparing and creating substantial datasets with a parallel and circulated calculation. In a similar rationale, proposed a foreseeing model dependent on C5.0 to gain huge data successfully between the foot issue gatherings and biomechanical parameters identified with manifestations. Likewise, proposed a methodology of grouping to fabricate an expectation display that can resolve the issue of huge information by utilizing Hadoop system and mahout to process and break down car crash.

**C4.5 Algorithm**: C4.5 is a standard algorithm for predicting decision rules in the form of DT. It used information gain ratio as a default criteria of choosing splitting attributes. The algorithm uses the function of the handset with a gain of entropy Split Info function to evaluate the attributes for each iteration. The algorithm needs to decide which split should be used to construct the tree. One option is to use the attribute with the highest purity measure that measured in terms of information value Info (D). C4.5 Algorithm utilizes entropy formula by giving a random variable that takes k values with probabilities P1, P2... Pk, the information value calculated with this following entropy Formula (1):

$$\text{Info(D)} = -\sum_{k=1}^{k} p_k \log_2 (P_k) \quad (1)$$

Where D refers to a specific data partition.
K is the number of class-values involving D in total.
Where D refers to a specific data partition.
K is the number of class-values involving D in total.
Pk is the probability of those class values occurring in K.

The expected information that is required by classification for a parameter $c_j$ ($j=1, 2... m$), is

$$E(c_j) = \sum_{k=1}^{k} P_k \, Info(D) \quad (2)$$

The information gain that is required by classification for a parameter cj (j=1, 2... m) is

$$Gain(c_j) = Info(D) - E(c_j) \quad (3)$$

For a parameter $c_j$ ($j=1, 2... m$), the split information that is required by classification is the same as information value Eq (1)

$$\text{SplitInfo(D)} = -\sum_{k=1}^{k} p_k \log_2 (P_k) \quad (4)$$

For a parameter $c_j$ ($j=1, 2... m$), the information gain ratio that is required by classification is

$$GainRatio(c_j) = \frac{Gain(c_j)}{SplitIno(c_j)} \quad (5)$$

The process of C4.5 algorithm30 is described in Figure 1.

```
Algorithm 1 C4.5(T)
Input:  training dataset T; attributes S.
Output:  decision tree Tree.
1:  if T is NULL then
2:      return  failure
3:  end if
4:  if S is NULL then
5:      return  Tree as a single node with most frequent class label in T
6:  end if
7:  if |S| = 1 then
8:      return  Tree as a single node S
9:  end if
10: set Tree = {}
11: for a ∈ S do
12:     set Info(a, T) = 0, and SplitInfo(a, T) = 0
13:     compute Entropy(a)
14:     for v ∈ values(a, T) do
15:         set T_{a,v} as the subset of T with attribute a = v
16:         Info(a, T)+ = |T_{a,v}|/|T_a| Entropy(a_v)
17:         SplitInfo(a, T)+ = -|T_{a,v}|/|T_a| log |T_{a,v}|/|T_a|
18:     end for
19:     Gain(a, T) = Entropy(a) - Info(a, T)
20:     GainRatio(a, T) = Gain(a,T)/SplitInfo(a,T)
21: end for
22: set a_{best} = argmax{GainRatio(a, T)}
23: attach a_{best} into Tree
24: for v ∈ values(a_{best}, T) do
25:     call C4.5(T_{a,v})
26: end for
27: return  Tree
```

*Figure 1.    Process of C4.5 algorithm.*

#### B. Apache Spark

Apache Spark8 is an open source huge information handling system worked around speed, and refined analytics. Created in UC Berkeley's AMP Lab, and publicly released in 2010 as an Apache venture. With the capabilities like in-memory, the execution can be a few times quicker than other huge information advancements.

#### Spark Architecture

Sparkle applications keep running on a group composed by flash setting in the principle program called driver program, the flash setting can associate with a few sorts of bunch man-agers, when associated flash obtain agents on hubs in the bunch, which are forms that run calculation and information stockpiling. Next, it sends the application to the agents, at last flash setting sends errands to the agent, Figure 2.

#### Spark Ecosystem

Flash give an extensive and brought together answer for oversee diverse enormous information use cases and necessities.
It is an option to Hadoop MapReduce, it contains extra libraries that are a piece of the Spark biological community and give extra abilities in enormous information investigation and machine-learning territories see Figure 3.
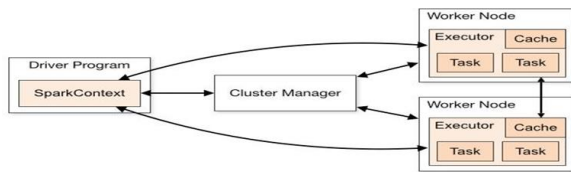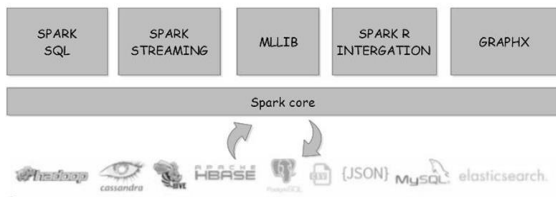
*Figure 2 Spark design.*



*Figure 3 Spark biological system.*

Apache Spark Run projects up to 100x quicker than Hadoop MapReduce in memory, or 10x quicker on the plate. Apache Spark has a RDD (Resilient Distributed Dataset) an accumulation of information things split into segments and put away in memory of laborer hubs of the group, and the Directed Acyclic Graph (DAG) a succession of calculations for each shaped on information .

**C.Proposed Approach**

To have a sufficient model for foreseeing mishap factors with regards to huge information, we think it is essential to adjust C4.5 calculation for dispersed figuring. The genius presented methodology includes three stages, the Figure 5 chows the full procedure dependent on SparkR with the pre-preparing of street mishap information, the DT where worked by utilizing C4.5 calculation, this proposed methodology is portrayed by the accompanying advances:

Pre-handling: In this progression, we allude to an ETL (Extraction Transformation Loading) device for getting ready and cleaning information identified with the street mishap by changing the information to an appropriate arrangement and choosing just certain sections to stack.

Decision rules extraction: In this progression, SparkR33 is utilized as a R bundle that gives a light-weight front end to utilize Apache Spark from R34. It comprises by Sparklyr that expert vides a dplyr interface to Spark Data Frames just as a R interface to Spark's conveyed Machine Learning (ML/H2O) pipelines.

Visualization: Data perception is the introduction of information in a pictorial or graphical configuration. It empowers leaders to see investigation exhibited outwardly, so they can get a handle on troublesome ideas or recognize new examples.

The working procedure of C4.5 calculation on Apache Spark is given in Figure 4

**Run C4.5**

**SparkContext:**

The constructor: new SparkContext(master, appName, [SparkHome]) is called to initialize SparkContext.

**Initialization**

Read and initialize attributes and their possible values from meta file

**RDD:**

The input training set is regarded as a RDD on Spark through csvFile(path, minSplits): RDD[String].

**FlatMap:**

Get a list through each input line, including:

1.<id+att+value+class, 1>

2.<id,1>

3.<total,1>

**ReduceBykey:**

Get the sum of the same key from the RDDs from flagMap.

**GenerateTree:**

Get the attribute that has the highest gain ratio in each node on current layers.

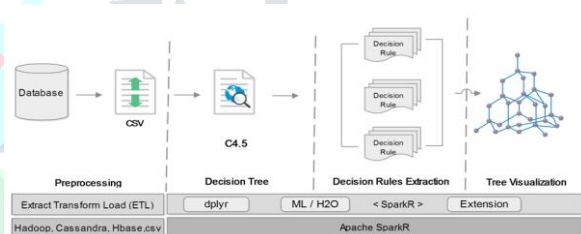*Figure 4.The working process of C4.5 on Apache Spark*



*Figure 5.The proposed approach based on Apache Spark and MLlib.*

## IV. CONCLUSION

This paper examines the issue of machine learning calculations on huge information through the street mishap analy-sister, which is unmistakably recognized by utilizing Apache Spark and C4.5 calculation to remove choice tenets from huge datasets. Along these lines, we discovered that Apache Spark, gave a quicker execution motor to disseminated preparing likewise gave a library to the machine learning calculations, called Machine Learning library (MLlib). Apache Spark asserted that it is a lot quicker than Hadoop MapReduce as it abuses the upsides of in-memory calculations which is especially increasingly valuable for iterative com-putations if there should be an occurrence of a machine learning calculations. We played out a few examinations on street mishap information to gauge the accelerate and scale up of implementa-tions of C4.5 calculations in Sparks' MLlib. We discovered much superior to anticipated outcomes for our analyses. The outcomes show that the proposed methodology is very adaptable and could give important informa-tion that can help the coordinations directors to enhance the exhibitions of transport quality and street security opti-mization. For further work, new approachs ought to be routed to process ongoing information by utilizing Apache

KAFKA the disseminated spilling stage, additionally the inte-gration of multi-criteria investigation will be helpful for the exactness of results.

## REFERENCES

[1] Fayyad UM, Piatetsky-Shapiro G, Smyth P. From data min-ing to knowledge discovery: An overview. Advances In Knowledge Discovery And Data Mining. AAAI Press/The MIT Press; 1996. p. 134.

[2] Agrawal R, Imielinski T, Swami A. Mining association rules between sets of items in large databases. In: Proceedings of ACM SIGMOD Conference on Management of Data (SIGMOD 1993); 1993. p. 207216. Crossref.

[3] Available from: http://hadoop.apache.org/

[4] Jeffrey D, Sanjay G. MapReduce: Simplified data process-ing on large clusters. Proceedings of the 6th Conference on Symposium on Operating Systems Design and Implementation, (OSDI'04); Berkeley, CA, USA: USENIX Association. 2004. p. 10-10.

[5] Saathoff BG, Hamid RA, Hill R, Staniforth A, Bayerl PS. Application of Big Data for National Security. Akhgar B, editor. Butterworth-Heinemann; 2015. p. 4.

[6] Chen Y, Li F, Fan J. Mining association rules in big data with NGEP Cluster Comput. 2015; 18: 577. Crossref.

[7] Jin S, Lin W, Yin H. Community structure mining in big data social media networks with MapReduce. Cluster Comput. 2015; 18:999. Crossref.

[8] Available from: http://spark.apache.org/

[9] Availablefrom: http://www.who.int/gho/road_safety/en/

[10] Kuhnert PM, Do KA, McClure R. Combining non-para-metric models with logistic regression an application to motor vehicle injury data. Computational Statistics and Data Analysis. 2000; 34(3):371–86. Crossref.