

UNIVERSITY OF PETROLEUM AND ENERGY STUDIES



BACHELOR OF TECHNOLOGY

Computer Science & Engineering

(specialization in Artificial Intelligence And Machine Learning)

A DATA MINING FRAMEWORK TO ANALYZE ROAD ACCIDENT DATA

PROJECT REPORT

Version No.	1.00
Authorized by	PIYUSH MALVIYA PRIYAL GUPTA SAKSHAM GARG RAHUL DHANOLA

Industry Mentor:	Guided By:
Sumit Shukla	Sujoy Chatterjee

S.NO	TOPIC
1	BACKGROUND
1.1	AIM
1.2	TECHNOLOGIES
1.3	HARDWARE ARCHITECTURE
1.4	SOFTWARE ARCHITECTURE
2	SYSTEM
2.1	REQUIREMENTS
2.1.1	FUNCTIONAL REQUIREMENTS
2.1.2	USER REQUIREMENTS
2.1.3	ENVIRONMENTAL REQUIREMENTS
2.2	DESIGN AND ARCHITECTURE
2.3	IMPLEMENTATION
2.4	TESTING
2.5	EVALUATION
3	SNAPSHOTS OF THE PROJECT
4	CONCLUSIONS
5	FURTHER DEVELOPMENT OR RESEARCH
6	REFERENCES

1. Background

In developed as well as developing countries, infrastructure development is one of the major investments by the government, while the safety of passengers on roads is of utmost importance. A road optimization during the construction or the maintenance phase requires that the engineers analyze all the parameters that play a crucial role in ensuring safety for the passengers and preventing accidents. One of the key objectives in accident data analysis is to identify the main factors associated with road accidents. Due to road accidents, a large number of lives are lost. From an analysis, it has been estimated that every year over 3,00,000 people die and 10 to 15 million people are injured due to road accidents in the entire world. There are a lot of vehicles driving on the roadway every day, and traffic accidents could happen at any time anywhere. Some accident involves a fatality, which means people die in that accident. As human beings, we all want to avoid accidents and stay safe. Traffic accidents have now earned india a dubious distinction; with nearly 140,000 deaths annually, the country has overtaken china to top the world in road fatalities. India is the only country in the world with more than 15 fatalities and 53 injuries every hour due to road crashes. The roads of india have not abated their contribution to traffic accident fatalities. The accident rate in india has been on an increase ever since the start of the century. Data mining analyses can help identify the major causes and help the transport authorities in improving safety requirements.

1.1 Aim

Our main aim is to study the states and the union territories of India against the contributing causes to facilitating road safety in the country. We are focused on taking the aid of clustering to group similar objects of this dataset to group regions based on vulnerability. The major problem in the analysis of accident data is its Heterogeneous nature. Thus, heterogeneity must be considered during the analysis of the data.

Road accident analysis aims to investigate the main factors that characterize an accident to understand patterns of behaviours and, consequently, to identify the appropriate countermeasures to adopt to avoid the accident. In this, we will use different data mining algorithms to analyze the data.

Different algorithms are applied to group the accident locations into clusters and mining techniques are used to characterize the locations. Most states of traffic management and information systems focus on data analysis. Python and Jupyter notebooks are mainly used. Since Python has a large number of libraries and packages, it has a very large ecosystem.

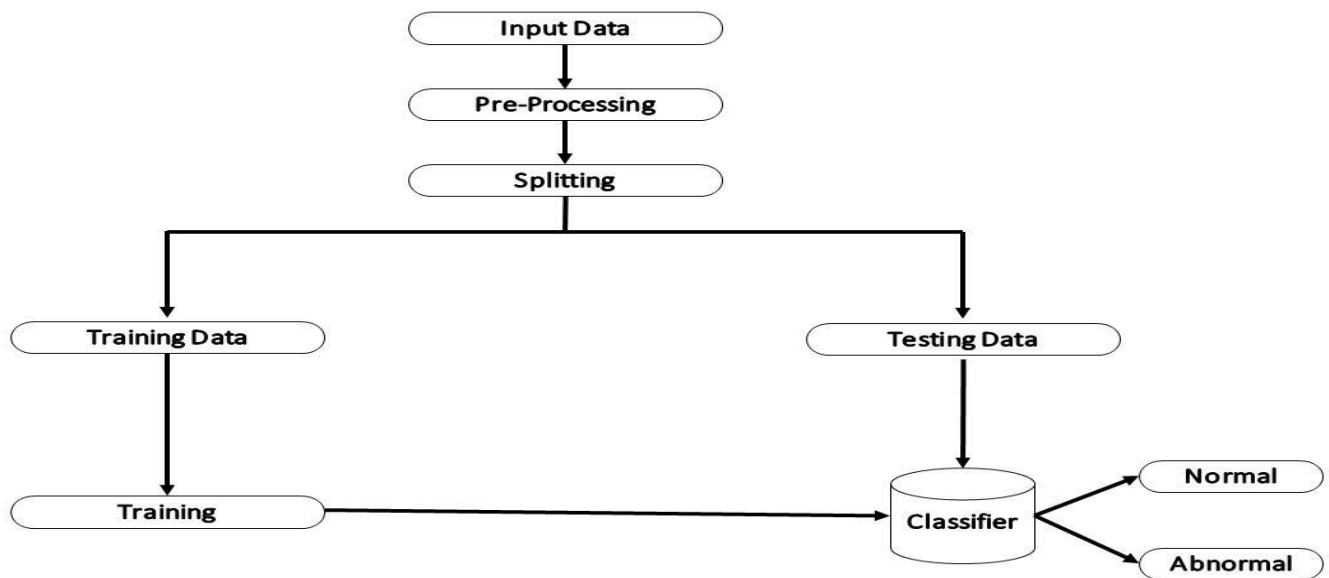
Python is used both in data scraping and in developing the server. Jupyter notebook is an open-source and web-based interactive environment for making notebook documents. The primary Jupyter online application and Jupyter python web server are the substances required for making a notebook

1.2 Technologies

In this, we will use different data mining algorithms to analyze the data. Different algorithms are applied to group the accident locations into clusters and mining techniques are used to characterize the locations. Most states of traffic management and information systems focus on data analysis. Python and Jupyter notebooks are mainly used.

Since Python has a large number of libraries and packages, it has a very large ecosystem. Python is used both in data scraping and in developing the server. Jupyter notebook is an open-source and web-based interactive environment for making notebook documents. The primary Jupyter online application and Jupyter python web server are the substances required for making a notebook.

1.3 Hardware Architecture

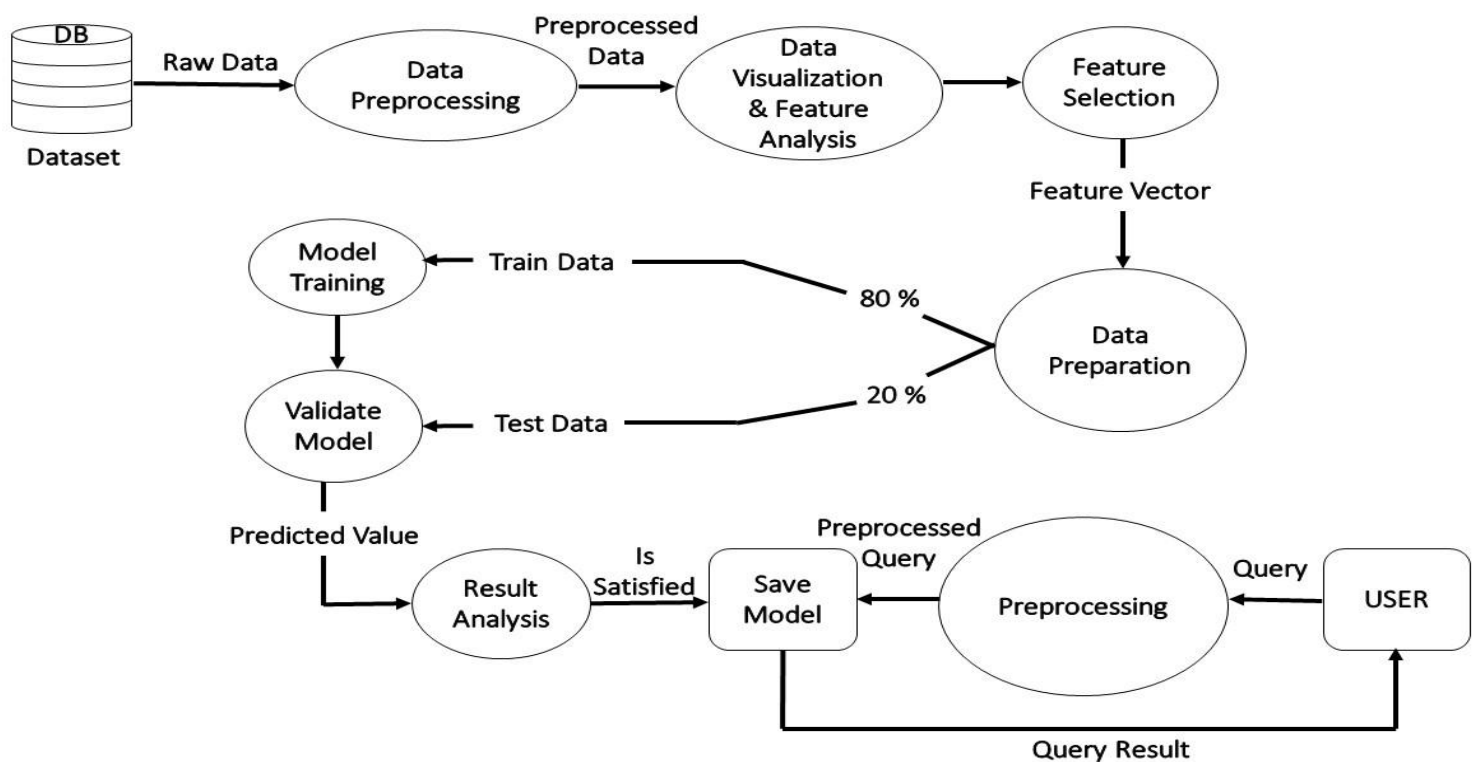


- ❖ **Data preprocessing:** Data preprocessing is one of the important tasks in data mining. Data preprocessing mainly deals with removing noise, handling missing values, and removing irrelevant attributes to make the data ready for analysis. In this step, we aim to preprocess the accident data to make it appropriate for the analysis.
- ❖ **Clustering algorithm:** There are several clustering algorithms in the literature. The objective of the clustering algorithm is to divide the data into different clusters or groups such that the objects within a group are similar to each other whereas objects in other clusters are different from each other. K means & Decision trees have been used in road accident analysis.
- ❖ **Association rules:** Association rule mining is a very popular data mining technique that extracts interesting and hidden relations between various attributes in a large data set. Association rule mining produces a set of rules that define the underlying patterns in the data set. The associativity of two characteristics of the accident is determined by the frequency of their occurrence together in the data set. A rule $A \rightarrow B$ indicates that if A occurs then B will also occur. Further association rules are generated from the frequent item sets and strong rules based on interestingness measures are taken for the analysis.

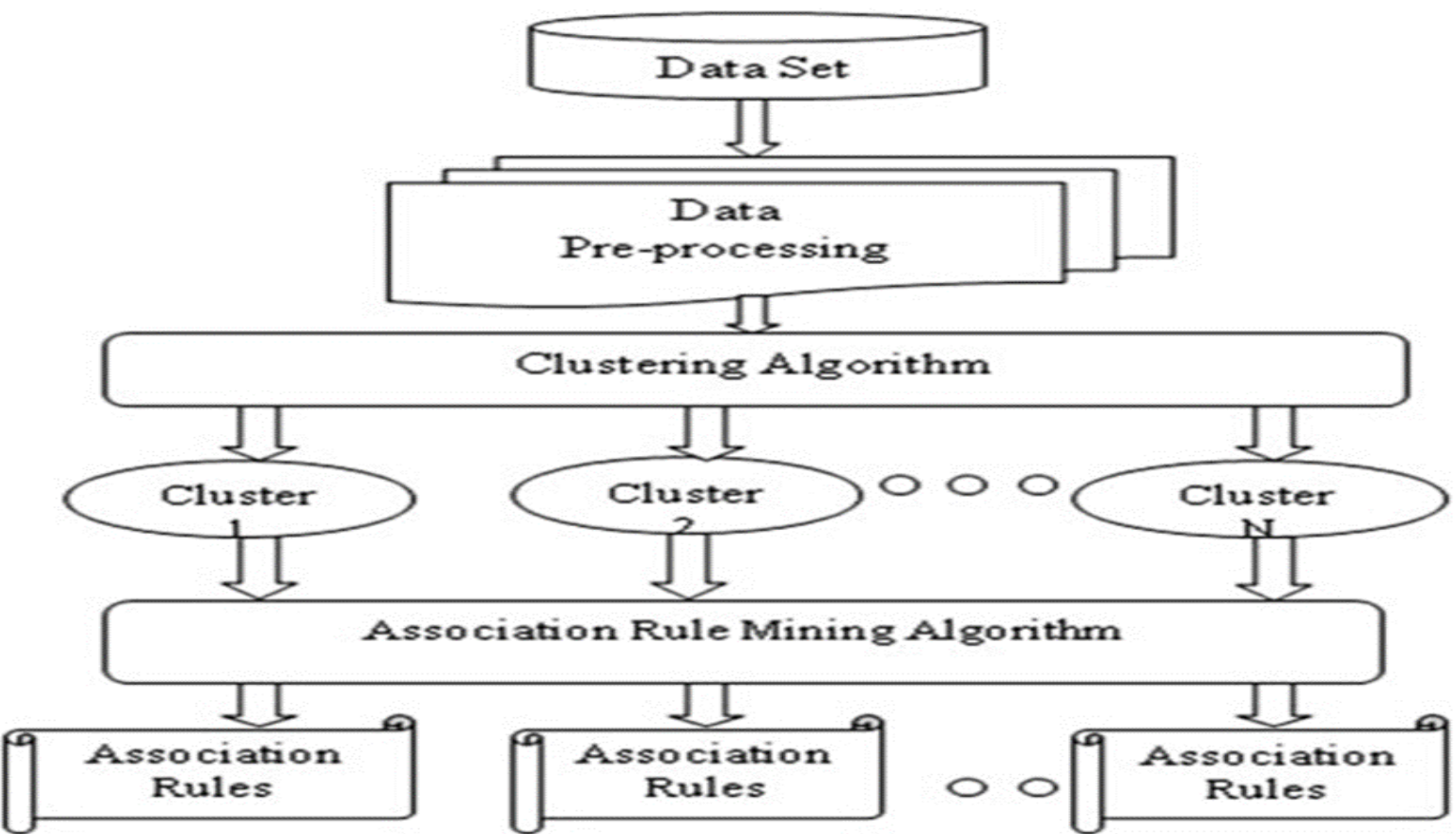
1.4 Software Architecture

This project can run on commodity hardware. The first part is the training phase which takes 10-15 mins and the second part is the testing part which only takes a few seconds to perform classification and calculate accuracy.

2. System



This system provides a well-analyzed data mining framework which indicates the accidents and their causes according to the authentic data obtained from data.gov.in. In this, we will use different data mining algorithms to analyze the data. Different algorithms are applied to group the accident locations into clusters and mining techniques are used to characterize the locations.



2.1 Requirements

2.1.1 Functional Requirement

Type	Name
Programming Language with Version	Python 3.5 & above
Operating System	Windows 7 & above. Linux based OS Mac OS

2.1.2 User Requirement: The user should have the required data before implementing the project.

Dataset 1: Road Accident Data Set – 2018

Rahul Dhanola

FileHomeInsertDrawPage LayoutFormulasDataReviewViewHelp

Tell me what you want to do

KU38

	KB	KC	KD	KE	KF	KG	KH	KI	KJ	KK	KL	KM	KN	KO	KP	KQ	KR	KS	KT	KU	KV
15	1718	920	7311	8231	39357	2	2491	4597	4952	51397	4388	7941	9127	10066	10061	5791	2250	1756	17	51397	
16	2	4	0	4	29958	5	1091	1616	3052	35717	3530	4928	5578	5880	6112	4266	2844	2456	123	35717	
17	4	16	24	40	298	27	93	71	139	601	53	100	107	108	149	62	13	9	0	601	
18	156	96	85	181	146	30	110	49	94	399	48	48	49	50	46	45	50	63	0	399	
19	5	1	10	11	39	35	3	7	4	53	5	5	7	16	7	8	5	0	0	53	
20	7	10	11	21	166	29	21	115	128	430	54	52	47	62	69	60	32	31	23	430	
21	827	1035	806	1841	7369	13	995	1775	1123	11262	1210	1648	1721	1996	2145	1136	576	751	79	11262	
22	387	194	94	288	3237	20	689	777	1725	6428	755	802	745	843	1220	845	411	431	376	6428	
23	1422	1208	1813	3021	19206	7	558	548	1431	21743	2369	3308	3523	4343	4392	2022	895	827	64	21743	
24	0	0	0	0	134	31	31	15	0	180	13	26	27	38	41	19	0	16	0	180	
25	250	12	948	960	50519	1	2247	6454	4700	63920	7657	9384	9478	11171	14238	6223	2450	3319	0	63920	
26	1340	608	4739	5347	17972	9	464	571	3223	22230	2372	3196	3442	4249	4485	2260	1043	1183	0	22230	
27	0	0	0	0	482	25	49	20	1	552	51	79	87	112	119	81	16	7	0	552	
28	89	161	14	175	694	24	3	30	741	1468	194	177	177	194	262	178	99	138	49	1468	
29	4814	4218	1895	6113	22803	6	5787	5194	8784	42568	5654	5998	5551	6081	5368						
30	0	239	117	356	6871	15	547	1361	3926	12705	1324	1936	1734	1756	1680						
31	0	0	0	0	254	28	0	0	0	254	35	38	40	48	52						
32	0	0	0	0	311	26	0	5	0	316	39	42	41	34	60						
33	2	0	0	0	47	34	0	6	27	80	5	12	8	14	26						
34	1	0	0	0	54	33	2	3	17	76	4	8	10	14	21						
35	373	286	1166	1452	3897	17	132	188	2298	6515	735	771	720	748	1044						
36	0	2	0	2	3	36	0	0	0	3	0	0	1	1	1						
37	34	111	173	284	1094	23	34	429	40	1597	188	244	213	256	297						
38																					
39																					
40																					

Road Accident Data Set - 2018

Ready

Accessibility: Unavailable

25°C

Rain to stop

<

Dataset 2: Road Accident Labelled Data

Road Accident Labelled Data - Excel

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

N577

00-300hrs - (Night)

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	States/UT	JUNCTION	VEHICLE A	HUMAN A	PERSON V	AREA	TYPE OF P	LOAD OF	TRAFFIC R	WEATHER	VEHICLE T	TYPE OF R	LICENSE	TIME							
2	Andhra Pr	T-Junction	Less than	18 Yrs	-Ma Drivers	Residentii	Urban	Normally	Over-Spe	Sunny/Cle	Pedestria	Straight R	License V	06-0900hrs - (Day)							
3	Andhra Pr	Y-Junction	5.1 - 10 Ye	18 Yrs	- Fe Passenger	Institutior	Rural	Overload;	Drunken	Rainy	Pedestria	Curved Rc	Learner's	09-1200hrs - (Day)							
4	Andhra Pr	Four arm	10.1 - 15 Y	18-25 Yrs	- Drivers	Market/C	Urban	Others	Driving on	Foggy & M	Bycycles -	Bridge	Without L	12-1500hrs - (Day)							
5	Andhra Pr	Staggered	> 15 Years	18-25 Yrs	- Passenger	Open Are	Rural	Normally	Jumping R	Hail/Sleet	Bycycles -	Culvert	License V	15-1800hrs - (Day)							
6	Andhra Pr	Round ab	Age not ki	25-35 Yrs	- Drivers	Residentii	Urban	Overload;	Use of Mo	Others	Two Whe	Pot Holes	Learner's	18-2100hrs - (Night)							
7	Andhra Pr	Staggered	Age not ki	25-35 Yrs	- Passenger	Institutior	Rural	Others	Over-Spe	Sunny/Cle	Two Whe	Steep Gra	Without L	21-2400hrs - (Night)							
8	Andhra Pr	T-Junction	5.1 - 10 Ye	35-40 Yrs	- Drivers	Market/C	Urban	Normally	Drunken	Rainy	Auto Rick	Ongoing R	License V	00-300hrs - (Night)							
9	Andhra Pr	Y-Junction	10.1 - 15 Y	35-40 Yrs	- Passenger	Open Are	Rural	Overload;	Driving on	Foggy & M	Auto Rick	Others	Learner's	03-600hrs - (Night)							
10	Andhra Pr	Four arm	> 15 Years	45-60 Yrs	- Drivers	Residentii	Urban	Others	Jumping R	Hail/Sleet	Cars & tax	Straight R	Without L	Unknown Time							
11	Andhra Pr	Staggered	Age not ki	45-60 Yrs	- Passenger	Institutior	Rural	Normally	Use of Mo	Others	Cars & tax	Curved Rc	License V	06-0900hrs - (Day)							
12	Andhra Pr	Round ab	Less than	60 Yrs	abo Drivers	Market/C	Urban	Overload;	Over-Spe	Sunny/Cle	Trucks/Loi	Bridge	Learner's	09-1200hrs - (Day)							
13	Andhra Pr	Others	5.1 - 10 Ye	60 Yrs	abo Passenger	Open Are	Rural	Others	Drunken	Rainy	Trucks/Loi	Culvert	Without L	12-1500hrs - (Day)							
14	Andhra Pr	Y-Junction	10.1 - 15 Y	Age not ki	Drivers	Residentii	Urban	Normally	Driving on	Foggy & M	Buses - M	Pot Holes	License V	15-1800hrs - (Day)							
15	Andhra Pr	Four arm	> 15 Years	Age not ki	Passenger	Institutior	Rural	Overload;	Jumping R	Hail/Sleet	Buses - Fe	Steep Gra	Learner's	18-2100hrs - (Night)							
16	Andhra Pr	Staggered	Age not ki	Age not ki	Passenger	Open Are	Rural	Others	Use of Mo	Others	Other Mo	Ongoing R	Without L	21-2400hrs - (Night)							
17	Andhra Pr	Staggered	Age not ki	Age not ki	Passenger	Open Are	Rural	Others	Use of Mo	Others	Other Mo	Others	Without L	00-300hrs - (Night)							
18	Arunachal	T-Junction	Less than	18 Yrs	-Ma Drivers	Residentii	Urban	Normally	Over-Spe	Sunny/Cle	Pedestria	Straight R	License V	06-0900hrs - (Day)							
19	Arunachal	Y-Junction	5.1 - 10 Ye	18 Yrs	- Fe Passenger	Institutior	Rural	Overload;	Drunken	Rainy	Pedestria	Curved Rc	Learner's	09-1200hrs - (Day)							
20	Arunachal	Four arm	10.1 - 15 Y	18-25 Yrs	- Drivers	Market/C	Urban	Others	Driving on	Foggy & M	Bycycles -	Bridge	Without L	12-1500hrs - (Day)							
21	Arunachal	Staggered	> 15 Years	18-25 Yrs	- Passenger	Open Are	Rural	Normally	Jumping R	Hail/Sleet	Bycycles -	Culvert	License V	15-1800hrs - (Day)							
22	Arunachal	Round ab	Age not ki	25-35 Yrs	- Drivers	Residentii	Urban	Overload;	Use of Mo	Others	Two Whe	Pot Holes	Learner's	18-2100hrs - (Night)							
23	Arunachal	Others	Less than	25-35 Yrs	- Passenger	Institutior	Rural	Others	Over-Spe	Sunny/Cle	Two Whe	Steep Gra	Without L	21-2400hrs - (Night)							
24	Arunachal	T-Junction	5.1 - 10 Ye	35-40 Yrs	- Drivers	Market/C	Urban	Normally	Drunken	Rainy	Auto Rick	Ongoing R	License V	00-300hrs - (Night)							

Snipping Tool

Snip copied to clipboard
Select here to mark up and share the image

25°C
Rain to stop

ENG US

1:38 PM
8/20/2022

Dataset 3: Final Road Accident Data

Final Road Accident Dataset - Excel

Rahul Dhanola

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

States/UTs

States/UT	JUNCTION	VEHICLE A	HUMAN A	PERSON V	AREA	TYPE OF P	LOAD OF	TRAFFIC R	WEATHER	VEHICLE T	TYPE OF R	LICENSE	TIME	ACCIDENT OCCURRENCE
Andhra Pr	T-Junctior	Less than	18 Yrs	-Ma Drivers	Residentii	Urban	Normally	Over-Spe	Sunny/Cle	Pedestria	Straight R	License V	06-0900hr	YES
Andhra Pr	Y-Junctior	5.1 - 10 Ye	18 Yrs	-Fe Passenger	Institutior	Rural	Overload	Drunken	F Rainy	Pedestria	Curved Rc	Learner's	09-1200hr	YES
Andhra Pr	Four arm	10.1 - 15 Y	18-25 Yrs	- Drivers	Market/C	Urban	Others	Driving on	Foggy & M	Bycycles -	Bridge	Without L	12-1500hr	YES
Andhra Pr	Staggered	> 15 Years	18-25 Yrs	- Passenger	Open Are	Rural	Normally	Jumping R	Hail/Sleet	Bycycles -	Culvert	License V	15-1800hr	YES
Andhra Pr	Round ab	Age not ki	25-35 Yrs	- Drivers	Residentii	Urban	Overload	Use of Mo	Others	Two Whe	Pot Holes	Learner's	18-2100hr	YES
Andhra Pr	Others	Less than	25-35 Yrs	- Passenger	Institutior	Rural	Others	Over-Spe	Sunny/Cle	Two Whe	Steep Gra	Without L	21-2400hr	YES
Andhra Pr	T-Junctior	5.1 - 10 Ye	35-40 Yrs	- Drivers	Market/C	Urban	Normally	Drunken	F Rainy	Auto Rick	Ongoing R	License V	00-300hrs	YES
Andhra Pr	Y-Junctior	10.1 - 15 Y	35-40 Yrs	- Passenger	Open Are	Rural	Overload	Driving on	Foggy & M	Auto Rick	Others	Learner's	03-600hrs	YES
Andhra Pr	Four arm	> 15 Years	45-60 Yrs	- Drivers	Residentii	Urban	Others	Jumping R	Hail/Sleet	Cars & tax	Straight R	Without L	Unknown	YES
Andhra Pr	Staggered	Age not ki	45-60 Yrs	- Passenger	Institutior	Rural	Normally	Use of Mo	Others	Cars & tax	Curved Rc	License V	06-0900hr	YES
Andhra Pr	Round ab	Less than	60 Yrs	abo Drivers	Market/C	Urban	Overload	Over-Spe	Sunny/Cle	Trucks/Lo	Bridge	Learner's	09-1200hr	YES
Andhra Pr	Others	5.1 - 10 Ye	60 Yrs	abo Passenger	Open Are	Rural	Others	Drunken	F Rainy	Trucks/Lo	Culvert	Without L	12-1500hr	YES
Andhra Pr	Y-Junctior	10.1 - 15 Y	35-40 Yrs	- Drivers	Residentii	Urban	Normally	Driving on	Foggy & M	Buses - M	Pot Holes	License V	15-1800hr	YES
Andhra Pr	Four arm	> 15 Years	Age not ki	Passenger	Institutior	Rural	Overload	Jumping R	Hail/Sleet	Buses - Fe	Steep Gra	Learner's	18-2100hr	YES
Andhra Pr	Staggered	Age not ki	Age not ki	Passenger	Market/C	Urban	Others	Use of Mo	Others	Other Mo	Ongoing R	Without L	21-2400hr	YES
Andhra Pr	Staggered	Age not ki	Age not ki	Passenger	Open Are	Rural	Others	Use of Mo	Others	Other Mo	Others	Without L	00-300hrs	YES
Arunachal	T-Junctior	Less than	18 Yrs	-Ma Drivers	Residentii	Urban	Normally	Over-Spe	Sunny/Cle	Pedestria	Straight R	License V	06-0900hr	NO
Arunachal	Y-Junctior	5.1 - 10 Ye	18 Yrs	- Fe Passenger	Institutior	Rural	Overload	Drunken	F Rainy	Pedestria	Curved Rc	Learner's	09-1200hr	NO
Arunachal	Four arm	10.1 - 15 Y	18-25 Yrs	- Drivers	Market/C	Urban	Others	Driving on	Foggy & M	Bycycles -	Bridge	Without L	12-1500hr	NO
Arunachal	Staggered	> 15 Years	18-25 Yrs	- Passenger	Open Are	Rural	Normally	Jumping R	Hail/Sleet	Bycycles -	Culvert	License V	15-1800hr	NO
Arunachal	Round ab	Age not ki	25-35 Yrs	- Drivers	Residentii	Urban	Overload	Use of Mo	Others	Two Whe	Pot Holes	Learner's	18-2100hr	NO
Arunachal	Others	Less than	25-35 Yrs	- Passenger	Institutior	Rural	Others	Over-Spe	Sunny/Cle	Two Whe	Steep Gra	Without L	21-2400hr	NO
Arunachal	T-Junctior	5.1 - 10 Ye	35-40 Yrs	- Drivers	Market/C	Urban	Normally	Drunken	F Rainy	Auto Rick	Ongoing R	License V	00-300hrs	NO

Final Road Accident Dataset

Ready Accessibility: Unavailable

25°C Rain to stop

ENG US 1:40 PM 8/20/2022

2.1.3 Environmental Requirement:

Type	Name
System Architecture	32 – bit or 64 – bit
Memory	8 GB 2400 MHz DDR4
Storage	500 GB SATA3 2.4 HDD
Central Processing Unit (CPU)	Intel or AMD - 2 GHz or Faster
Graphical Processor Unit (GPU)	2 GB Nvidia Graphic Processor

2.2 Implementation

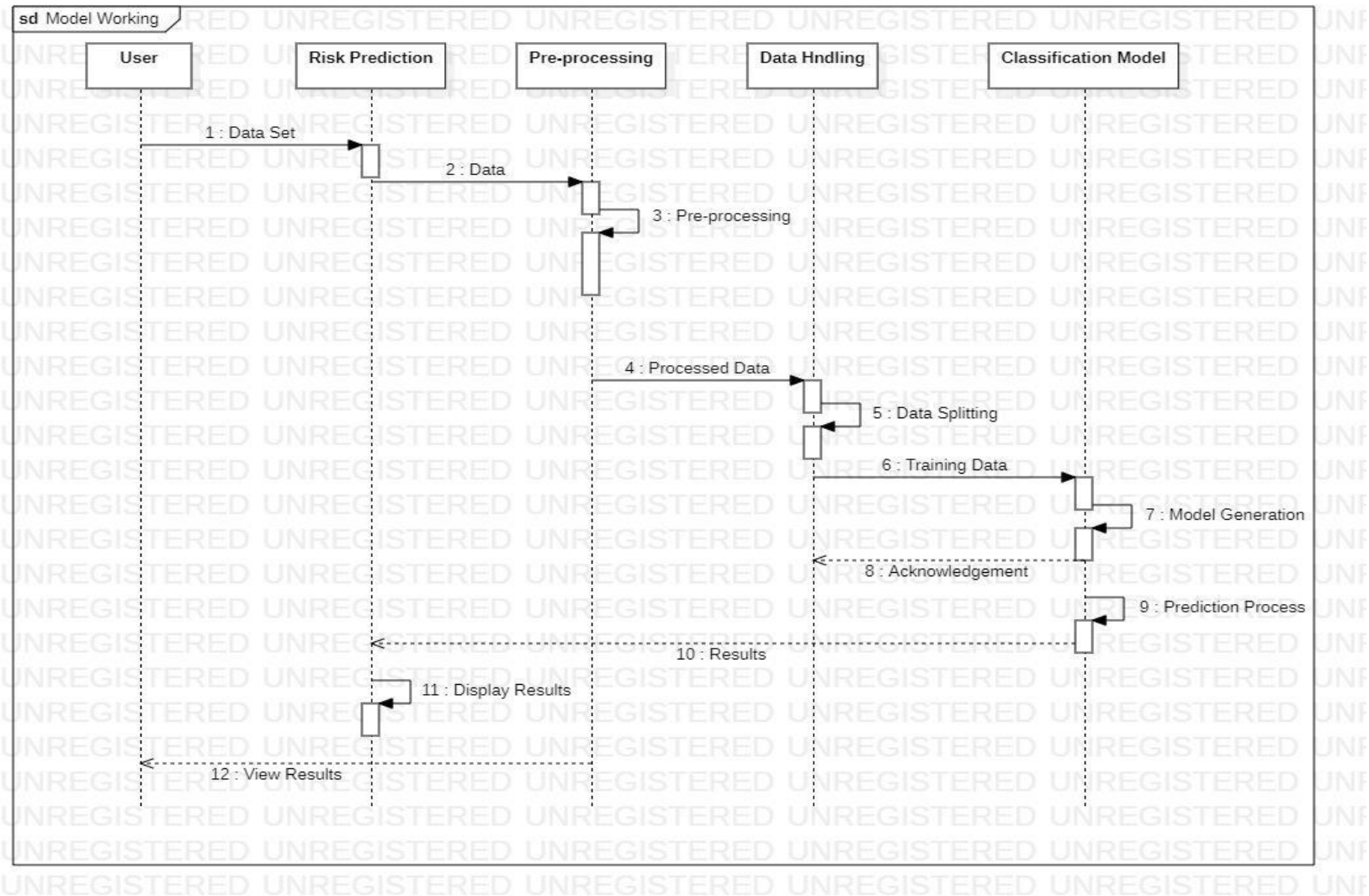


Fig: Sequence Diagram of System

2.3 Testing

We will develop and train the machine learning model from scratch and subsequently, a testing framework will also be made for the classification module. We will use many classification datasets to test the model. We will also document the performance and accuracy of our model.

Test Data1:

```
In [103] testData1 = {
    'States/UTs': ['Kerala'],
    'JUNCTION': ['Four arm Junction'],
    'VEHICLE AGE': ['> 15 Years'],
    'HUMAN AGE AND SEX': ['45-60 Yrs- Male'],
    'PERSON WITHOUT SAFETY PRECAUTIONS': ['Drivers'],
    'AREA': ['Residential Area'],
    'TYPE OF PLACE': ['Urban'],
    'LOAD OF VEHICLE': ['Others'],
    'TRAFFIC RULES VIOLATION': ['Jumping Red Light'],
    'WEATHER': ['Hail/Sleet'],
    'VEHICLE TYPE AND SEX': ['Cars & taxies Vans & LMV - Male'],
    'TYPE OF ROAD': ['Straight Road'],
    'LICENSE': ['Without Licence'],
    'TIME': ['Unknown Time']
}

In [104] # Now the test values are compared with the Column Codes and store them
for col in testData1:
    code = [columnCodes[''.join(testData1[col])]]
    testData1[col] = code
print(testData1)

testDataFrame = pd.DataFrame.from_dict(testData1)
print(testDataFrame)

{'States/UTs': [17], 'JUNCTION': [0], 'VEHICLE AGE': [2], 'HUMAN AGE AND SEX': [9], 'PERSON WITHOUT SAFETY PRECAUTIONS': [0], 'AREA': [3], 'TYPE OF PLACE': [1], 'LOAD OF VEHICLE': [1], 'TRAFFIC RULES VIOLATION': [2], 'WEATHER': [1], 'VEHICLE TYPE AND SEX': [7], 'TYPE OF ROAD': [7], 'LICENSE': [2], 'TIME': [8]}
   States/UTs  JUNCTION  VEHICLE AGE  HUMAN AGE AND SEX  \
0           17         0           2             9      \
0  PERSON WITHOUT SAFETY PRECAUTIONS  AREA  TYPE OF PLACE  LOAD OF VEHICLE  \
0           0           3           1             1      \
0  TRAFFIC RULES VIOLATION  WEATHER  VEHICLE TYPE AND SEX  TYPE OF ROAD  \
0           2           1             7             7      \
0  LICENSE  TIME
0           2     8

In [105] predictionValue = trainedModel.predict(testDataFrame)
if (predictionValue == 1):
    print("Yes, there is a chance of Accident")
else:
    print("No, it is safe")
Yes, there is a chance of Accident
```

Test Data2:

```
In [106] testData2 = {
    'States/UTs': ['Chhattisgarh'],
    'JUNCTION': ['Staggered Junction'],
    'VEHICLE AGE': ['> 15 Years'],
    'HUMAN AGE AND SEX': ['45-60 Yrs- Male'],
    'PERSON WITHOUT SAFETY PRECAUTIONS': ['Passengers'],
    'AREA': ['Residential Area'],
    'TYPE OF PLACE': ['Rural'],
    'LOAD OF VEHICLE': ['Others'],
    'TRAFFIC RULES VIOLATION': ['Jumping Red Light'],
    'WEATHER': ['Hail/Sleet'],
    'VEHICLE TYPE AND SEX': ['Other Motor Vehicles - Female'],
    'TYPE OF ROAD': ['Others'],
    'LICENSE': ['Without Licence'],
    'TIME': ['Unknown Time']
}

In [107] # Now the test values are compared with the Column Codes and store them
for col in testData2:
    code = [columnCodes[''.join(testData2[col])]]
    testData2[col] = code
print(testData2)
testDataFrame = pd.DataFrame.from_dict(testData2)

{'States/UTs': [6], 'JUNCTION': [3], 'VEHICLE AGE': [2], 'HUMAN AGE AND SEX': [9], 'PERSON WITHOUT SAFETY PRECAUTIONS': [1], 'AREA': [3], 'TYPE OF PLACE': [0], 'LOAD OF VEHICLE': [1], 'TRAFFIC RULES VIOLATION': [2], 'WEATHER': [1], 'VEHICLE TYPE AND SEX': [8], 'TYPE OF ROAD': [1], 'LICENSE': [2], 'TIME': [8]}

In [108] predictionValue = trainedModel.predict(testDataFrame)
if (predictionValue == 1):
    print("Yes, there is a chance of Accident")
else:
    print("No, it is safe")
Yes, there is a chance of Accident
```

Test Data3:

```
In [109_] testData3 = {
    'States/UTs': ['Puducherry'],
    'JUNCTION': ['Staggered Junction'],
    'VEHICLE AGE': ['> 15 Years'],
    'HUMAN AGE AND SEX': ['45-60 Yrs- Male'],
    'PERSON WITHOUT SAFETY PRECAUTIONS': ['Passengers'],
    'AREA': ['Residential Area'],
    'TYPE OF PLACE': ['Rural'],
    'LOAD OF VEHICLE': ['Others'],
    'TRAFFIC RULES VIOLATION': ['Jumping Red Light'],
    'WEATHER': ['Hail/Sleet'],
    'VEHICLE TYPE AND SEX': ['Other Motor Vehicles - Female'],
    'TYPE OF ROAD': ['Others'],
    'LICENSE': ['Without Licence'],
    'TIME': ['Unknown Time']
}
```

```
In [110_] # Now the test values are compared with the Column Codes and store them
for col in testData3:
    code = [columnCodes[''.join(testData3[col])]]
    testData3[col] = code
print(testData3)
testDataFrame = pd.DataFrame.from_dict(testData3)

{'States/UTs': [26], 'JUNCTION': [3], 'VEHICLE AGE': [2], 'HUMAN AGE AND SEX': [9], 'PERSON WITHOUT SAFETY PRECAUTIONS': [1], 'AREA': [3], 'TYPE OF PLACE': [0], 'LOAD OF VEHICLE': [1], 'TRAFFIC RULES VIOLATION': [2], 'WEATHER': [1], 'VEHICLE TYPE AND SEX': [8], 'TYPE OF ROAD': [1], 'LICENSE': [2], 'TIME': [8]}
```

```
In [111_] predictionValue = trainedModel.predict(testDataFrame)
if (predictionValue == 1):
    print("Yes, there is a chance of Accident")
else:
    print("No, it is safe")

No, it is safe
```

Test Data4:

```
In [112_] testData4 = {
    'States/UTs': ['Lakshadweep'],
    'JUNCTION': ['Staggered Junction'],
    'VEHICLE AGE': ['10.1 - 15 Years'],
    'HUMAN AGE AND SEX': ['45-60 Yrs- Male'],
    'PERSON WITHOUT SAFETY PRECAUTIONS': ['Drivers'],
    'AREA': ['Institutional Area'],
    'TYPE OF PLACE': ['Rural'],
    'LOAD OF VEHICLE': ['Others'],
    'TRAFFIC RULES VIOLATION': ['Jumping Red Light'],
    'WEATHER': ['Others'],
    'VEHICLE TYPE AND SEX': ['Other Motor Vehicles - Female'],
    'TYPE OF ROAD': ['Others'],
    'LICENSE': ['License Valid Permanent'],
    'TIME': ['00-300hrs - (Night)']
}
```

```
In [113_] # Now the test values are compared with the Column Codes and store them
for col in testData4:
    code = [columnCodes[''.join(testData4[col])]]
    testData4[col] = code
print(testData4)
testDataFrame = pd.DataFrame.from_dict(testData4)

{'States/UTs': [18], 'JUNCTION': [3], 'VEHICLE AGE': [0], 'HUMAN AGE AND SEX': [9], 'PERSON WITHOUT SAFETY PRECAUTIONS': [0], 'AREA': [0], 'TYPE OF PLACE': [0], 'LOAD OF VEHICLE': [1], 'TRAFFIC RULES VIOLATION': [2], 'WEATHER': [1], 'VEHICLE TYPE AND SEX': [8], 'TYPE OF ROAD': [1], 'LICENSE': [1], 'TIME': [0]}
```

```
In [114_] predictionValue = trainedModel.predict(testDataFrame)
if (predictionValue == 1):
    print("Yes, there is a chance of Accident")
else:
    print("No, it is safe")

No, it is safe
```

2.4 Evaluation

We will develop and train the machine learning model from scratch and subsequently, a testing framework will also be made for the classification module. We will use many classification datasets to test the model. We will also document the performance and accuracy of our model.

Whenever we build Machine Learning models, we need some form of metric used for the measurement of the goodness of the model. Bear in mind that the “goodness” of the model could have multiple interpretations, but generally when we speak of it in a Machine Learning context, we are talking of the measure of a model's performance on new instances that weren't a part of the training data.

Some common intrinsic metrics to evaluate the system are as follows:

- **Accuracy**

The accuracy of a Machine Learning classification algorithm is one way to measure how often the algorithm classifies a data point correctly. Accuracy is the number of correctly predicted data points out of all the data points.

$$\text{Accuracy} = \frac{\text{Number of correct predictions}}{\text{Total number of prediction}}$$

- **Precision**

Precision (also called positive predictive value) is the fraction of relevant instances among the retrieved instances

$$\begin{aligned} \text{Precision} &= \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}} \\ &= \frac{\text{True Positive}}{\text{Total predicted positive}} \end{aligned}$$

- **Recall**

Recall measures the proportion of actual positive labels correctly identified by the model.

$$\begin{aligned} \text{Recall} &= \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}} \\ &= \frac{\text{True Positive}}{\text{Total Actual positive}} \end{aligned}$$

3 Snapshots of the Project

Snapshot 1:

A Data Mining Framework To Analyze Road Accident Data - Road Accident Data Analysis



Introduction

Road accidents are uncertain and unpredictable incidents and their analysis requires the knowledge of the factors affecting them. Fatalities and injuries resulting from road traffic accidents are a major and growing public health problem in India. Every week nearly 2,650 people get killed and 9,000 get injured due to traffic accidents.

Traffic accidents have now earned India a dubious distinction; with nearly 140,000 deaths annually, the country has overtaken China to top the world in road fatalities. India is the only country in the world which faces more than 15 fatalities and 53 injuries every hour as a consequence of road crashes.

The major problem in the analysis of accident data is its Heterogeneous nature. Thus, heterogeneity must be considered during analysis of the data. Road accident analysis aims to investigate the main factors that characterize an accident to understand patterns or behaviors and, consequently, to identify the appropriate countermeasures to adopt to avoid the accident.

Importing the Libraries

```
In [1]: # Importing the necessary libraries
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
```

Reading the Dataset



```
In [2]: # Reading the csv file
data = pd.read_csv(
    '../input/road-accident-analysis-data/Road Accident Analysis Data/Analysis/Road Accident Data Set - 2018.csv',
    index_col=0)
```

Snapshot 2:

Exploratory Data Analysis (EDA)

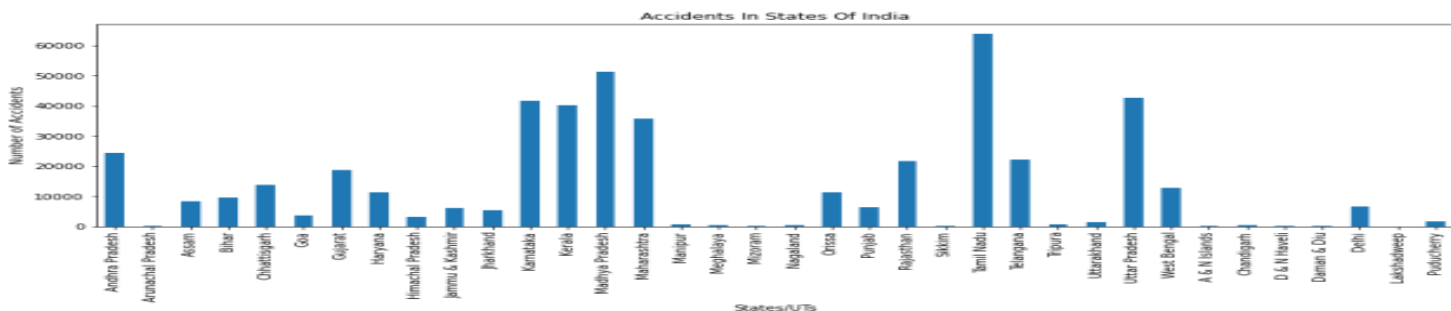


In statistics, exploratory data analysis is an approach of analyzing data sets to summarize their main characteristics, often using statistical graphics and other data visualization methods. A statistical model can be used or not, but primarily EDA is for seeing what the data can tell us beyond the formal modeling or hypothesis testing task.

Total Number Of Accidents Occurred across The Indian States

```
In [8]: # Total Number Of Accidents Occurred across The Indian States
data['Total Accidents'].plot.bar(figsize=(15, 5),
    title="Accidents In States Of India")
plot.ylabel('Number of Accidents')
```

Out[8]: Text(0, 0.5, 'Number of Accidents')



Snapshot 3:

Model - 2 => Using Random Forest

```
In [99]: RF = RandomForestClassifier()
RF.fit(accidentData_train, target_train)
target_pred_rf = RF.predict(accidentData_test)
print('Accuracy score: {0:0.2f}'.format(
    metrics.accuracy_score(target_test, target_pred_rf)))
```

Accuracy score: 0.57

We have achieved a accuracy of 55%.

Model - 3 => Using Decision Tree

```
In [100]: DT = DecisionTreeClassifier()
DT.fit(accidentData_train, target_train)
target_pred_dt = DT.predict(accidentData_test)
print('Accuracy score: {0:0.2f}'.format(
    metrics.accuracy_score(target_test, target_pred_dt)))
```

Accuracy score: 0.91

We have achieved a accuracy of 94%.

Model - 4 => Using Decision Tree - Using AdaBoostClassifier

```
In [101]: # Create adaboost classifier object
adaboostClassifier = AdaBoostClassifier()

# Train Adaboost Classifier
trainedModel = adaboostClassifier.fit(accidentData_train, target_train)

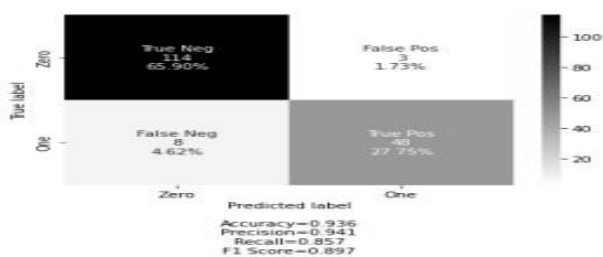
# Predict the response for test dataset
target_pred_ada = trainedModel.predict(accidentData_test)

# Printing the Accuracy
print("Accuracy: {0:0.2f}".format(
    metrics.accuracy_score(target_test, target_pred_ada)))
```

Accuracy: 0.94

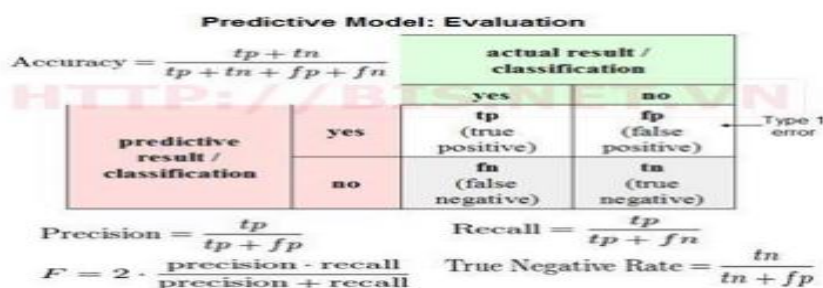
We have achieved a accuracy of 95%.

Snapshot 4:



Evaluating The Model

Whenever we build Machine Learning models, we need some form of metric used for the measurement of the goodness of the model. Bear in mind that the "goodness" of the model could have multiple interpretations, but generally when we speak of it in a Machine Learning context, we are talking of the measure of a model's performance on new instances that weren't a part of the training data.



4 *Conclusions*

This framework brings us into an analysis of road accident data which indicates to us the major causes and conditions of an accident where it is important to understand the basic probabilities where we get the maximum accident depending upon the region as well as the factors like over speeding, drink and drive and other major and minute factors to inculcate the accident and also to improve or prevent them from happening in all situations depending on weather conditions and regions where the accidents take place majorly.

This framework also helps to analyze the fatal rates to indicate the level of the cause and accidents, using clustering and algorithms we can conclude the probabilities and accuracy of the model that what will be the expected rate of an accident to take place and further what can be done to prevent and take effective measure to ensure the safety and health of a driver or a suspect of the accident, this also brings us an understanding of how is it important to inculcate the best way to prevent an accident on whatever the condition is like a natural cause of weather or traffic timings and regions and localities also this brings us to a better way to drive and be safe.

This project can help an organization which can implement and measure better equipment and all the aspects to prevent and help suspect to save his/her life in any minimal or major accident taking place across the Indian region also this can help a government department of road safety to understand and analysis root cause and their solutions to take immediate action plan to suppose the betterment of each citizens road safety.

5 *Further Development Or Research*

In this project further, there can be analyzed based on all parameters including datasets of weather conditions, road health and safety measures, road traffic, traffic signals and all other geographical situations.

The optimistic approach for the framework can result in an analysis that can help build an integrated solution to prevent major accidents that can be very useful to save lives as it can indicate factors and aspects that can improve the road safety measures and their use in a global level to ensure the better road traffic and management, also this can lead to the betterment of government policies on road safety and also in constructing the roads based on the better placements of every equipment to ensure the better approach and safety as well as precautions.

6 *References*

- [1] A. Priyanka and K. Sathiyakumari, "A comparative study of classification algorithm using accident data", International Journal of Computer Science & Engineering Technology (IJCSET)
- [2] Karlaftis M, Tarko A. Heterogeneity considerations in accident modelling. *Accid Anal Prev.* 1998;30(4):425–33.
- [3] Article(online) Available: <https://www.hindustantimes.com/mumbai-news/india-had-most-deaths-in-road-accidents-in-2019-report/story-pikRXxsS4hptNVvf6J2g9O.html>
- [4] Road Traffic Accidents in India: Issues and Challenges, Sanjay Kumar Singh <https://www.sciencedirect.com/science/article/pii/S2352146517307913>
- [5] Ng KS, Hung WT, Wong WG. An algorithm for assessing the risk of traffic accidents. *J Saf Res.* 2002;33:387–410.
- [6] T. Dipo, I Akomolafe and Akinbola Olutayo, "Using Data Mining Technique to Predict Cause of Accident and Accident Prone Locations on Highways", *American Journal of Database Theory and Application* 2012
- [7] Article by S. Nagendra Babu “A Data Mining Framework to Analyze Road Accident Data using Map Reduce Methods CCMF and TCAMP Algorithms”