# DATA WAREHOUSING MSIS 5663 TERM PROJECT – SPRING 2019

## Team 7

Roshith Elangovan- A20161234 Krithika Ganesh Kumar- A20162108 Anudeep Subraveti- A20168818

# **CONTENT**

1.	Introdu	iction	2
2.	Establi	shing Connection	2
3.	Creatin	ng data source view	2
4.	Creatin	ng named calculations	3
5.	Creatin	ng dimensions	4
6.	Creatin	ng and deploying cube	5
7.	Creatin	ng hierarchies	6
8.	Creatin	ng Partitions and Aggregations	10
9.	Creatin	ng and Executing MDX Queries	11
10.	Data M	fining	19
11.	Predict	for and target variables	19
12.	Buildir	ng mining models	21
13.	Creatin	ng mining structure	21
14.	Buildir	ng models	22
	a.	Decision tree	22
	b.	Results	23
	c.	Neural networks	24
	d.	Results	25
	e.	Logistic regression	26
	f.	Results	27
15.	Compa	urison of results	28
	a.	Lift Chart	28
	b.	Classification matrix	29
16.	Cross -	validation	30
17.	Summa	ary of results	37
18.	Conclu	asion	38

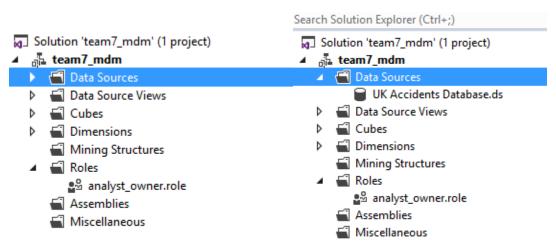
## **Introduction**

The dataset given is "UK\_Accidents\_database", is used in relational database format. It consists of 4,650,859 records and four dimensions namely Date, Vehicles, Causalities and Accidents.

The following steps have been followed:

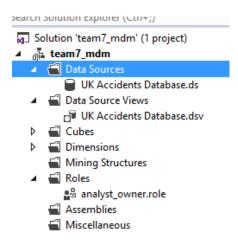
#### **Step 1: Establishing Connection:**

- A data Connection is made to UK\_Accidents\_database on the server (stwssbsql01.ad.okstate.edu)
- We selected NewProject>Import from Server from Analysis Server under Business Intelligence
- We added the data source

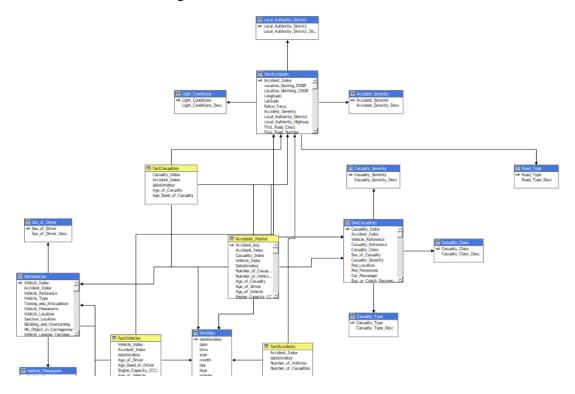


## **Step 2: Creating Data Source View:**

 We created a Data Source View of the UK\_Accidents\_database which will serve as a disconnected data source for the project



We checked the data source view diagram



# **Step 3: Creating Named Calculations:**

We created the following named calculations:

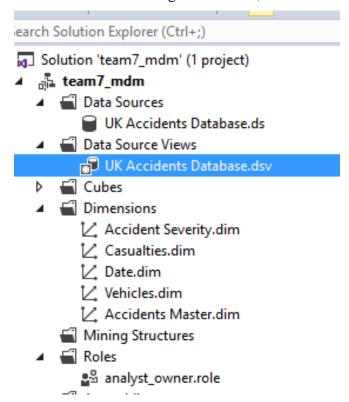
Named Calculation	Description
Police_Officer_Attended_Desc	Description of Police officer who attended the
	scene

Junc_Control_Desc	Description of the Junction at which the
	accident took place
Journey_Purpose_Desc	Journey of the passenger in the vehicle during
	the time of accident
DateofAccident	Date at which the accident took place in the
	format Month Date, Year

# **Step 4: Creating Dimensions:**

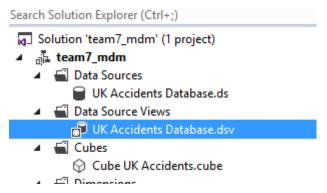
Next step is to create Dimensions.

We created the following Dimensions,



# **Step 5: Creating and deploying the cube:**

Since the dimensions are intact, next step in the process is to create cube.



## **Step 6: Creating calculated measures:**

We created the following calculated measures:

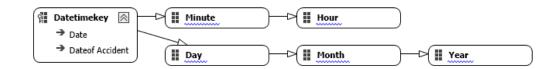
Measures	Description	Formula
[Average Age of Casualties]	This measures gives the	[Measures].[Age Of
	average age of the casualties	Casualty]/[Measures].[Number
		Of Casualties]
[Driver to Casualties Ratio]	This gives the ratio between	[Measures].[Age Of
	driver's age and casualty	Vehicle]/[Measures].[Number
		Of Vehicles]
Minimum Number of	This gives the minimum	
Casualties	number of casulaties	
Engine_Capacity_CC	This gives the distinct count	
Distinct Count	of Engine capacity that will	
	help in determining the	
	fatality of the accident	

## **Step 7: Creating Hierarchies:**

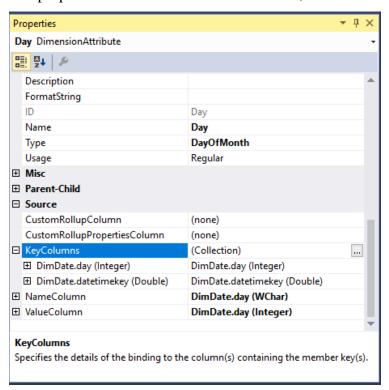
We created hierarchies in two dimensions:

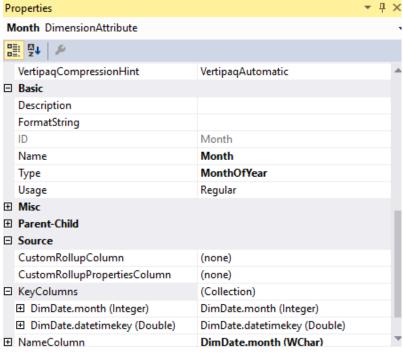
#### (i) Date Dimension:





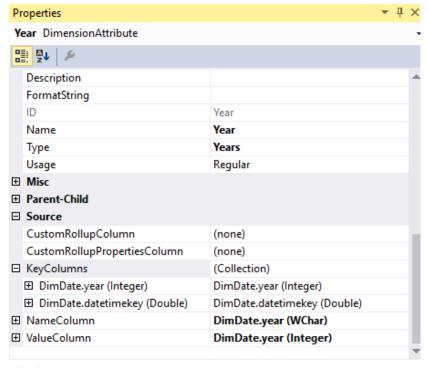
The properties of the attributes are as follows,





#### KeyColumns

Specifies the details of the binding to the column(s) containing the member key(s).

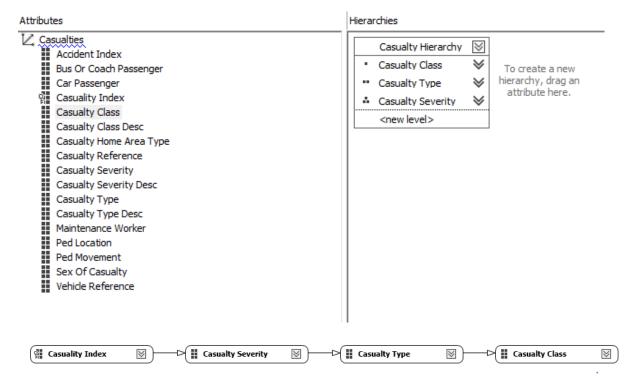


#### KeyColumns

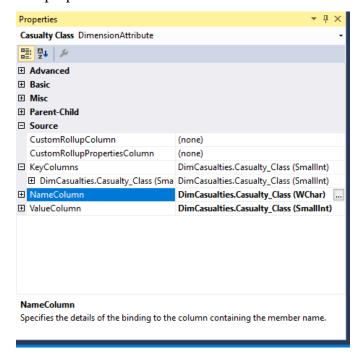
Specifies the details of the binding to the column(s) containing the member key(s).

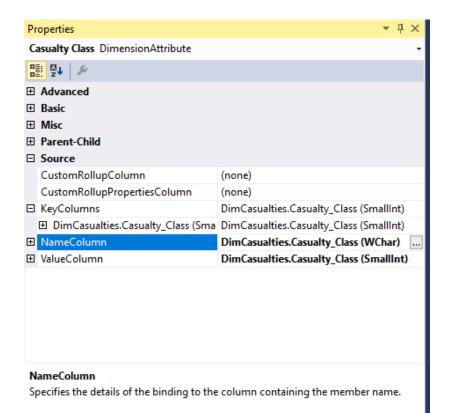
#### (ii) Casualties Dimension:

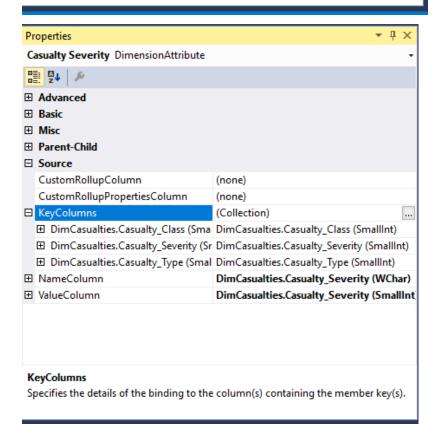
We created casualties hierarchies in casualties dimension



The properties of the attributes are as follows:





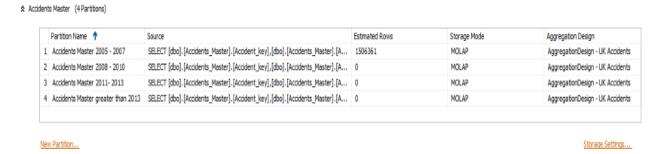


## **Step 8: Creating Partitions and Aggregations:**

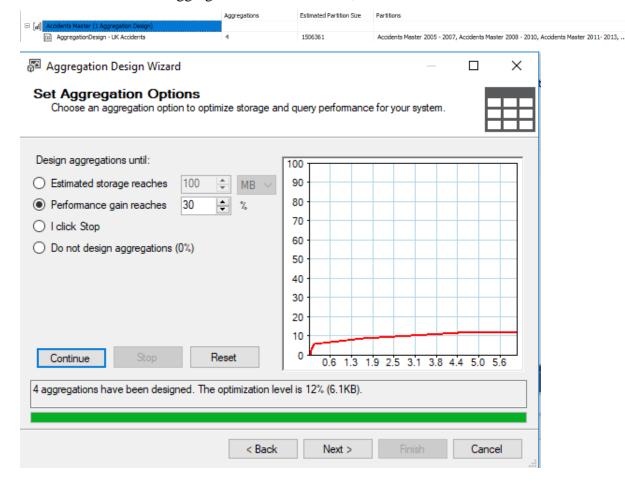
The main purpose of partitions and aggregations are to speedup the process of querying as we know data warehouse deals with large data and speeding up queries is vital.

We created partitions in Accident Master table based on years.

Since there are 11 years, 2005-2011. We decided to create 4 partitions containing 3 years each.

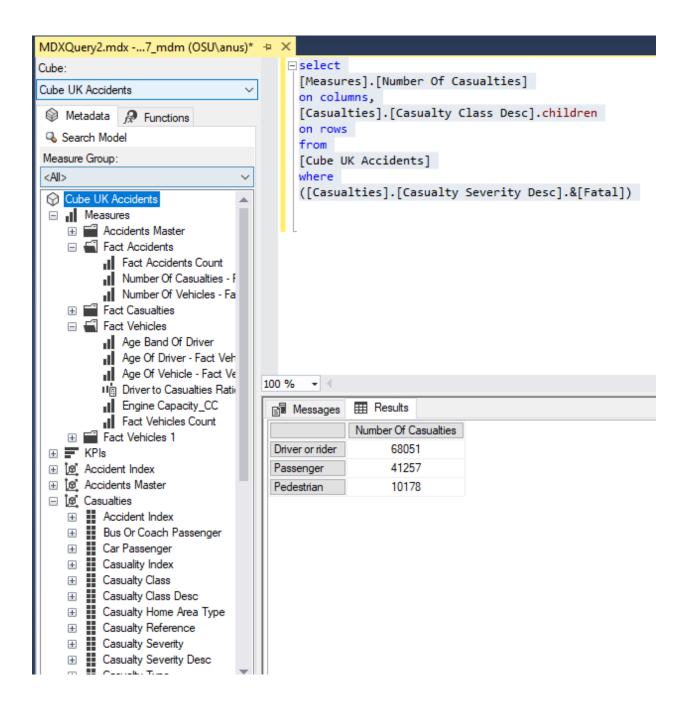


We decided to do 30% aggregation as shown below,



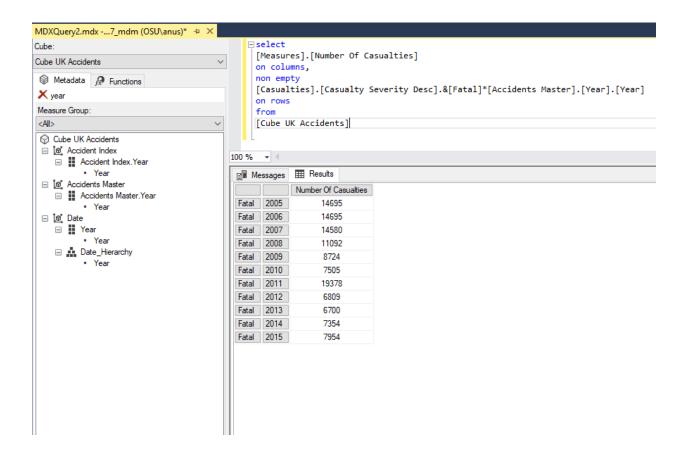
## **Step 9:Designing and Executing MDX Queries:**

1. What are the total number of casualties in each class?



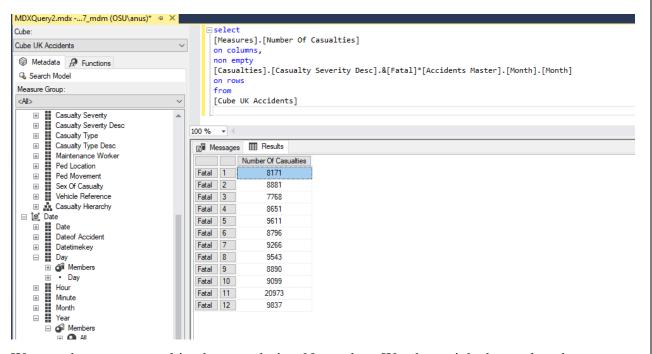
Driver or Rider are the most affected in casualties than other classes.

2. Number of fatal accidents in each year?



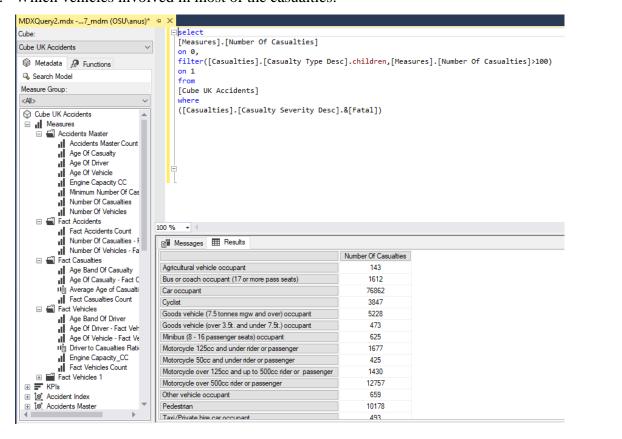
From the query we see that, 2005 and 2006 have the most number of casualties.

3. Total number of casualties in each month of all years



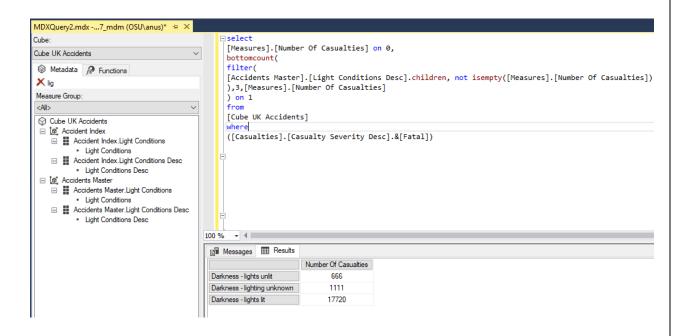
We see that most casualties happen during November, Weather might have played an important role in these accidents.

4. Which vehicles involved in most of the casualties:



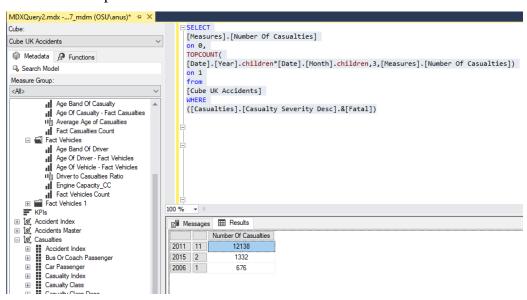
We see that car occupants are the most affected by casualties

5. Which light conditions have the most casualties?



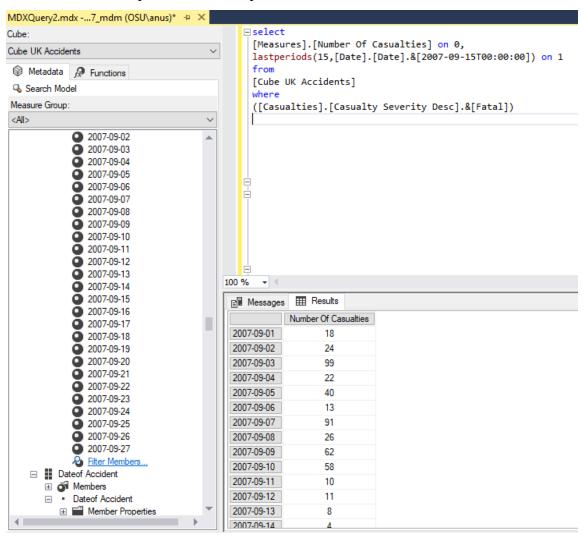
We can understand from the query that dark roads with lights lit suffer the most number of casualties.

6. What are the top 3 months in which the casualties occur?

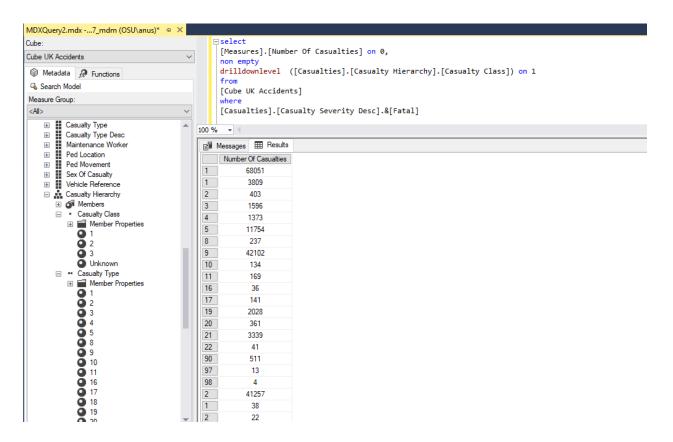


Casualties occur the most in November, followed by February and January.



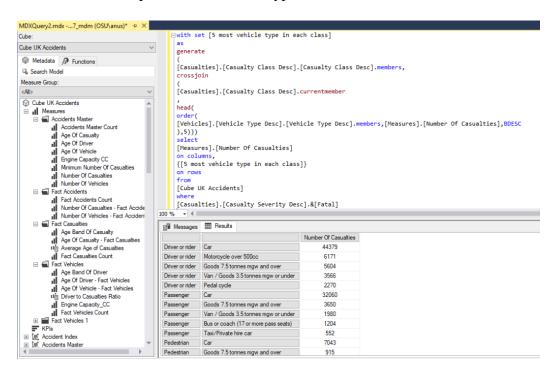


8. Drilldown the casualties hierarchy where the casualty severity is fatal

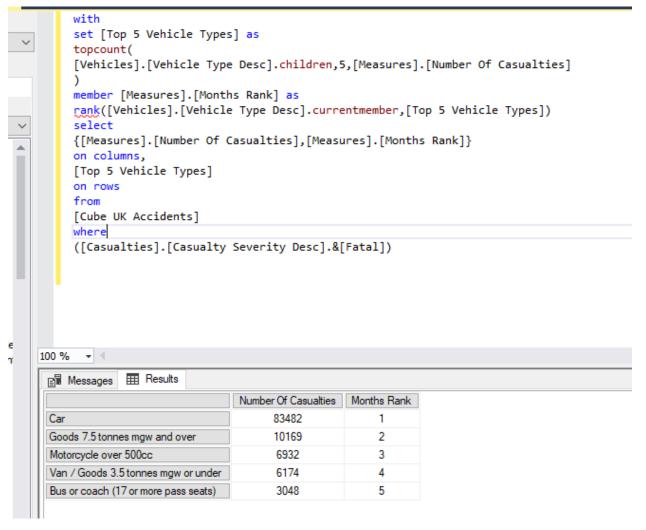


We see that in level one, the descendants are from 1 to 98 followed by the hierarchy level 2 with number of fatalities.

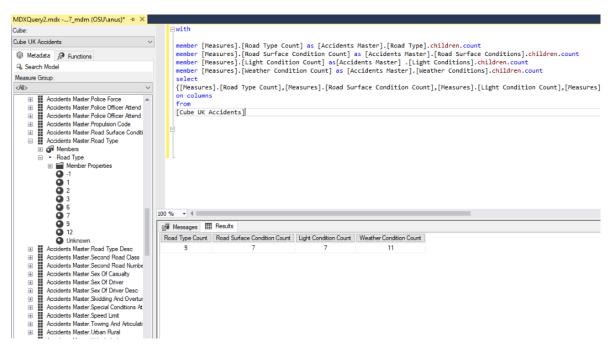
9. What are the top 5 vehicles in each types that suffer fatal accidents



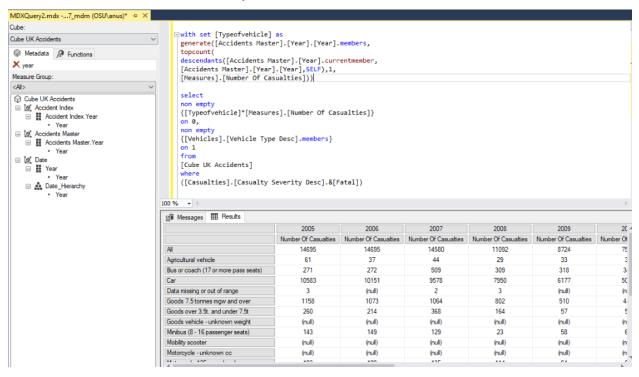
10. Rank each vehicles according to the number of fatalities



11. Display the count of Road type, Road surface condition, Light condition, Weather condition to know the types



12. Number of fatalities that occur each year in each vehicle



## **Data Mining:**

As we know that data mining is process of knowledge discovery from the large databases and datasets that can be used for further analysis in intelligent decision making. Without proper data mining and knowledge discovery, there is no purpose of maintaining large databases and Enterprise warehouses with large setting and operating costs. Data mining focuses on specific machine learning methods and procedures to empower predictive modeling and pattern finding hence enabling intelligent decision making.

With SQL Server Management Studio, the mining models can be created using Microsoft specific language called Data Mining Extensions abbreviated as (DMX) which will be used for creating mining structures, mining models, to train, test and evaluate the built models.

In the project, we have used UK Accidents Database. The predictive models will be built to predict the Accident severity. This will be predicted based on various independent variables which will be discussed in later part.

Since the predictor variable is Accident Severity, we have chosen to use

Decision Tree: As it performs better for both classification and regression problems

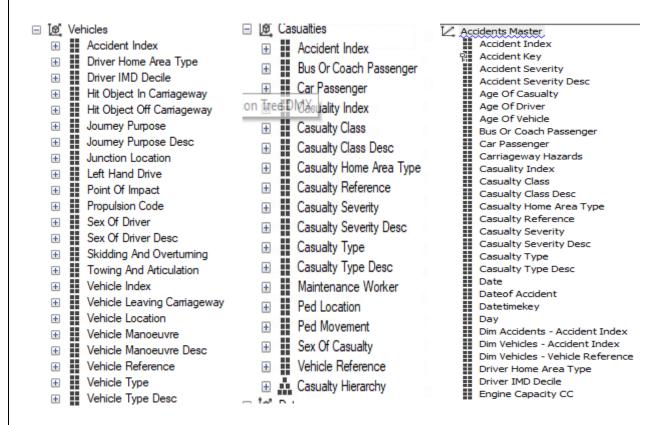
Logistic Regression: As logistic regression is the most sought classification problem for categorical variables

Neural Networks: We chose Neural Networks as we believe Neural Networks to be better for both supervised and unsupervised learning in classification problems.

### **Predictor and Target variables:**

As we have chosen the Target variable as accident severity, we thought it would be interesting to know which type of accident lead to serious injury.

We built the model using SSMS, considering the following tables,



As using all attributes can lead to underperformance and overfitting of the model. We chose to use several variables that were considered helpful in predicting the accident severity.

```
[Accident Index],
[Accident Severity] PREDICT,
[Carriageway Hazards],
[First Road Class],
[Junction Control],
[Junction Detail],
[Light Conditions],
[Ped Cross Physical],
[Road Surface Conditions],
[Road Type],
[Special Conditions At Site],
[Speed Limit],
[Urban Rural],
[Weather Conditions],
[Sex of Casualty],
[Journey Purpose],
[Vehicle Type]
```

We chose the above variables to predict the accident severity.

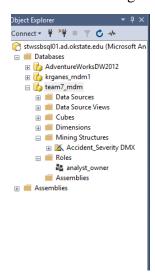
## **Building Mining Models:**

#### **Creating Mining Structure:**

We started with creating the Accident-Severity DMX

```
CREATE MINING STRUCTURE [Accident_Severity DMX]
[Accident Index] LONG KEY,
[Accident Severity] LONG DISCRETE,
[Carriageway Hazards] TEXT DISCRETE,
[First Road Class] TEXT DISCRETE,
[Junction Control] TEXT DISCRETE,
[Junction Detail] TEXT DISCRETE,
[Light Conditions] TEXT DISCRETE,
[Ped Cross Physical] TEXT DISCRETE,
[Road Surface Conditions] TEXT DISCRETE,
[Road Type] TEXT DISCRETE,
[Special Conditions At Site] TEXT DISCRETE,
[Speed Limit] TEXT DISCRETE,
[Urban Rural] TEXT DISCRETE,
[Weather Conditions] TEXT DISCRETE,
[Sex of Casualty] TEXT DISCRETE,
[Journey Purpose] TEXT DISCRETE,
[Vehicle Type] TEXT DISCRETE
)
WITH HOLDOUT (30 PERCENT)
```

The created mining structure is shown below,



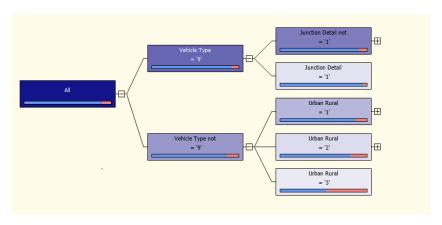
#### **Building Models:**

#### **Decision Trees:**

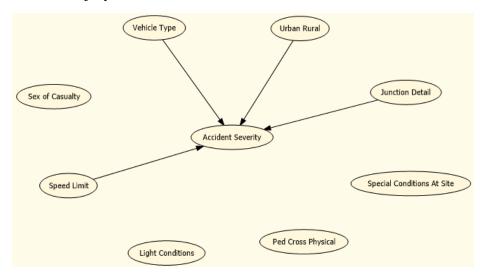
```
ALTER MINING STRUCTURE [Accidents DMX]
ADD MINING MODEL [Decision Tree DMX]
( [Accident Index],
[Accident Severity] PREDICT,
[Carriageway Hazards],
[First Road Class],
[Junction Control],
[Junction Detail],
[Light Conditions],
[Ped Cross Physical],
[Road Surface Conditions],
[Road Type],
[Special Conditions At Site],
[Speed Limit],
[Urban Rural],
[Weather Conditions],
[Sex of Casualty],
[Journey Purpose],
[Vehicle Type] )
USING Microsoft_Decision_Trees
```

```
INSERT INTO MINING STRUCTURE [Accidents DMX]
( [Accident Index], [Accident Severity],
[Carriageway Hazards], [First Road Class],
[Junction Control], [Junction Detail],
[Light Conditions], [Ped Cross Physical],
[Road Surface Conditions],
[Road Type], [Special Conditions At Site], [Speed Limit], [Urban Rural],
[Weather Conditions], [Sex of Casualty],
[Journey Purpose], [Vehicle Type] )
OPENQUERY([UK Accidents Database],
'SELECT TOP 100000 DimAccidents.Accident Index,
Accident_Severity,
Carriageway_Hazards,
First Road Class,
Junction_Control,
Junction_Detail,
Light_Conditions,
Ped_Cross_Physical,
Road_Surface_Conditions,
Road_Type,
Special_Conditions_At_Site,
Speed_Limit, Urban_Rural,
Weather_Conditions, Sex_of_Casualty,
Journey_Purpose, Vehicle_Type FROM DimAccidents,
DimCasualties, DimVehicles
where DimAccidents.Accident_Index = DimCasualties.Accident_Index AND
DimAccidents.Accident Index = DimVehicles.Accident Index AND
Accident_Severity != 1 AND
Carriageway_Hazards != -1 AND
First_Road_Class != 6 AND
Junction_Control != -1 AND
Junction_Detail != -1 AND
Light Conditions != -1 AND
Ped Cross Physical != -1 AND
Road_Surface_Conditions != -1 AND
Road_Type != -1 AND
Road_Type != 9 AND
Special_Conditions_At_Site != 9 AND
Weather_Conditions != -1 AND
Sex of Casualty != -1 AND
Vehicle_Type != -1 AND
Journey_Purpose != 15 AND
Journey_Purpose != -1 AND
Journey_Purpose != 6 ')
```

#### **Results:**



We see from the decision the results that vehicle type= 9 which is car impacts the severity of the accident. If a person, who is a male and is traveling in car with a speed greater than 60 is prone to severe injury.



The important model variables from decision tree determining Accident Severity are Speed limit, vehicle type, urban rural and junction detail.

#### **Neural Networks:**

```
//Neural Networks Model
ALTER MINING STRUCTURE [Accidents DMX]
ADD MINING MODEL [Neural Networks DMX]
( [Accident Index],
[Accident Severity] PREDICT,
[Carriageway Hazards],
[First Road Class],
[Junction Control],
[Junction Detail],
[Light Conditions],
[Ped Cross Physical],
[Road Surface Conditions],
[Road Type],
[Special Conditions At Site],
[Speed Limit],
[Urban Rural],
[Weather Conditions],
[Sex of Casualty],
[Journey Purpose],
[Vehicle Type] )
USING Microsoft_Neural_Network
```

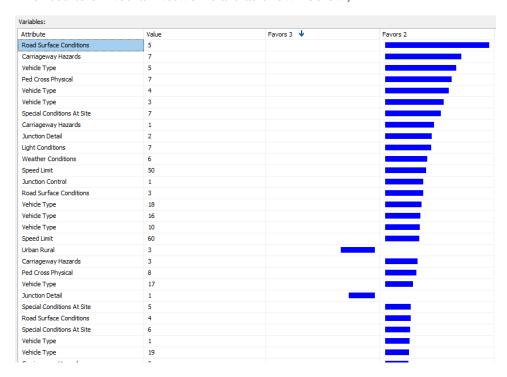
```
INSERT INTO MINING STRUCTURE [Accidents DMX]
( [Accident Index], [Accident Severity],
[Carriageway Hazards], [First Road Class],
[Junction Control], [Junction Detail],
[Light Conditions], [Ped Cross Physical],
[Road Surface Conditions],
[Road Surface Conditions],
[Speed Limit], [Urban Rural],
[Weather Conditions], [Sex of Casualty],
[Journey Purpose], [Vehicle Type])

OPENQUERY([UK Accidents Database],

'SELECT TOP 100000 DimAccidents.Accident_Index,
Accident_Severity,
Carriageway.Hazards,
First_Road_Class,
Junction_Control,
Junction_Detail,
Light_Conditions,
Ped_Cross_Physical,
Road_Surface_Conditions,
Road_Type,
Special_Conditions, At_Site,
Speed_Limit, Urban_Rural,
Weather_Conditions, Sex_of_Casualty,
Journey_Purpose, Vehicle Type FROM DimAccidents,
DimCasualties, DimWehicles
where DimAccidents.Accident_Index = DimCasualties.Accident_Index AND
DimAccidents.Accident_Index = DimVehicles.Accident_Index AND
Accident_Severity != 1 AND
Carriageway_Hazards != -1 AND
Junction_Control != -1 AND
Junction_Control != -1 AND
Junction_Control != -1 AND
Road_Type != -1 AND
Road_Typ
```

#### **Results:**

The results of Neural Networks are as shown below,



We can see that Road surface conditions being , followed by vehicle load on the carriage way and vehicle type stand as most important factors in determining the severity of the accident being serious.

#### **Logistic Regression:**

```
//Logistic Regression Model
ALTER MINING STRUCTURE [Accidents DMX]
ADD MINING MODEL [Logistic Regression DMX]
( [Accident Index],
[Accident Severity] PREDICT,
[Carriageway Hazards],
[First Road Class],
[Junction Control],
[Junction Detail],
[Light Conditions],
[Ped Cross Physical],
[Road Surface Conditions],
[Road Type],
[Special Conditions At Site],
[Speed Limit],
[Urban Rural],
[Weather Conditions],
[Sex of Casualty],
[Journey Purpose],
[Vehicle Type] )
USING Microsoft_Logistic_Regression
```

```
INSERT INTO MINING STRUCTURE [Accidents DMX]
( [Accident Index], [Accident Severity],
[Carriageway Hazards], [First Road Class],
[Junction Control], [Junction Detail],
[Light Conditions], [Ped Cross Physical],
[Road Surface Conditions],
[Road Type], [Special Conditions At Site],
[Speed Limit], [Urban Rural],
[Weather Conditions], [Sex of Casualty],
[Journey Purpose], [Vehicle Type] )
OPENQUERY([UK Accidents Database],
'SELECT TOP 100000 DimAccidents.Accident_Index,
Accident_Severity,
Carriageway_Hazards,
First_Road_Class,
Junction_Control,
Junction_Detail,
Light Conditions,
Ped_Cross_Physical,
Road_Surface_Conditions,
Road_Type,
Special_Conditions_At_Site,
Speed_Limit, Urban_Rural,
Weather_Conditions, Sex_of_Casualty,
Journey_Purpose, Vehicle_Type FROM DimAccidents,
DimCasualties, DimVehicles
where DimAccidents.Accident_Index = DimCasualties.Accident_Index AND
DimAccidents.Accident_Index = DimVehicles.Accident_Index AND
Accident_Severity != 1 AND
Carriageway_Hazards != -1 AND
First_Road_Class != 6 AND
Junction_Control != -1 AND
Junction_Detail != -1 AND
Light Conditions != -1 AND
Ped_Cross_Physical != -1 AND
Road_Surface_Conditions != -1 AND
Road_Type != -1 AND
Road_Type != 9 AND
Special Conditions At Site != 9 AND
Weather_Conditions != -1 AND
Sex_of_Casualty != -1 AND
Vehicle_Type != -1 AND
Journey_Purpose != 15 AND
Journey_Purpose != -1 AND
Journey_Purpose != 6 ')
```

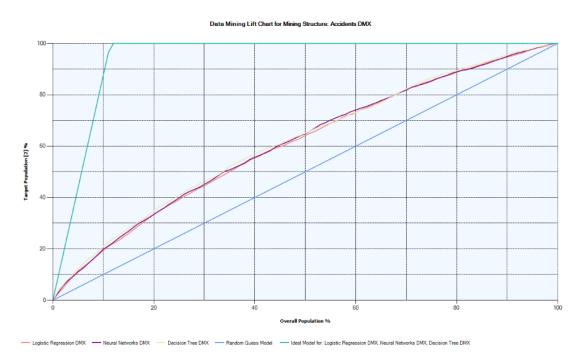
#### **Results:**

ttribute	Value	Favors 3 🔱	Favors 2
ehide Type	16		
Veather Conditions	6		1
ehide Type	18		•
load Surface Conditions	5		1
pecial Conditions At Site	2		1
load Surface Conditions	3		1
ed Cross Physical	7		1
Carriageway Hazards	7		1
ehide Type	5		1
Veather Conditions	3		1
ehide Type	4		1
Carriageway Hazards	6		1
ight Conditions	7		1
ehide Type	3		1
arriageway Hazards	1		I
ehicle Type	10		1
unction Control	1		1
unction Detail	2		1
pecial Conditions At Site	4		1
pecial Conditions At Site	3		1
ight Conditions	5		1
Irban Rural	3		1
pecial Conditions At Site	5		1
Veather Conditions	9		1
peed Limit	50		1

We see from the screenshot above that vehicle type and surface conditions play an important role in accident severity to be fatal . Weather conditions, road surface conditions play important role in determining the accident to be serious.

## **Comparison of Results:**

#### **Lift Chart:**



Series, Model	Score	Target population	Predict probability
Logistic Regression DMX	0.64	44.55%	13.14%
Neural Networks DMX	0.65	45.10%	12.33%
Decision Tree DMX	0.65	46.21%	12.96%
Random Guess Model		30.00%	
Ideal Model for: Logistic Regression D		100.00%	

The lift score of all the models are nearly the same. We chose the population percentage of 30%. This can be interpreted as follows:

The logistic model will correctly identify 44.55% of the people who are seriously injured in the entire population. To identify the severe injury among the people, we would use query to retrieve cases with a predict probability of 13.14%

The Neural Network will correctly identify 45.10% of the people who are seriously injured in the entire population. To identify the severe injury among the people, we would use query to retrieve cases with a predict probability of 12.33%

The Decision Tree will correctly identify 46.21% of the people who are seriously injured in the entire population. To identify the severe injury among the people, we would use query to retrieve cases with a predict probability of 12.96%

From the score, the performance of Neural Network and Decision Tree are better but from the prediction probability it can be concluded that Decision Tree performs the best.

#### **Classification Matrix:**

Counts for Logistic Regression DMX on Accident Severity			
	Predicted	3 (Actual)	2 (Actual)
	3	26534	3405
	2	47	14
Counts for Neural Networks DMX on Accident Severity			
	Predicted	3 (Actual)	2 (Actual)
	3	26579	3419
	2	2	0
Counts for Decision Tree DMX on Accident Severity			•
	Predicted	3 (Actual)	2 (Actual)
	3	26575	3373
	2	6	46

S.No	Model Name	Correct Classification Accuracy
1	Decision Tree	0.887
2	Logistic Regression	0.884
3	Neural Networks	0.885

Classification model proves that Decision Tree has better accuracy.

# **Cross Validation:**

A 10 fold cross validation on injury severity=2 with a threshold of 0.3 was performed.

Logistic Regres	sion DMX			
Partition	<b>Partition Size</b>	Test	Measure	Value
Index				
1	6999	Classification	True Positive	26
2	7000	Classification	True Positive	37
3	7001	Classification	True Positive	35
4	7000	Classification	True Positive	46
5	7000	Classification	True Positive	42
6	7000	Classification	True Positive	47
7	7000	Classification	True Positive	35
8	7000	Classification	True Positive	50
9	7000	Classification	True Positive	41
10	7000	Classification	True Positive	39
			Average	39.8001
			Standard Deviation	6.6751
1	6999	Classification	False Positive	82
2	7000	Classification	False Positive	99
3	7001	Classification	False Positive	61
4	7000	Classification	False Positive	94
5	7000	Classification	False Positive	77
6	7000	Classification	False Positive	101
7	7000	Classification	False Positive	86
8	7000	Classification	False Positive	104
9	7000	Classification	False Positive	98
10	7000	Classification	False Positive	115
			Average	91.6997
			Standard Deviation	14.7113
1	6999	Classification	True Negative	6112
2	7000	Classification	True Negative	6096
3	7001	Classification	True Negative	6134
4	7000	Classification	True Negative	6101
5	7000	Classification	True Negative	6118
6	7000	Classification	True Negative	6094
7	7000	Classification	True Negative	6109
8	7000	Classification	True Negative	6091
9	7000	Classification	True Negative	6097
10	7000	Classification	True Negative	6080

			Average	6103.2003
			Standard Deviation	14.6483
1	6999	Classification	False Negative	779
2	7000	Classification	False Negative	768
3	7001	Classification	False Negative	771
4	7000	Classification	False Negative	759
5	7000	Classification	False Negative	763
6	7000	Classification	False Negative	758
7	7000	Classification	False Negative	770
8	7000	Classification	False Negative	755
9	7000	Classification	False Negative	764
10	7000	Classification	False Negative	766
			Average	765.2999
			Standard Deviation	6.7534
1	6999	Likelihood	Log Score	-0.3476
2	7000	Likelihood	Log Score	-0.3485
3	7001	Likelihood	Log Score	-0.3469
4	7000	Likelihood	Log Score	-0.3453
5	7000	Likelihood	Log Score	-0.3472
6	7000	Likelihood	Log Score	-0.3508
7	7000	Likelihood	Log Score	-0.349
8	7000	Likelihood	Log Score	-0.3461
9	7000	Likelihood	Log Score	-0.3499
10	7000	Likelihood	Log Score	-0.3485
			Average	-0.348
			Standard Deviation	0.0016
1	6999	Likelihood	Lift	0.0093
2	7000	Likelihood	Lift	0.0083
3	7001	Likelihood	Lift	0.0102
4	7000	Likelihood	Lift	0.0116
5	7000	Likelihood	Lift	0.0097
6	7000	Likelihood	Lift	0.006
7	7000	Likelihood	Lift	0.0079
8	7000	Likelihood	Lift	0.0107
9	7000	Likelihood	Lift	0.0069
10	7000	Likelihood	Lift	0.0083
			Average	0.0089
			Standard Deviation	0.0016

1	6999	Likelihood	Root Mean Square Error	0.1418
2	7000	Likelihood	Root Mean Square Error	0.1448
3	7001	Likelihood	Root Mean Square Error	0.1423
4	7000	Likelihood	Root Mean Square Error	0.148
5	7000	Likelihood	Root Mean Square Error	0.1419
6	7000	Likelihood	Root Mean Square Error	0.1462
7	7000	Likelihood	Root Mean Square Error	0.1414
8	7000	Likelihood	Root Mean Square Error	0.1461
9	7000	Likelihood	Root Mean Square Error	0.1471
10	7000	Likelihood	Root Mean Square Error	0.149
			Average	0.1449
			Standard Deviation	0.0027
<b>Neural Network</b>	s DMX			
Partition Index	Partition Size	Test	Measure	Value
1	6999	Classification	True Positive	22
2	7000	Classification		
3		N 461551111C6111C11	True Positive	25
-	7001	+	True Positive True Positive	25 15
<b>  4</b>	7001 7000	Classification	True Positive	15
4 5	7000	Classification Classification	True Positive True Positive	15 42
5		Classification	True Positive	15
	7000 7000	Classification Classification Classification	True Positive True Positive True Positive	15 42 39
5 6 7	7000 7000 7000	Classification Classification Classification Classification	True Positive True Positive True Positive True Positive	15 42 39 34
5 6	7000 7000 7000 7000	Classification Classification Classification Classification Classification	True Positive True Positive True Positive True Positive True Positive	15 42 39 34 20
5 6 7	7000 7000 7000 7000 7000	Classification Classification Classification Classification Classification Classification	True Positive	15 42 39 34 20 34
5 6 7 8 9	7000 7000 7000 7000 7000 7000 7000	Classification Classification Classification Classification Classification Classification Classification Classification	True Positive	15 42 39 34 20 34 33 33 29.6999 8.2953
5 6 7 8 9 10	7000 7000 7000 7000 7000 7000	Classification Classification Classification Classification Classification Classification Classification Classification	True Positive Average Standard	15 42 39 34 20 34 33 33 29.6999 8.2953
5 6 7 8 9 10	7000 7000 7000 7000 7000 7000 7000	Classification	True Positive Average Standard Deviation	15 42 39 34 20 34 33 33 29.6999 8.2953
5 6 7 8 9 10	7000 7000 7000 7000 7000 7000 7000	Classification	True Positive Average Standard Deviation False Positive	15 42 39 34 20 34 33 33 29.6999 8.2953

5	7000	Classification	False Positive	64
6	7000	Classification	False Positive	68
7	7000	Classification	False Positive	47
8	7000	Classification	False Positive	51
9	7000	Classification	False Positive	75
10	7000	Classification	False Positive	97
			Average	67.9995
			Standard	10 1020
			Deviation	18.1938
1	6999	Classification	True Negative	6123
2	7000	Classification	True Negative	6121
3	7001	Classification	True Negative	6157
4	7000	Classification	True Negative	6100
5	7000	Classification	True Negative	6131
6	7000	Classification	True Negative	6127
7	7000	Classification	True Negative	6148
8	7000	Classification	True Negative	6144
9	7000	Classification	True Negative	6120
10	7000	Classification	True Negative	6098
	<u>.</u>	·	Average	6126.9005
			Standard	10 2127
			Deviation	18.2127
1	6999	Classification	False Negative	783
2	7000	Classification	False Negative	780
3	7001	Classification	False Negative	791
4	7000	Classification	False Negative	763
5	7000	Classification	False Negative	766
6	7000	Classification	False Negative	771
7	7000	Classification	False Negative	785
8	7000	Classification	False Negative	771
9	7000		False Negative	772
10	7000	Classification	i aise negative	, , <u>~</u>
1	7000	Classification	False Negative	772
			False Negative	772 775.4001
			False Negative Average	772
1			False Negative Average Standard	772 775.4001
1 2	7000	Classification	False Negative Average Standard Deviation	772 775.4001 8.476
1	7000 6999	Classification	False Negative Average Standard Deviation Log Score	772 775.4001 8.476 -0.3493
1 2	7000 6999 7000	Classification  Likelihood  Likelihood	False Negative Average Standard Deviation Log Score Log Score	772 775.4001 8.476 -0.3493 -0.3507
1 2 3	7000 6999 7000 7001	Classification  Likelihood Likelihood Likelihood	False Negative Average Standard Deviation Log Score Log Score Log Score	772 775.4001 8.476 -0.3493 -0.3507 -0.3488
1 2 3 4	7000 6999 7000 7001 7000	Likelihood Likelihood Likelihood Likelihood	False Negative Average Standard Deviation Log Score Log Score Log Score Log Score	772 775.4001 8.476 -0.3493 -0.3507 -0.3488 -0.3472

11				
8 9	7000	Likelihood	Log Score	-0.3461
	7000	Likelihood	Log Score	-0.35
10	7000	Likelihood	Log Score	-0.3503
			Average	-0.3492
			Standard	0.0015
			Deviation	0.0013
1	6999	Likelihood	Lift	0.0076
2	7000	Likelihood	Lift	0.0062
3	7001	Likelihood	Lift	0.0083
4	7000	Likelihood	Lift	0.0096
5	7000	Likelihood	Lift	0.0076
6	7000	Likelihood	Lift	0.0051
7	7000	Likelihood	Lift	0.0079
8	7000	Likelihood	Lift	0.0107
9	7000	Likelihood	Lift	0.0069
10	7000	Likelihood	Lift	0.0066
			Average	0.0076
			Standard	0.0016
			Deviation	0.0016
1	6999	Likelihood	Root Mean	0.1371
T	0999	Likeiiiiood	Square Error	0.13/1
2	7000	Likelihood	Root Mean	0.1374
2	7000	Likelinood	Square Error	0.13/4
3	7001	Likelihood	Root Mean	0.1277
5			Square Error	0.12//
4	7000	Likelihood	Root Mean	0.1464
T	7000	LIKEIIIIOOU	Square Error	0.1101
5	7000	Likelihood	Root Mean	0.1373
5	7000	LIKCIII 1000	Square Error	0.13/3
6	7000	Likelihood	Root Mean	0.1368
O	7000	LIKCIII 1000	Square Error	0.1500
7	7000	Likelihood	Root Mean	0.1308
,	7000	LIKCIII 1000	Square Error	0.1500
8	7000	Likelihood	Root Mean	0.135
	, 550	Lincilliood	Square Error	0.133
9	7000	Likelihood	Root Mean	0.1403
	, 555	Littelli 100d	Square Error	311 103
10	7000	Likelihood	Root Mean	0.1428
	, 555	Enteniiood	Square Error	
			Average	0.1372
			Standard	0.0051
			Deviation	0.0002
<b>Decision Tre</b>	ee DMX			

Partition Index	Partition Size	Test	Measure	Value
1	6999	Classification	True Positive	12
2	7000	Classification	True Positive	8
3	7001	Classification	True Positive	5
4	7000	Classification	True Positive	8
5	7000	Classification	True Positive	54
6	7000	Classification	True Positive	7
7	7000	Classification	True Positive	10
8	7000	Classification	True Positive	9
9	7000	Classification	True Positive	49
10	7000	Classification	True Positive	43
			Average	20.4999
			Standard Deviation	18.6829
1	6999	Classification	False Positive	1
2	7000	Classification	False Positive	1
3	7001	Classification	False Positive	0.000e+000
4	7000	Classification	False Positive	1
5	7000	Classification	False Positive	117
6	7000	Classification	False Positive	1
7	7000	Classification	False Positive	1
8	7000	Classification	False Positive	1
9	7000	Classification	False Positive	117
10	7000	Classification	False Positive	116
			Average Standard	35.6
	leans.	- ICI	Deviation	53.072
1	6999	Classification	True Negative	6193
2	7000	Classification	True Negative	6194
3	7001	Classification	True Negative	6195
4	7000	Classification	True Negative	6194
5	7000	Classification	True Negative	6078
6	7000	Classification	True Negative	6194
7	7000	Classification	True Negative	6194
8	7000	Classification	True Negative	6194
9	7000	Classification	True Negative	6078
10	7000	Classification	True Negative	6079
			Average	6159.3
			Standard Deviation	53.0077
1	6999	Classification	False Negative	793

2	7000	Classification	False Negative	797
3	7001	Classification	False Negative	801
4	7000	Classification	False Negative	797
5	7000	Classification	False Negative	751
6	7000	Classification	False Negative	798
7	7000	Classification	False Negative	795
8	7000	Classification	False Negative	796
9	7000	Classification	False Negative	756
10	7000	Classification	False Negative	762
			Average Standard Deviation	784.6001 18.7681
1	6999	Likelihood	Log Score	-0.3426
2	7000	Likelihood	Log Score	-0.3425
3	7001	Likelihood	Log Score	-0.3458
4	7000	Likelihood	Log Score	-0.3423
5	7000	Likelihood	Log Score	-0.3451
6	7000	Likelihood	Log Score	-0.3442
7	7000	Likelihood	Log Score	-0.3411
8	7000	Likelihood	Log Score	-0.3408
9	7000	Likelihood	Log Score	-0.3472
10	7000	Likelihood	Log Score	-0.3438
			Average Standard Deviation	-0.3435 0.002
1	6999	Likelihood	Lift	0.0143
2	7000	Likelihood	Lift	0.0144
3	7001	Likelihood	Lift	0.0113
4	7000	Likelihood	Lift	0.0145
5	7000	Likelihood	Lift	0.0117
6	7000	Likelihood	Lift	0.0127
7	7000	Likelihood	Lift	0.0158
8	7000	Likelihood	Lift	0.016
9	7000	Likelihood	Lift	0.0096
10	7000	Likelihood	Lift	0.013
			Average Standard Deviation	0.0133 0.0019
1	6999	Likelihood	Root Mean Square Error	0.1245
2	7000	Likelihood	Root Mean Square Error	0.1255

3	7001	Likelihood	Root Mean Square Error	0.1247
4	7000	Likelihood	Root Mean Square Error	0.1244
5	7000	Likelihood	Root Mean Square Error	0.1589
6	7000	Likelihood	Root Mean Square Error	0.1241
7	7000	Likelihood	Root Mean Square Error	0.1247
8	7000	Likelihood	Root Mean Square Error	0.124
9	7000	Likelihood	Root Mean Square Error	0.1563
10	7000	Likelihood	Root Mean Square Error	0.1555
			Average	0.1343
			Standard Deviation	0.0149

Comparing Root Mean Square error, Lift and log score for all the three models,

S.No	Model Name	RMSE	Lift Score	Log Score
1	Decision Tree	0.1343	0.0133	-0.3435
2	Logistic Regression	0.1449	0.0089	-0.348
3	Neural Networks	0.1372	0.0076	-0.3492

We see from the comparison that Decision Tree performed better in terms of RMSE, Lift and Log score.

## **Summary of Results:**

Among the three models Decision tree performed better in terms of predicting the serious injury in accidents.

Hierarchy of variable importance in predicting the accident severity is as follows,

- 1) Vehicle Type
- 2) Junction Details
- 3) Urban rural
- 4) Speed Limit

If a person, who is a male and is traveling in car with a speed greater than 60 is prone to severe injury.

The Decision Tree will correctly identify 46.21% of the people who are seriously injured in the entire population. To identify the severe injury among the people , we would use query to retrieve cases with a predict probability of 12.96%

## **Conclusion:**

The decision tree made sensible prediction in determining the serious injury factors. People shouldn't over speed in a car to avoid accident.