

Data cleaning in Excel file

Each attribute in each dataset have been thoroughly studied and determined if it is needed or not in related with the objective of this project. Attributes which are used for administrative purpose like vehicle plate, driver license, Municipality code, Area code, Rout No and other irrelevant variables are removed from the dataset that reduced the attributes by significant number. Then, formatting for some attributes have been changed, for example, birth year of the driver “03-10_1970 ”was changed into 1970 and name of some attributes are changed as well to understand easily, for example, RD_DIV to ROAD TYPE, CV BODY to VEHICLE BODY. More than 90% of the attributes data entries are coded (01, 02,88, 99) and understanding each code helps for the data cleaning purpose, for example, 88 = “OTHER” and 99 = “UNKNOWN”. Before the data set was imported to SAS EM, all UNKNOWNs were changed to NULL to make it easy for later cleaning in SAS EM. Some codes are written in decimal, for instance, in vehicle information Body Type code 27.88 = “Snowmobile” changed into 27. Categorical Target Variable (Report Type) was changed into numeric factor by assigning “1” for occurrence of fatal/injury accident and “0” for car damage in planning of using logistic regression as a binary outcome (Appendix 1)

In most of reported accidents there were more than one person involved in the car accident, for example, children or passengers who got killed or injured and registered on personal information dataset with the same Report No. Here, the goal of the project is to identify the cause of the car crash due to the condition of the driver, and all passengers are excluded from the analysis. This was

identified by “Person Type” attribute (Driver, Occupant, Pedestrian). In this project exclusion of Fatal/injury accident on Occupant and Pedestrian can be considered as a limitation.

Creating of a Variable

Data set combination Techniques and Preparing SAS TABLE

The three data sets were imported from Excel file to MS-SQL Database and combined by using INNER JOINT SCRIPT with a common variable of Report No (Appendix 4). After merging the datasets, 26,660 rows and 32 variables were exported to excel file and then imported again into SAS EM library to create SAS Table. The property of SAS Table was set up by assigning the Role of each attribute (Target, Input, ID) and choosing the correct level (Binary, Nominal, Interval). For example, some categorical attributes like County Codes, Accident Time were labeled as interval Variable and edited to Nominal. Out of 32 variables, four of them are interval variables (Speed limit, Reference _Distance, Longitude, and Latitude), Report Type (Target variable) is a Binary variable and all the rest are Nominal (Appendix 5). Advanced Option was used to set up missing percentage Threshold (Default was used where variables with more than 50% of missing value were rejected) and 100 class levels count Threshold was used to avoid rejection of some variables with many classification, for example, car make, county Code, Vehicle body made, accident time, Date of Birth). Out of the total Reported accident 0.27% of them are fatal/injury with Prior probability of 0.2644 for fatal/injury and

0.7336 for car damage. In order to increase the chance of catching rare events adjusted equal prior probability (0.5) was applied and a decision weight of 0.99 for fatal/injury and 0.01 of car damage was used in a sense that fatal/injury event is 99 times more important than car damage. Finally, 25% of the observation are considered to create the Data Source (SAS Table).

Initial Exploration was used to visualize the input data helps to observe possible pattern and check missing value before building a model. The result of the exploration showed the percentage of missing value, Min, Max, Mean, SD, Skewness, and Kurtosis values. The Min and Max values helps to see if there is any outliers, for example speed limit with a min value of 0 and max value of 75 and Reference distance with min and max value of 0 and 75 are found in acceptable range Four variables had a missing value of more than 50% and rejected from being used in the analysis. These variables are Drug Test, Movement code, occ-seat with a missing value of 99%, 96% and 96 %, respectively. The Target variable, Accident Time, Airbag Deployment, Longitude, Latitude, Speed limit have 0 missing value and the rest have missing value with the range of 1.89% -12.62%.

The Stat Explore Node also showed A chi-square plot orders the top 20 variables by their chi-square statistics and a variable worth plot orders the input variables by their importance in predicting target variable (Appendix 6). The stat Explore also provided the summary statistics for class and interval variables. Among the four interval variables, Distance has the highest SD which is more than 100 (146) with high variability (Appendix 7). Variable Clustering Node was run to see the correlation between interval variables and found out that there is high correlation between them, and all are kept in the analysis (Appendix 8).

Replacing and Imputation

Before replacing unknown value and outliers by missing value, data partition (Training dataset of 60%, Validation data set of 30% and Test dataset of 10%) was done. Interval Variables don't have any missing value. With Replacement Node, all class variable and interval variables with unknown values are replaced with missing value (Appendix 8). Then Impute Node has been run to replace all missed class value with the most occurrence event (Mode). For example, Airbag has 400 missing value and imputed with a value of code value of 1 because 63.69% of Airbag deployment was responded on the questioner as code 1. The second class variable with the highest missing value is Alcohol Test code and imputed with code 0 as 92.75% of response were code 0. In general, for bimodal variables, the highest mode was taken as an imputed value for missing value. (Appendix 8).

Creating Additional Variables/ Data Transformation

Rate of Accident

A variable name “Rate of Accident” has been created from Population size in each state in Maryland and total number of accident in each state. Rate of Accident data was calculated by dividing the number of fatal/injury accident by the total number of population in each state. The same thing was done for car damage in each state and divided by the total number of accident to get the rate. This was done by assuming in all state starting age of driving, socioeconomic status and demographic behaviors of the people are the same. However, taking assumption and creating this variable is much more important than counting the number of accident in each state and deliver this statistics to safety planner, insurance company and health organizations. For example, in Baltimore the population is much greater than any other county of the state and getting more car crash than any small populated county. This variable especially very helpful for visualization.

Variable Transformation

Variable transformation for skewed distribution is needed to fulfill the assumption of Normality to build Logistic Regression model. After the missing values are imputed, variable Transformation was done by selecting Transformation Node from SAS EM. The distribution of all variables are shown on (Appendix 5) and the type of transformation is selected according to the distribution of the variable. Latitude, Longitude, Speed limit, Date of Birth have close normal distribution. Some of the variables have pick at the left

or right and flat in the middle which means most of the values are concentrated at the end. For example, person condition has the highest pick at the right (code 1), movement of the person, Airbag deployment, Alcohol test are positively skewed. All the rest of the variables have also negative or positive skewness which violates Normality and needs transformation. In order to approach Normality the log transformation can be used to make highly skewed distributions less skewed. This can be valuable both for making patterns in the data more interpretable and for helping to meet the assumptions of inferential statistics.

The best thing about SAS EM, whenever transformation or imputation are conducted, it keeps the original variable and create a new feature. However, the status of the original variables are changed into rejection. After transformation, the new feature will have less standard deviation or variability and pretty much enough to be representative.

Data Analysis

Tableau is going to be used for visualization which is one of the most important part of data analysis to present the data in an understandable and visually appealing format. In this project, various types of visualization will be provided for end users. For example, fatal /injury data by county, Airbag deployment, age, name of road, in addition, the most interesting part of tableau is providing geospatial information from longitude and latitude variables. The below screenshot showed the exact place where accident was occurred, and the red dots indicates the place where fatal/injury accident occurred. In the map, on Stone Rd in 2016 there were more fatal/ injury than other places. This helps to Maryland safety Highway office to plan and implement on the right

time to reduce accident. More than two variables can be included to get more detailed information about the accident like the time, the driver condition, the type of road and so on. In general, Tableau does a lot of visualization and presenting reality with a simple bar graph, pie chart, line graph, maps, so any stakeholders, business organizations, individual can grasp the information very easily.

SAS EM is the right tool for statistics analysis and in this project, it will be used for model development. SAS EM is the best as it does all the data preparation, cleaning, transformation, model development, evaluation and comparison.

Data Visualization

Visualize large amounts of complex data is easier than presenting spreadsheets or reports as it is easy way to convey concepts in a universal manner. In this project various visualization are prepared in related with fatal/injury car crash in Maryland to answer basic questions like in which county is the highest death due to Car crash? What age groups are mainly involved in Car accident? What is the impact of Air bag performance at the time of accident on the rate of fatal/injury accident and other similar questions?

Photos - Fatal_injury Car Crash Report by Airbag Deployment.png

See all photos

Add to a creation



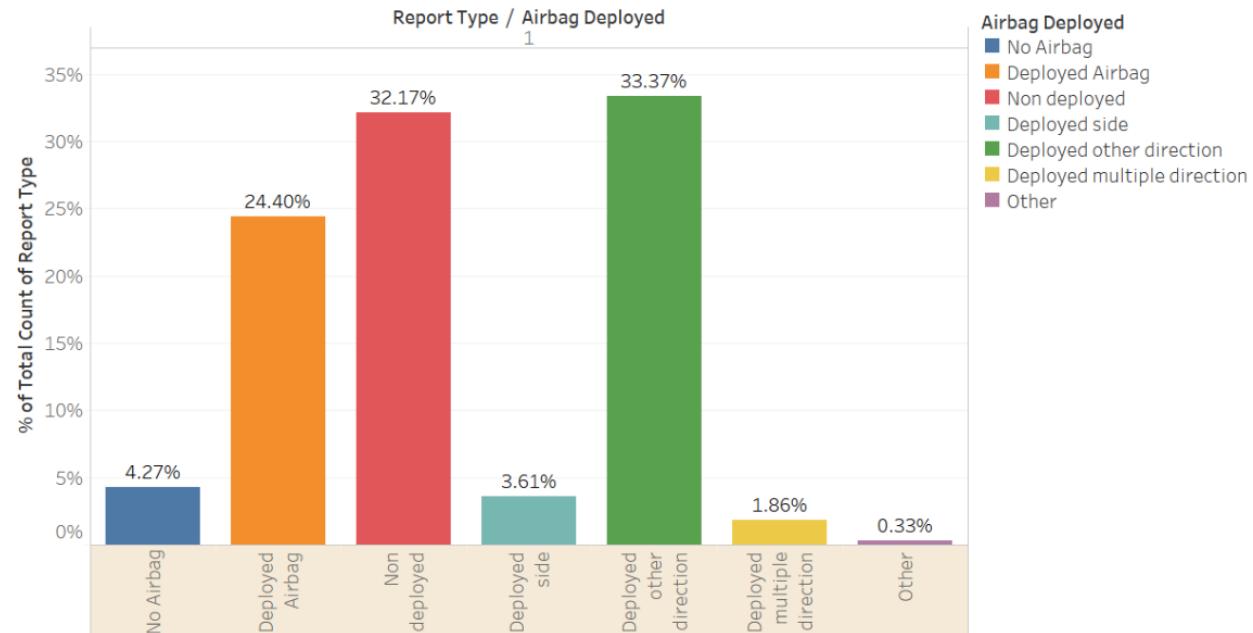
Edit & Create

Share



...

For Destroyed Car Damage the Number of Fatal/injury Car Crash by Airbag deployment



% of Total Count of Report Type for each Airbag Deployed broken down by Report Type. Color shows details about Airbag Deployed. The data is filtered on Damage Code, which keeps Destroyed. The view is filtered on Report Type and Airbag Deployed. The Report Type filter keeps 1. The Airbag Deployed filter excludes Null. Percents are based on each row of the table.



The purpose of the air bag is to provide a cushion between the occupants and the vehicle's interior. For air bags to be effective they must be fully inflated in a short amount of time, before the occupants contact them (Jesse, 2018). The above Visualization 1 showed the Percentage of fatal-injury car crash by Air bag deployment performance at the time when the car is damaged. One of the car accident tragedy happens when the car is destroyed. The question is what was the performance of the Air bag at the time of this tragedy? Was the Air bag deployed on the right time and on the right direction? or was the Air bag failed? Or was the Air bag deployed in different direction? Here, it can be possible to see the effect of Air bag deployment by selecting only one type of Car damage which is "Destroy".

The row and the column of the bar graph shows the percentage of fatal/injury and the performance of Air bag deployment. As it can be seen from the graph, out of the total fatal/injury Car crash in knowing that the Car was "destroyed", 33.37% of the fatal/injury was contributed by poor performance of Air bag with a code of "Deployed in other direction". The second highest fatal/injury (32.1%) was contributed again by another type of poor performance of Air bag with a code of "non deployed at all". In contrast, the lowest fatal/injury car crash was happened with Air bag performance of code "Deployed multiple direction "which is only 1.86% of fatal/injury happened out of the total death of car destroy. The second lowest death was showed when the Air bag was deployed in one direction with a percentage of 3.61%.

However, one of the columns of the Air bag deployment (Air bag deployed) showed unexpected result of fatal/injury car crash with a contribution of 24.4% which was considered the highest as far as the air bag was deployed. The reason might be different in such

a way that the performance is not only measured because It was deployed. However, other questions like does the Air bag deployed on time? Or does the Airbag deploy with a perfect inflation? Moreover, other reasons also affect the performance.

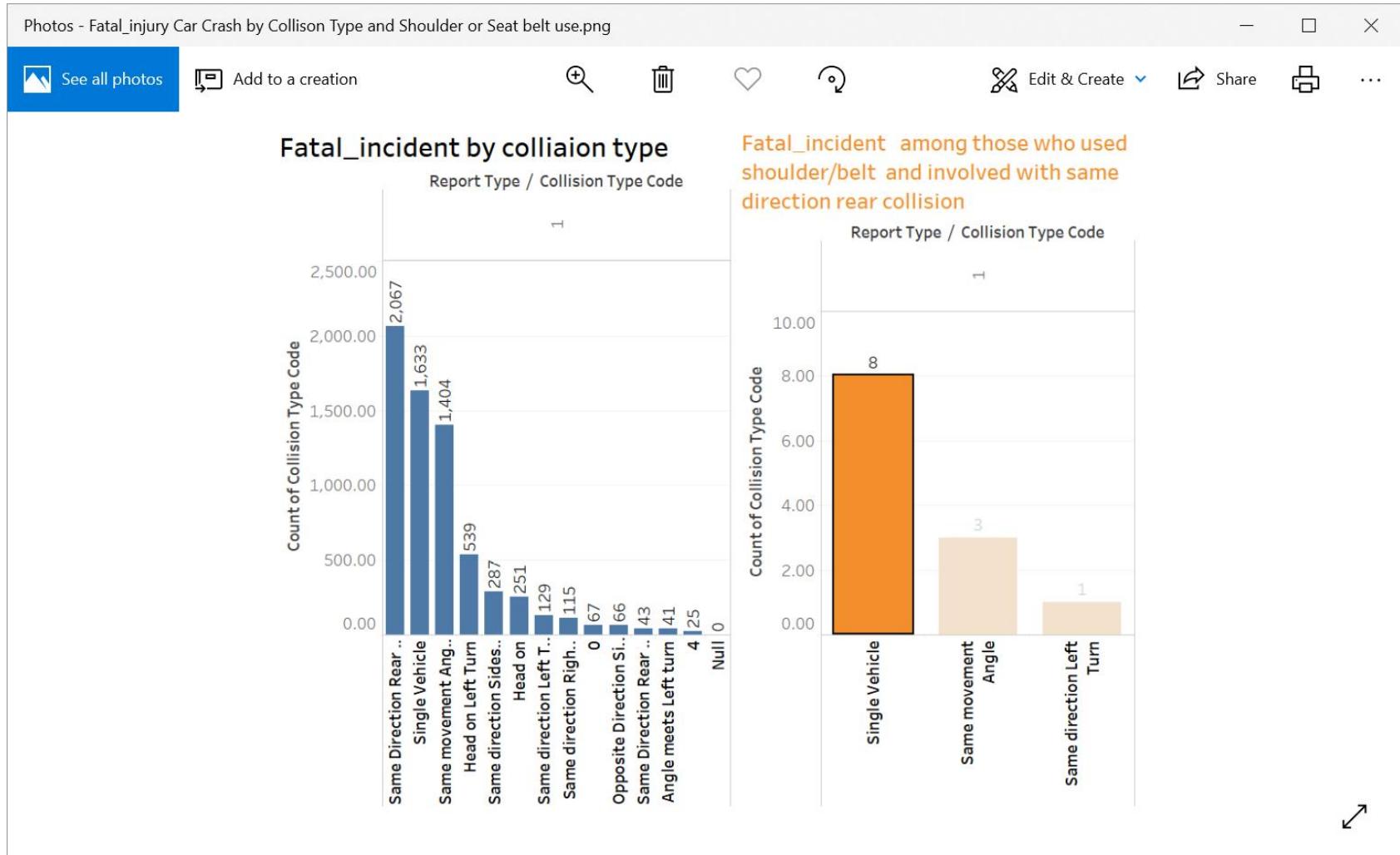
In general, the insights gained from the visualization is that the performance of Air bag deployment plays a significant role in saving the life of people even at the time of big accident tragedy. Good Performance of Air bag is a combination effect of various factors such as deployment time, the size of the inflation, and the direction. Form this result Car Manufactures can get good insight to assess or evaluate the performance of their car by combining the result and their long-time experience in the domain. Car insurance company also get some knowledge about the type of car their customers are driving in calculating the risk of selling their insurance with respect to the type of car with the history of Airbag deployment performance. Least but not last, each and everyone who wants to buy a car search up the performance history of the car.

Earlier, in this project, the performance of Air bag deployment was proposed to be one of the important factors to determine the fatal/injury Car Crash. The result from the visualization also supported the scope of this project which confirms to move forward and develop models by taking Air bag performance as one of the most effective variables to train the model. However, the gap of information of the

One of the lessons that was gained from this result is that detailed information is the base for the quality of analysis. For example, the result from the column of Airbag deployment with a code of "Airbag deployed" resulted a high fatal/injury percentage. This

implies that the question must go further in order to address the problem of the data quality by going deeper and asking questions like the time the Air bag deployed during the accident or the size of the inflation to protect the occupant.

Data Visualization 2: Fatal/injury Car Crash by Collision Type and using Safety equipment



Using safety equipment like Lap or shoulder belt are the most effective way to prevent fatal/injury Car crash. The above visualization displayed two bar graphs. The first one with blue color showed the number of the fatal/injury Car crash with the type of collision. The result revealed that the highest fatal/injury Car Crash was happened when the type of Collision was “same direction with rear end” with 2067 death. The second highest death was happened with a collision of a car with some other object like building, pole or something external other than another car with 1633. The third and fourth highest number of deaths appeared with a collision type of same movement angle and head on left turn with 1405 and 519, respectively. In contrary, the lowest death was observed with same direction sideswipe and angle meet left turn with a total fatal of 43 and 41, respectively.

In order to see the effect of wearing lap/shoulder belt on fatal/ injury Car Crash, the second bar graph was done by taking the number of fatal /injury car crash and experience of wearing safety equipment and filtered by “lap/shoulder” and applied to all type of collusion. Surprising, there was no one who died due to “same direction with rear end Collision” accident because of the safety equipment. The column was reduced into three with the lowest death in all. Over all only 13 people died even though they wear lap/shoulder belt. The total number of deaths by “single car” crash was 1633, among these deaths only 8 people with belt were died. This result proved that even though the accident is serious if a person wear lap/shoulder belt, the chance of being killed by accident is very rare. “Five children were killed, and two adults were injured in a Maryland car crash, officials said the children were not wearing seat belt, but the driver and the passenger were wearing seat belt and were saved” CNN, 2019).

The insight gained from this result is seat belt prevent fatal car crash in a higher chance. The result of the visualization helps to add the attribute “safety equipment” as one of important variable to train the various model (logistics, SVM) to increase the accuracy of predicting fatal/ injury car crash in Maryland.

From the point of business objectives, SHSP get knowledge from the result to plan and implement to approach their goal in reduction of death due to car crash. This information also helps Stakeholders like MVA, and other partners can depend on the result to act on the rule of seat belt. For example, car and life insurance can get insight how wearing seat belt can save life of their customer and indirectly reduce the expense of their business. In this aspect, insurance company donate every year to SHSP to the implementation of the rule such as dispatching more officers on the road for the enforcement of seat belt, “Click it or Ticket”.

Data Visualization 3: Fatal/injury Car Crash by County and by Road Main

See all photos

Add to a creation



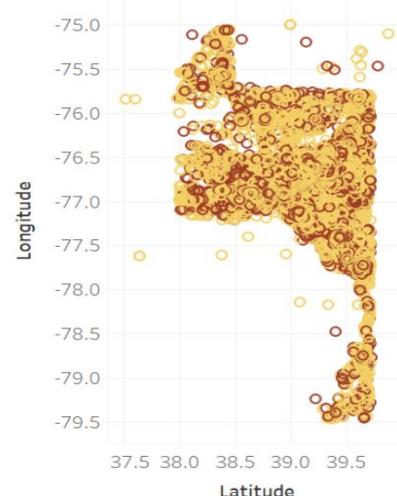
Edit & Create

Share



...

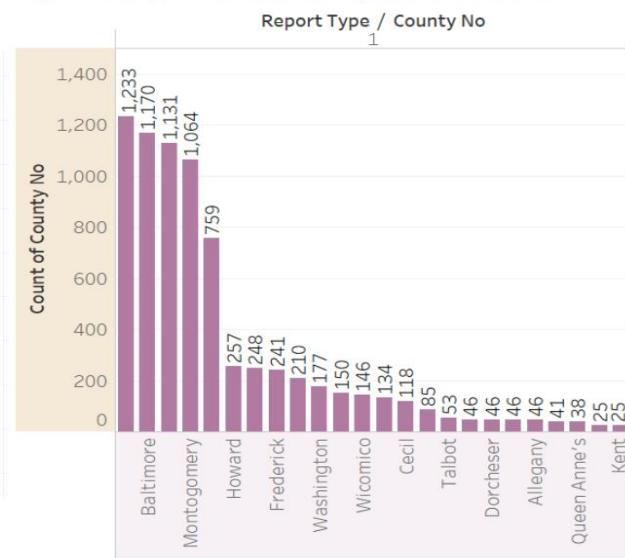
Number of fatal/injury and property Damage due to car crash by geographical location and Marinroad name

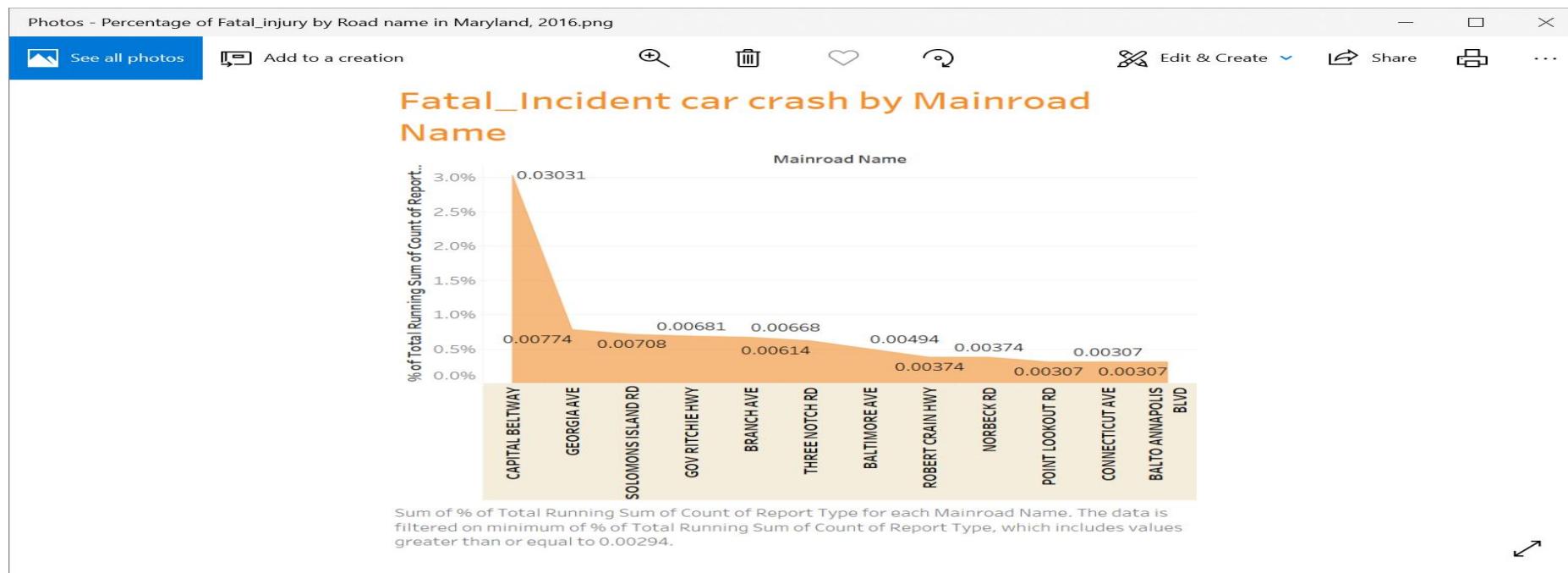


Sum of Report Type



Number of fatatl/injury and car damage by county name in Maryland in 2016





The above dashboard displayed two visualization. The Left side shows geographical location of Maryland Roads with fatal/injury accident. On the map, each red dot indicates fatal accident on that specific road which is more red color on one spot means more accident was occurred on that road. In order to see clearly the name of the road with high frequency of accident, another visualization was added below the dashboard which shows the name of Maryland Roads with highest death. Capital Belt Road contributed 3% of

fatal/injury Car Crash in Maryland in the year of 2016. The next highest death or injury happened on the road of Georgia Ave and Solomon Island Rd, with 0.77% and 0.7% respectively. Many fatal/injury car crashes were also observed on the road of Branch Ave, Three North, Baltimore Ave, Robert Ave.

On the Dash board, the second visualization on the right side revealed the distribution of fatal/injury accident by county. Baltimore was the highest with a death of 4087 followed by Baltimore City with 1233 fatal/injury car crash. The next most county with highest death/injury were Prince George's, Montgomery, Anne Arundel, Howard and Harford with death/injury of 1131, 1964, 759, 257, and 248, respectively.

The insight gained from the above result with respect to the project scope, road name can be one of the important variables in predicting the fatal/injury Car crash. Earlier, this attribute was not proposed as a significant variable and now it will be considered in the variable choice in building various models.

From the business objective point of view, the name of the county and road name provided the place where most fatal/ injury accidents occurred. Organization like SHSP who is responsible for planning and implanting car crash in Maryland can use this insight as a starting point to make further research why more accident happened on these places and finding out the reason helps to take the right decision and action.

Proposed Visualization 1: Driver condition by Report type of fatal/injury

The first proposed visualization is displaying driver condition by fatal/injury accident. Driver condition states the physical and mental condition before the accident. The condition code includes “Apparently Normal”, “Had been drinking”, “using drugs”, “Physical Defects”, “other Handicaps”, “Apparently Asleep”, “Emotional depressed or angry”, “influenced by medication or drug”. Alcohol or drug impaired driver is more exposed for fatal/injury accident because the mind of abnormal person can’t process information on time. For example, drunk person is not able to keep enough distance from the vehicle in front of him or press the break on time to avoid any Collision which can be prevented. The same thing with apparently Asleep driver or fatigued person can easily engaged into accident because of failure to control himself.

The number of fatal/injury accident helps SHSP to make the necessary enforcement and rule to save the life of the driver and others. In addition, if this visualization is sorted by age and sex, it will provide good estimate to focus the rule and regulation on these types of people. Car insurance and life insurance can make risk analysis of taking these group of people as a customer or set the sell of their insurance.

Proposed Visualization 2: Road Type by Fatal/injury

The type of the road in Maryland contribute for fatal/injury accident. Looking into the distribution of Road type by fatal/injury gives a lot of insights. The type of roads in the data description are “Two-way Trafficway”, one-way Trafficway”, “Two-way Divided “, “Two way not divided” and many. The display helps the responsible organizations to reconstruct the road to reduce accident.

Predictive Model on Fatal/injury car crash in Maryland, 2016

Intended predictive model is the one which classify the target variable with a higher accuracy rate of sensitivity and overall accuracy. Classification is one of the major problems that we solve while working on standard business problems across industries (Lalit, 2015). In this project, 12 models will be developed including Ensemble model and compared based on ROC value of each models. Before developing the model, the imbalanced Target variable was handled.

Handling Imbalanced Target variable

The prior probability of the target variable is 0.26 (Fatal/injury accident) and 0.74 (Car damage) which is not deadly imbalanced. However, in order to accomplish the business objective which is reducing death due to car accident towards zero, increasing the proportion of the rare event(fatal/injury) would be the best way by using equal sample size in both groups. This was done by using Sample Node in SAS EM and selecting stratified sampling method with equal sample size from both group of target variable, 5241 for each group.

The dataset allocation was done by partitioning the data into training, validation and test dataset with a proportion of 70%, 30% and 10% of the whole sampled dataset. Training set is used for learning to fit the parameters of the classifier and Validation set use to tune the parameters of a classifier and set to find the “optimal” number of hidden which uses to validate the model. Finally, the model was tested with a test dataset.

Altogether, 13 models have been built and compared by ROC value. The top three models are HP Forest, SVM (Active set with pol 2) and Neural Network. The below screenshot showed each step which was used to build the model. In addition, the tables showed the calculated total accuracy rate, sensitivity, specificity and ROC value for Training, Validation and Test dataset. Detailed models properties will be explaining for the three champion models

Models Development

Enterprise Miner - ACC02272019

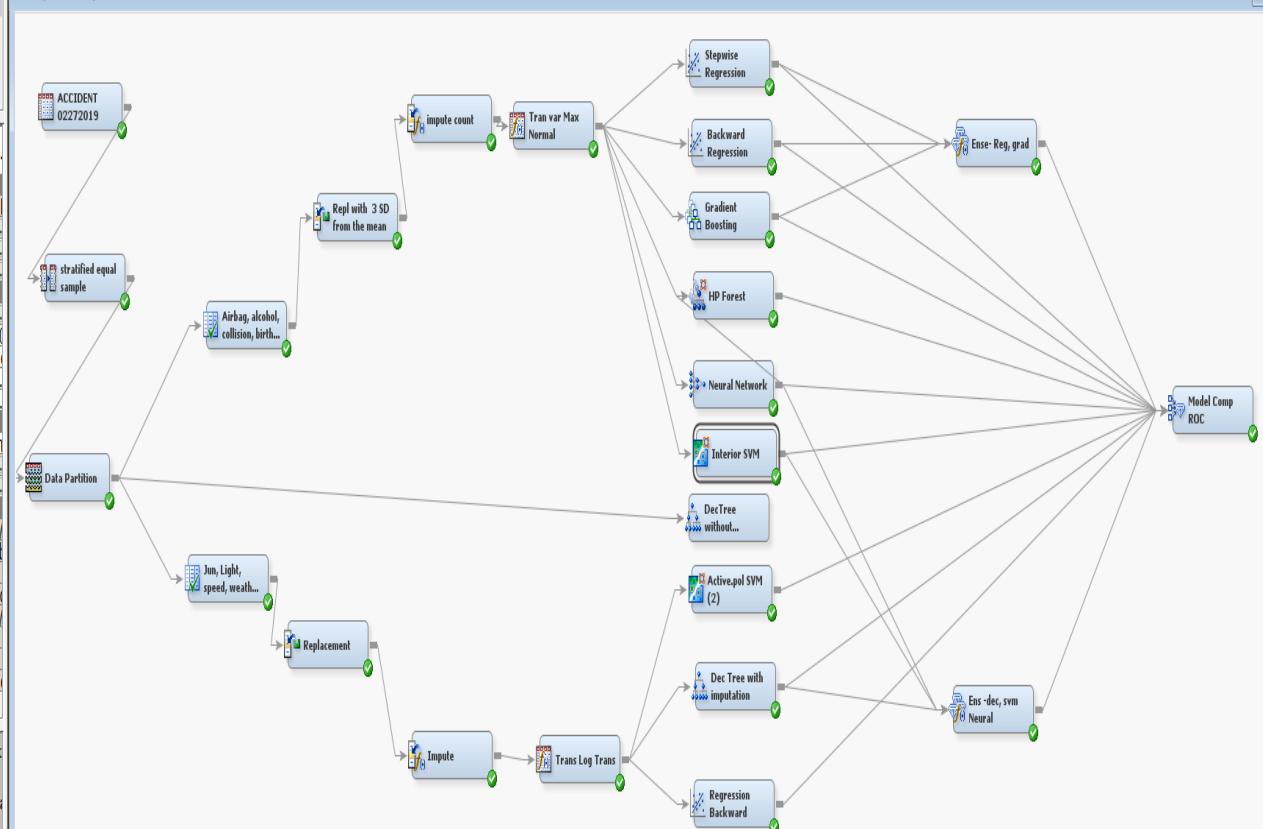
- X

File Edit View Actions Options Window Help



Sample Explore Modify Model Assess Utility Credit Scoring HPDM Applications Text Mining Time Series

Equal Sample size



100%

Diagram Log

Desktop

Calculated value of Sensitivity, Specificity and overall Accuracy Rate for Training Dataset

Calculated value of Sensitivity, Specificity and overall Accuracy Rate for Validation Dataset

ROC 04132019.xlsx - Read-Only - Excel

RUTH GEMECHU (876073)

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

Clipboard Font Alignment Number Styles Cells Editing

AutoSave (Off)

F29

Models	TP	TN	FP	FN	Sensitivity	Specificity	Total Accuracy Rate
Reg. Stepwise	954	1090	407	544	0.636849132	0.728122912	0.682470785
Reg. Backward	982	1172	325	514	0.656417112	0.782899132	0.719679252
Gradiant Boosting	960	958	539	538	0.640854473	0.63994656	0.640400668
HP Forest	925	1232	265	573	0.617489987	0.822979292	0.720200334
Neural Network	980	1188	309	518	0.654205607	0.793587174	0.723873122
SVM Interior	905	1222	275	593	0.604138852	0.816299265	0.710183639
SVM(pol 2, sigmoid)	971	966	531	527	0.648197597	0.645290581	0.646744574
Decision Tree	804	1329	168	694	0.536715621	0.887775551	0.712186978
Ense 1(Step. Reg, Gr)	987	1089	408	511	0.658878505	0.72745491	0.693155259

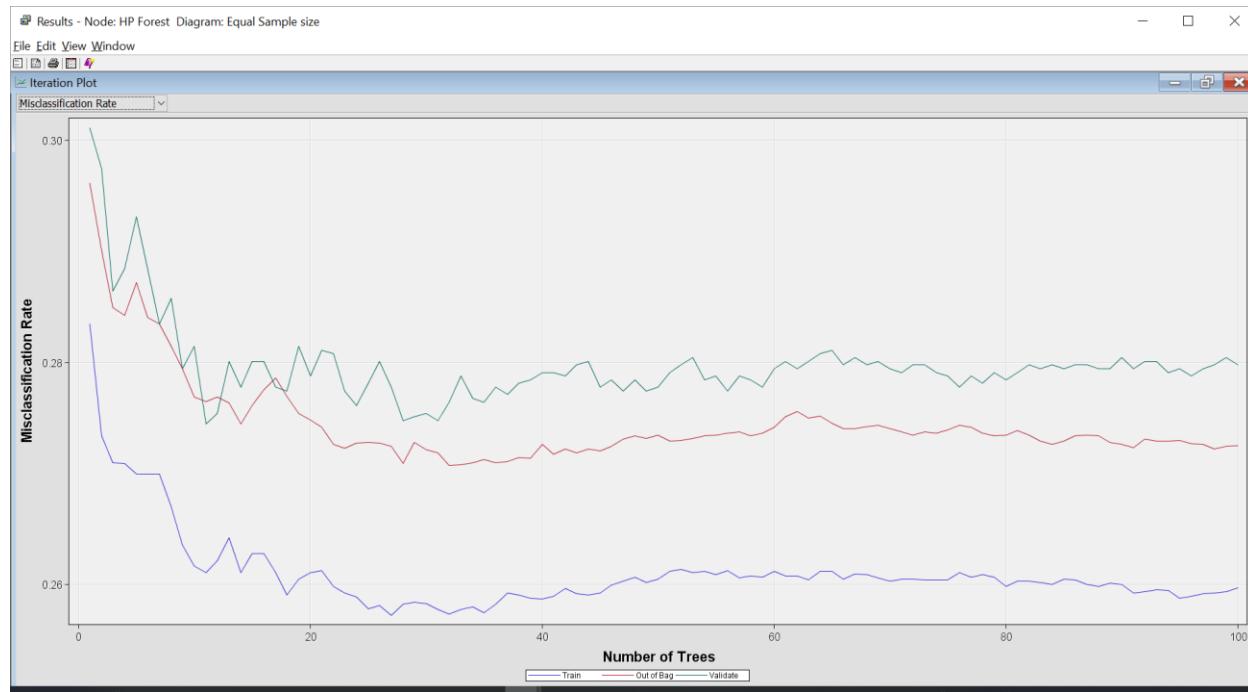
Description of the Three Champion models

First Champion Model: HP Forest

Random forest algorithm is a supervised classification algorithm. The algorithm creates the forest with several trees. The accuracy of the model increases as the number of trees increases. In contrary, average square error decreases as the number of trees increase.

In this project, HP Random was built with 100 trees (Maximum number of trees). The proportion of observations in each tree was 60% of the total observations and the sample was taken with replacement which implies that one observation (one incident report) can have a chance of to be include in more than one tree or in the same tree. The maximum variable used at the start of the tree was 25 and Loss Reduction was used to select important variables with a p- value of 0.05 significant level that is if no association meet this threshold the node will not be split. Based on these properties, the HP Forest model was built. As it can be seen from the below

screenshot, misclassification rate the datasets (Train, Validation, and Test) decreases as the number of trees increased in reverse the average square error decreased as the trees increases.



Before HP model was developed, all unknowns and outliers (3 SD away from the mean) have been replaced by missing value and then all missing class values were imputed by count. Max Transformation was applied for all skewed variable. The classification table for Training and Validation data have been shown below.

Classification Table for Training data set

Outcome	Predicted fatal/injury	Predicted Car Crash	Total
Actual fatal/injury	3351	1890	5241
Actual Car Crash	833	4409	5242
Total	4992	5491	10483

$$\text{Total Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 0.74 \sim 74\%$$

$$\text{Sensitivity} = (\text{TP} / (\text{TP} + \text{FN})) = 0.639 \sim 64\%$$

$$\text{Specificity} = (\text{TN} / (\text{TN} + \text{FN})) = 0.841 \sim 84\%$$

Specificity is predicting fatal/injury accident that showed how much of the events are caught by the model. Out of the total car crash, 3351 or 64% of them are predicted as fatal/injury car crash and 1890 or 36% of them are predicted as car damage. Specificity was predicted 84% which is higher than specificity. According to the model the most important variables are Airbag deployment, Alcohol, Type of collision, collision, condition of the driver, Going direction, junction, light code.

Classification Table for Validation data set

Outcome	Predicted fatal/injury	Predicted Car Crash	Total
Actual fatal/injury	953	545	1498
Actual Car Crash	414	1093	1507
Total	1367	1638	3005

$$\text{Total Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 0.68 \sim 68\%$$

$$\text{Sensitivity} = (\text{TP} / (\text{TP} + \text{FN})) = 0.64 \sim 64\%$$

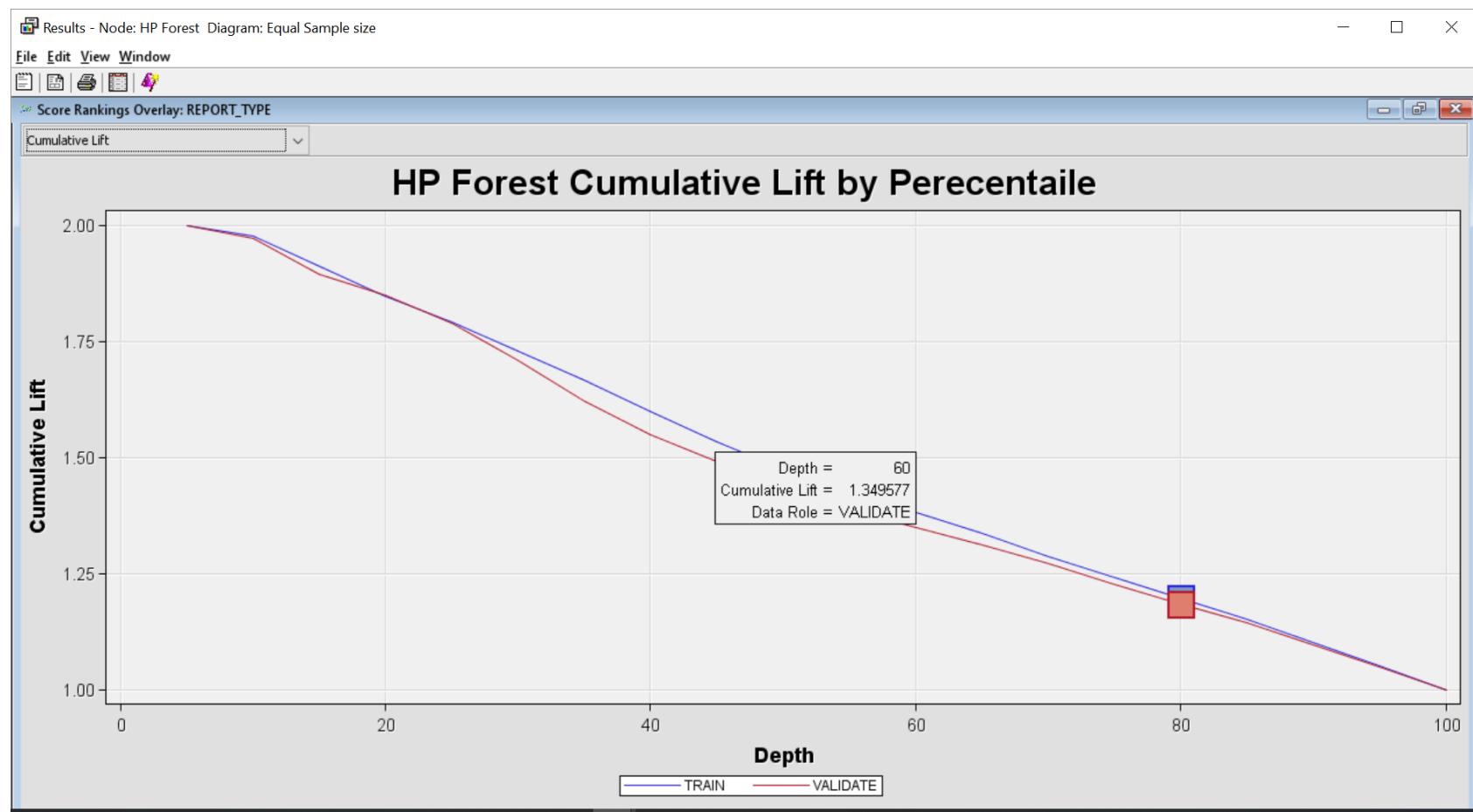
Specificity = $(TN / (TN+FN)) = 0.73 \sim 73\%$

HP Forest model predict 953 fatal/injury car crash correctly with overall accuracy rate of 64% and predict Car Crash correctly with a specificity rate of 0.73 or 73%.

When the accuracy rate of Training dataset was compared with the validation dataset, there is a difference of 0.03 rate

Cumulative Lift for HP Random Forest

Cumulative Lift chart is another evaluating model effectiveness, the x-axis shows the percentile and the y-axis shows lift. The horizontal line intersecting the y-axis at 1 is the chance of identifying the fatal/injury at random is 10% out of the total observation without using model. For example, the Depth 60 and cumulative Lift 1.349, showed that model can capture 60 percentiles (top 60) with a cumulative lift of 1.349 (about 67% accuracy) which is 6 times as many as if no model is used or the lift is above the baseline (1.349). The cumulative lift chart visually shows the advantage of using a predictive model than predicting without the model and do random sample.



Second model: SVM with Pol 2, sigmoid 2

SVM has a story of success and most popular for nonlinear distribution. One of the reasons for the success is using different optimization method with the option of interior point and active point such as kernel, polynomial degree with is suitable for the distribution.

In this project, SVE with Optimal Method of Active point set of pols 2 with max iteration time of 25 was used. The most important variables with a chi-square value of 3.6 and above were selected to proceed the model to reduce the time of processing and complexity of the model.

SVM Training model classification Table

Outcome	Predicted fatal/injury	Predicted Crash	Car	Total

Actual fatal/injury	3366	1875	5241
Actual Car Crash	685	4557	5242
Total	4051	6432	10483

Total Accuracy = $(TP + TN) / (TP+TN+FP+FN)$ = 0.76 ~ 76%

Sensitivity = $(TP/TP+FN)$ = 0.64 ~ 64%

Specificity = $(TN/ (TN+FN))$ = 0.87 ~87%

In the training dataset, SVM predicted only 64% of fatal/injury accident correctly. However, Car Crash Damage was classified correctly with accuracy rate of 87%. The overall accuracy rate was 76% where 7923 reports were classified correctly out of 10,483.

SVM Validation Data set Classification Table

Outcome	Predicted fatal/injury	Predicted Crash	Car	Total

Actual fatal/injury	971	527	1498
Actual Car Crash	531	966	1497
Total	1502	1493	2995

Total Accuracy = $(TP + TN) / (TP + TN + FP + FN)$ = 0.65 ~65%

Sensitivity = $(TP / (TP + FN))$ = 0.65 ~ 65%

Specificity = $(TN / (TN + FN))$ = 0.65 ~65 %

In the validation dataset, SVM predicted Sensitivity, Specificity and Total Accuracy rate with a rate of 65%. When the accuracy rate of the Training and Validation dataset was compared, the accuracy rate for sensitivity is 1% greater than the training dataset.

The Third Model: Neural Network

Neural networks are a class of parametric models that can capture a variety of nonlinear relationships between a setoff predictor and a target variable than can a traditional logistic regression model. Building a neural network involves two main stages. First, one must define the network configuration or structure; and then iteratively trains the model based on the given network structure (Kechen et al, na).

In this project, Neural Network was used after all class variable with missing values are replaced by count and transformed with Log Transformation method to reduce the skewness of the distribution. Model selection criteria was used for classification.

Neural Network classification Table for Training dataset

Outcome	Predicted fatal/injury	Predicted Car Crash	Total
Actual fatal/injury	3538	1703	5241
Actual Car Crash	1421	3821	5241
Total	4959	5524	

$$\text{Total Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 70\%$$

$$\text{Sensitivity} = (\text{TP} / (\text{TP} + \text{FN})) = 68\%$$

$$\text{Specificity} = (\text{TN} / (\text{TN} + \text{FN})) = 73\%$$

Neural Network classification Table for Validation dataset

Outcome	Predicted fatal/injury	Predicted Crash	Total
Actual fatal/injury	980	518	1498
Actual Car Crash	309	1188	1497
Total	1289	1706	2995

Total Accuracy = $(TP + TN) / (TP+TN+FP+FN)$ = 72%

Sensitivity = $(TP/TP+FN)$ = 65 %

Specificity = $(TN/ (TN+FN))$ = 79 %

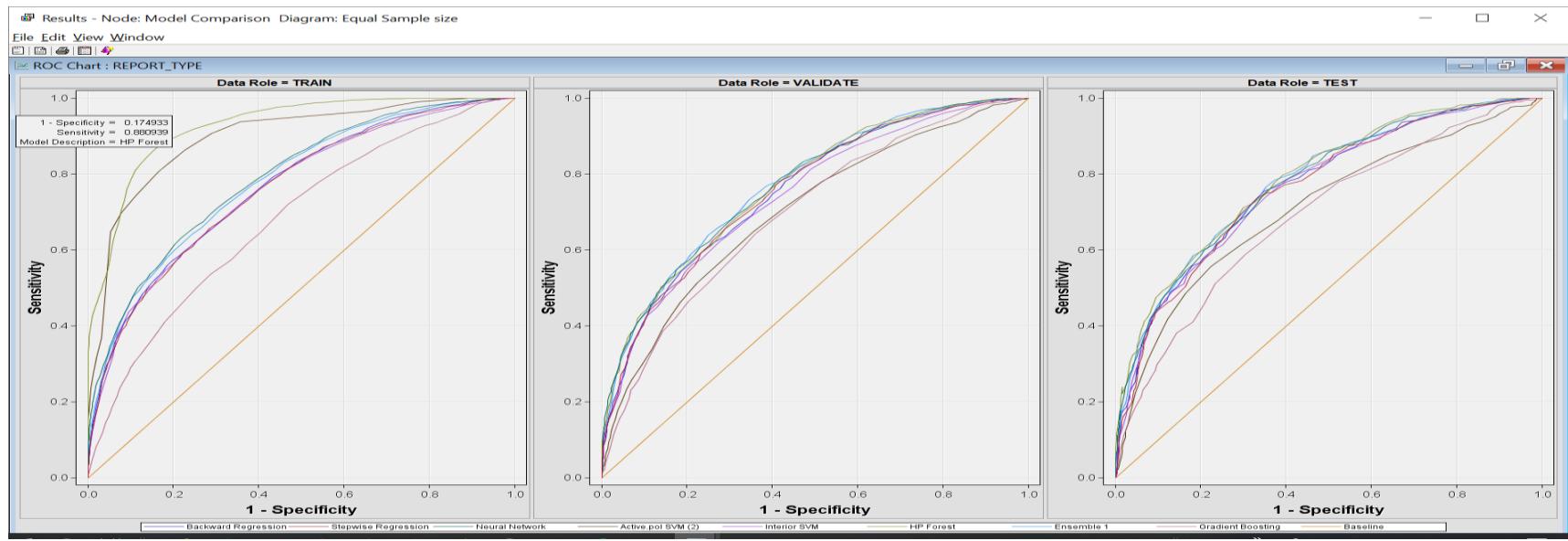
Neural Network model predicted fatal/injury car crash with accuracy rate of 68% on the training dataset and 65% on the validation dataset which shows that a difference of 3% overfit.

Comparison of Models

ROC Index and classification chart to compare the champion model

The area under the ROC (Receiver Operating Characteristic) can be used as a criterion to measure the test's discriminative ability.

The graph is plotted with true positive rate (sensitivity) as a Y axis and 1-Sensitivity as x-axis which is classifying false negative as true positive. As the Y axis increases, the sensitivity increases, in contrast as the x-axis increase the false positive increase. The goal is to find the optimum point where the model gets the max sensitivity and min false positive. Each point on the ROC curve represents a sensitivity/(1-specificity) pair correspond to a decision threshold. In this project, 0.05 cut off value was used. As it can be seen from the ROC curve and the summary table, the champion model is HP forest with ROC value of 0.86, 0.81, and 0.80 for training, validation, and test dataset. The second champion model is SVM and the ROC values for training, validation and test dataset are 0.84, 0.785 and 0.788 respectively. The third champion is Neural Network, with ROC value of 0.825 for training, 0.814 for validation and 0.80 for test dataset. In general, all three has ROC value above 80 which is catching fatal/injury car crash significantly. Moreover, all three of the models, revealed the most important variables which are Airbag deployment, Alcohol Test, collision, Car Damage code, and safety equipment.



Selection Editor-WORK.EMOUTFIT

Selected Model	Predecessor Node	Model Node	Model Description	Target Variable	Target Label	Selection Criterion: Train: Roc Index	Train: Average Squared Error	Train: Divisor for ASE	Train: Maximum Absolute Error	Train: Sum of Frequencies	Train: Root Average Squared Error	Train: Sum of Squared Errors	Train: Frequency of Classified Cases	Train: Miscl Rate
Yes	HPDMForest	HPDMForest	HP Forest	REPORT_TYPE	REPORT_TYPE	0.86	0.1573232109	20966	0.9260375215	10483	0.3966399008	3298.4384392	10483	0
No	HPSVM2	HPSVM2	Active.pol SVM (2)	REPORT_TYPE	REPORT_TYPE	0.843	0.2137876049	20966	0.6880321492	10483	0.4623717173	4482.2709253	10483	0
No	Neural	Neural	Neural Network	REPORT_TYPE	REPORT_TYPE	0.825	0.168841431	20966	0.998704329	10483	0.4109031894	3539.9294428		0
No	Neural2	Neural2	Neural Network (2)	REPORT_TYPE	REPORT_TYPE	0.825	0.168841431	20966	0.998704329	10483	0.4109031894	3539.9294428		0
No	Ensmbl2	Ensmbl2	Ensemble (2)	REPORT_TYPE	REPORT_TYPE	0.821	0.1777448761	20966	0.9957801071	10483	0.421598003	3726.5990731	10483	0
No	Reg2	Reg2	Backward Regression	REPORT_TYPE	REPORT_TYPE	0.818	0.1725638706	20966	0.9973655688	10483	0.4154080772	3617.9741108		0
No	Reg3	Reg3	Regression Backward	REPORT_TYPE	REPORT_TYPE	0.818	0.1725638706	20966	0.9973655688	10483	0.4154080772	3617.9741108		0
No	Reg	Reg	Stepwise Regression	REPORT_TYPE	REPORT_TYPE	0.817	0.1728156785	20966	0.9973343134	10483	0.4157110517	3623.2535149		0
No	Ensmbl	Ensmbl	Ensemble 1	REPORT_TYPE	REPORT_TYPE	0.816	0.1761692463	20966	0.9149526439	10483	0.4197252033	3693.5644174	10483	0
No	HPSVM	HPSVM	Interior SVM	REPORT_TYPE	REPORT_TYPE	0.813	0.2058364397	20966	0.9999996286	10483	0.4536920097	4315.5667948	10483	0
No	Tree	Tree	Decision Tree	REPORT_TYPE	REPORT_TYPE	0.768	0.1852128907	20966	0.9886363636	10483	0.4303636726	3883.1734665		0
No	Boost	Boost	Gradient Boosting	REPORT_TYPE	REPORT_TYPE	0.742	0.2079405958	20966	0.8247777085	10483	0.4560050393	4359.6825321		0

ROC Table: HP Forest, SVM (pol deg 2), and Neural Network are the three-champion model based on ROC value

The screenshot shows an Excel spreadsheet titled "ROC 04132019.xlsx". The table has the following data:

Models	Train ROC Index	Valid ROC Index	Test ROC Index
Reg. Stepwise	0.817	0.809	0.809
Reg. Backward	0.818	0.809	0.809
Gradient Boosting	0.742	0.761	0.749
HP Forest	0.86	0.818	0.819
Neural Network	0.825	0.814	0.809
SVM Interior	0.813	0.81	0.806
SVM(pol 2, sigmoid)	0.843	0.785	0.788
Decision Tree	0.768	0.805	0.767
Ense 1(Step. Reg, Gra	0.816	0.812	0.81

Conclusion

What is the best model? The quality of the model shouldn't be only measured by accuracy rate, but the percentage of overfitting or underfitting, simplicity, easily understandable by decision maker or layman and convenience to implement do matter the most.

According to the result of the ROC chart and cumulative Lift the first champion model is HP forest with a ROC value of 0.86, 0.818, 0.819 for the training, validation and test data respectively. Even though HP Forest is the champion model, most of the models came up with highest ROC value and when these models are compared with HP Forest, the difference is very low. As it can be seen from the above ROC Table, out of 9 models, only Decision Tree (ROC = 0.76) and Gradient Boosting (0.74) have ROC value less than 0.80. These implies that only ROC value can't determine the best model, instead assessing classification matrix and combine the overall result helps to select the working model.

As the classification matrix indicates, Logistic Regression with backward selection method takes the first rank in classifying sensitivity (fatal/injury car crash) with a rate of 0.667, Specificity of 0.79 and overall accuracy rate of 0.72 for training dataset. In contrary, HP forest predicted fatal/injury with a rate of 0.639 (sensitivity), 0.84 (specificity) and 0.74(overall accuracy rate) with a difference of 3%. So, based on the combination effect of ROC value and classification matrix, I personally, suggest that Logistic reg with back ward method is a better model than HP Forest.

When these two models are compared with respect to overfitting, which is when a model learns the detail and noise in the training dataset in such a way that it affects the performance of the model, again Logistic Reg is leading the way with a minimum overfitting of 0.011, (0.667 -0.656). In contrast, HP Forest has overfitting with a value of 0.038 (0.639-0.61). Thus, Logistic Regression is a better model in performing with unknown dataset than HP Forest.

Finally, the two models are compared with respect to complexity, interpretability and ease of implementation. As it is known, one of the strength parts of HP Forest was building the model with 100 trees and a max variable of 25, then dropping the weak variables with loss reduction and posterior probability. Even though, the approach is very systematic in dropping redundant variable or insignificant variables by using loss reduction variable importance with valid Gini or Valid Margin values, it is very complicated to interpret the result. Logistic regression with back ward methods used simple likelihood ratio test for Null Hypothesis of Beta value with the input value of effect, Df, chi-square test and probability. In this regard, both models identified the same type of variables, including their codes, to determine the most significant variables. These are Airbag dep, Alcohol test, driver condition, damage code, and safety equipment.

In general, in all aspect of statistical measurement except ROC value, Logistic regression with backward method showed good performance than Hp forest and can be the best models to implement for the reduction of fatal/ injury accident in Maryland. In addition, the sensitivity is below 80% in all models which was the goal of this project to predict as min accuracy rate. However, to compensate the drawback, cost of classification was considered by applying more weight for death or injury event than car damage accident. I believe it is true and acceptable among the decision makers to reduce the risk of death by giving appropriate attention with respect to time, money and knowledge to come to accomplish the vision of ZERO death in Maryland due to car crash.

Computing cost of classification

Why do decision makers use cost classification? Cost classification would be the best solution to bring to management's attention certain costs that are considered more crucial than others. In this project, it is absolutely recommended to impose different weight for fatal/injury accident and car damage because of the consequence of fatal/injury on individual life, family, society and county economy. The below cost matrix was used to calculate the cost of each occurrence based on the Logistic Regression with backward selection method.

Cost classification based on the Logistic Regression backward

Cost Matrix		Predicted class	
Actual class	C(I/J)	+	-
+	-1(3499)	100(1742)	
	1(1107)	0(3821)	

- The cost of predicting fatal/injury car crash incorrectly weight 100 times more than predicting car damage as fatal. Based on this the cost classification will be = $-3499 + 1107 + 174200 + 0 = 170,808$. This implies that fatal/injury car crash is very serious which needs to be addressed seriously.

Logistic Regression with backward selection method classification Table for Training dataset

Outcome	Predicted fatal/injury	Predicted Car Crash	Total
Actual fatal/injury	3499	1742	5241
Actual Car Crash	1107	3821	4928
Total	4959	5524	10229

$$\text{Total Accuracy} = (\text{TP} + \text{TN}) / (\text{TP} + \text{TN} + \text{FP} + \text{FN}) = 72\%$$

$$\text{Sensitivity} = (\text{TP} / (\text{TP} + \text{FN})) = 667\%$$

Specificity = $(TN / (TN + FN)) \times 100$ = 78%

Reference

Brandon. W (2019): Crashes more than twice as common on I-83 as MD Highways.

MHSP (2017): Maryland Highway safety office: Retrieved from

http://www.mva.maryland.gov/safety/_docs/FFY2017MarylandHSP.pdf

MDT (2018): Maryland Department of Transportation: Retrieved from <http://www.mva.maryland.gov/safety/mhso/index.htm>

Six Reasons why data science project fail: Retrieved from <https://medium.com/@ODSC/6-reasons-why-data-science-projects-fail-6240bf9326f6>

Miller and Zois (2016): Car Accident statistics in Maryland: Retrieved from <https://www.millerandzois.com/car-accident-statistics.html>

<http://www.mva.maryland.gov/safety/mhso/Maryland-Traffic-Safety-Data.htm>

<http://www.mva.maryland.gov/safety/mhso/index.htm>

https://www.cio.com.au/article/401353/how_define_scope_project/

<https://rapidbi.com/the-difference-between-goals-objectives/>

<https://www.pinderplotkin.com/er-vs-urgentcare/>

https://www.iii.org/sites/default/files/docs/pdf/auto_rates_wp_092716-62.pdf

http://www.marylandconsumers.org/penn_station/folders/issues/auto/insurance/Auto_Insurance_Brief_with_Infographics_for_website.pdf

<https://www.hindawi.com/journals/mpe/2013/547904/>

United States Census Bureau. B01001 SEX BY AGE, 2017 American Community Survey 5-Year Estimates. U.S. Census Bureau, American Community Survey Office. Web. 6 December 2018. <http://www.census.gov/>.

http://www.mva.maryland.gov/_resources/docs/MarylandSHSP_2016-2020-Final.pdf

<http://jesshampton.com/2013/11/11/model-evaluation-explaining-the-cumulative-lift-chart/>

Appendix 1

VARIABLE	CATEGORY	DATA TYPE
SEX_CODE	F=FEMALE M=MALE U=UNKNOWN	Nominal
CONDITION_CODE	U- UNKNOWN 00=NOT APPLICABLE 01=APPARENTLY NORMAL 02=HAD BEEN DRINKING 03=USING DRUGS 04=PHYSICAL DEFECTS 05=OTHER HADICAPS 06=III 07=FATIGUED FAINTED 08=APPARENTLY ASLEEP 09=EMOTIONAL DEPRESSED ANGRY DISTURBED	Nominal

	10=INFLUENCED BY Medications and/or DRUGS AND /OR ALCOHOL 88=OTHER HADICAPS 99UNKNOWN	
OCC-SEAT-POST- CODE	00=NOT APPLICABLE 01=IN VEHICLE 02=CENTER FRONT SEAT 03=RIGHT FRONT SEAT 04=LEFT REAR/MOTORCYCLE PASSENGER 05=CENTER REAR SEAT 06=RIGHT REAR SEAT 07=OTHER IN VEHICLE 08=CARGO AREA 09=UNENCLOSED CARGO AREA 15=TRAILER UNIT 88= OTHER 99= UNKNOWN	Nominal
ALCHOL _TEST _CODE	00=NOT APPLICABE 01=TEST REFUSED 02=POSITIVE PRELIMINARY TEST 03=EVIDENCE TEST GIVEN 88=OTHE R 99=UNKNOWN	Nominal

DRUG_TEST RESULT_CODE	00=NOT APPLICABE 01=TEST REFUSED 02=POSITIVE PRELIMINARY TEST 03=EVIDENCE TEST GIVEN 88=OTHE R 99=UNKNOWN	Nominal
EQUIP-PROOB-CODE		Nominal
SAF-EQUIP-CODE		Nominal
AIRBAG_DEPLOYED		Nominal
DATA OF BIRTH		Nominal
EMS-UNIT LABEL		Nominal
LIGHT CODE	00=NOT APPLICABEL 01=DAYLIGHT 03=DARK LIGHT ON 04=DARK NO LIGHTS 05=DAWN 06=DUSK 07=DARK-UNKNOWN 88=OTHER 99=UNKNOWN	Nominal
COUNTY_NO	01=ALLEGANY 02=ANNE ARUNDEL 03=BALTIMORE 04=CALVERT 05=CAROLINE 06=CARROLL 07=CECIL 08=CHARLES 09=DORCHESTER	Nominal

	10=FREDERICK 11=GARRETT 12=HARFORD 13=HOWARD 14=KENT 15=MONTGOMERY 16=PRINCE GEORGE'S 17=QUEEN ANNE'S 18=ST. MARY'S 19=SOMERSET 20=TALBOT 21=WASHINGTON 22=WICOMICO 23=WORCESTER 24=BALTIMORE CITY	
JUNCTION_TYPE_CODE		Nominal
		Nominal
		Nominal
COLLISION_TYPE_CODE		Nominal
SURFACE_CODNITION_CODE		Nominal
RD_COND_CODE		Nominal
RD_DIV_CODE		Nominal
WEATHER_CODE		Nominal
ACC_TIME		Nominal
RTE-NO		Nominal

APPENDIX 2 JOINED DATASET WITH MS-SQL

JOINED ACCIDENT 2 EDITED.xlsx.xlsx - Excel

RUTH GEMECHU (876073)

The screenshot shows a Microsoft Excel spreadsheet titled "JOINED ACCIDENT 2 EDITED.xlsx.xlsx". The ribbon menu is visible at the top, showing tabs like File, Home, Insert, Draw, Page Layout, Formulas, Data, Review, View, and Help. The Home tab is selected. The status bar indicates "RUTH GEMECHU (876073)". The main area displays a table with the following columns and data:

	S	T	U	V	W	X	Y
1	RD_DIV_CODE	REPORT_TYPE	WEATHER_CODE	ACC_TIME	DISTANCE	LATITUDE	LONGITUDE
2		4	1	7	15:37	50	40
3			0	6	15:26		40
4		4	0	0	06:56	0	40
5		4	1	6	06:57	100	40
6			0	6	14:39	15	40
7		1	0	7	10:36	20	40
8		1	0	7	07:48	30	40
9		3	0	10	12:22	90	40
10		4	1	6	17:10	0	40
11		4	1	7	19:25	75	39
12		4	1	6	15:45	0	39
13			0	6	18:48		40
14		4	0	6	19:10	50	39
15		3	1	5	22:57	56	40
16		4	0	6	19:40	75	39
17		4	0	12	17:15	50	40
18		3	1	6	17:05	0	40

Sheet1 | Sheet2 | Sheet3 | +

APPENDIX 3 COMBINING THE THREE DATASETS WITH MS-SQL INNER JOIN WITH COMMON VARIALBE OF REPORT NO

Microsoft SQL Server Management Studio

File Edit View Query Project Debug Tools Window Community Help

New Query | Object Explorer | SQL Server Object Explorer | Task List | Object Explorer Details |

Object Explorer

SEBASTIAN\ALIC_DR (SQL Server)

Databases

- System Database
- ACC
- Adventure Database
- Tables
- Views
- Synonyms
- Programmability
- Service Broker
- Storage
- Security
- Data670
- meki
- Database Diagrams
- Tables
- System Tables
- dbo.crash\$
- dbo.crashperse\$
- dbo.crashvehicl\$
- dbo.Person\$
- dbo.Sheet1\$
- dbo.Sheet2\$
- dbo.Sheet3\$
- dbo.Vehicle\$
- Views
- Synonyms
- Programmability
- Service Broker
- Storage
- Security
- mirror (Restoring...)
- Security
- Server Objects
- Replication
- Management

SQLQuery3.sql - SEBASTIAN\meki (54)" | SQLQuery2.sql - SEBASTIAN\meki (58)" | Object Explorer Details |

```
# Joining the three dataset crash, vehicle and person
select * from dbo.Person$ 
join dbo.crash$ ON dbo.Person$.REPORT_NO = dbo.crash$.REPORT_NO
join dbo.[Vehicle $] ON dbo.[Vehicle $].REPORT_NO = dbo.Person$.REPORT_NO
```

Results Messages

REPORT_NO	HARM_EVENT_CODE	DAMAGE_CODE	MOVEMENT_CODE	VEH_YEAR	VEH_MAKE	VEH_MODEL	GOING_DIRECTION_CODE	BODY_TYPE_CODE	SPEED_LIMIT	AREA_DAMAGE_CODE_IMP1	AREA_DAMAGE_CODE2
AA04040021	1	99	3	NULL	UNKNOWN	UNKNOWN	E	NULL	50	6	NULL
AA04040022	2	1	11	2001	DODGE	CARAVAN	S	2	10	6	NULL
AA073000N	1	4	2	2010	TOYOTA	COROLLA	E	2	40	5	5
AA073000N	1	3	3	2014	TOY	COROLLA	E	2	40	12	12
AA0790011	9	4	1	2009	MITS	4S	S	2	30	1	12
AA0790014	1	3	12	NULL	UNKNOWN	UNKNOWN	S	NULL	25	5	5
AA0790015	2	4	1	2006	HONDA	4S	N	2	25	12	11
AA0106001F	16	2	3	2008	CHEV	4S	N	17	40	5	NULL
AA0107005Y	1	4	1	1998	TOYOTA	CAMRY	W	2	40	12	12
AA01070061	1	4	1	2007	CHEVROLET	MALIBU	E	2	40	12	11
AA01070062	1	4	1	2003	HYUNDAI	ELANTRA	W	2	40	5	5
AA01070063	2	3	13	2006	CHEVROLET	4DR	E	2	0	10	9
AA01070065	1	3	7	NULL	UNKNOWN	UNKNOWN	W	99	40	4	5
AA01070066	3	3	2	1995	FORD	RANGER	S	20	25	12	12

Query executed successfully.

SEBASTIAN\ALIC_DR (10.50 SP2) | SEBASTIAN\meki (54) | meki | 00:00:04 | 26967 rows

Ln 8 Col 1 Ch 1 INS

APPENDIX 4: 1 INITIAL EXPLORATION

Variables - Ids

(none) ▾ not Equal to ▾ Label

Missing Bytes Statistics

Apply Reset

Columns: Name Role Level Report Order Drop Lower Limit Upper Limit Type Format Informat Length Number of Levels Percent Missing Minimum Maximum Mean Standard Deviation

Name	Role	Level	Report	Order	Drop	Lower Limit	Upper Limit	Type	Format	Informat	Length	Number of Levels	Percent Missing	Minimum	Maximum	Mean	Standard Deviation
ACC_TIME	Input	Nominal	No	No	No	-	-	Character	\$6.	\$6.	6	13	0	-	-	-	
AIRBAG_DEPLOYED	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	7	6.140515	-	-	-	
ALCOHOL_TEST_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	5	4.38154	-	-	-	
ARM_DAMAGE_CODE_MP1	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	13	4.00001	-	-	-	
BODY_TYPE_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	29	6.405674	-	-	-	
COLLISION_TYPE_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	17	14.09187	-	-	-	
CONDITION_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	12	10.77578	-	-	-	
COUNT_ID	Rejected	Nominal	No	No	No	-	-	Numeric	BEST.		8	1	0.003725	1	24	12.30196	
DAMAGE_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	6	6.272023	-	-	-	
DATE_OF_BIRTH	Input	Nominal	No	No	No	-	-	Numeric	YEAR4.0		8	-	8.202293	-	-	-	
DISTANCE	Input	Interval	No	No	No	-	-	Numeric	BEST.		8	-	9.666455	0	986.31	64.82719	
DRUG_TESTRESULT_CODE	Input	Nominal	No	No	No	-	-	Character	\$1.	\$1.	1	4	99.33142	-	-	141.6894	
END_STOP_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	11	13.0004	-	-	-	
GONG_DIRECTION_CODE	Input	Nominal	No	No	No	-	-	Character	\$1.	\$1.	1	5	8.127591	-	-	-	
JUNCTION_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	11	11.07085	-	-	-	
LATITUDE	Input	Interval	No	No	No	-	-	Numeric	BEST.		8	0	37.52082	39.87399	39.14282	0.321627	
LONGITUDE	Input	Nominal	No	No	No	-	-	Character	\$2.	\$2.	2	61	13.94345	-	-	-	
LIGHT_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	7	1.93845	-	-	-	
LONGITUDE	Input	Interval	No	No	No	-	-	Numeric	BEST.		8	-	0	-79.469	-75	-76.7492	
MOVEMENT_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	15	96.96709	-	-	0.449798	
MOVEMENT_CODE_L1	Rejected	Nominal	No	No	No	-	-	Numeric	BEST.		8	-	5.5220025	0	20.03	4.512317	
ODOMETER_POS_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	12	98.6272	-	-	4.644554	
PERSON_CONDITION	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	24	1.008479	-	-	-	
RD_COND_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	9	10.79072	-	-	-	
RD_DIV_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	-	11.86643	0	5.01	2.357285	
REPORT_ID	Input	Nominal	No	No	No	-	-	Character	\$12.	\$12.	12	-	-	-	-	-	
REPORT_TYPE	Target	No	No	No	No	-	-	Numeric	BEST.		8	2	0	-	-	-	
SAP_EQUIP_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	14	12.8525	-	-	-	
SEX_CODE	Input	Nominal	No	No	No	-	-	Character	\$1.	\$1.	1	3	8.034214	-	-	-	
SPEED_LIMIT	Input	Interval	No	No	No	-	-	Numeric	BEST.		8	-	0	0	75	34.01076	
SLRF_COND_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	10	10.47035	-	-	15.81744	
WEATHER_CODE	Input	Nominal	No	No	No	-	-	Numeric	BEST.		8	11	1.438016	-	-	-	

Explore OK Cancel

APPENDIX 4:2

Results - Node: ACC 2 EDITED03072019 Diagram: Data preparation 03202019

File Edit View Window

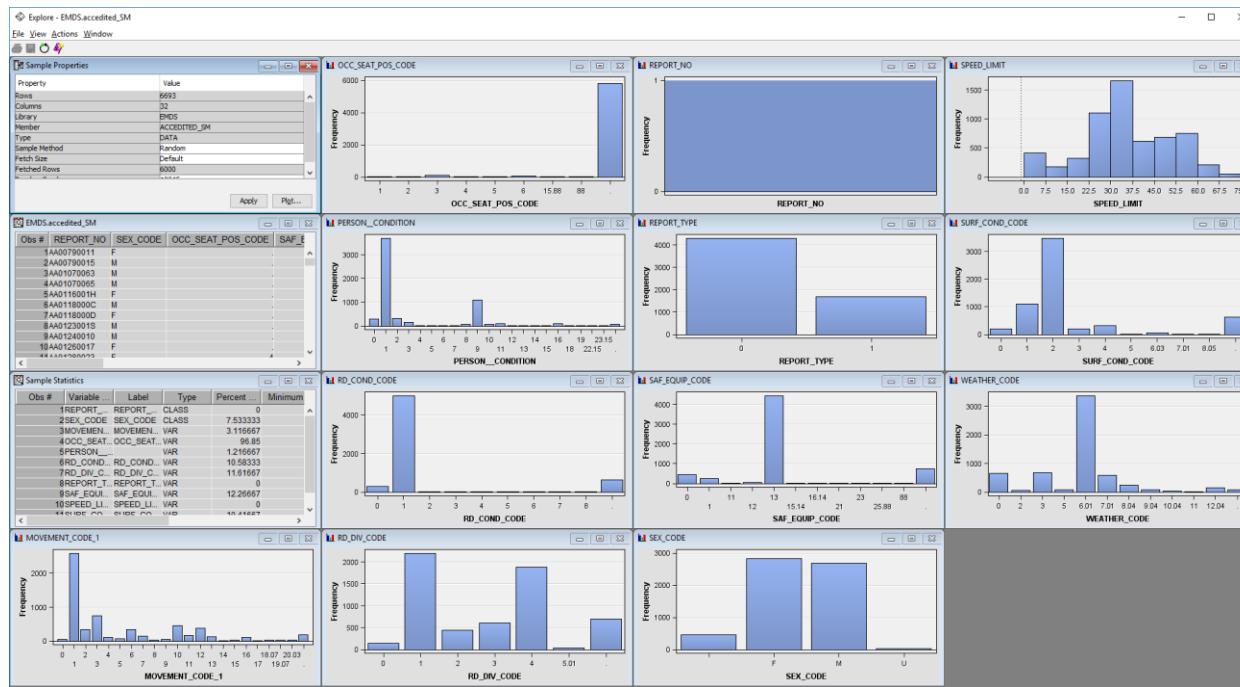
Statistics Table

Variable Name	Role	Measurement Level	Type	Number of Levels	Percent Missing	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	
ACC TIME	INPUT	NOMINAL	C	513	0	5.976393	
AIRBAG DEPLO...	INPUT	NOMINAL	N	5	4.213357	
ALCOHOL TES...	INPUT	NOMINAL	N	5	3.200508	
ARMED WEAP...	INPUT	NOMINAL	N	29	6.036157	
BODY TYPE...	INPUT	NOMINAL	N	17	14.8364	
COLLISION_TY...	INPUT	NOMINAL	N	12	10.2794	
CONDITION_CO...	INPUT	NOMINAL	N	24	0.016111	
COUPLE NO.	INPUT	NOMINAL	N	6	6.036157	
DAMAGE_CODE	INPUT	NOMINAL	N	89	7.485433	
DATE OF BIRTH	INPUT	NOMINAL	N	9	9.26341	0	950	67.17112	146.4317	3.034248	9.779442	
DISTANCE	INPUT	INTERVAL	O	4	9.573	
DRUG TEST REJECT...	INPUT	NOMINAL	N	6	11.62513	
GOING DIRECTI...	INPUT	NOMINAL	C	5	7.769311	
JUNCTION CODE	INPUT	NOMINAL	N	11	10.42881	
LATITUDE	INPUT	INTERVAL	O	43	0	37.98582	39.87399	39.14291	0.320504	-0.94885	0.772228	
LICENSE_STAT...	INPUT	NOMINAL	C	7	12.49066	
LIGHT_CODE	INPUT	NOMINAL	N	1	1.897505	
LONGITUDE	INPUT	INTERVAL	N	13	96.9371	0	-79.4458	-75.0544	-76.7535	0.45391	-0.89222	7.23114
MOVEMENT_C...	INPUT	NOMINAL	N	21	3.212311	
MOVEMENT_C... REJECTED	INPUT	NOMINAL	N	8	96.83251	
OCC SEAT PO...	INPUT	NOMINAL	N	21	1.120574	
PERSON CON...	INPUT	NOMINAL	N	9	10.3898	
ROAD CODE	INPUT	NOMINAL	N	6	11.41491	
ROAD CODE	INPUT	NOMINAL	N	11	1.419393	
REPORT NO.	INPUT	NOMINAL	C	2	0	
REPORT_TYPE	TARGET	BINARY	N	11	12.23667	
SANCTION CO...	INPUT	NOMINAL	N	3	7.410728	
SEX_CODE	INPUT	NOMINAL	C	0	0	75	34.33139	15.67391	-0.10923	-0.26592	.	
SPEED LIMIT	INPUT	INTERVAL	N	9	10.21963	
SURF_COND_C...	INPUT	NOMINAL	N	11	1.419393	
WEATHER_CODE	INPUT	NOMINAL	N	11	1.419393	

1.89

Appendix 5 Initial Exploration





APPENDIX 6 STAT EXPLORE NODE WITH CHI-SQUARE PLOT AND WORTH

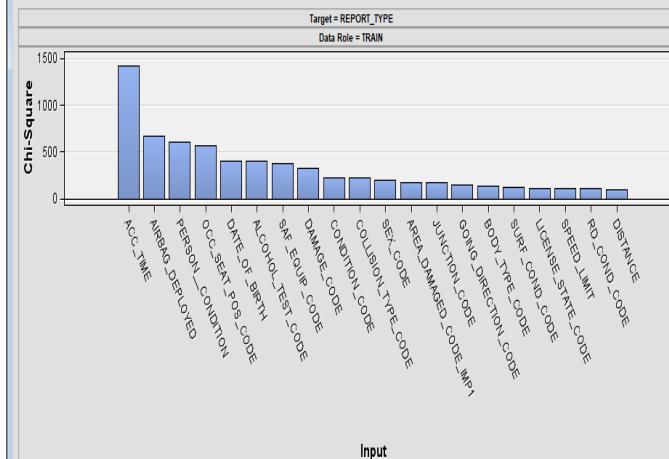
Results - Node: StatExplore Diagram: Data preparation 03202019

File Edit View Window



Chi-Square Plot

Chi-Square



Input

Output

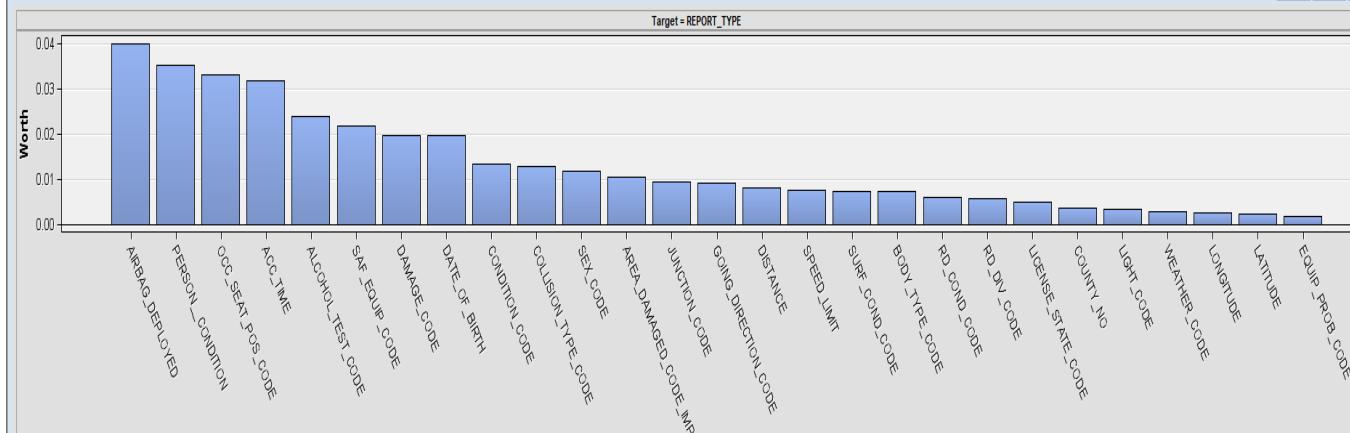
Variable Summary

Role	Measurement Level	Frequency Count
ID	NOMINAL	1
INPUT	INTERVAL	4
INPUT	NOMINAL	23
REJECTED	NOMINAL	3
TARGET	BINARY	1

Variable Levels Summary
(maximum 500 observations printed)

Variable	Role	Frequency Count

Variable Worth



Variable

Appendix 5

Class Variable Summary Statistics

```
File Edit View Window
□ □ | □
Output
34
35
36 Class Variable Summary Statistics
37 (maximum 500 observations printed)
38
39 Data Role=TRAIN
40
41          Number
42          Role  Variable Name   Role  Levels  Missing  Mode  Mode2  Percentage  Mode2 Percentage
43
44
45 TRAIN ACC_TIME    INPUT    513      0  11:00  0.71  17:00  0.71
46 TRAIN AIRBAG_DEPLOYED INPUT    8      400  1  63.69  2  11.64
47 TRAIN ALCOHOL_TEST_CODE INPUT    6      282  0  92.75  .  4.21
48 TRAIN AREA_TYPE_CODE_INP1 INPUT   16      265  2  40.0  6  16.73
49 TRAIN BODY_TYPE_CODE    INPUT   30      415  2  65.74  23.08  11.21
50 TRAIN COLLISION_TYPE_CODE INPUT   19      993  3  25.24  17  24.79
51 TRAIN CONDITION_CODE    INPUT   13      688  1  77.20  .  10.28
52 TRAIN DAMAGE_CODE     INPUT    7      404  4  37.25  3  24.00
53 TRAIN DIASTOLIC_BP    INPUT   90      901  .  90.0  2011  2.01
54 TRAIN DRUG_TESTRESULT_CODE INPUT    5      6554  .  99.42  U  0.27
55 TRAIN EQUIP_Prob_CODE  INPUT    7      945  1  67.98  0  19.20
56 TRAIN GOING_DIRECTION_CODE INPUT   12      520  N  27.27  S  25.16
57 TRAIN JUNCTION_CODE    INPUT   12      698  1  31.45  2  24.47
58 TRAIN LIGHT_CODE       INPUT   44      856  M  76.59  14  11.49
59 TRAIN LIGHT_CODE       INPUT   9      127  1  55.39  3  25.56
60 TRAIN MOVEMENT_CODE   INPUT   14      6408  .  96.94  51  0.99
61 TRAIN OCC_HEAT_POS_CODE INPUT   9      6481  .  96.83  3  1.64
62 TRAIN PECULIAR_CONDITION INPUT   22      75  1  61.87  9  18.02
63 TRAIN PECUL_COND_CODE  INPUT   10      693  1  69.41  .  14.30
64 TRAIN RL_DIV_CODE      INPUT   7      764  1  36.86  4  31.24
65 TRAIN SAF_EQUIP_CODE   INPUT   12      819  13  74.39  .  12.24
66 TRAIN SEC_CODE         INPUT   4      498  F  46.81  M  45.18
67 TRAIN SURF_CODE        INPUT   10      684  5  56.46  1  19.45
68 TRAIN SURF_CODE        INPUT   12      96  6.01  56.71  3  11.41
69 TRAIN REPORT_TYPE      TARGET   2      0  72.03  1  27.97
70
71
72 Distribution of Class Target and Segment Variables
73 (maximum 500 observations printed)
74
75 Data Role=TRAIN
76
77          Frequency
78          Data  Variable  Role  Level  Count  Percent
79
80
81 TRAIN REPORT_TYPE TARGET  0      4821  72.0305
82 TRAIN REPORT_TYPE TARGET  1      1872  27.9695
83
84
85
86 Interval Variable Summary Statistics
87 (maximum 500 observations printed)
88
```

Interval Variable Summary Statistics

```

Results - Node: StatExplore Diagram: Data preparation 03202019
File Edit View Window
Output
70
71
72 Distribution of Class Target and Segment Variables
73 (maximum 500 observations printed)
74
75 Data Role=TRAIN
76
77
78 Data Variable Frequency
79 Role Name Role Level Count Percent
80
81 TRAIN REPORT_TYPE TARGET 0 4821 72.0305
82 TRAIN REPORT_TYPE TARGET 1 1872 27.9695
83
84
85 Interval Variable Summary Statistics
86 (maximum 500 observations printed)
87
88 Data Role=TRAIN
89
90
91 Variable Standard Non
92 Mean Deviation Missing Minimum Median Maximum Skewness Kurtosis
93
94 DISTANCE INPUT 67.17112 146.4317 6973 620 0 0.5 950 3.034240 9.779442
95 LATITUDE INPUT 39.14291 0.320504 6993 0 37.96582 39.22344 39.87399 -0.94885 0.772228
96 LONGITUDE INPUT -76.7553 0.45391 6903 0 -79.4459 -76.7194 -75.6758 -0.89222 7.23114
97 SPEED_LIMIT INPUT 54.53159 15.67591 6993 0 0 35 75 -0.10923 -0.46092
98
99
100 Class Variable Summary Statistics by Class Target
101 (maximum 500 observations printed)
102
103 Data Role=TRAIN Variable Name=ACC_TIME
104
105
106 Target Number
107 Target of Mode Mode2
108 Level Levels Missing Mode Percentage Mode2 Percentage
109
110 REPORT_TYPE 0 513 0 09:30 1.06 15:30 1.06
111 REPORT_TYPE 1 513 0 09:30 0.96 09:30 0.68
112 _OVERALL_ 513 0 11:00 0.72 17:00 0.71
113
114
115 Data Role=TRAIN Variable Name=AIRBAG_DEPLOYED
116
117 Target Number
118 Target of Mode Mode2
119 Level Levels Missing Mode Percentage Mode2 Percentage
120
121 REPORT_TYPE 0 8 170 1 70.79 0 11.10
122 REPORT_TYPE 1 8 230 1 45.41 2 19.60
123 _OVERALL_ 8 400 1 63.69 2 11.64
124

```

Appendix 6 Variable Clustering to see correlation between interval variable



Appendix 7 Replacement Editor

Replacement Editor-WORK.OUTCLASS

Variable	Formatted Value	Replacement Value	Frequency Count	Type	Character Unformatted Value	Numeric Value
.F_EQUIP_CODE	88		5N		.	88
.F_EQUIP_CODE	15.14		4N		.	15.14
.F_EQUIP_CODE	16.14		3N		.	16.14
.F_EQUIP_CODE	25.88		3N		.	25.88
.F_EQUIP_CODE	23		2N		.	23
.F_EQUIP_CODE	_UNKNOWN_	DEFAULT	N		.	
X_CODE	F		1881C	F	.	
X_CODE	M		1800C	M	.	
X_CODE			312C		.	
X_CODE	U		21C	U	.	
X_CODE	_UNKNOWN_	DEFAULT	C		.	
IRF_COND_CODE	2		2362N		.	2
IRF_COND_CODE	1		725N		.	1
IRF_COND_CODE	.		423N		.	
IRF_COND_CODE	4		207N		.	4
IRF_COND_CODE	0		129N		.	0
IRF_COND_CODE	3		128N		.	3
IRF_COND_CODE	6.03		35N		.	6.03
IRF_COND_CODE	5		3N		.	5
IRF_COND_CODE	7.01		2N		.	7.01
IRF_COND_CODE	_UNKNOWN_	DEFAULT	N		.	
EATHER_CODE	6.01		2272N		.	6.01
EATHER_CODE	0		452N		.	0
EATHER_CODE	3		447N		.	3
EATHER_CODE	7.01		377N		.	7.01
EATHER_CODE	8.04		166N		.	8.04
EATHER_CODE	12.04		85N		.	12.04
EATHER_CODE	.		64N		.	
EATHER_CODE	9.04		49N		.	9.04
EATHER_CODE	5		46N		.	5

< >

OK Cancel

Appendix 8

Replacing unknown with missing value or mode for class variable

```

Results - Node: Replacement Diagram: Data preparation 03202019
File Edit View Window
Output
28
29 Replacement Values for Class Variables
30
31
32 Variable Formatted Value Type Unformatted Value Numeric Replacement Value Label
33
34
35 ACC_TIME Unknown C . . _blank_ ACC_TIME
36 AIRBAG_DEPLOYED Unknown N . . AIRBAG_DEPLOYED
37 ALCOHOL_TEST_CODE Unknown N . . ALCOHOL_TEST_CODE
38 AREA_DAMAGE_CODE_IMPL Unknown N . . AREA_DAMAGE_CODE_IMPL
39 BODY_TYPE_CODE Unknown N . . BODY_TYPE_CODE
40 COLLISION_TYPE_CODE Unknown N . . COLLISION_TYPE_CODE
41 CONDITION_CODE Unknown N . . CONDITION_CODE
42 COUNTY_NO Unknown N . . COUNTY_NO
43 DAMAGE_CODE Unknown N . . DAMAGE_CODE
44 DATE_OF_BIRTH Unknown N . . DATE_OF_BIRTH
45 EQUIP_PROD_CODE Unknown N . . EQUIP_PROD_CODE
46 EQUIPMENT_CONDITION_CODE Unknown C . . _blank_ EQUIPMENT_CONDITION_CODE
47 JUNCTION_CODE Unknown N . . JUNCTION_CODE
48 LICENSE_STATE_CODE Unknown C . . _blank_ LICENSE_STATE_CODE
49 LIGHT_CODE Unknown N . . LIGHT_CODE
50 OC_CHEAT_POC_CODE Unknown N . . OC_CHEAT_POC_CODE
51 POC_CONDITION Unknown N . . POC_CONDITION
52 RD_COND_CODE Unknown N . . RD_COND_CODE
53 RD_DIV_CODE Unknown N . . RD_DIV_CODE
54 REPORT_TYPE Unknown N . . REPORT_TYPE
55 SAF_EQUIP_CODE Unknown N . . SAF_EQUIP_CODE
56 SEAT_CODE Unknown C . . _blank_ SEAT_CODE
57 SURF_COND_CODE Unknown N . . SURF_COND_CODE
58 WEATHER_CODE Unknown N . . WEATHER_CODE
59
60
61 * Report Output
62
63
64
65
66
67
68 Replacement Counts
69
70 Obs Variable Role Label Train Validation Test
71
72 1 ACC_TIME INPUT ACC_TIME 1603 968 325
73 2 AIRBAG_DEPLOYED INPUT AIRBAG_DEPLOYED 0 0 0
74 3 ALCOHOL_TEST_CODE INPUT ALCOHOL_TEST_CODE 0 0 0
75 4 AREA_DAMAGE_CODE_IMPL INPUT AREA_DAMAGE_CODE_IMPL 0 0 0
76 5 BODY_TYPE_CODE INPUT BODY_TYPE_CODE 0 0 0
77 6 COLLISION_TYPE_CODE INPUT COLLISION_TYPE_CODE 0 0 0
78 7 CONDITION_CODE INPUT CONDITION_CODE 0 0 0
79 8 COUNTY_NO INPUT COUNTY_NO 0 0 0
80 9 DAMAGE_CODE INPUT DAMAGE_CODE 0 0 0

```

```
Results - Node: Replacement (2) Diagram: Data preparation 03202019
File Edit View Window
Output
73 REP_SURF_COND_CODE Unknown N . 2 Replacement: SURF_COND_CODE
74 REP_WEATHER_CODE Unknown X . 6.01 Replacement: WEATHER_CODE
75
76
77 *-----*
78 * Report Output
79 *-----*
80
81
82
83
84 Replacement Counts
85
86 Obs Variable Label Role Train Validation Test
87
88 1 DISTANCE DISTANCE INPUT 187 108 33
89 2 LATITUDE LATITUDE INPUT 15 15 2
90 3 LONGITUDE LONGITUDE INPUT 72 32 14
91 4 REP_ACT_TIME replacement: ACC_TIME INPUT 0 0 0
92 5 REP_AIRBAG_DEPLOYED replacement: AIRBAG_DEPLOYED INPUT 0 0 0
93 6 REP_ALCOHOL_TEST_CODE replacement: ALCOHOL_TEST_CODE INPUT 0 0 0
94 7 REP_ANIMAL_CODE replacement: ANIMAL_CODE_IMP1 INPUT 0 0 0
95 8 REP_AVG_VEHICLE_CODE replacement: AVG_VEHICLE_CODE INPUT 0 0 0
96 9 REP_COLLISION_TYPE_CODE replacement: COLLISION_TYPE_CODE INPUT 0 0 0
97 10 REP_CONDITION_TYPE_CODE replacement: CONDITION_CODE INPUT 0 0 0
98 11 REP_COUNTY_CD replacement: COUNTY_CD INPUT 0 0 0
99 12 REP_DRIVER_CODE replacement: DRIVER_CODE INPUT 0 0 0
100 13 REP_DATE_OF_BIRTH replacement: DATE_OF_BIRTH INPUT 0 0 0
101 14 REP_EQUIP_PROD_CODE replacement: EQUIP_PROD_CODE INPUT 0 0 0
102 15 REP_GOING_DIRECTION_CODE replacement: GOING_DIRECTION_CODE INPUT 0 0 0
103 16 REP_JUNCTION_CODE replacement: JUNCTION_CODE INPUT 0 0 0
104 17 REP_LIC_PLATE_STATE_CODE replacement: LIC_PLATE_STATE_CODE INPUT 0 0 0
105 18 REP_LIGHT_CODE replacement: LIGHT_CODE INPUT 0 0 0
106 19 REP_OCC_SEAT_POS_CODE replacement: OCC_SEAT_POS_CODE INPUT 0 0 0
107 20 REP_PERSON_CONDITION replacement: PERSON_CONDITION INPUT 0 0 0
108 21 REP_RD_COND_CODE replacement: RD_COND_CODE INPUT 0 0 0
109 22 REP_RD_DVY_CODE replacement: RD_DVY_CODE INPUT 0 0 0
110 23 REP_EQUIP_TYPE replacement: EQUIP_TYPE TARGET 0 0 0
111 24 REP_SAF_EQUIP_CODE replacement: SAF_EQUIP_CODE INPUT 0 0 0
112 25 REP_SEX_CODE replacement: SEX_CODE INPUT 0 0 0
113 26 REP_SURF_COND_CODE replacement: SURF_COND_CODE INPUT 0 0 0
114 27 REP_WEATHER_CODE replacement: WEATHER_CODE INPUT 0 0 0
115 28 SPEED_LIMIT SPEED_LIMIT INPUT 0 0 0
116
```

Results - Node: Impute (2) (Diagram: Data preparation 03202019)						
File Edit View Window						
Output						
Number of Observations						
46						
Inputs						
50	Variable Name	Method	Imputed Variable	Impute Value	Role	Measurement Level
51	AIRBAG_DEPLOYED	COUNT	IMP_AIRBAG_DEPLOYED	1	INPUT	NOMINAL AIRBAG_DEPLOYED
52	ALCOHOL_TEST_CODE	COUNT	IMP_ALCOHOL_TEST_CODE	0	INPUT	NOMINAL ALCOHOL_TEST_CODE
53	AREA_DAMAGE_CODE_IMP1	COUNT	IMP_AREA_DAMAGE_CODE_IMP1	12	INPUT	NOMINAL AREA_DAMAGE_CODE_IMP1
54	BODY_TYPE_CODE	COUNT	IMP_BODY_TYPE_CODE	2	INPUT	NOMINAL BODY_TYPE_CODE
55	COLLISION_TYPE_CODE	COUNT	IMP_COLLISION_TYPE_CODE	3	INPUT	NOMINAL COLLISION_TYPE_CODE
56	CONDITION_CODE	COUNT	IMP_CONDITION_CODE	1	INPUT	NOMINAL CONDITION_CODE
57	DAMAGE_CODE	COUNT	IMP_DAMAGE_CODE	4	INPUT	NOMINAL DAMAGE_CODE
58	DEATH_INJURY_CODE	MEAN	IMP_DEATH_INJURY_CODE	8108.7055180	INPUT	INTERVAL DEATH_INJURY_CODE
59	DISTANCE	MEAN	IMP_DISTANCE	65.15489549	INPUT	INTERVAL DISTANCE
60	EQUIP_PROD_CODE	COUNT	IMP_EQUIP_PROD_CODE	1	INPUT	NOMINAL EQUIP_PROD_CODE
61	GOING_DIRECTION_CODE	COUNT	IMP_GOING_DIRECTION_CODE	N	INPUT	NOMINAL GOING_DIRECTION_CODE
62	JUNCTION_CODE	COUNT	IMP_JUNCTION_CODE	1	INPUT	NOMINAL JUNCTION_CODE
63	LICENCE_STATE_CODE	COUNT	IMP_LICENCE_STATE_CODE	MD	INPUT	NOMINAL LICENCE_STATE_CODE
64	LIGHT_CODE	COUNT	IMP_LIGHT_CODE	1	INPUT	NOMINAL LIGHT_CODE
65	PERSON_CONDITION	COUNT	IMP_PERSON_CONDITION	1	INPUT	NOMINAL PERSON_CONDITION
66	RD_COND_CODE	COUNT	IMP_RD_COND_CODE	1	INPUT	NOMINAL RD_COND_CODE
67	RD_DIV_CODE	MEAN	IMP_RD_DIV_CODE	2.3797551365	INPUT	INTERVAL RD_DIV_CODE
68	SEX_CODE	COUNT	IMP_SEX_CODE	15	INPUT	NOMINAL SEX_CODE
69	SURF_COND_CODE	COUNT	IMP_SURF_COND_CODE	F	INPUT	NOMINAL SURF_COND_CODE
70	WEATHER_CODE	COUNT	IMP_WEATHER_CODE	2	INPUT	NOMINAL WEATHER_CODE
71						
72						
73						
74						
75						
76						
77	Variable Distribution Training Data					
78						
79	Number of Missing	Number of Variables	Percent of Variables			
80	for TRAIN					
81	Obs					
82						
83	1	625	1	4.76190		
84	2	530	1	4.76190		
85	3	512	1	4.76190		
86	4	486	1	4.76190		
87	5	461	1	4.76190		
88	6	430	1	4.76190		
89	7	405	1	4.76190		
90	8	424	1	4.76190		
91	9	423	1	4.76190		
92	10	378	1	4.76190		
93	11	317	1	4.76190		
94	12	315	1	4.76190		
95	13	312	1	4.76190		
96	14	265	1	4.76190		
97	15	256	1	4.76190		
98	16	243	1	4.76190		
99	17	189	1	4.76190		
100	18	159	1	4.76190		
101	19	159	1	4.76190		
102	20	159	1	4.76190		

Results - Node: Impute (2) Diagram: Data preparation 03/2019

File Edit View Window

Imputation Summary

Variable Name	Impute Method	Imputed Variable	Impute Value	Role	Measurement Level	Label	Number of Missing for TRAIN
AIRBAG_DEPLOYED	COUNT	IMP_AIRBAG_DEPLOYED	1	INPUT	NOMINAL	AIRBAG_DEPLOYED	243
ALCOHOL_TEST_CODE	COUNT	IMP_ALCOHOL_TEST_CODE	0	INPUT	NOMINAL	ALCOHOL_TEST_CODE	169
AREA_DAMAGE_CODE_IMP1	COUNT	IMP_AREA_DAMAGE_CODE_IMP1	12	INPUT	NOMINAL	AREA_DAMAGE_CODE_IMP1	159
BODY_TYPE_CODE	COUNT	IMP_BODY_TYPE_CODE	2	INPUT	NOMINAL	BODY_TYPE_CODE	250
COLLISION_TYPE_CODE	COUNT	IMP_COLLISION_TYPE_CODE	3	INPUT	NOMINAL	COLLISION_TYPE_CODE	823
CONDITION_CODE	COUNT	IMP_CONDITION_CODE	1	INPUT	NOMINAL	CONDITION_CODE	430
DAMAGE_CODE	COUNT	IMP_DAMAGE_CODE	4	INPUT	NOMINAL	DAMAGE_CODE	265
DATE_OF_BIRTH	MEAN	IMP_DATE_OF_BIRTH	8108785188	INPUT	INTERVAL	DATE_OF_BIRTH	317
DISTANCE	MEAN	IMP_DISTANCE	65.15489549	INPUT	INTERVAL	DISTANCE	378
EQUIP_PROB_CODE	COUNT	IMP_EQUIP_PROB_CODE	1	INPUT	NOMINAL	EQUIP_PROB_CODE	512
GOING_DIRECTION_CODE	COUNT	IMP_GOING_DIRECTION_CODE	N	INPUT	NOMINAL	GOING_DIRECTION_CODE	315
ARMED_CODE	COUNT	IMP_ARMED_CODE	1	INPUT	NOMINAL	ARMED_CODE	424
LICENCE_STATE_CODE	COUNT	IMP_LICENCE_STATE_CODE	MD	INPUT	NOMINAL	LICENCE_STATE_CODE	520
LIGHT_CODE	COUNT	IMP_LIGHT_CODE	1	INPUT	NOMINAL	LIGHT_CODE	81
PERSON_CONDITION	COUNT	IMP_PERSON_CONDITION	1	INPUT	NOMINAL	PERSON_CONDITION	54
RD_COND_CODE	COUNT	IMP_RD_COND_CODE	1	INPUT	NOMINAL	RD_COND_CODE	426
RD_DIV_CODE	MEAN	IMP_RD_DIV_CODE	2.3797551365	INPUT	INTERVAL	RD_DIV_CODE	461
SAF_EQUIP_CODE	COUNT	IMP_SAF_EQUIP_CODE	13	INPUT	NOMINAL	SAF_EQUIP_CODE	486
SEX_CODE	COUNT	IMP_SEX_CODE	F	INPUT	NOMINAL	SEX_CODE	312
SURF_COND_CODE	COUNT	IMP_SURF_COND_CODE	2	INPUT	NOMINAL	SURF_COND_CODE	423
WEATHER_CODE	COUNT	IMP_WEATHER_CODE	6.01	INPUT	NOMINAL	WEATHER_CODE	64

Results - Node: Transform Variables Diagram: Data preparation 03202019

File Edit View Window

Transformations Statistics

Source	Method	Variable Name	Formula	Number of Levels	Non Missing	Mising	Minimum	Maximum	Mean	Standard Deviation	Skewness	Kurtosis	Label
Input	Original	IMP_ALCOHOL_TEST_		5	0	0	Imputed: ALCOHOL_T...
Input	Original	IMP_ALCOHOL_DEPLO_		7	0	0	Imputed: ALCOHOL_DE...
Input	Original	IMP_AREA_DAMAGEDE_		15	0	0	Imputed: AREA_DAMAGE...
Input	Original	IMP_BODY_TYPE_CO_		29	0	0	Imputed: BODY_TYPE_C...
Input	Original	IMP_CONDITION_CD_		17	0	0	Imputed: CONDITION_C...
Input	Original	IMP_CONDITION_CD_		12	0	0	Imputed: CONDITION_C...
Input	Original	IMP_DAMAGE_CODE_		6	0	0	Imputed: DAMAGE_C...
Input	Original	IMP_EQUIP_PROB_C_		5	0	0	Imputed: EQUIP_PROB...
Input	Original	IMP_EQUIP_PROB_C_		5	0	0	Imputed: EQUIP_PROB...
Input	Original	IMP_JUNCTION_CODE_		11	0	0	Imputed: JUNCTION_C...
Input	Original	IMP_LICENSE_STATE_		36	0	0	Imputed: LICENSE_S...
Input	Original	IMP_LIGHT_CODE_		7	0	0	Imputed: LIGHT_CODE...
Input	Original	IMP_MODEL_CD_		19	0	0	Imputed: MODEL_CD...
Input	Original	IMP_RD_COND_CODE_		9	0	0	Imputed: RD_COND...
Input	Original	IMP_SAFEQUIP_CD_		11	0	0	Imputed: SAF_EQUIP...
Input	Original	IMP_SEX_CODE_		3	0	0	Imputed: SEX_CODE...
Input	Original	IMP_VEHICLE_CD_		8	0	0	Imputed: VEHICLE_CD...
Input	Original	IMP_WEATHER_CODE_		10	0	0	Imputed: WEATHER_C...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_AIRBAG_DEP_	Dummy	2	0	0	IMP_AIRBAG_DEPLO...
Output	Computed	T1_IMP_ALCOHOL_TE_	Dummy	2	0	0	IMP_ALCOHOL_TEST...

Output

```

25 (maximum 500 observations printed)
26
27
28 Input Name      Rule   Level    Name      Level  Formula
29
30 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED1  BINARY Dummy
31 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED2  BINARY Dummy
32 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED3  BINARY Dummy
33 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED4  BINARY Dummy
34 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED5  BINARY Dummy
35 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED6  BINARY Dummy
36 IMP_AIRBAG_DEPLOYED INPUT  NOMINAL TI_IMP_AIRBAG_DEPLOYED7  BINARY Dummy
37 IMP_ALCOHOL_TEST_CODE INPUT  NOMINAL TI_IMP_ALCOHOL_TEST_CODE1  BINARY Dummy
38 IMP_ALCOHOL_TEST_CODE INPUT  NOMINAL TI_IMP_ALCOHOL_TEST_CODE2  BINARY Dummy
39 IMP_ALCOHOL_TEST_CODE INPUT  NOMINAL TI_IMP_ALCOHOL_TEST_CODE3  BINARY Dummy
40 IMP_ALCOHOL_TEST_CODE INPUT  NOMINAL TI_IMP_ALCOHOL_TEST_CODE4  BINARY Dummy
41 IMP_ALCOHOL_TEST_CODE INPUT  NOMINAL TI_IMP_ALCOHOL_TEST_CODE5  BINARY Dummy
42 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
43 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
44 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
45 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1A  BINARY Dummy
46 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1B  BINARY Dummy
47 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1C  BINARY Dummy
48 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
49 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
50 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy
51 IMP_AREA_DAMAGE_CODE_IMP1 INPUT  NOMINAL TI_IMP_AREA_DAMAGE_CODE_IMP1  BINARY Dummy

```

Results - Node: Backward Regression Diagram: Equal Sample size

File Edit View Window

Output

```
130
131
132      Likelihood Ratio Test for Global Null Hypothesis: BETA=0
133
134      -2 Log Likelihood      Likelihood
135      Intercept    Intercept &      Ratio
136      Only        Covariates   Chi-Square   DF      Pr > ChiSq
137
138      14532.524     10667.944    3864.5795    41      <.0001
139
140
141      Type 3 Analysis of Effects
142
143      Wald
144      Effect          DF      Chi-Square   Pr > ChiSq
145
146      TI_G_AIRBAG_DEPLOYED1    1      160.0701    <.0001
147      TI_G_AIRBAG_DEPLOYED2    1      0.0446      0.8328
148      TI_G_AIRBAG_DEPLOYED3    0      0.0000      .
149      TI_G_ALCOHOL_TEST_CODE1  1      60.5049    <.0001
150      TI_G_ALCOHOL_TEST_CODE2  1      17.7006    <.0001
151      TI_G_ALCOHOL_TEST_CODE3  0      0.0000      .
152      TI_G_COLLISION_TYPE_CODE1 1      62.0844    <.0001
153      TI_G_COLLISION_TYPE_CODE2 1      63.8871    <.0001
154      TI_G_COLLISION_TYPE_CODE3 1      0.0702      0.7911
155      TI_G_COLLISION_TYPE_CODE4 1      5.6492      0.0175
156      TI_G_COLLISION_TYPE_CODE5 0      0.0000      .
157      TI_G_CONDITION_CODE1     1      112.9630    <.0001
158      TI_G_CONDITION_CODE2     1      98.4284    <.0001
159      TI_G_CONDITION_CODE3     1      15.6780    <.0001
160      TI_G_CONDITION_CODE4     1      4.3590      0.0368
161      TI_G_CONDITION_CODE5     0      0.0000      .
162      TI_G_DAMAGE_CODE1       1      68.8384    <.0001
163      TI_G_DAMAGE_CODE2       1      13.5934      0.0002
```

Results - Node: Backward Regression Diagram: Equal Sample size

File Edit View Window

Output

169	TI_G_JUNCTION_CODE1	1	3.3469	0.0673
170	TI_G_JUNCTION_CODE2	1	3.2194	0.0728
171	TI_G_JUNCTION_CODE3	1	1.1413	0.2854
172	TI_G_JUNCTION_CODE4	0	0.0000	.
173	TI_G_LIGHT_CODE1	1	0.3753	0.5402
174	TI_G_LIGHT_CODE2	1	0.0791	0.7786
175	TI_G_LIGHT_CODE3	1	0.0766	0.7819
176	TI_G_LIGHT_CODE4	1	1.3115	0.2521
177	TI_G_LIGHT_CODE5	0	0.0000	.
178	TI_G_SAF_EQUIP_CODE1	1	188.6822	<.0001
179	TI_G_SAF_EQUIP_CODE2	1	5.4720	0.0193
180	TI_G_SAF_EQUIP_CODE3	0	0.0000	.
181	TI_G_SEX_CODE1	1	195.7405	<.0001
182	TI_G_SEX_CODE2	1	185.5497	<.0001
183	TI_G_SEX_CODE3	0	0.0000	.
184	TI_G_SPEED_LIMIT1	1	4.2640	0.0389
185	TI_G_SPEED_LIMIT2	1	1.2388	0.2657
186	TI_G_SPEED_LIMIT3	1	0.0257	0.8726
187	TI_G_SPEED_LIMIT4	1	0.4536	0.5006
188	TI_G_SPEED_LIMIT5	1	0.7352	0.3912
189	TI_G_SPEED_LIMIT6	0	0.0000	.
190	TI_G_SURF_COND_CODE1	1	3.8906	0.0486
191	TI_G_SURF_COND_CODE2	1	1.8970	0.1684
192	TI_G_SURF_COND_CODE3	1	3.5824	0.0584
193	TI_G_SURF_COND_CODE4	0	0.0000	.
194	TI_G_WEATHER_CODE1	1	1.5243	0.2170
195	TI_G_WEATHER_CODE2	1	1.0186	0.3129
196	TI_G_WEATHER_CODE3	1	1.0584	0.3036
197	TI_G_WEATHER_CODE4	1	0.0192	0.8897
198	TI_G_WEATHER_CODE5	1	0.0742	0.7853
199	TI_G_WEATHER_CODE6	0	0.0000	.
200				
201				
202				

Analysis of Maximum Likelihood Estimates