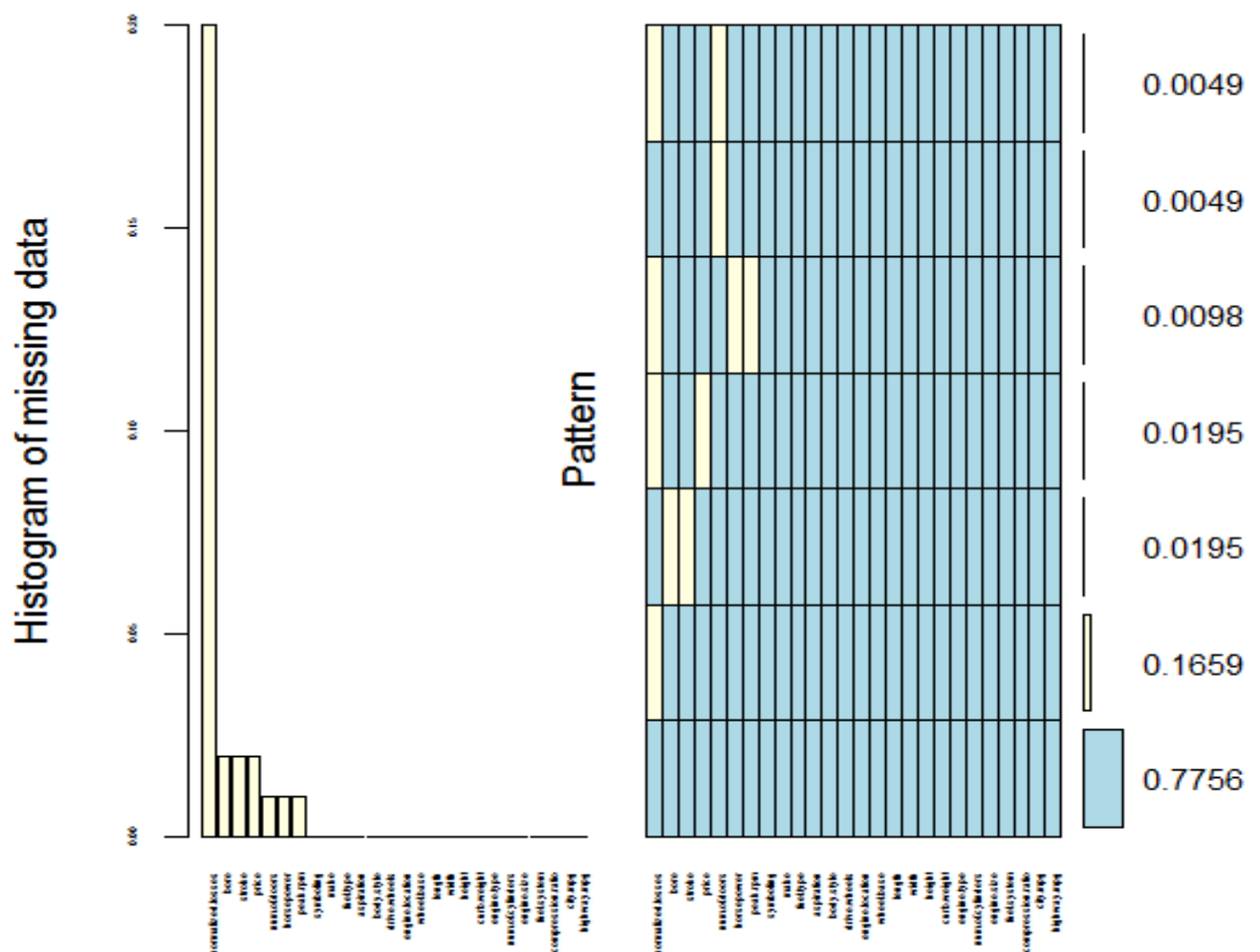


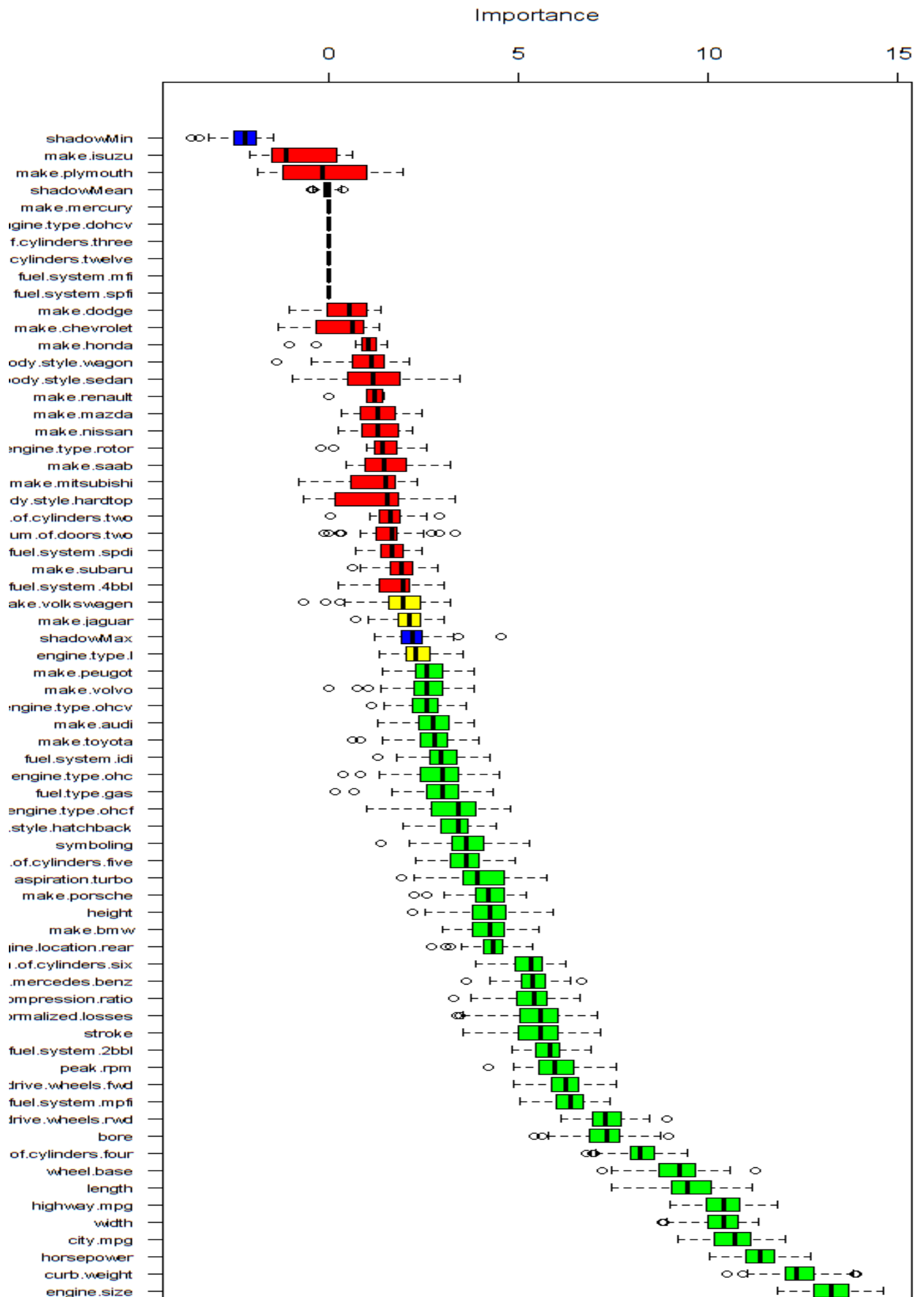
## Output Summary:-

The dataset consists of the automobile details and the task is to build a model for price prediction. **Caret** package has been mainly used, the implementation of **Boruta** to find important features has also been shown.

- **Data Loading and exploration** – The Data is loaded in a way to treat unwanted strings as NAs. The structure, summary and missing values are checked. The following plot is observed for the missing values:



- **Preprocessing** – The data is scaled, centered and the missing values are imputed using the KNN algorithm with the help of caret. Categorical variables are then converted to numerical using the one-hot encoding
- **Feature selection using Boruta** – Boruta package is used to find the important features out of several features present in the dataset. The following plot is obtained:

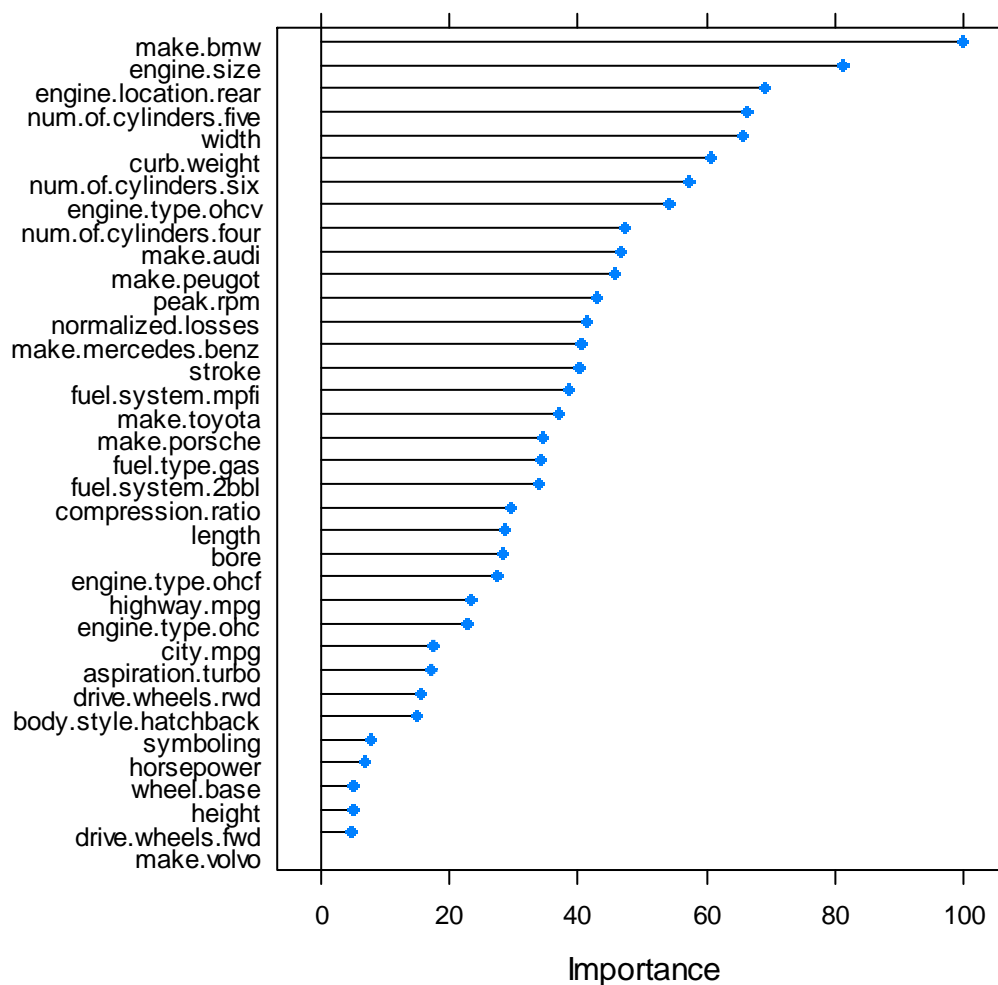


Using the `getSelectedAttributes` function, the important features are selected for model building.

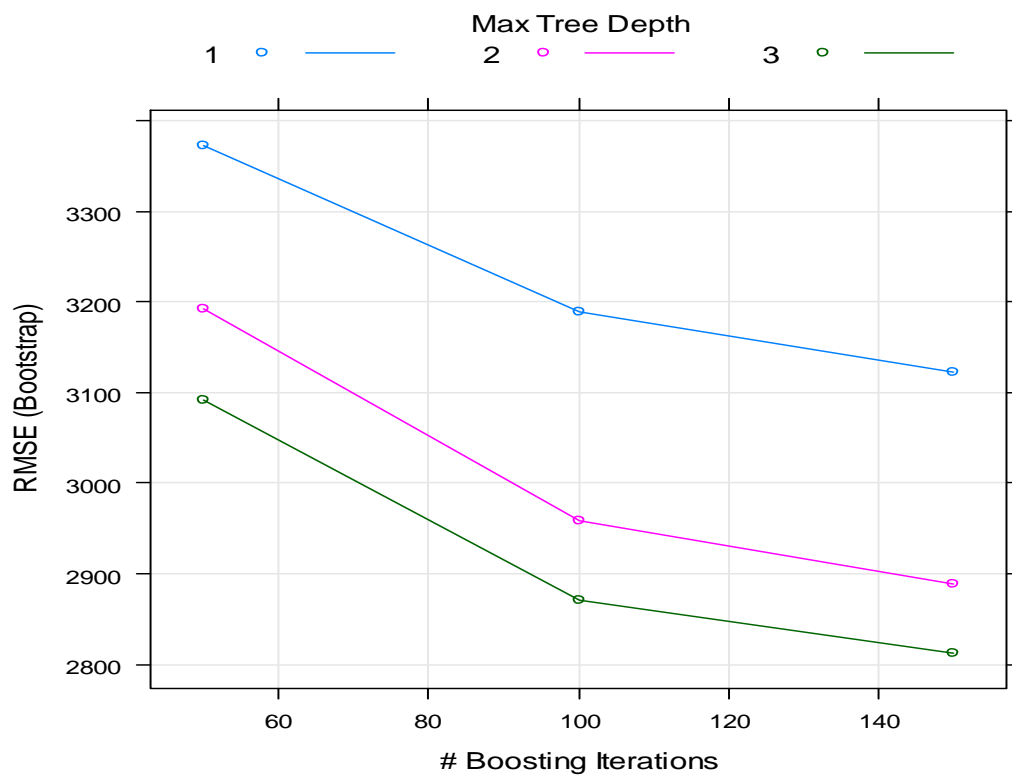
- **Splitting the data into train and test set** – The dataset is split into training and testing in the ratio of 8:2 using `createDataPartition` function.
- **Model Training** – Model training is performed on the training dataset using the `train` function from the `caret` package. Three algorithms – Linear model, GLM and NNET has been used to perform predictions. The following plots have been obtained for the three algorithms :

**Output Plot for LM :**

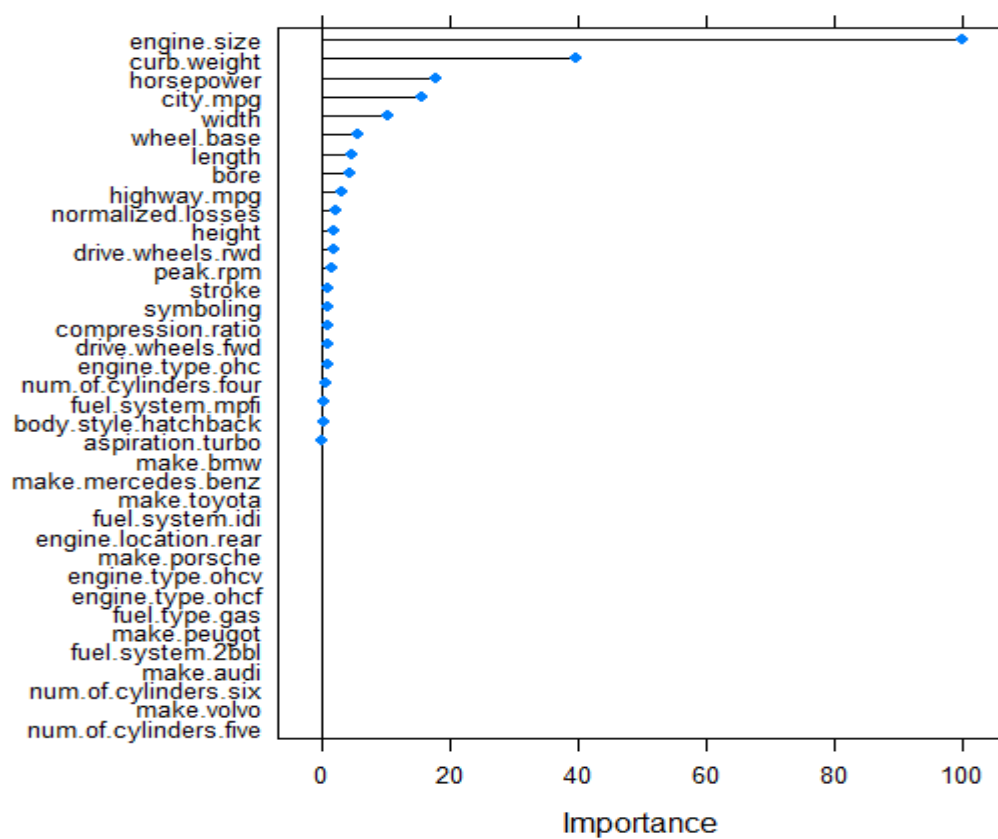
**LM - Variable Importance**



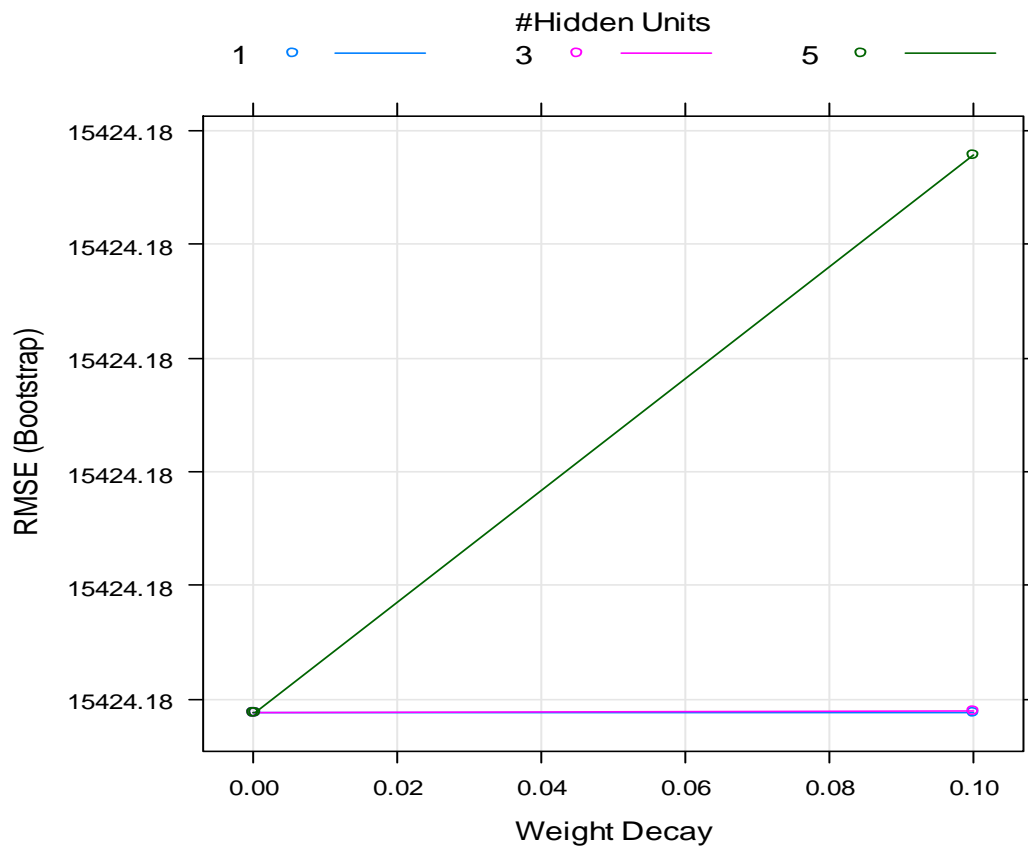
### GBM Plots :



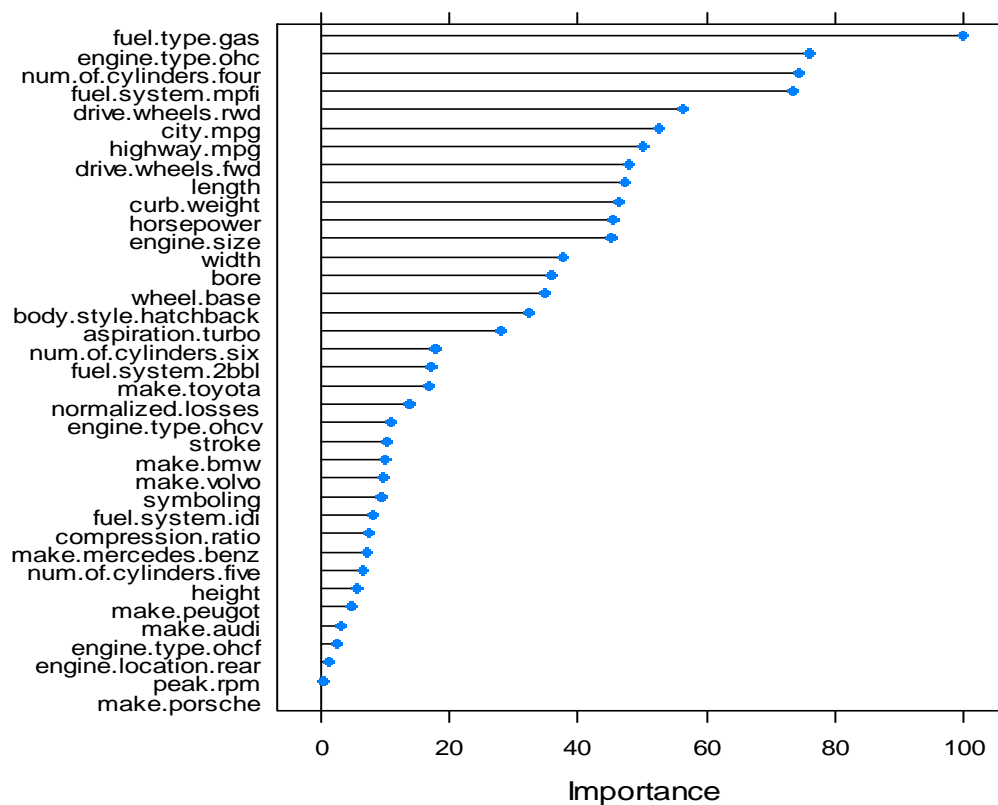
### GBM - Variable Importance



### NNET Output Plots :



### NNET - Variable Importance

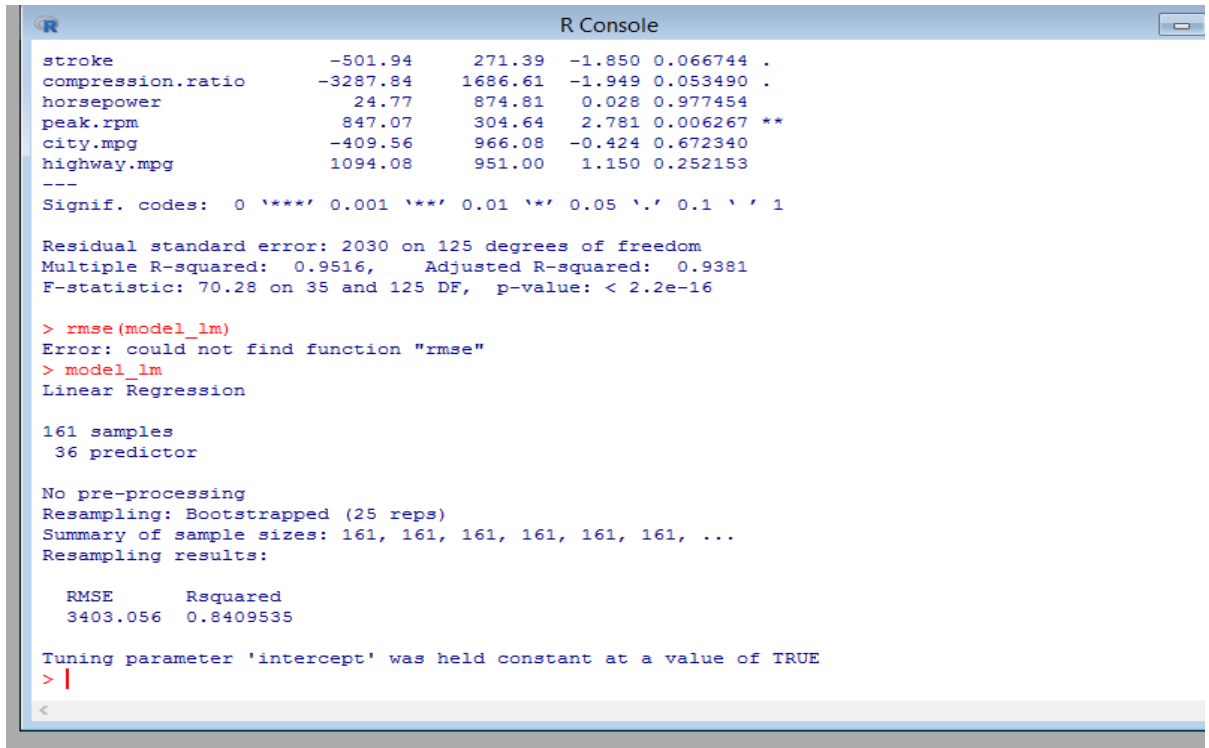


- **Price Prediction & Model Evaluation** – The prices can be predicted using any of the above used algorithms. Linear Regression model has been used to compute the required errors:

**Mean Absolute Error = mean(abs(predictions - testSet\$price))**

**MAE = 2178.836**

**RMSE = 3403.056**



```

R Console
stroke           -501.94      271.39   -1.850 0.066744 .
compression.ratio -3287.84     1686.61   -1.949 0.053490 .
horsepower        24.77       874.81    0.028 0.977454
peak.rpm          847.07      304.64    2.781 0.006267 **
city.mpg          -409.56      966.08   -0.424 0.672340
highway.mpg       1094.08      951.00    1.150 0.252153
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2030 on 125 degrees of freedom
Multiple R-squared:  0.9516,    Adjusted R-squared:  0.9381 
F-statistic: 70.28 on 35 and 125 DF,  p-value: < 2.2e-16

> rmse(model_lm)
Error: could not find function "rmse"
> model_lm
Linear Regression

161 samples
 36 predictor

No pre-processing
Resampling: Bootstrapped (25 reps)
Summary of sample sizes: 161, 161, 161, 161, 161, 161, ...
Resampling results:

      RMSE      Rsquared
3403.056  0.8409535

Tuning parameter 'intercept' was held constant at a value of TRUE
> |

```

**Relative Absolute Error = (mean(abs(predictions- testSet\$price))) / (mean(abs(mean(testSet\$price)-testSet\$price)))**

**RAE = 0.3782024**

**Relative Squared Error = (mean((predictions- testSet\$price)^2))/(mean((mean(testSet\$price)-testSet\$price)^2))**

**RSE = 0.1866436**

**Coefficient of Determination(  $R^2$  ) = 0.8409535**