

# Prediction of the SGEMM kernel runtime

Kashish Kothari



# Data

- Data consist of 250000 instances of the runtime.
- It has 14 independent variable describing different parameters of SGEMM kernel.
- 4 independent variables are Numerical.
- 10 independent variables are categorical.
- Dependent variable is Run\_time in Millisecond which is numerical.

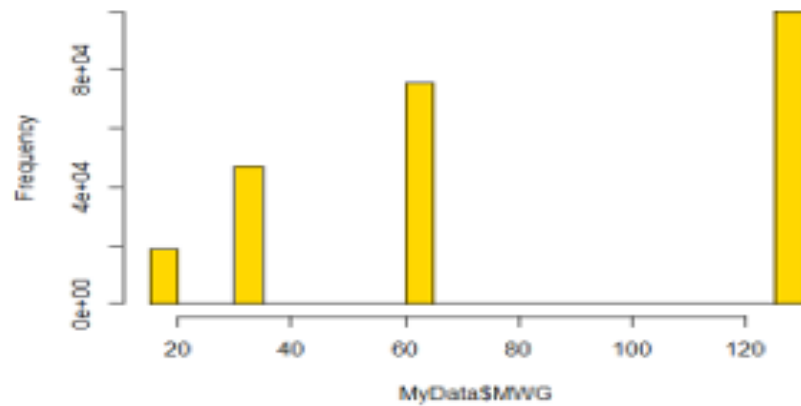


# Question

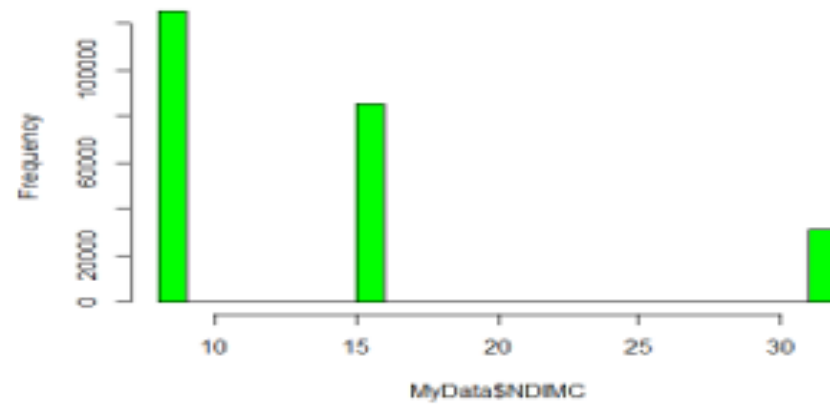
- My aim is to predict the Run time based on the given Variable values for the SGEMM Kernel
- So my question was to find out whether there is Linear Relationship between explanatory and dependent Variables.
- I wanted to select regression algorithm based on linear and non linear relationships in the data.
- Whether I should select parametric or non parametric approach.

# Exploratory Analysis

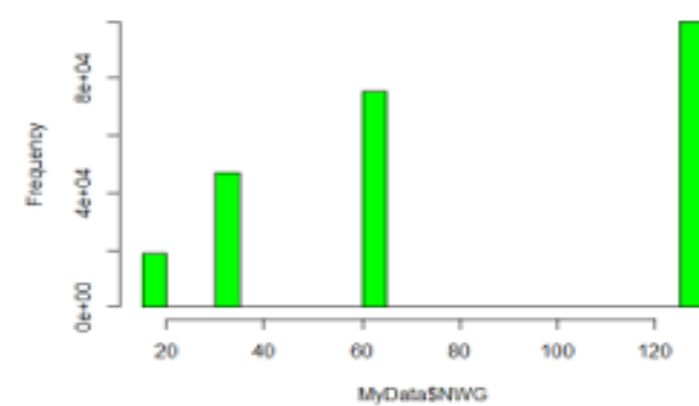
Histogram of MyData\$MWG



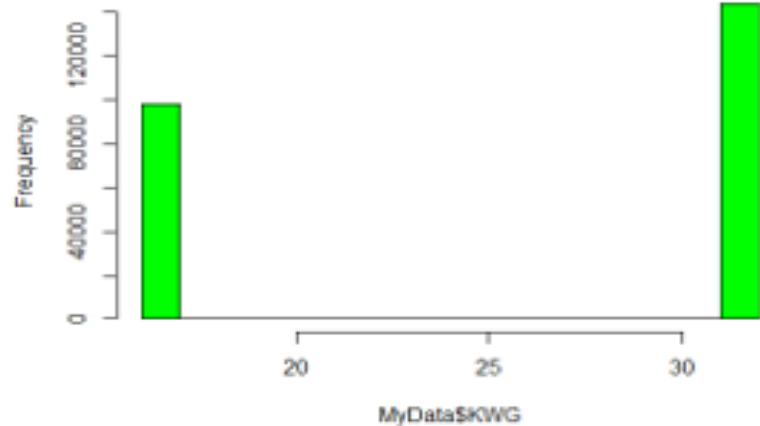
Histogram of MyData\$NDIMC



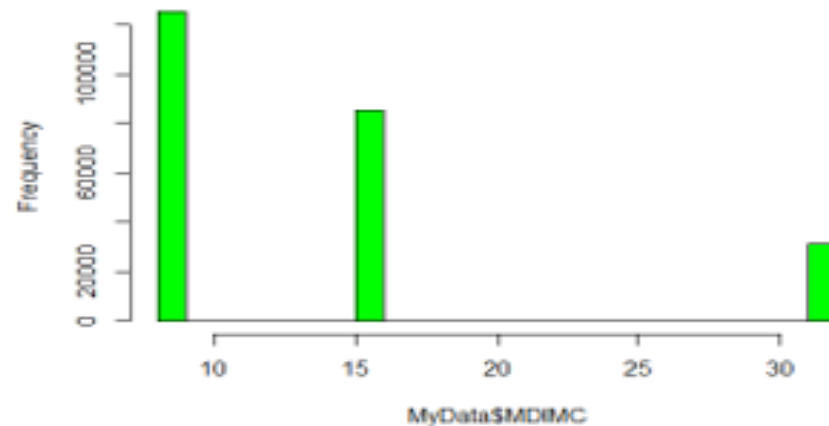
Histogram of MyData\$NWG



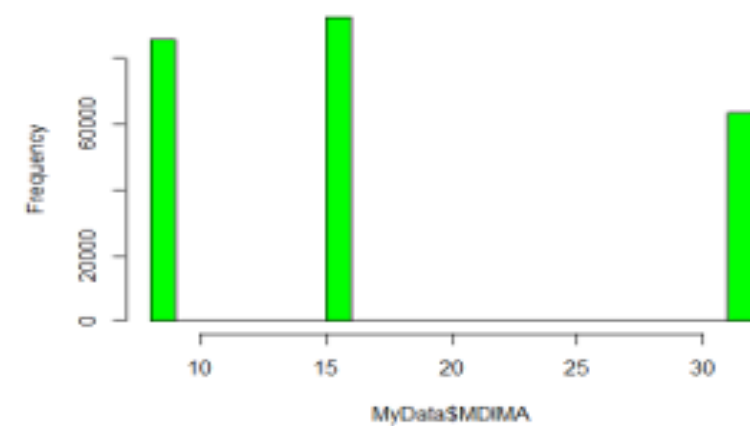
Histogram of MyData\$KWG



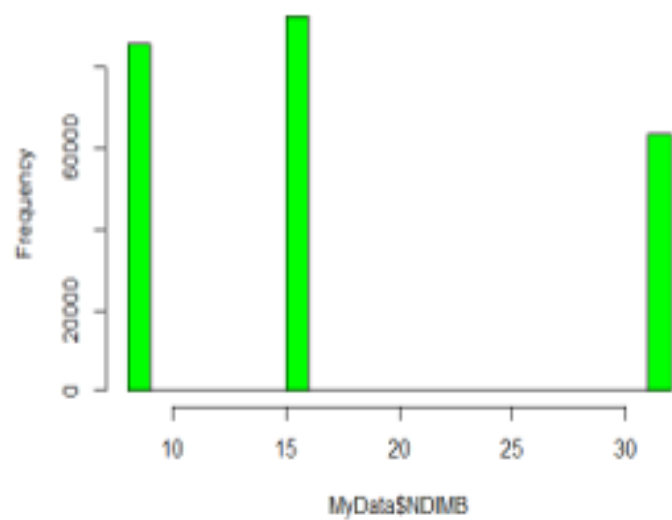
Histogram of MyData\$MDIMC



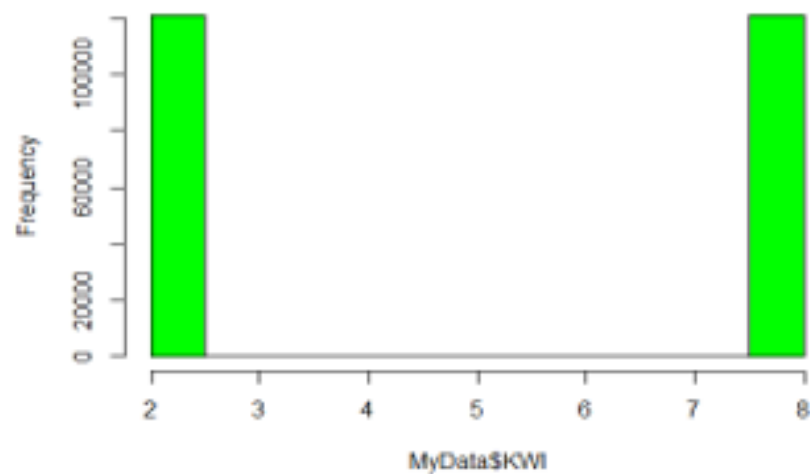
Histogram of MyData\$MDIMA



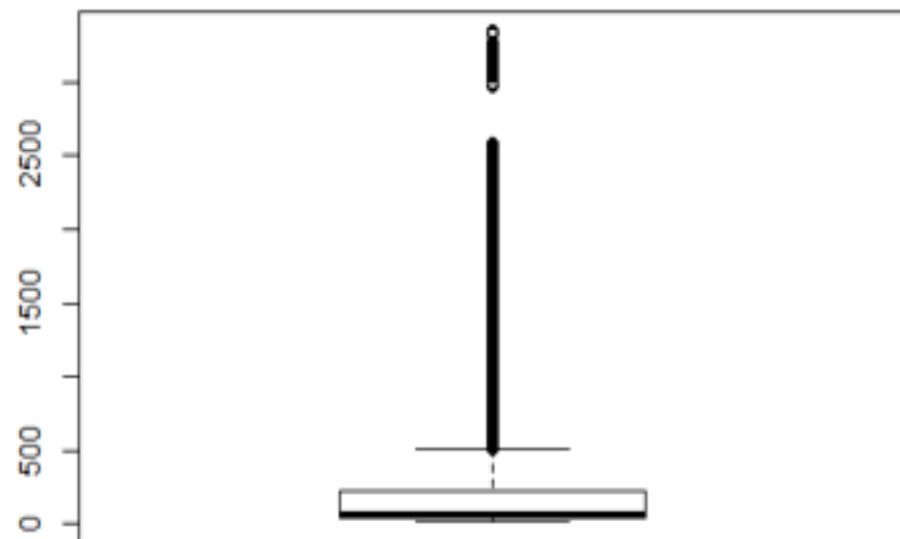
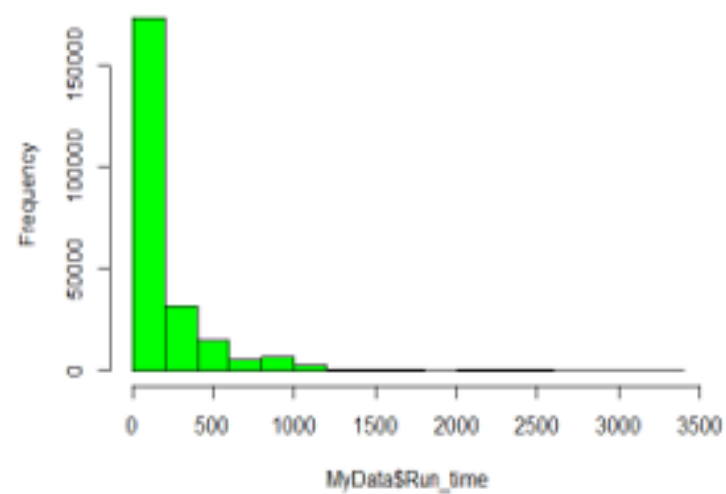
Histogram of MyData\$NDIMB



Histogram of MyData\$KWI



Histogram of MyData\$Run\_time

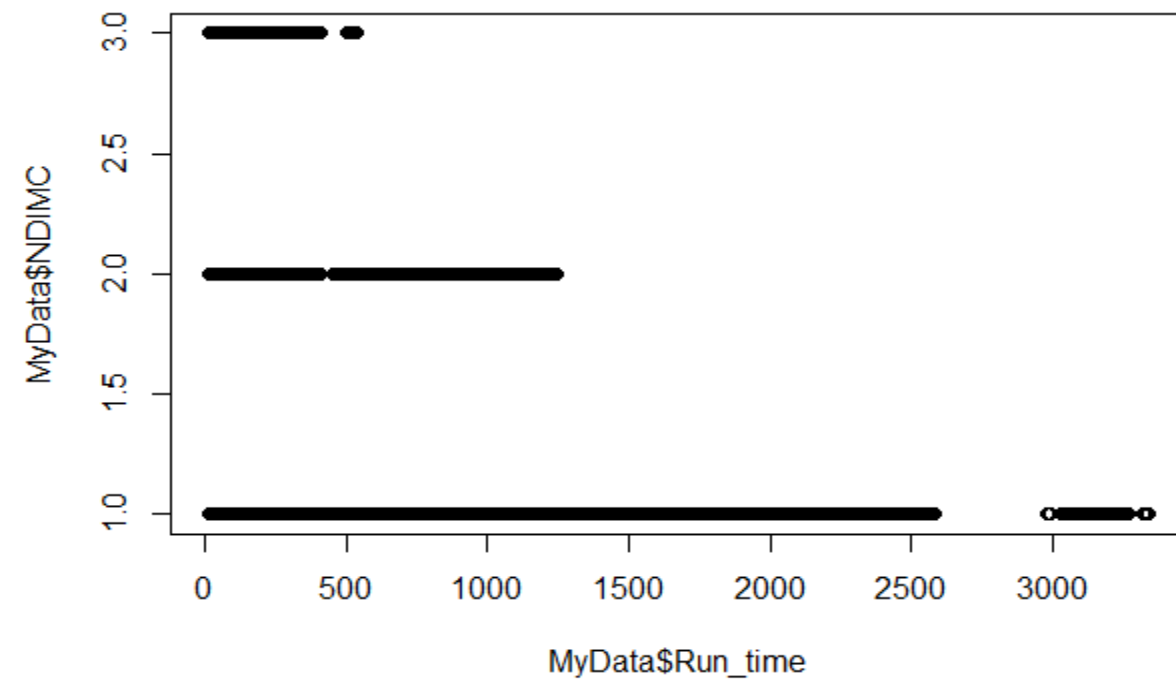
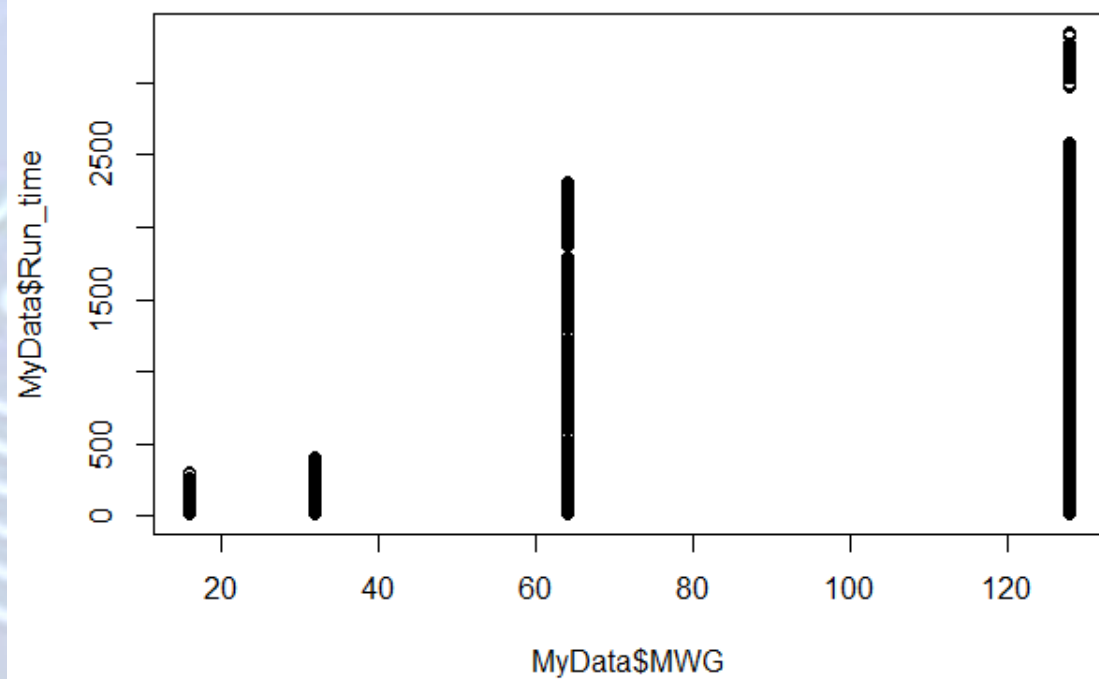




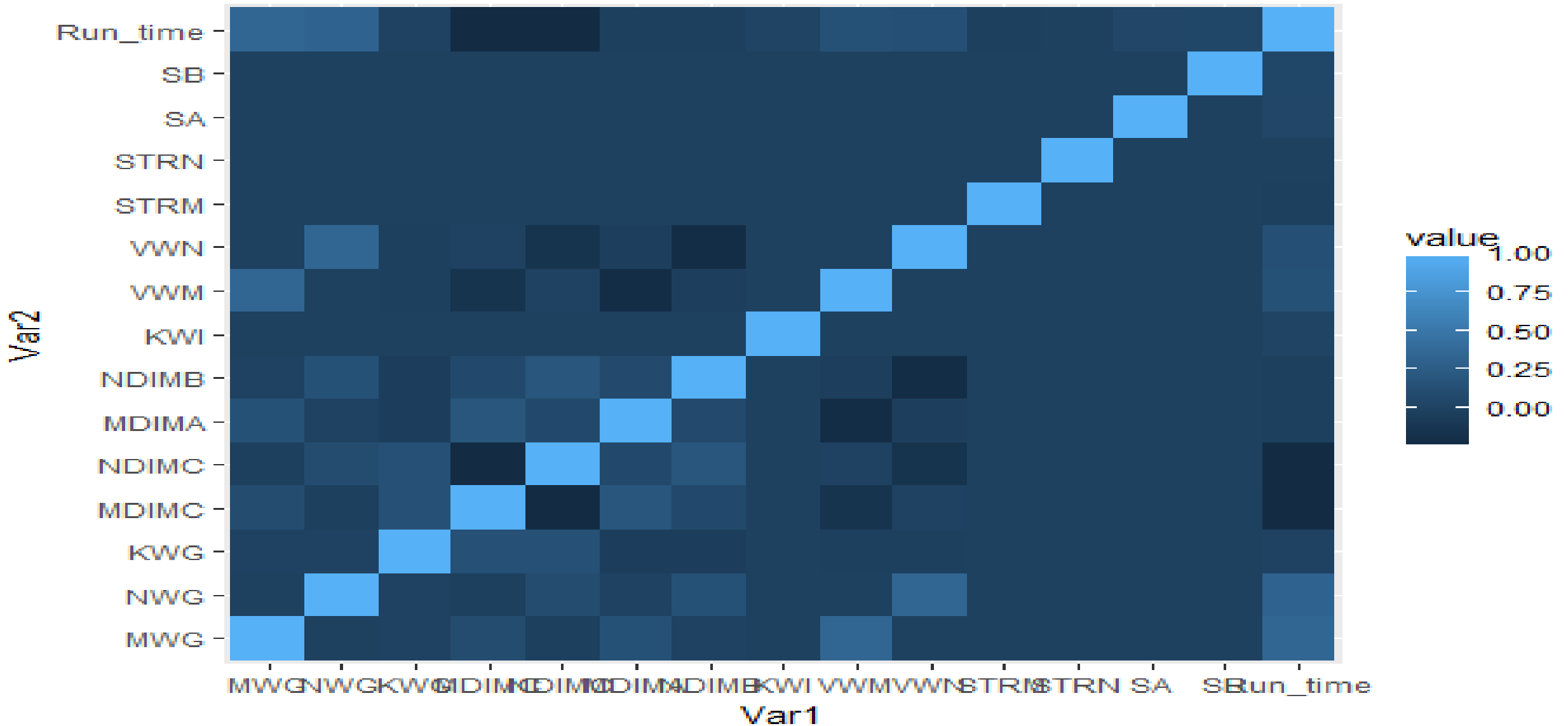
# Potential Problems to be checked before applying Linear Regression

- Non-linearity of the response-predictor relationships.
- Correlation of error terms.
- Non-constant variance of error terms.
- Outliers.
- High-leverage points.
- Collinearity.

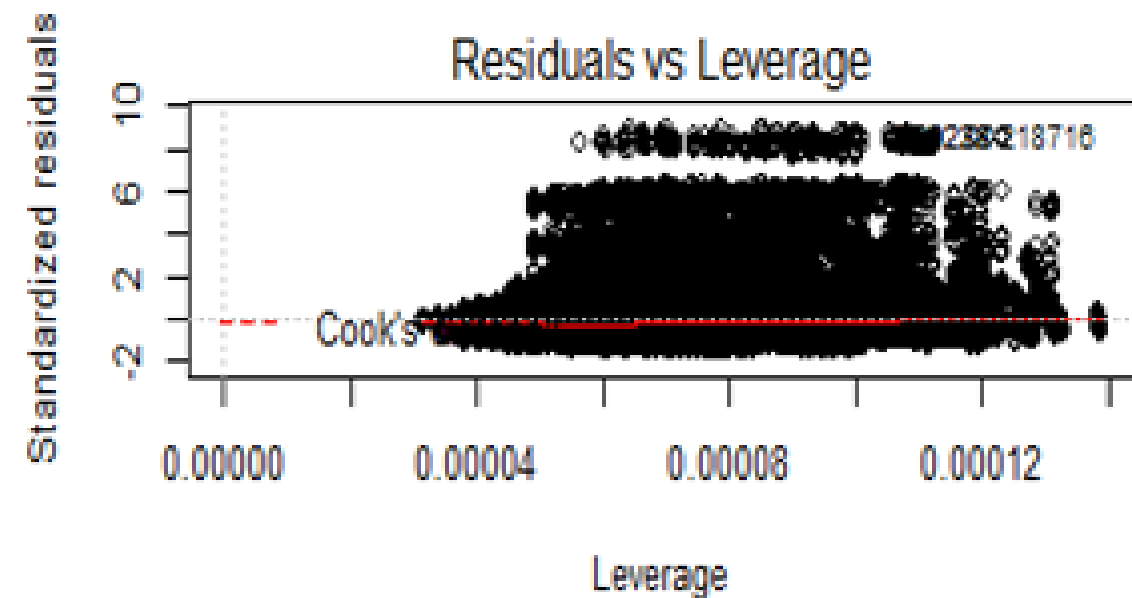
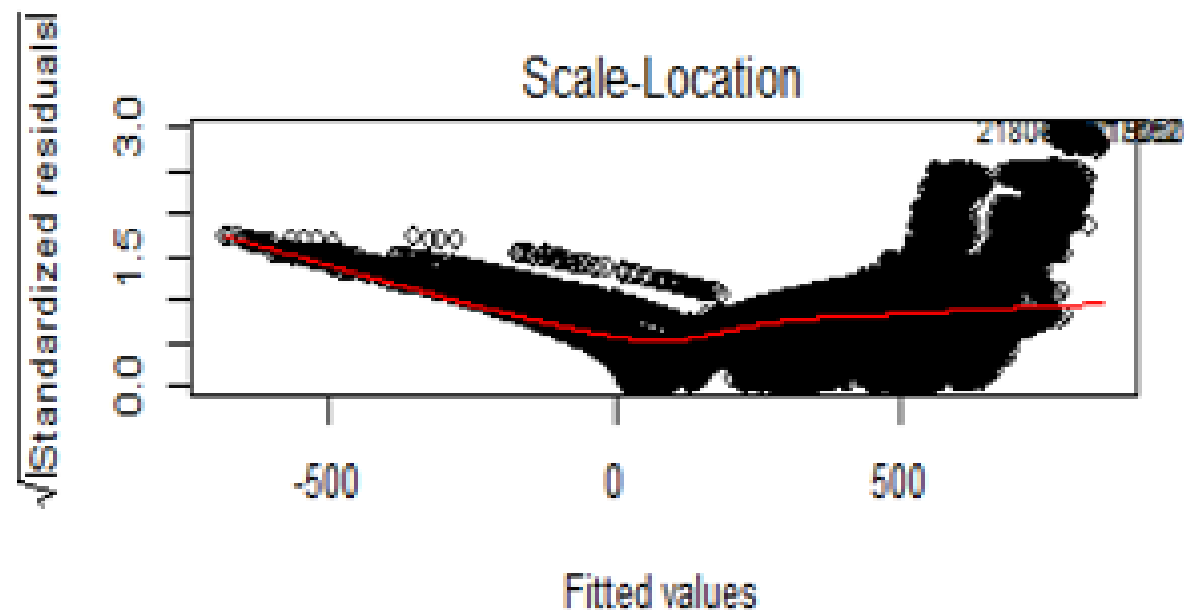
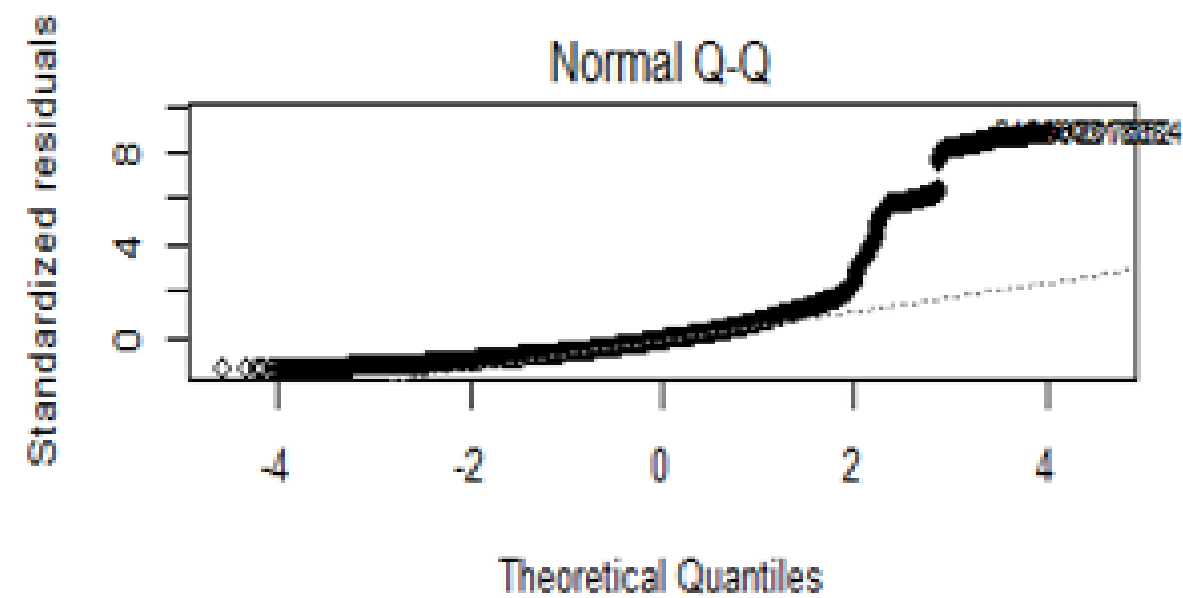
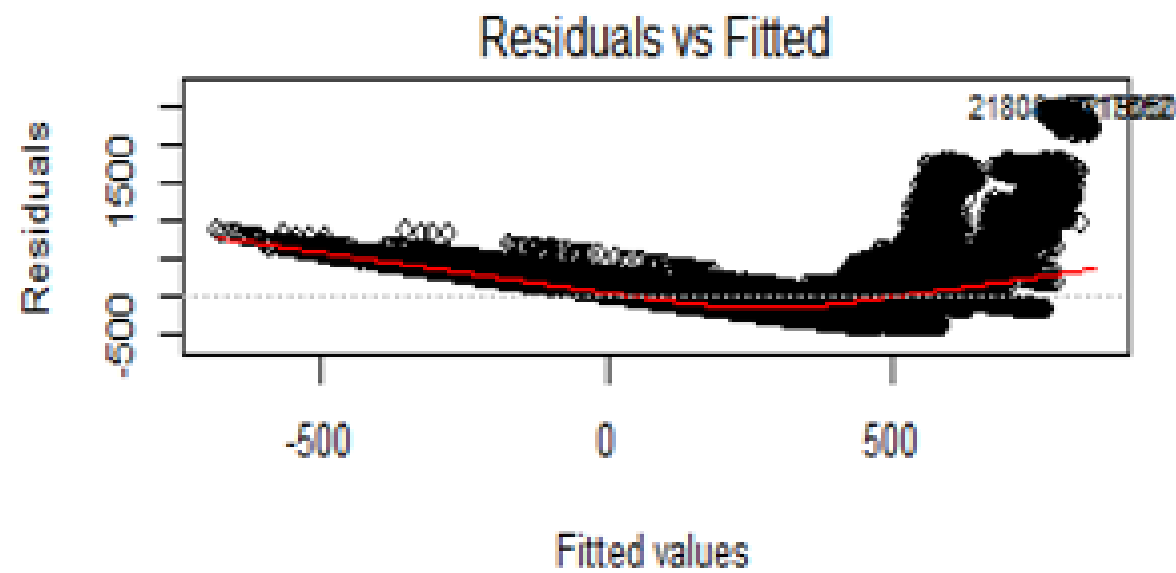
# Some Linear Association



# Checking for Collinearity







# What went Wrong?

- Linearity Assumptions were violated.
- Normal Distribution Conditions were Violated
- As there were nonlinear distribution in data, I opted for Non linear regression method such as GAM.
- I selected GAM because even though in non linear space it is easy for interpretation which is important.
- I applied smoothing spline to 4 numerical variables.



# Why GAM Did not work

- As majority of Independent variables are categorical and rest of the variables are also integers with only 4 unique values.
- Therefore predictors are not able to cover the variance of the predicted variable.
- As the distribution of the data is not known and majority of variables are categorical .



# Non Parametric Approach

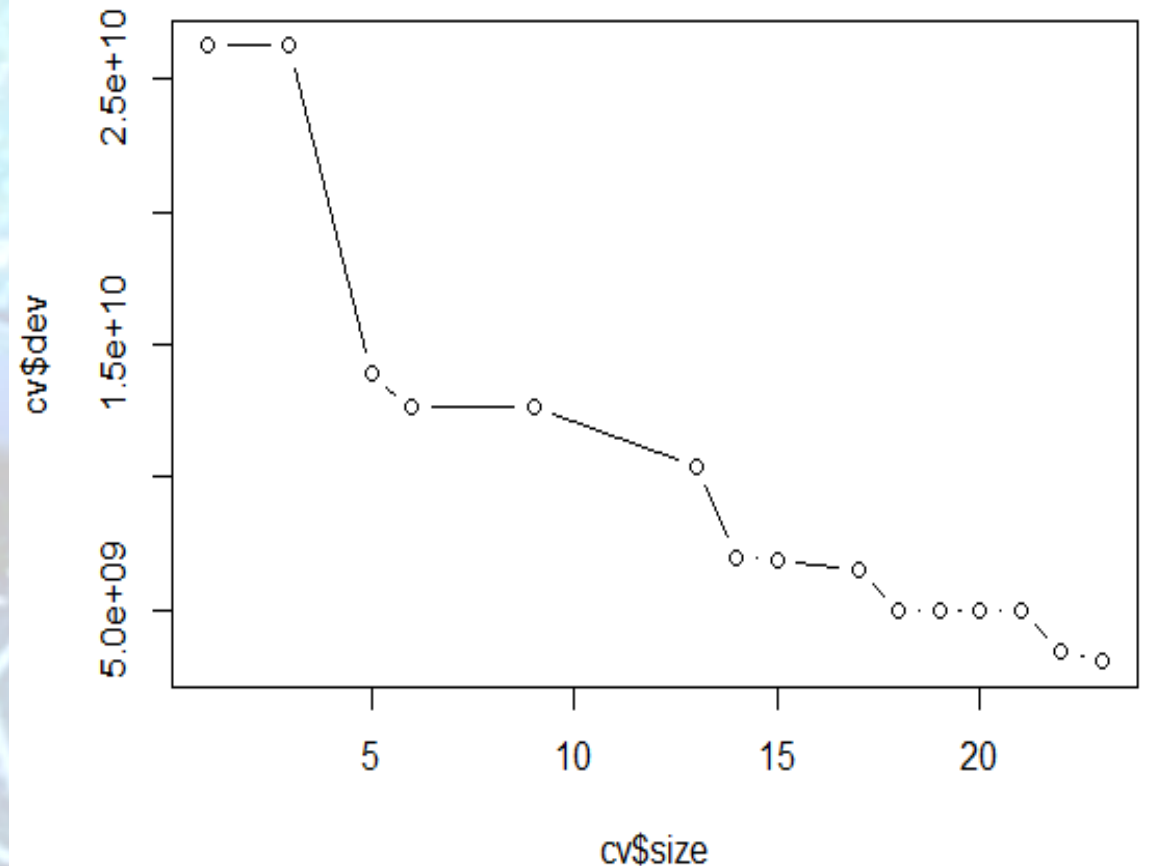
- I have selected Regression tree because I wanted to predict non linearity in the data.
- Regression tree is non linear model which can segregate predictor space into simple regions
- They are easy to interpret
- They generally work better when majority of independent variables are categorical
- There is less effort in data preparation and feature selection.
- They are robust to outliers and leverage points.

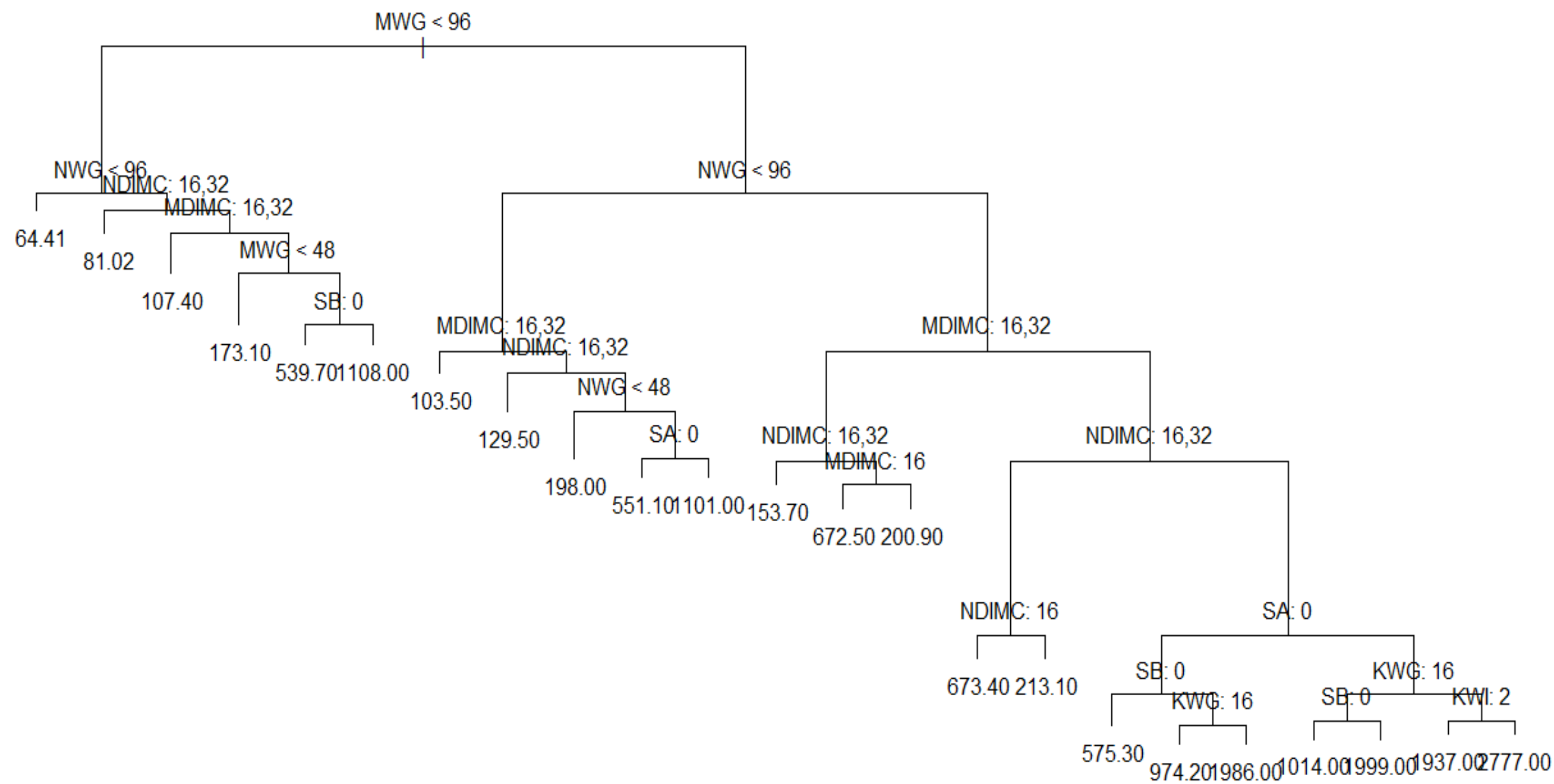


# Regression Tree

- I have used cross validation for selecting the effective size of the tree.
- Tree with 23 terminal nodes has lowest error .
- Therefore this tree was used for the prediction

```
Regression tree:  
tree(formula = Run_time ~ ., data = training)  
variables actually used in tree construction:  
[1] "MWG" "NWG" "NDIMC" "MDIMC" "SB" "SA" "KWG" "KWI"  
Number of terminal nodes: 23  
Residual mean deviance: 16270 = 3.145e+09 / 193300  
Distribution of residuals:  
      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.  
-1031.00  -50.19   -22.36    0.00   28.10  1272.00
```

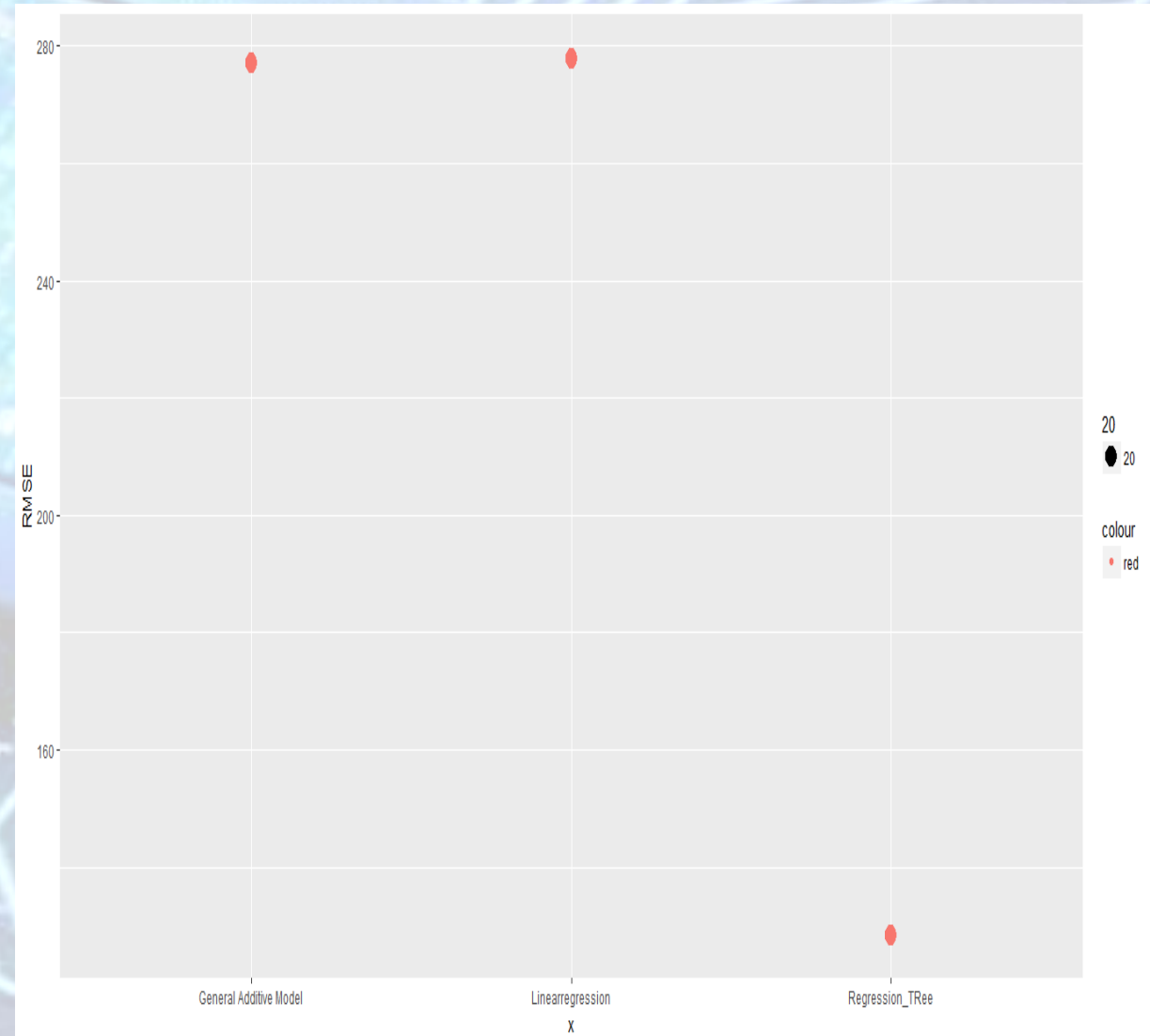






# Result of prediction on testing dataset

- Regression Tree performed better in this data with lowest RMSE compare to other two algorithm.
- There was no difference in the RMSE when I applied Linear Model and GAM model on test data.



A stylized illustration of a computer chip on a circuit board. The chip is a square component with a grid of pins on its sides, mounted on a blue circuit board with white traces. The top surface of the chip is a lighter blue and features a grid of binary code (0s and 1s) in a golden-yellow color. The text "Thank You!" is written in a black, sans-serif font across the center of the chip. The background is a soft, out-of-focus view of the circuit board's traces and other components.

Thank You!