

Regression Analysis Project Work on “ Boston House Data”

Summited by:-

1) **Abhilash Shingh**

2) **Ankit Gupta**

Project Work Duration:-Feb'2020-April'2020

INTRODUCTION

In this Project work we, we worked on BOSTON HOUSING data which represent the housing values in various suburbs of Boston. First, we perform Exploratory Data Analysis (EDA), in which we checked for missing value observation, Categorical variables and made Boxplots

In the Boston housing dataset, consider the variable 'MEDV' as the response, and the remaining variables, except 'ZN', 'CHAS' and 'RAD', as predictors. Standardize the regressors (not the response). Randomly split the data into 2 parts: one part containing 100 samples (say, X_{test} , y_{test} : test data), and the other part containing the remaining samples (say, X_{train} , y_{train} : training data).

First of all, We with the leverage points and outliers. We chose threshold according to data so too much points not get dropped. Now, compute Cook's distance, DFFIT, DFBETA and COVRATIO. Hence, find the influential points and delete those. Plot Cook's distance, DFFIT and COVRATIO.

Now we Deal with Curvatures. First check whether the residuals show any pattern when plotted against the individual regressors. Check the difference between Augmented Partial Residual (APR) and Component Plus Residual (CPR) plots. Then Choose a suitable transformation for removing pattern.

After that deal with heteroscedasticity and Normality which are assumption for linear regression model. For checking Heteroscedasticity Breush-Pagan test and check Non-Normality using QQ-plot. For removing Non-Normality apply Box-Cox transformation. And then finally apply Shapiro test for normality checking.

We compute the RMSE of model at every step. And at We fit a model on final data and compute test error.

Data Structure:

The Boston data set has 506 rows and 14 columns. This data set contains the following columns:

CRIM :- per capita crime rate by town.

ZN :- proportion of residential land zoned for lots over 25,000 sq.ft.

INDUS :- proportion of non-retail business acres per town.

CHAS :- Charles River dummy variable (= 1 if tract bounds river; 0 otherwise).

NOX :- nitrogen oxides concentration (parts per 10 million).

RM :- average number of rooms per dwelling.

AGE :- proportion of owner-occupied units built prior to 1940.

DIS :- weighted mean of distances to five Boston employment centres.

RAD :- index of accessibility to radial highways.

TAX :- full-value property-tax rate per 10,000.

PTRATIO :- pupil-teacher ratio by town. **BLACK** :- $1000(B_k - 0.63)^2$ where B_k is the proportion of blacks by town. **LSTAT** :- lower status of the population (percent).

We have these 13 Attributes as regressors(x_j). Which also have 12 Attributes as of continuous form and one is of binary valued Attribute(CHAS)

MEDV :- median value of owner-occupied homes(in 10000s dollars). This variable is taken as response variable (Y).

Dealing with Leverage point and influential points.

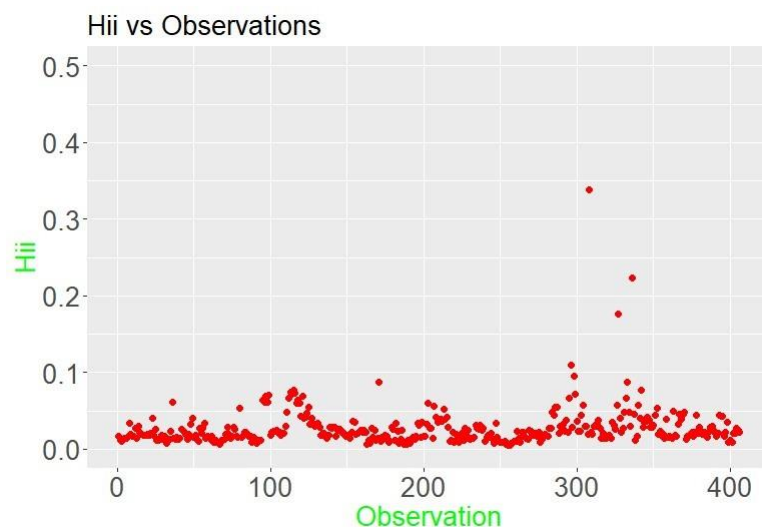
Leverage points

To find leverage point, we usually focus on the diagonal elements h_{ii} of the H matrix, which may be written as

$$h_{ii} = \mathbf{x}_i' (X'X)^{-1} \mathbf{x}_i$$

where x_i' is the i row of X matrix. The Hat matrix diagonal is a standardized measure of the distance of the i th observation from the centre of the \mathbf{x} space, thus the large hat diagonal reveals observations that are potentially influential because they are remote in \mathbf{x} space from the rest of the sample. We assume that any observation for which hat diagonal is greater than $\frac{2p}{n} + 0.2$ is considered a leverage point.

Graph btw h_{ii} and Observations are given below



For detecting influential points we will use Cook's Distance, DFBETA, DFFIT and COVRATIO.

Cook's Distance:

Formula for calculating Cook's Distance

$$D_i = \frac{(\hat{\beta}_i - \hat{\beta})' X' X (\hat{\beta}_i - \hat{\beta})}{pMS_{res}}$$

Point with large values of D_i have considerable influence on the least squares estimates $\hat{\beta}$. We usually consider points for which

$$D_i > 1$$

to be influential.

DFBETA:

The first of these is a statistics that indicates how much the regression coefficient $\hat{\beta}_j$ changes, in standard deviation units, if the i th observation is deleted. It is denoted by

$$\mathbf{DFBETAS}_{j,i} = \frac{\hat{\beta}_j - \hat{\beta}_{j(i)}}{\sqrt{S_i^2 C_{jj}}}$$

where C_{jj} is the j th diagonal element of $(X'X)^{-1}$ and $\hat{\beta}_{j(i)}$ is the j th regression coefficient computed without use of i th observation. A large value of $\mathbf{DFBETAS}_{j,i}$, indicates that observation i has considerable influence on the j th regressor coefficient. We use a cut off of $\frac{2}{\sqrt{n}}$ i.e., if $|\mathbf{DFBETAS}_{j,i}| > \frac{2}{\sqrt{n}}$ then i th observation warrants examination.

DFFITS:

We may also investigate the deletion influence of the i th observation on the predicted or fitted value. This is done by \mathbf{DFFITS}_i

$$\mathbf{DFFITS}_i = \sqrt{\frac{|\mathbf{DFBETAS}_i|}{1 - h_{ii}})} t_i$$

where h_{ii} is the diagonal elements of hat matrix and t_i is the R-Student. We will use the fact that any observation for which $|\mathbf{DFFITS}_i| > 2\sqrt{\frac{p}{n}}$ warrants attention, here p denotes the rank of hat matrix.

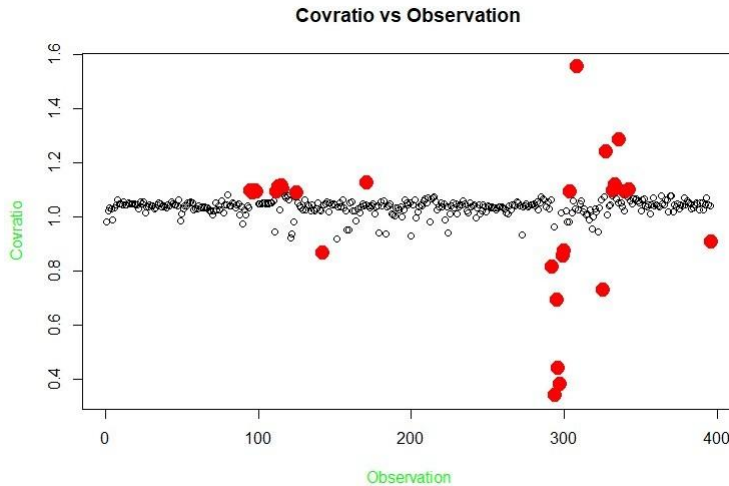
COVRATIO:

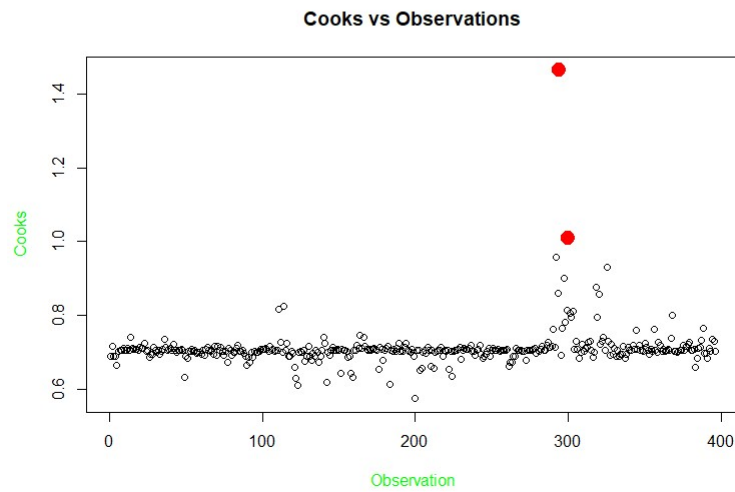
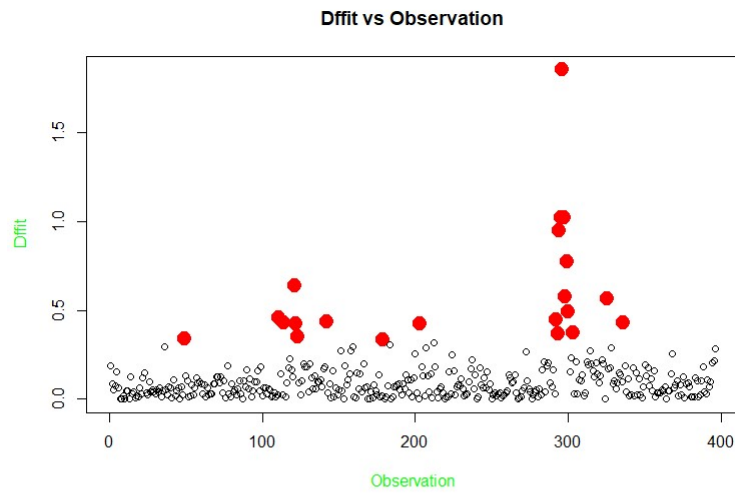
To express the role of the i th observation on the precision of estimation, we could define

$$\mathbf{COVRATIO}_i = \frac{(S_i^2)^p}{MS_{Res}^p} \left(\frac{1}{1 - h_{ii}} \right)$$

We will use fact that if $|\mathbf{COVRATIO}_i - 1| > 3\sqrt{\frac{p}{n}}$ then i th point should be considered as influential Point.

Now we plot COVRATIO , DFFITS , Cooks Distance

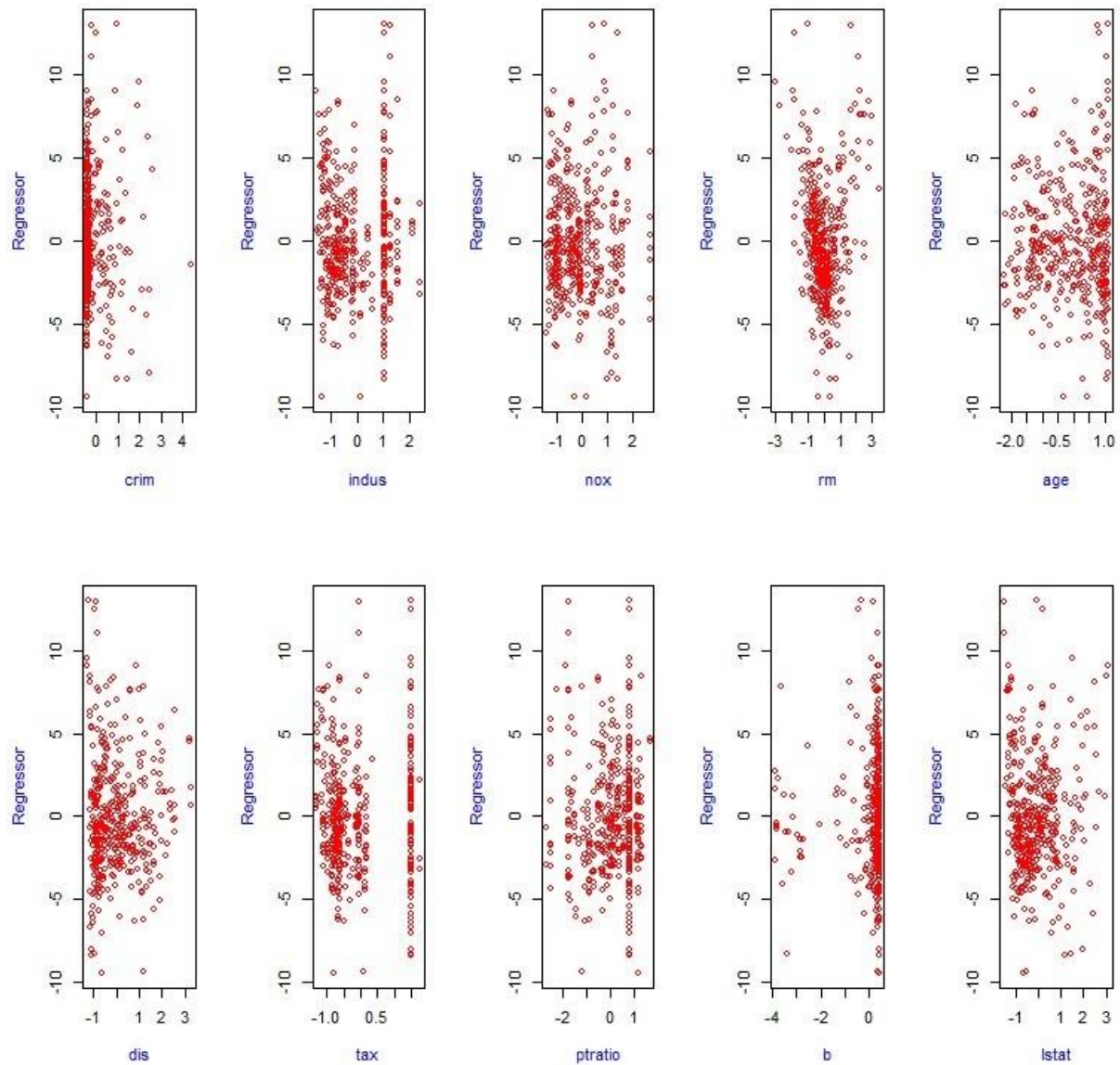




In this graph Red Points denotes the influential point.

Value of RMSE before deleting the leverage Point and influential Point is 4.549473 and after deleting the leverage point and influential point is 4.051916
Hence value of RMSE is decreases after deleting the influential point and leverage point.

Dealing with Curvatures

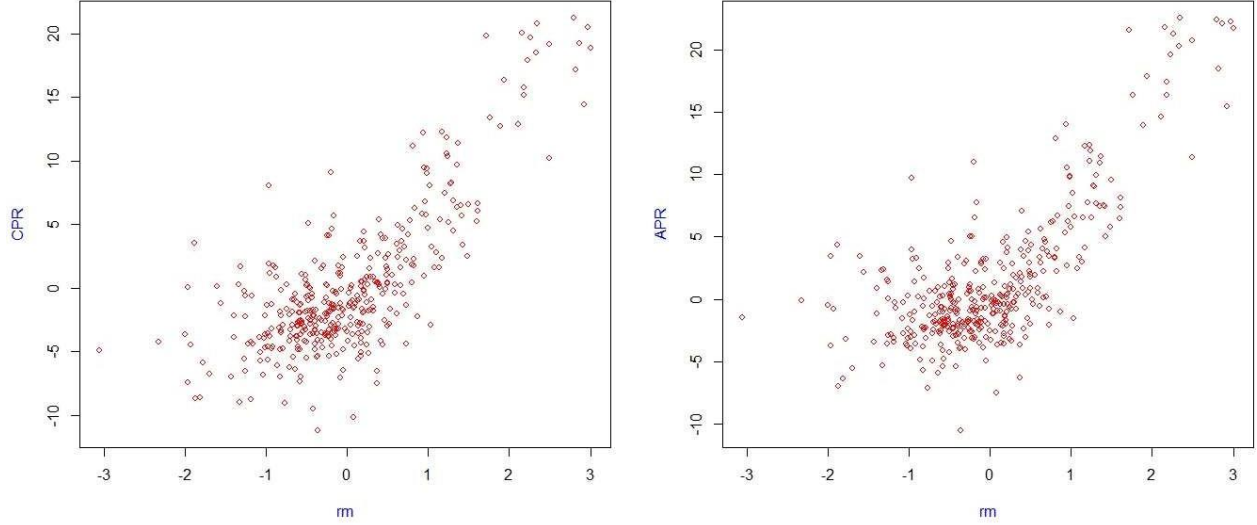


We first observe the residual vs regressor plot for each regressor.

Fig 1: Residual plots

The residual plot of rm regressor shows some U-shaped pattern. We claim that rm may not be linearly related with the response.

We now plot the CPR and APR plot of regressor(attached at the last of the pdf). There is no such major difference between APR and CPR plots of regressors ,However a close scrutiny reveals that there is some difference for regressor rm (attached below)



We now observe the CPR plot of rm to choose suitable transformation.

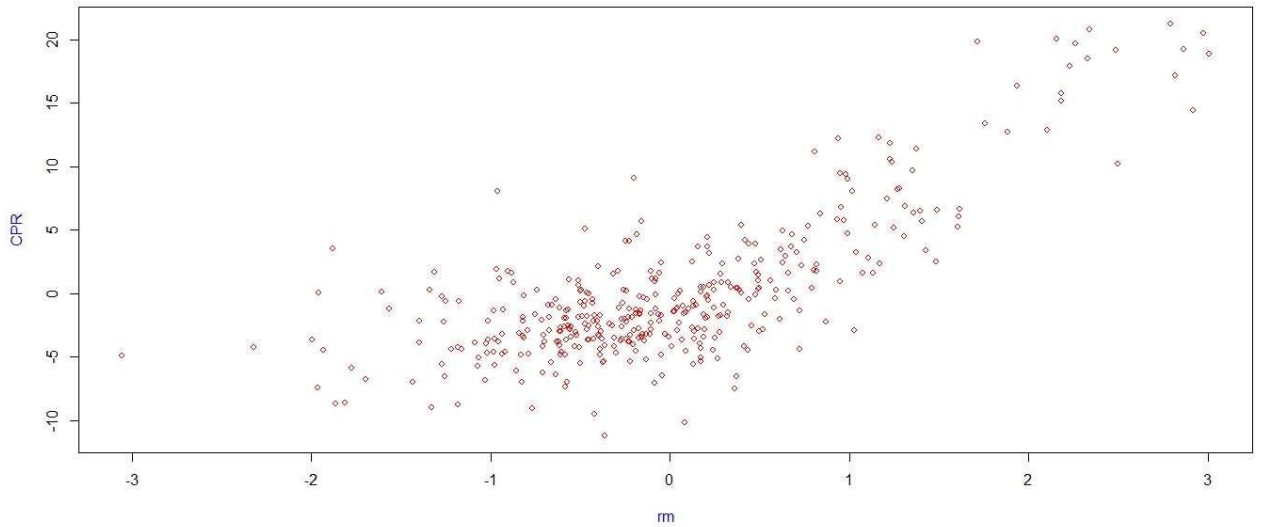


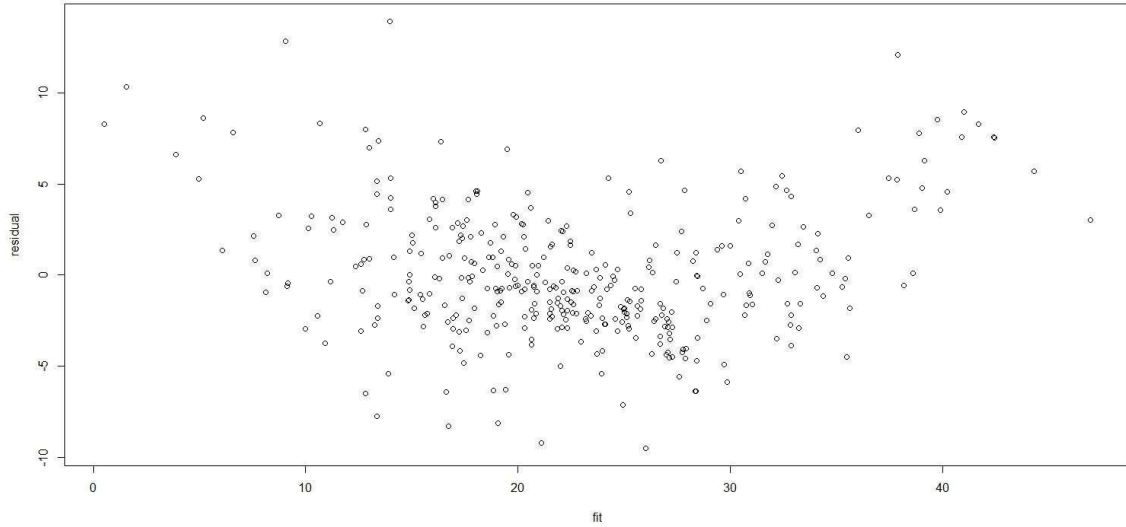
Fig 1: CPR plot of rm

Observing the graph we use the transformation $g(x) = \beta_0 e^{\beta_1 x_{rm}}$. Taking log on both sides of the model we get $\log(y_{Train}) = \log(\beta_0) + \beta_1 x_{Train}$, which is a linear model. We now estimate β_0 and β_1 by regressing $\log(y_{Train} + c)$ on $\log(\beta_0) + \beta_1 x_{rm}$. We get $\beta_0 = 20.831$ and $\beta_1 = 0.2978335$.

We now calculate correlation coefficient between y_{Train}, x_{rm} and $y_{Train}, g(x)$
Here $corr(y_{Train}, x_{rm}) = 0.7922$ and $corr(y_{Train}, g(x)) = 0.82$. As the correlation creases we use this transformation on the regresor rm and fit the regression model with the transformed rm. We now calculate the RMSE again. The RMSE comes out to be 15.72489, which is higher than the RMSE obtained without transformation. As the RMSE increases we do not take this transformation and keep the rm regressor as it was.

Dealing with the heteroscedasticity

We first observe the plot below:



From the graph of residual vs fitted value we see that the graph is U-shaped. Again from the graph of residual vs regressor in 2(a) we see that the graph of regressor 'rm' shows a U-shaped pattern.

Only one regressor shows some pattern. So, we choose $d_i^* = \alpha_0 + \alpha_1 Z_1 + \epsilon^*$
Now we choose: $Z_1 = C_0 + C_1 X + C_2 X^2$ where x corresponds to 'rm'

$$= \frac{cov^2(d_i^*, d_i^*)}{var(\hat{d}_i^*) var(d_i^*)}$$

we reject H_0 at 5% level of significance if $Q_{obs} \geq \chi_{0.05,1}^2$. Here $Q = 66.36$ which is greater than $\chi_{0.05,1}^2 = 3.84$. So we reject H_0 .

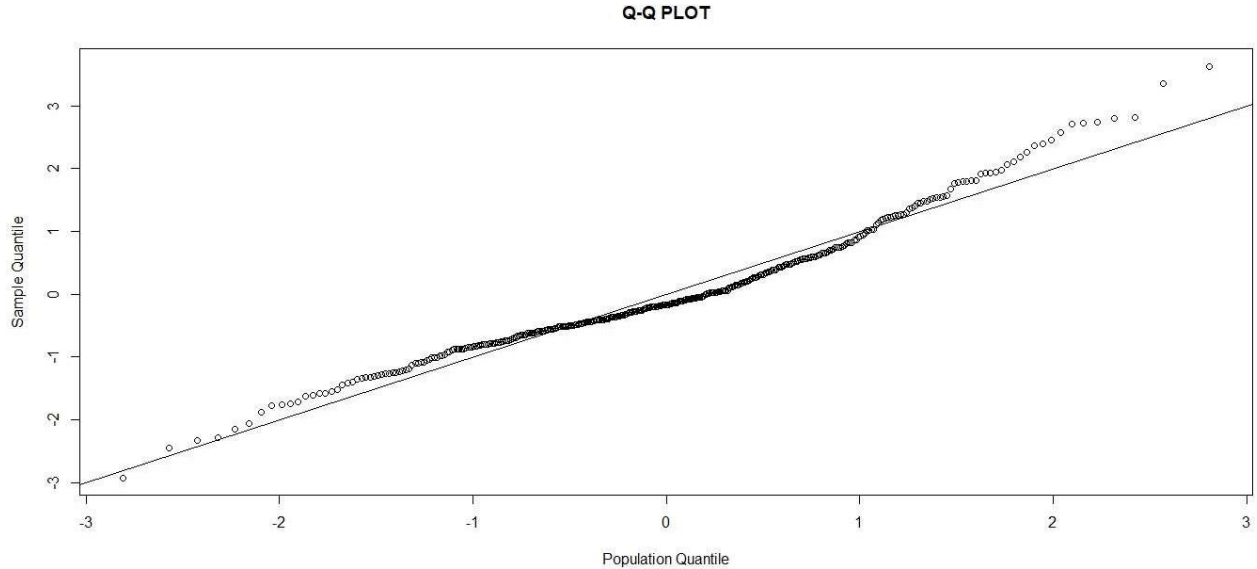
Here we take: $\sigma_i^2 = \alpha_0 + \alpha_1 Z_1$

Let Σ be the dispersion matrix with diagonals as $(\sigma_1^2, \dots, \sigma_n^2)$. Here we see that same as d_i^* . We find $\Sigma = \text{diag}(\sigma_1^2, \dots, \sigma_n^2)$. So, $\hat{\beta} = (X^T \Sigma^{-1} X)^{-1} X^T \Sigma^{-1} Y$. So be changed and hence d_i^* . So, $\hat{\alpha}$ will be changed. We will continue updating $\hat{\alpha}$ and till there is significant difference between two consecutive values of $\hat{\alpha}$ and $\hat{\beta}$.

We now calculate RMSE with the updated β after the process converges. It is found that the RMSE comes out to be 5.013884, which is greater than the previous one. So, we suspect that our chosen $h(a, \alpha, \beta)$ is not appropriate.

Checking For normality and Box-Cox Transformation

Now we draw QQ-plot of the R-student residuals against it's population distribution. Here we consider R-student residual quantile as sample quantile and y-training data-quantile as population quantile.



From this graph, it is visible that distribution of y-training data is deviated from our normality assumption.

4.b To fix the problem, we do box-cox transformation to normalize the Ytrain data.

Box-Cox Transformation:

Formula of Box-Cox Transformation

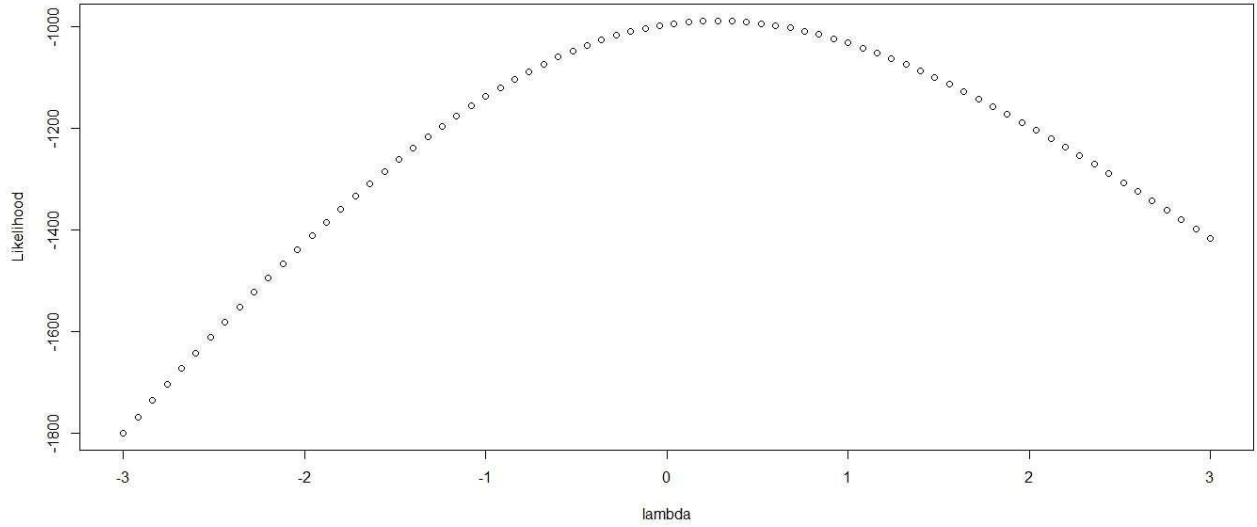
$$y(\lambda) = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0; \\ \log y, & \text{if } \lambda = 0. \end{cases}$$

To estimate the parameters of the box-cox transformation, we need the value of λ for which the expression :

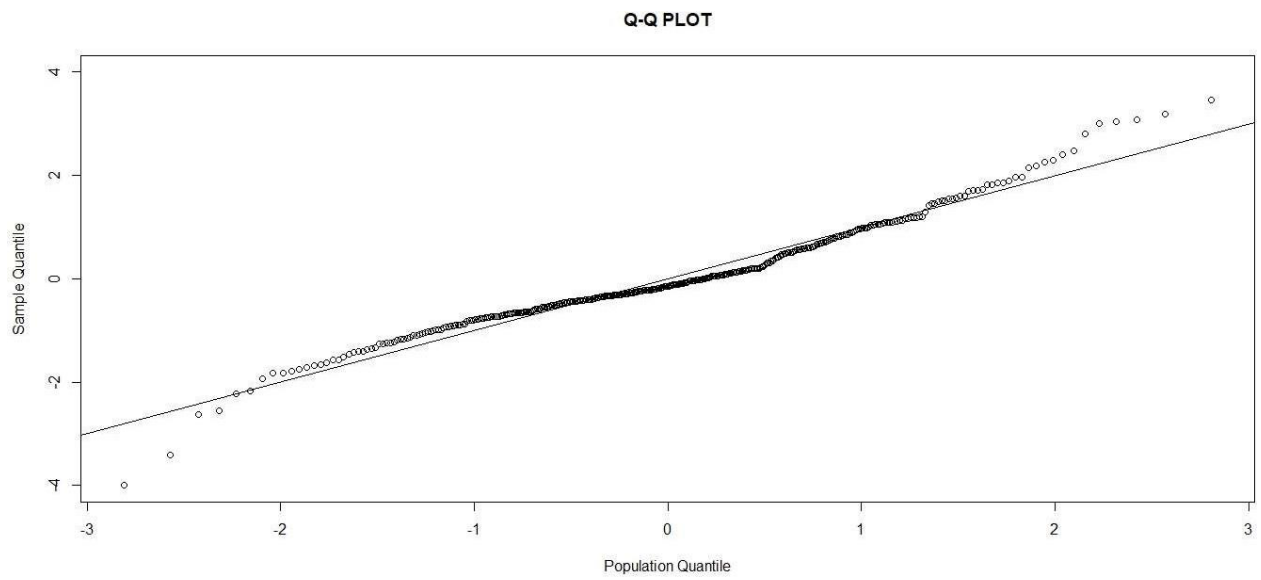
$$L(\lambda, \hat{\beta}_{MLE}^\lambda, \hat{\sigma}_{MLE}^{2,\lambda}) = -\frac{n \log 2\pi}{2} - \frac{n \log \frac{RSS_\lambda}{n}}{2} - \frac{n}{2} + n(\lambda - 1) \log(G)$$

is maximized.

We compute the likelihood function for different λ values and plot it against λ to find the maximum value of likelihood function



Now we transform the y-training data through box-cox transformation and again fit regression model with the transformed y and recompute the R-student residuals. Again we draw QQ-plot of the new R-student residuals against it's population distribution. We obtain the graph as follows :



We observe from the graph that, the transformed data now more close to the normality assumption.