

Boston House Price Analysis
(A Regression Analysis on Boston House Price Data Using R)

Group Members

Avik Biswas

Rasamanjari Nandan

Shamriddha De

Som Paul

Sreya Sarkar

Report

The present analysis deals with predicting the median value of a home (**MEDV**) using 10 predictors, *viz.*, **CRIM**, **INDUS**, **NOX**, **RM**, **AGE**, **DIS**, **TAX**, **PTRATIO**, **B**, **LSTAT**, which have been described in the Boston Housing Dataset. For an effective analysis, we have standardized the 10 regressors. Further, we have divided the entire dataset containing 506 observations into a test set of 100 observations selected at random, and a training set of remaining 406 observations.

To the training data, we have fit a linear regression model of the form

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (0.1)$$

through the method of ordinary least squares. Using this regression equation, we have predicted the values of the response MEDV for the test set and tried to determine the extent to which these values agree with the actual values provided to us. The measure used for this purpose is called Root Mean Squared Error, mathematically expressed as

$$RMSE = ||\mathbf{y}_{test} - \hat{\mathbf{y}}_{test}||/\sqrt{n}. \quad (0.2)$$

For the present pair of training and test data, the value of RMSE has been obtained as **4.6894**.

In the instant problem, we are essentially interested to check whether certain assumptions of OLS indeed hold or not for the model.

1 Dealing with Leverage Points and Outliers

At first, we shall try to identify and subsequently discard the leverage points in the training data. To that effect, we consider the diagonal elements h_{ii} 's of the projection matrix $X(X^T X)^{-1}X^T$. An observation would be declared a leverage point if its corresponding value of h_{ii} exceeds a certain threshold. Theoretically, a cutoff of $2p/n$ serves the purpose, where p is the number of regression coefficients in the linear regression model (in this case, $p = 11$, including the intercept term) and n is the number of observations (to begin with, we have $n = 406$). For the given problem, we obtain a cutoff equal to 0.0542, which suspects 34 leverage points. However, this number seems to be large. In order to avoid too many detections of leverage points, we increase the cutoff by 0.04. In fact, a plot of h_{ii} versus i (see Figure 1.1) suggests that the cutoff 0.0542 is indeed quite conservative, and an enhancement of the same by 0.025 may seem appropriate. Besides, if we consider the RMSE after excluding the 34 observations, we observe that it increases to 4.9206. Accordingly, we deem fit to proceed with a threshold of $(0.0542 + 0.04)$, that is, **0.0942**. Based on this threshold, we would discard 5 observations (which is small enough), marginally reducing the RMSE to **4.6805**. The number of observations consequently becomes 401.

Once we have dropped the leverage points, our interest shifts to detection of influential points. Towards that, we would employ four measures - Cook's Distance, *DFFITs*, *COVRATIO* and *DFBETAS*, whose working formulae and theoretical cutoffs are mentioned below.

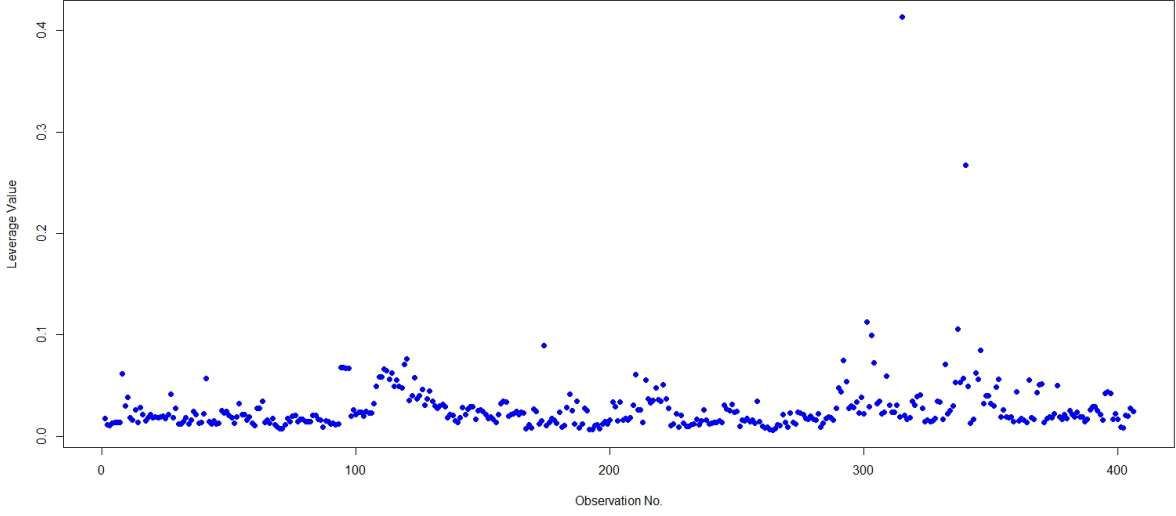


Figure 1.1: Plot of Leverage Values (h_{ii})

- **Cook's Distance:**

$$D_i = \frac{\|\hat{\mathbf{y}}_{(i)} - \hat{\mathbf{y}}_i\|^2}{pMS_{Res}} \quad (1.1)$$

Here, $\hat{\mathbf{y}}_{(i)}$ denotes the predicted response vector obtained by fitting the regression model excluding the i -th observation. Typically, $D_i > 1$ suspects the i -th observation to be influential.

- **DFFITS:**

$$DFFITS_i = \frac{\hat{y}_{(i)} - \hat{y}_i}{\sqrt{h_{ii}S_{(i)}^2}} \quad (1.2)$$

In this case, $\hat{y}_{(i)}$ is the predicted value of y_i based on the model excluding the i -th observation, while $S_{(i)}^2$ indicates the MSE of that model. Mostly, the i -th observation may be detected as an influential point if the absolute value of $DFFITS_i$ exceeds $2\sqrt{p/n}$. Accordingly, for the given problem, the cutoff becomes 0.3312.

- **COVRATIO:**

$$COVRATIO_i = \left(\frac{S_{(i)}^2}{MS_{Res}} \right)^p \left(\frac{1}{1 - h_{ii}} \right) \quad (1.3)$$

In most of the situations, an observation is termed as influential if $|COVRATIO - 1| > 2p/n$. That is, in our case, a value of $COVRATIO$ outside the interval $[0.9177, 1.0823]$ might be enough for a point to be influential.

- **DFBETAS:**

$$DFBETAS_{j,i} = \frac{\hat{\beta}_j - \hat{\beta}_{(i)j}}{\sqrt{C_{jj}S_{(i)}^2}}, \quad (1.4)$$

where $C = (X^T X)^{-1}$, and $\hat{\beta}_{(i)}$ is the estimate of the vector of regression coefficients computed by excluding the i -th observation. A cutoff of $2/\sqrt{n}$ for $|DFBETAS_{j,i}|$ generally serves our purpose, and the i -th observation may be highly influential if the value exceeds the cutoff. In view of the value of n in our case, the cutoff turns out to be 0.0999.

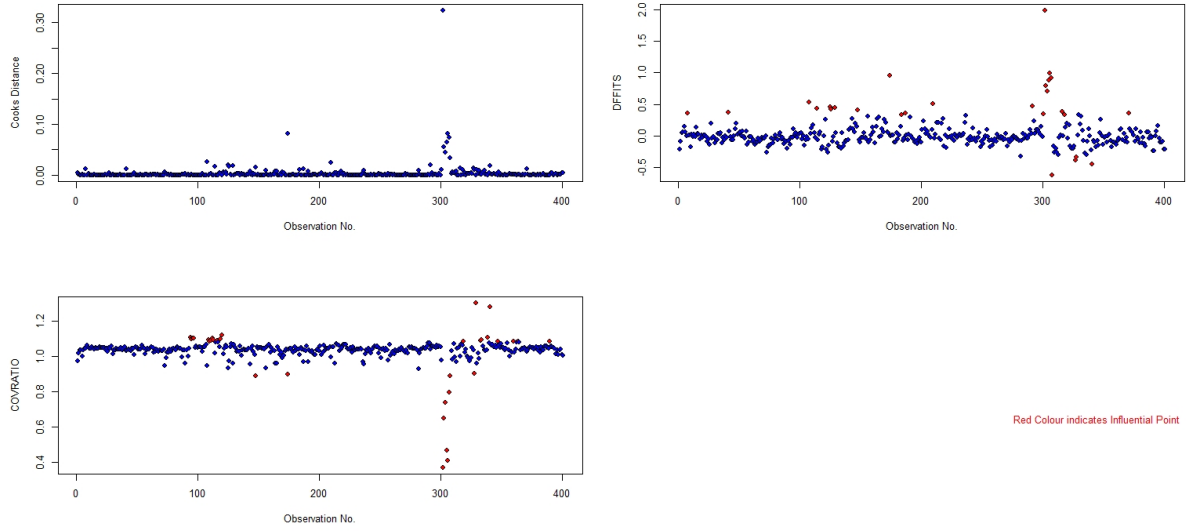


Figure 1.2: Plots of Cook's Distance, DFFITS and COVRATIO, with influential points being suspected using usual cutoffs

In this regard, we would take note of a certain fact about DFBETAS. It can be observed that there are 11 DFBETAS values corresponding to each observation. For our consideration, if at least 60% of those exceed the aforesaid cutoff for DFBETAS, the observation is suspected to be influential. Further, for an overall conclusion, we decide a criterion that a point would be discarded as influential if at least three of the measures suspect it to be so. It must be noted that there may be other plausible criteria. However, we have chosen our benchmark because of its simplicity.

Analyzing the given data yields no suspected influential points using Cook's Distance and 6 using DFBETAS. On the other hand, DFFITS and COVRATIO suspect 27 and 32 points to be influential. Certainly, these figures are quite large. Although only 5 points are declared as influential in overall, the large number of suspicions through DFFITS and COVRATIO warrant to be checked. Besides, through the plots (see Figure 1.2) of Cook's Distance, DFFITS and COVRATIO, we find that that the cutoffs for DFFITS and COVRATIO deserve leniency. In view of this, and also taking help from the plots, we modify the cutoffs as mentioned below, followed by the plots of these measures (Figure 1.3).

- Cook's Distance: 1 (unchanged)
- DFFITS: 0.8312 (increased by 0.5)
- COVRATIO: [0.8177, 1.1823] (higher and lower cutoffs are respectively increased and decreased by 0.1)
- DFBETAS: 0.0999 (unchanged)

Resultantly, the Cook's Distance, DFFITS, COVRATIO and DFBETAS suspect 0, 5, 8 and 6 observations respectively as influential. Observably, these figures are corroborative. In the totality, we would discard 4 observations, reducing the number of observations n to 397. Looking into the RMSE, we observe that it has further decreased to **4.5846**. Accordingly, we can wrap up the analysis of influential points.

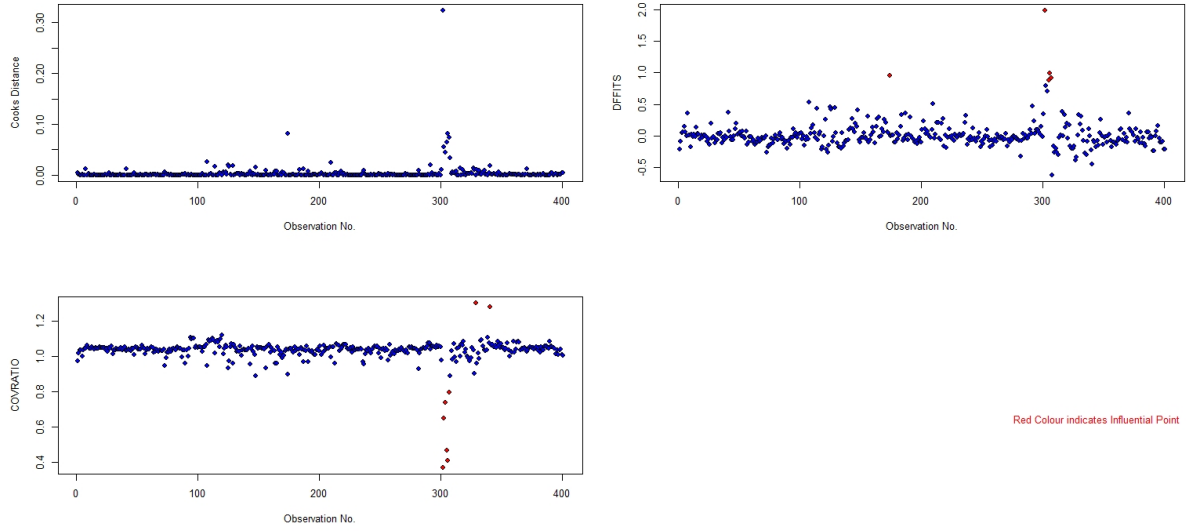


Figure 1.3: Plots of Cook's Distance, DFFITS and COVRATIO classifying the observations as influential or not using the modified cutoffs

2 Dealing with Non-Linearity

After dropping the unusual observations in the previous section, we now proceed to analyze the linearity of relationship between the response and the regressors. In this regard, we first plot the residuals against the individual regressors (Figure 2.1).

It can be observed that there is a U-shaped pattern in the residual plot corresponding to the regressor **RM**, which roughly indicates a possible presence of non-linearity. A slightly decreasing pattern is visible in the residual plot against **DIS** as well. The remaining plots seem to be devoid of any pattern.

To substantiate the presumptions based on residual plots, we take Component plus Residual (CPR) and Augmented Partial Residual (APR) plots into consideration. At this juncture, let us briefly revisit the theoretical aspects of these plots. In CPR plot, we plot $(e_{y|X} + \hat{\beta}_j x_j)$, the portion of response explained linearly by the regressor x_j only, against x_j . The APR plot, on the other hand, depicts the portion of response explained by RM only through a quadratic function. To clarify more, the APR diagram plots $(e_{y|X^*} + \hat{\beta}_j x_j + \hat{\gamma}_{jj} x_j^2)$ against RM, where X^* is obtained by augmenting x_j^2 as a column to the original design matrix X . In case the response y is linearly related to x_j , the CPR and APR plots should not differ significantly.

In the present scenario (see Figure 2.2), there appears to be a marked difference between the two plots for RM and a marginal difference for LSTAT. For each of the other predictors, including DIS, the APR and CPR plots look fairly similar. As such, we may assume that non-linearity, if any, is caused due to LSTAT and RM only.

To get an overview of the nature of non-linear relationship, we focus on the partial residual plots corresponding to RM and LSTAT in Figure 2.3. In this case, we at first obtain the partial residual vector $(e_{y|X_{-j}})$ by regressing the response on all the predictors except x_j . Next, we regress x_j on the remaining predictors to obtain another set of residuals given by $e_{x_j|X_{-j}}$. The partial residual diagram plots $e_{y|X_{-j}}$ against $e_{x_j|X_{-j}}$ to indicate the actual nature of relationship between y and x_j . This is so because the two sets of residuals ultimately filter out the effects of other regressors from both y and x_j .

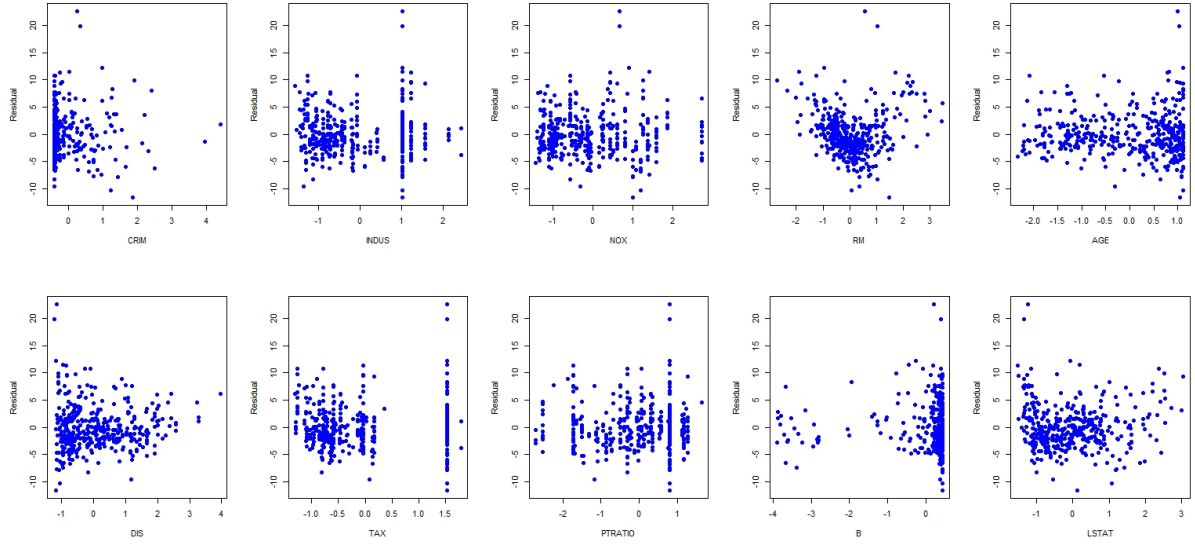


Figure 2.1: Residual Plots (against regressors)

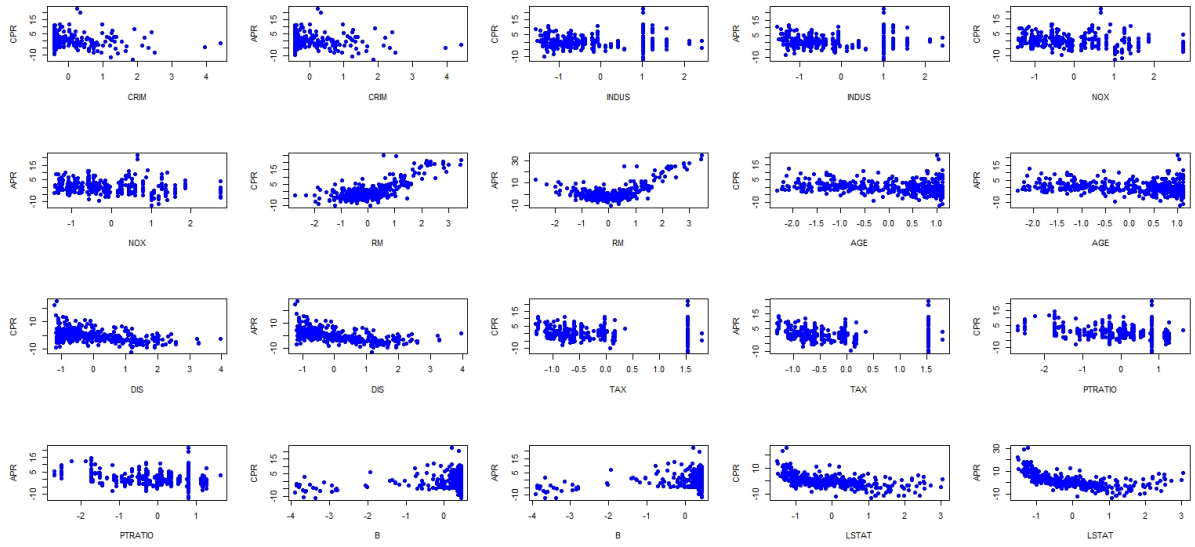
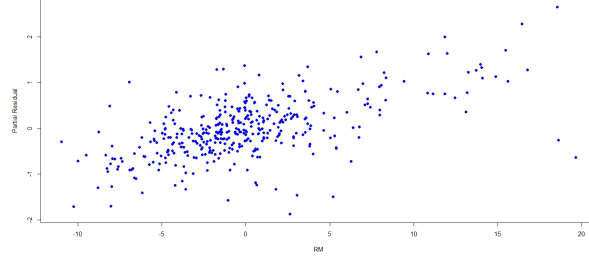
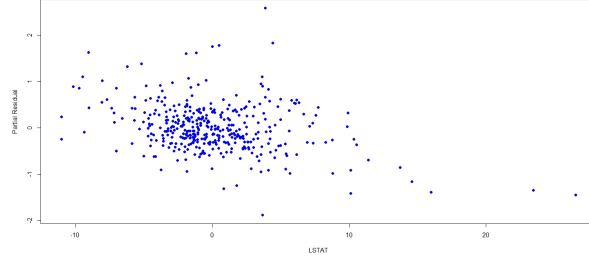


Figure 2.2: APR and CPR Plots for each regressor



(a) For RM



(b) For LSTAT

Figure 2.3: Partial Residual Plots for RM and LSTAT

Once we are done with the plotting, we observe that the non-linearity is not quite evident between RM and MEDV. Therefore, we do not deem fit to transform RM. The regressor LSTAT, on the other hand, seems to affect MEDV through a decreasing non-linear function which somewhat mimics an exponential function of the type b^x . We accordingly run a model of the nature

$$y_i = \gamma^{x_i} \delta_i, \quad i = 1, 2, \dots, n,$$

through the data, taking \mathbf{y} as MEDV, \mathbf{x} as LSTAT and δ_i to be a multiplicative error term. Taking logarithm on both sides, this model can be reconstructed as a linear regression model of $\log(\text{MEDV})$ on LSTAT as follows, presuming that the model satisfies the usual assumptions.

$$\mathbf{y}^* = \gamma^* \mathbf{x} + \boldsymbol{\delta}^*$$

We obtain the OLS estimate of $\hat{\gamma}^*$ as -0.517, and accordingly, $\hat{\gamma} = \exp \hat{\gamma}^* = 0.596$. Resultantly, LSTAT gets transformed into $\log(\text{LSTAT})$.

Let us check whether this transformation at all improves the linearity. Before transforming, correlation coefficient between MEDV and LSTAT was -0.756, while upon transformation, value turns out to be 0.816. Thus, there is a slight increase in the magnitude of correlation coefficient, which suggests that the extent of linearity improves marginally (but certainly not significantly) after transformation.

Further, let us compute the RMSE to check whether the transformation results in an overall decrease in RMSE. The value of RMSE is obtained as 13.5845, which is almost thrice the RMSE obtained before the transformation. Hence we refrain from considering the transformed variable in the Model (0.1) and retain the regressor LSTAT in its original form.

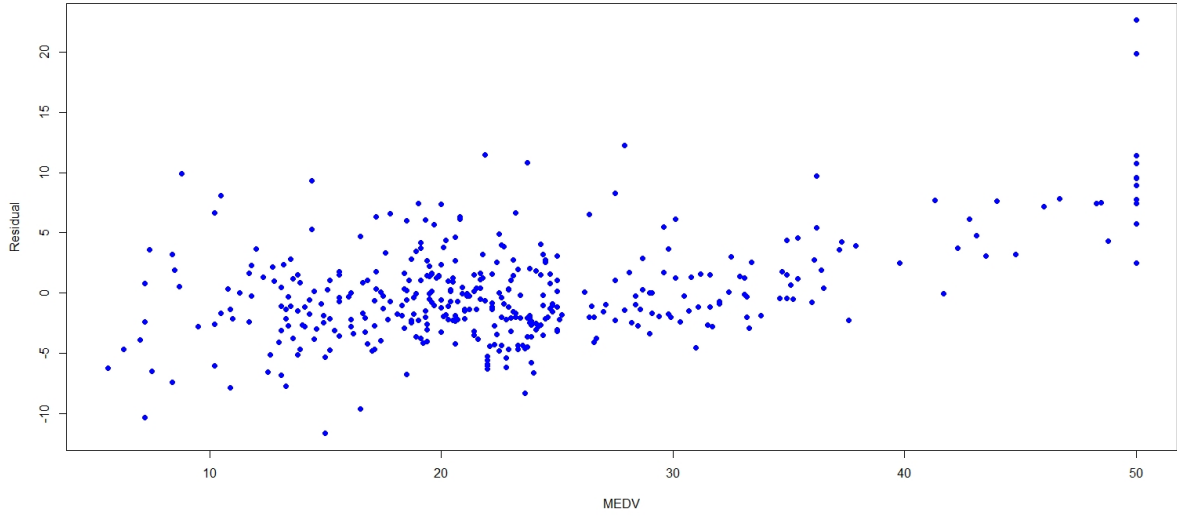


Figure 3.1: Residual Plots (against fitted values)

3 Dealing with Heteroscedasticity

The focus now shifts to the assumption of homoscedasticity of errors. We begin by plotting the residuals against the fitted values of the responses (Figure 3.1), as well as against the individual regressors (see Figure 2.1 in the previous section).

In the residual plot against fitted values, we find a nearly linear pattern. Besides, in the previous section, we have already mentioned that patterns are visible in residual plots against RM, DIS and LSTAT. Suffice to say, heteroscedasticity may have crept into the data. To look deeper into this aspect, let us try to identify the regressors with respect to which the variance (approximated by the square of residual e_i^2) is not constant. On plotting e_i^2 against the regressors (Figure 3.2), we observe that the estimated variances vary with respect to CRIM, AGE, DIS and B. Accordingly, these predictors are likely to cause heteroscedasticity.

To affirm our contentions, we carry out the Breusch-Pagan test to check the validity of the null hypothesis of homoscedasticity. At first, we compute $d_i^* = ne_i^2 / \sum_i e_i^2$. Next we find the appropriate variables z_1, z_2, \dots, z_k (which are nothing but some of the regressors, or their suitable functions) for which the following relation approximately holds,

$$d_i^* = \alpha_0 + \alpha_1 z_{i1} + \dots + \alpha_k z_{ik} + \epsilon_i^*, \quad (3.1)$$

with ϵ_i^* 's being i.i.d. standard normal variates. Here, the appropriate z_j 's are CRIM, AGE, DIS and B.

The test statistic is accordingly computed using

$$Q = nR_z^2 = n \frac{\text{cov}^2(d_i^*, \hat{d}_i^*)}{\text{var}(d_i^*)\text{var}(\hat{d}_i^*)}, \quad (3.2)$$

which is distributed (asymptotically) as χ_k^2 under the null hypothesis. The null hypothesis of homoscedasticity is rejected at level α if the observed value of Q exceeds $\chi_{\alpha,k}^2$. In our case, we obtain $Q_{obs} = 17.83$, which becomes large enough to reject the null hypothesis. Hence we can restate the regression model in

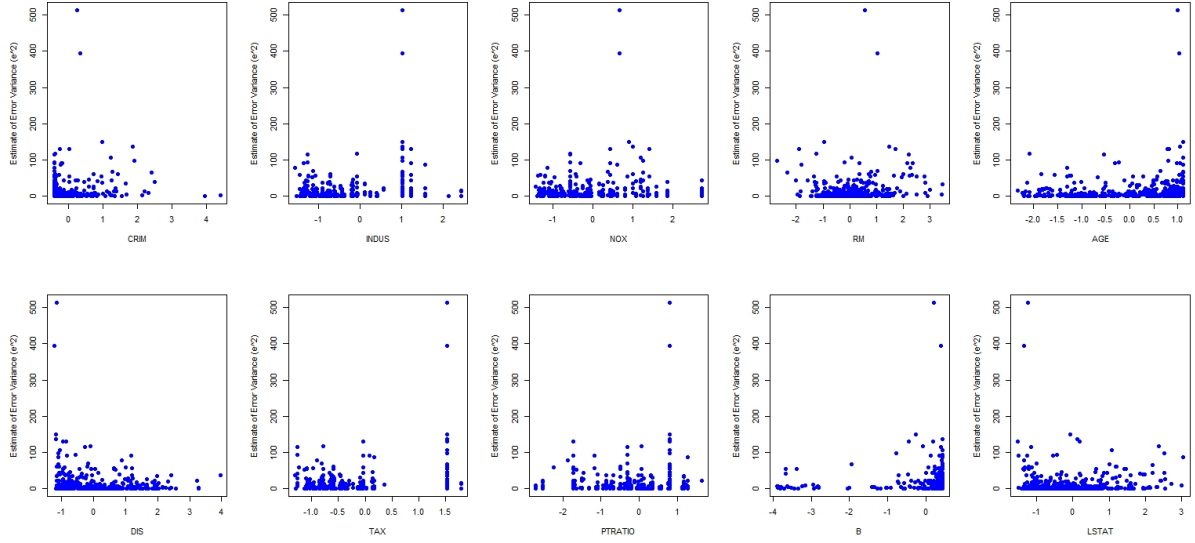


Figure 3.2: Plotting residual squares against regressors to identify the regressors causing heteroscedasticity

the form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (3.3)$$

in which $\boldsymbol{\epsilon} \sim N_n(\mathbf{0}, \sigma^2 \boldsymbol{\Sigma})$, $\boldsymbol{\Sigma}$ being a positive definite diagonal matrix with at least two unequal diagonal elements. Needless to say, $\boldsymbol{\Sigma}$ is unknown as well.

Our next task is to estimate the regression coefficient $\boldsymbol{\beta}$ in the presence of heteroscedasticity. We at first consider $h(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\beta})$ to be of the nature

$$h(\mathbf{Z}, \boldsymbol{\alpha}, \boldsymbol{\beta}) = \exp(\alpha_0 + \sum_{j=1}^k \alpha_j z_{ij}). \quad (3.4)$$

. We must note that such a choice of h would ensure its values to be positive. Further, taking logarithm on both sides of (3.4) yields a linear function on the right hand side, which is quite simple to handle. Following this, we take recourse to the prescribed algorithm with a slight change. To update $\hat{\alpha}_{(I)}$, we minimize $S = \sum_i (\log e_i^2 - \log h(\mathbf{Z}, \hat{\boldsymbol{\alpha}}_{(I-1)}, \hat{\boldsymbol{\beta}}_{(I-1)}))^2$. The change occurs due to the fact that e_i^2 has been used to approximate σ_i^2 which is further estimated through the function h . Subsequently, we update $\hat{\boldsymbol{\beta}}_{(I)}$ and $\hat{\boldsymbol{\Sigma}}_{(I)}$ as mentioned.

Once the algorithm converges, we compute the RMSE based on the new estimate of $\boldsymbol{\beta}$ and obtain the value to be **4.7464**, which is slightly larger than the one obtained prior to this section (the previous value was 4.5846). Thus the nascent estimate of regression coefficient, although takes into account the heteroscedasticity, marginally compromises with the prediction performance. In view of this, it would not be appropriate to transform the original model (0.1) to get rid of the heteroscedasticity which appears to have crept in.

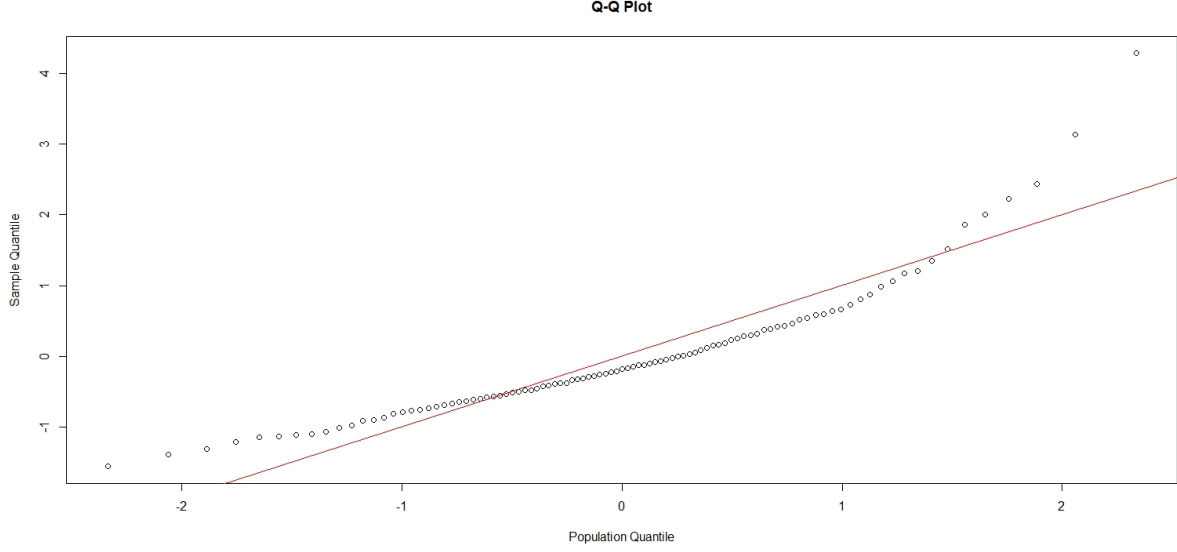


Figure 4.1: QQ-plot based on R-student residuals of original responses

4 Dealing with Non-Normality

The final aspect which shall be dealt in the present situation is the validity of the normality assumption of the error (and subsequently the response) in Model (0.1). A graphical idea of such may be fathomed through a QQ-plot (quantile-quantile plot), which plots sample quantiles against population quantiles. If the underlying distribution is symmetric, we can expect the points in the plot to lie on a straight line passing through the origin with unit slope. For our purpose, we would consider the sample quantiles of the R-student residuals with respect to the corresponding population quantiles in the QQ-plot. Mathematically, R-student residuals are defined as

$$r_i^* = \frac{e_i}{\sqrt{S_{(i)}^2(1 - h_{ii})}}, \quad (4.1)$$

which, under usual regression assumptions, follows a t-distribution with $(n - p)$ degrees of freedom. Therefore, in the QQ-plot, we consider the population quantiles of a $t_{(n-p)}$ distribution on the horizontal axis.

The resulting QQ-plot is shown in Figure 4.1. We observe that the points largely deviate from the desired straight line, thereby implying a departure of errors from normality. More precisely, the points form a curve which is concave upward, indicating a positive skew in the underlying distribution of the R-student residuals. To normalize the data, we undertake Box-Cox transformation, which is given by

$$u_\lambda = \begin{cases} \frac{y^\lambda - 1}{\lambda}, & \text{if } \lambda \neq 0 \\ \log y, & \text{if } \lambda = 0 \end{cases} \quad (4.2)$$

The parameter λ is estimated using the method of profile likelihood. In this method, we first fix a value λ_0 of λ , and estimate the nuisance parameters β and σ^2 by maximizing the likelihood function $L(\lambda_0, \beta, \sigma^2)$. Doing so provides the estimates as $\hat{\beta}_{\text{MLE}} = (X^T X)^{-1} X^T \mathbf{u}_\lambda$ and $\hat{\sigma}_{\text{MLE}}^2 = \frac{RSS_\lambda}{n}$, where

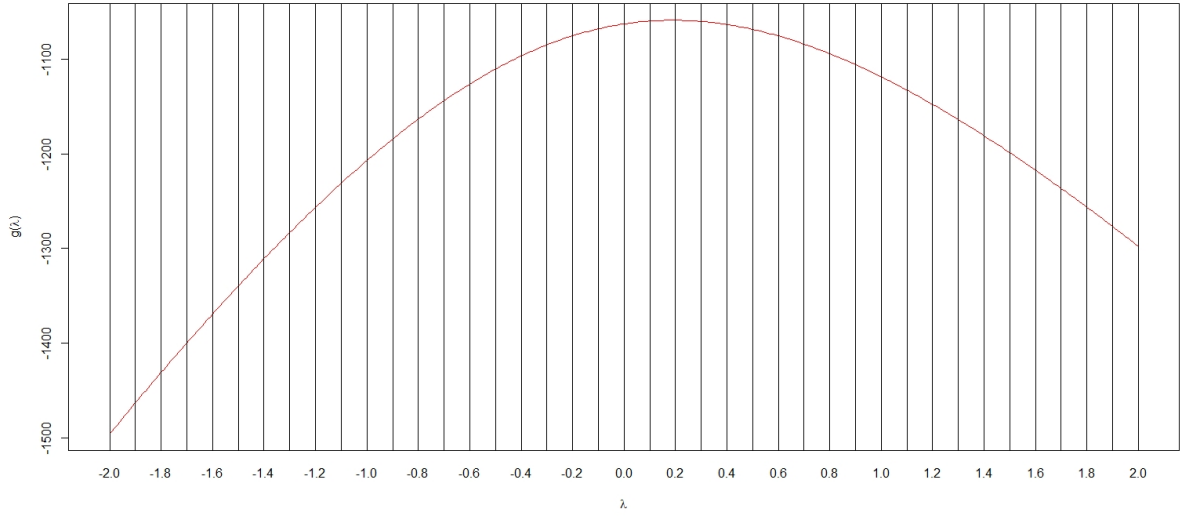


Figure 4.2: Maximizing $g(\lambda)$ graphically w.r.t. λ

$RSS_\lambda = \|\mathbf{u}_\lambda - X\hat{\boldsymbol{\beta}}_{\text{MLE}}\|^2$. Subsequently, the estimate of λ is estimated by maximizing the function

$$g(\lambda) = L(\lambda, \hat{\boldsymbol{\beta}}_{\text{MLE}}, \hat{\sigma}_{\text{MLE}}^2) = -\frac{n}{2} \log(2\pi) - \frac{n}{2} \log \frac{RSS_\lambda}{n} - \frac{n}{2} + n(\lambda - 1) \log G, \quad (4.3)$$

G being the geometric mean of the responses y_1, y_2, \dots, y_n . Naturally, it is cumbersome to maximize this function analytically, and as such we employ graphical method.

From the graph of $g(\lambda)$ versus λ in Figure 4.2, we observe that the maxima occurs at $\lambda = 0.2$. Therefore, undertake the Box-Cox transformation by putting $\lambda = 0.2$ in (4.2). To check the validity of normality assumption on the transformed response, we draw the QQ-plot of R-Student residuals of transformed variables. From the plot (Figure 4.3), we can conclude that the extent of normality in the

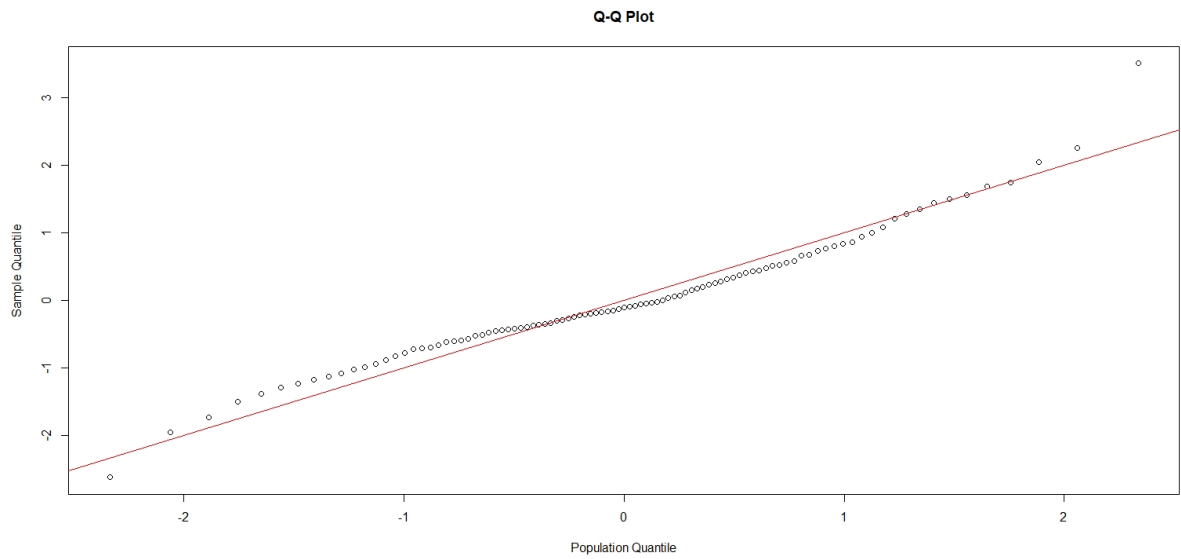


Figure 4.3: QQ-plot based on R-student residuals of transformed responses

data has improved. However, normality assumption cannot be said to be perfectly true. It is because

the actual distribution seems to have a slightly thicker tail than the desired distribution, as evident from the plot.