

DSC 423: DATA ANALYSIS AND
REGRESSION

FINAL PROJECT

BOSTON HOUSING DATASET

Swathi BABU ID: 2061690



TABLE OF CONTENTS

S. No.	Topics	Pg. No.
1	Introduction	2
2	Methodology	2
3	Exploratory Data Analysis	3
4	Full model	5
5	Log transformation	8
6	Quadratic Polynomial Regression	11
7	Comparing two models	17
8	The final model	17
9	Conclusion	18
10	References	19
11	Appendix	20

I. INTRODUCTION

The objective of this report was to predict the median value of houses occupied by the owners using SAS. The dataset used for this analysis is the Boston Housing Data collected by Harrison D. and Rubinfield D.L. who used it in the paper "Hedonic prices and the demand for clean air", Journal of Environmental Economics and Management, Vol. 5, Issue 1, 1978. [1]

The dataset that has been used here is a subset of the original set and it has 375 observations and 14 variables. There is one dependent variable, medv and 13 independent variables. One objective is to find the variables that would be a part of the model that is predicted.

To understand the dataset and how the variables affect each other, looking at articles regarding median value of houses helped at setting things in perspective. It helped in understanding what the important variables are and how we would incorporate it in the model. For example, one of the independent variables, dis, which is the weighted distance from five employment centers in Boston, is very important when looking for places. Proximity to work opportunities drastically improves the condition of living in a place. So, it is kept in mind that this is an important variable and must be included in the model. [2]

II. METHODOLOGY

To predict the dependent variable, medv, median value of owner-occupied homes, we use linear regression and build a model with significant independent variables. As there are multiple independent variables to predict one dependent variable, we use multiple linear regression. [3]

The dataset is imported under the data file "boston" using the infile function and using procedure print the first 10 observations of the imported dataset is checked (Refer Appendix A, fig. 1). Before proceeding with the process of model prediction, the dataset is analyzed first. Relationship between variables between dependent variable and independent variables and between independent variables themselves are looked through and further analysis is done.

The regression model is constructed after that. This report uses two methods to predict the model and concludes with one model based on various parameters. The two methods exercised are fitting the model with log transformation and through quadratic polynomial regression. The final model is obtained by comparing all these models. The results are further explained in detail under each section in this report.

III. EXPLORATORY DATA ANALYSIS

There are 13 independent variables – crime, zn, indus, chas, nox, rm, age, dis, rad, tax, ptratio, minor and lstat. Variables crime, zn, indus, nox, rm, age, dis, rad, tax, ptratio, minor and lstat are numeric variables. Variable chas is categorical binary variable that has values 0 and 1. Although there are no dummy variables added in the code, there is a binary variable. There aren't any variables that go well together enough to create interaction variables.

The important descriptive of all the variables are checked using the procedure means. The mean, minimum, maximum, standard deviation, standard error, first quartile, second quartile or median, third quartile, confidence intervals and mode (Refer Appendix A, fig. 2). The crime variable has a minimum of 0.00632, mean of 0.7182351, maximum of 18.4982 and Q3 is 0.537. This shows that 75% of the observations have crime rate below 0.537 and the maximum crime rate is much further away than the Q3. There could be some issues such as outliers on the right side there. Similar effect can be seen in zn, indus, rad and minor. For the zn variable, the minimum is 0, maximum is 100, mean is 15.334, Q2 is 0 and Q3 is 22. So, in this case most observations might have the value 0 for zn variable. For indus variable, minimum is 0.46, maximum is 25.65 and Q3 is 10.59. There are chances of outliers on the right side. The rad variable has minimum 1 and maximum 24. 75% of the observations have rad variable 5 or less than 5. The maximum is larger here as well. Chance of outliers on the right side can be identified from that. The variable minor has minimum 70.8, maximum 369.9 and Q3 is 396.06. This shows that there might be outliers on the left side of the graph. Also, from the standard deviations, the variables zn, age, tax and minor look like they would add a lot to the model since they have high standard deviations compared to the other variables. But that could also be because these variables have a higher range than the other variables i.e., they are not normalized. [3] Moreover, there are no missing values in the dataset that was used.

To check the distribution of data with respect to all the variables, histograms are plotted using procedure univariate. From the crime histogram, it can be seen that the graph is skewed to the right with possible outliers on the right side and it is unimodal (Refer Appendix A, fig. 3). The zn is also skewed to the right and is unimodal (Refer Appendix A, fig. 4). The indus variable is skewed to the right and is also unimodal (Refer Appendix A, fig. 5). Chas is a categorical binary variable and there are more observations with without track bounded rivers than with them (Refer Appendix A, fig. 6). The nox variable is slightly skewed to the right due to outliers and it is unimodal (Refer Appendix A, fig. 7). The rm variable is slightly skewed to the left due to some outliers and it is also unimodal (Refer Appendix A, fig. 8). The age variable is skewed to the left and it is unimodal (Refer Appendix A, fig. 9). The variable dis is skewed to the right, has some outliers on the right and is unimodal (Refer Appendix A, fig. 10). The variable rad is also skewed to the right due to some observations that have extreme values (outliers) and it is unimodal (Refer Appendix A, fig. 11). The tax variable is unimodal, skewed to the right and has some outliers on the right side of the graph (Refer Appendix A, fig. 12). The ptratio variable is skewed to the left, has outliers on the left side of the graph and it is also unimodal (Refer Appendix A, fig. 13). The minor is extremely skewed to the left, has outliers on the left side and is unimodal as well (Refer Appendix A, fig. 14). The lstat is skewed to the right, has outliers on the right and

is unimodal (Refer Appendix A, fig. 15). The medv is skewed to the right (Refer Appendix A, fig. 16).

Moving on to boxplots, boxplots and histograms convey similar information but boxplots are clearer on the outliers whereas histograms are not, but histograms are clearer on modal distribution of the graphs. There are many outliers on the right side with respect to crime variable. More than 15 observations higher crime rate than the estimated maximum (Refer Appendix A, fig. 17). There are 9 outliers on the right side for the variable zn. As already seen, most of the observations, around 240 of observations have the value 0 for zn (Refer Appendix A, fig. 18). The indus variable has 2 outliers on the right side (Refer Appendix A, fig. 19). Chas variable is a categorical variable so no point in checking univariate boxplot (Refer Appendix A, fig. 20). The nox variable has 3 outliers on the right side (Refer Appendix A, fig. 21). For rm variable, there are more than 5 outliers on the right side and there are exactly 3 on the left side (Refer Appendix A, fig. 22). The age variable has no outliers (Refer Appendix A, fig. 23). The variable dis has one outlier on the right side (Refer Appendix A, fig. 24). The rad variable has 3 outliers on the right side and 2 on the left (Refer Appendix A, fig. 25). The tax variable has 1 outlier on the right (Refer Appendix A, fig. 26). There are no outliers for ptratio (Refer Appendix A, fig. 27). There are more than 15 outliers on the left side for minor variable (Refer Appendix A, fig. 28). There are around 8 outliers on the right for lstat variable (Refer Appendix A, fig. 29). The boxplot for medv has also been plotted (Refer Appendix A, fig. 30). Boxplot is an important analysis step that tells the user a lot about the data that are crucial for further steps. [3]

Scatterplots are an important way to look at linearity in the relationships between dependent variable medv and the independent variables. The variable crime seems to slightly linear with respect to the dependent variable medv with possibly very low negative correlation (Refer to Appendix A, fig. 31). The zn variable seems to have a linear relationship with medv variable. They seem to have a low positive correlation (Refer to Appendix A, fig. 32). The indus variable also has a slight linear relation with respect to medv and also has low negative correlation (Refer to Appendix A, fig. 33). Chas variable is a categorical binary variable so no point of checking for linearity (Refer to Appendix A, fig. 34). The nox variable also only a slight linear relationship with medv and has low negative correlation (Refer to Appendix A, fig. 35). The rm clearly has a linear relationship with medv and they have high positive correlation (Refer to Appendix A, fig. 36). The age has a very subtle linear relationship with variable medv and has low negative correlation (Refer to Appendix A, fig. 37). The dis variable also has very slight linear relationship with medv and has very low negative correlation (Refer to Appendix A, fig. 38). The rad variable which is actually an index of values 1 to 24 is discrete and not continuous. Stating that, it has low linear relationship and has low positive correlation (Refer to Appendix A, fig. 39). The tax variable also has slight linear relationship with medv and has low negative correlation (Refer to Appendix A, fig. 40). The ptratio has a linear relationship with medv and has moderate positive correlation (Refer to Appendix A, fig. 41). The minor variable also has low linearity with respect to medv and has low positive correlation (Refer to Appendix A, fig. 42). The lstat variable has a good linear relationship with medv and has moderate negative correlation (Refer to Appendix A, fig. 43).

Another sure way of analyzing the linear relationship between the variables is Pearson's correlation coefficients. High correlation values between independent variables can also cause issues like multicollinearity etc.... Looking at the correlation table, the dependent variable medv has low positive correlation with zn, chas, rad, minor, has high positive correlation with rm, has low negative correlation with crime, indus, nox, age, dis, tax, has moderate negative correlation with ptratio, lstat. The variables crime and rad, tax have high positive correlation. Dis and zn seem to have moderate positive correlation. Dis and indus seem to have moderate negative correlation. Nox and indus also seem to have high positive correlation. Dis and nox have high negative correlation. Age and nox have moderate positive correlation. Lstat and rm have moderate negative correlation. Age and dis have high negative correlation. Rad and tax have high positive correlation. (Refer to Appendix A, fig. 44) Another important point to note is that there aren't any independent variables that have correlation as high as +/- 0.9 which means that there probably aren't any variables that are collinear. [3]

IV. FULL MODEL

The full model with all 13 independent variables was modeled using the procedure regression in SAS to predict the medv variable. All the parameters, assumptions and model diagnostics are checked for the full model.

The model obtained has an R^2 value of 0.7801 and adjusted R^2 of 0.7722 which means that 78.01% of the variance of medv variable can be explained by the 13 independent variables and also 77.22% of the variance of medv variable is explained by the 13 independent variables crime, zn, indus, chas, nox, rm, age, dis, rad, tax, ptratio, minor and lstat. The error RMSE value is 4.18606 which is pretty low.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 is that no independent variable has an effect on the dependent variable medv i.e., all $\beta_i = 0$. The alternate hypothesis H_1 is that there is at least one independent variable that has an effect on the dependent medv variable i.e., at least one $\beta_i \neq 0$. On running the goodness of fit test, we get F-value as 98.53 which is not that high but p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables crime, zn, indus, chas, rm, rad and minor have a positive effect on the dependent variable medv and the variables nox, dis, age, tax, ptratio, lstat have a negative effect on medv. The effect can be expanded into if the per capita crime rate of the town increase by 1%, then the median value of the house will increase by $(1.28883 * 1000) = 1288.83\$$. This can be done for all the variables in the model. Logically speaking, increase in crime rate might not lead to increase in home value so there could be some inconsistency here.

Based on standardized beta estimates, rm is the independent variable with the highest influence on medv followed by lstat, dis, crime, nox, rad, tax, ptratio, zn, minor, indus, age

and finally chas. This is estimated using the absolute value of the standardized beta estimates. Variable rm had the highest value followed by the other variables in the order mentioned above. [3]

The model equation is $medv = 12.63324 + 1.28883crime + 0.03579zn + 0.07362indus + 0.12124chas - 17.35683nox + 6.41932rm - 0.00650age - 1.18892dis + 0.39406rad - 0.01565tax - 0.71652ptratio + 0.01167minor - 0.042714lstat$, where chas = 1 if tract bounds river, 0 otherwise.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that the variables chas, age and indus are insignificant as their p-values are greater than 0.05 (Refer to Appendix B, table 1).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption can be checked using scatterplots. It was addressed in the scatterplot section of data exploration. The assumption is satisfied. The constant variance and independency assumption go hand in hand. The residual plot for variable crime has a cone kind of shape and hence does not satisfy the constant variance and independency assumptions (Refer to Appendix B, fig. 2). The residual plot for zn variable does not have randomly scattered points. It does not satisfy the constant variation and independency assumptions (Refer to Appendix B, fig. 3). The indus variable seems to have a tunnel shape and hence does not satisfy constant variance and independency (Refer to appendix B, fig. 4). The assumptions need not be checked for chas variable as it is categorical (Refer to Appendix B, fig. 5). The nox variable also has a tunnel shape. Therefore, the assumptions constant variance and independency are not satisfied (Refer to Appendix B, fig. 6). Rm has a U-shape (Refer to Appendix B, fig. 7). Age has a tunnel shape (Refer to Appendix B, fig. 8). Variable dis also has a tunnel shape (Refer to Appendix B, fig. 9). They do not satisfy the constant variance and independence assumption. Variable rad has index values so they are discrete ad do not satisfy the assumptions (Refer to Appendix B, fig. 10). The variables pat, ptratio and minor have a tunnel shape. So, they do not satisfy the assumptions either (Refer to Appendix B, fig. 11,12,13). The variable lstat has a U -shape (Refer to Appendix B, fig. 14). The predicted student residual plot also has a U-shape (Refer Appendix B, fig. 15). The normal probability plot seems to have a slight S shape but it is almost a straight line (Refer to Appendix B, fig. 16). [3]

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. All the variables have variation inflation factor (VIF) lesser than 10. So, there is no issue of collinearity which was observed through correlation values as well (Refer to Appendix B, table 1). There are outliers and influential points for this model. Observations that are both outliers and influential points are observation number 215, 357, 365, 366, 369, 372, 373 and 375 (Refer to Appendix B, fig. 1). There are many more influential points but those are left in the model as removing them will lead to too much loss of data. Moving forward, we will be removing the above-mentioned observations from the boston dataset and create a new one without these- "boston_new" which is what will be used for the rest of the steps. [3]

This model is a moderate model with 78% accuracy but there are assumptions that are not met and many outliers. Proceeding to running the model again on “boston_new” dataset after removing the outliers.

The model built after removing outliers have R^2 value of 0.8599 and adjusted R^2 of 0.8547 and that means that 85.99% of the variance of medv is explained by the independent variables and 85.47% of the variance of medv is explained by them. This is better than the previous model with almost 5% increase in accuracy. The error RMSE is 3.25627 which has also decreased from the previous model.

The model null hypothesis and alternate hypothesis remains the same and it can be seen here that F-value is 166.65, it has increased as well and the p-value is less than 0.05. Hence, the null hypothesis is rejected. So, this model is good and has at least one independent variable that influences the dependent variable medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables crime, zn, indus, chas, rm, rad and minor have a positive effect on the dependent variable medv and the variables nox, dis, age, tax, ptratio, lstat have a negative effect on medv. The effect can be expanded into if the per capita crime rate of the town increase by 1%, then the median value of the house will increase by $(1.725577 * 1000) = 1725.77\$$. This can be done for all the variables in the model.

Based on standardized beta estimates, rm is the independent variable with the highest influence on medv followed by crime, dis, lstat, nox, tax, ptratio, rad, age, zn, minor, indus, and chas. This is estimated using the absolute value of the standardized beta estimates (Refer to Appendix B, table 2).

The model equation is $medv = -8.01046 + 1.72577 \text{ crime} + 0.02307 \text{zn} + 0.05365 \text{indus} + 1.08971 \text{chas} - 12.53887 \text{nox} + 8.51415 \text{rm} - 0.02453 \text{age} - 0.90588 \text{dis} + 0.21037 \text{rad} - 0.01464 \text{tax} - 0.57630 \text{ptratio} + 0.01385 \text{minor} - 0.24693 \text{lstat}$, where chas = 1 if tract bounds river, 0 otherwise.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that the variables chas and indus are insignificant as their p-values are greater than 0.05 (Refer to Appendix B, table 2).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption can be checked using scatterplots. It was addressed in the scatterplot section of data exploration. The assumption is satisfied. The constant variance and independency assumption go hand in hand. The residual plots haven't changed from the previous model much. The assumptions constant variance and independency are not satisfied. The normal probability plot looks like a straight line. So, the assumption of normality is satisfied (Refer to Appendix B, fig. 18 – 32).

Using the VIF, r, influence add-ons in the model, the model diagnostics is assessed. All the variables have variation inflation factor (VIF) lesser than 10. So, there is no issue of

collinearity in this model as well. There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 162, 182, 365 and 367 (Refer to Appendix B, fig. 17). Removing these points does not add much to the R². There are many more influential points but those are left in the model as removing them will lead to too much loss of data here as well.

This model with 86% accuracy is better than the previous one but still not the best one as the assumptions constant variance and independency are violated. To resolve this, two methods are undertaken here as mentioned in the methodology section.

V. LOG TRANSFORMATION

A new dataset “boston_log” from the “boston_new” dataset was created with the log variable i.e., the variable logmedv which is log of medv is added as the independent variable in the new dataset. The full model with all 13 variables is run to predict logmedv rather than medv.

The model obtained has an R² value of 0.8777 and adjusted R² of 0.8732 which means that 87.77% of the variance of logmedv variable can be explained by the 13 independent variables and also 87.32% of the variance of logmedv variable is explained by them. This is not a huge increase in the values but this model is better. The error RMSE value is 0.11013 which is a decrease from the previous model RMSE.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H₀ is that no independent variable has an effect on the dependent variable logmedv i.e., all $\beta_i = 0$. The alternate hypothesis H₁ is that there is at least one independent variable that has an effect on the dependent logmedv variable i.e., at least one $\beta_i \neq 0$. On running the goodness of fit test, we get F-value as 194.96 which has also increased from the previous model and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on logmedv. [4]

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables crime, zn, indus, chas, rm, rad and minor have a positive effect on the dependent variable logmedv and the variables nox, dis, age, tax, ptratio, lstat have a negative effect on logmedv. The effect can be expanded into if the per capita crime rate of the town increase by 1%, then the median value of the house will increase by ($e^{0.04567} - 1$) *100% i.e., 4.67%. This can be done for all the variables in the model. [4]

Based on standardized beta estimates, rm is the independent variable with the highest influence on logmedv followed by lstat, dis, nox, crime, ptratio, tax, rad, age, minor, zn, indus and chas. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix C, table 1)

The model equation is $\log(\text{medv}) = 2.51620 + 0.04567\text{crime} + 0.00074642\text{zn} + 0.00227\text{indus} + 0.03477\text{chas} - 0.57152\text{nox} + 0.25083\text{rm} - 0.00092026\text{age} - 0.03512\text{dis} + 0.01240\text{rad} - 0.00012\text{lstat}$

$0.00057326\text{tax} - 0.02534\text{ptratio} + 0.00060553\text{minor} - 0.01516\text{lstat}$, where $\text{chas} = 1$ if tract bounds river, 0 otherwise.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that the variables chas and indus are insignificant as their p-values are greater than 0.05 in this model also(Refer to Appendix C, table 1).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption can be checked using scatterplots. The assumption is satisfied. The variables student residual plots are more random around the origin now than the previous model especially the student residuals vs predicted values plot has significant improvement with the transformation (Refer to Appendix C, fig. 2-16). [4]

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. All the variables have variation inflation factor (VIF) lesser than 10. So, there is no issue of collinearity (Refer to Appendix C, table 1). There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 8, 182 and 367 (Refer to Appendix C, fig. 1). The outliers are not removed as they don't add much R² to the model. Removing influential points could lead to too much loss in data. This model with 87% accuracy is a pretty good model and it is better than the previous one because it satisfies all the assumptions and model diagnostics.

After looking at the model, the dataset is split into training and testing through the holdout method with 80% training and 20% testing with seed value 317380 using procedure surveyselect and a new dataset is created for the training called "boston_all_log" where a new variable "new_medv" is created to assign the logmedv value to it if it part of the training. There are 294 observations for training and 73 for testing.

Further, two model selection methods Stepwise and backward selection is run on the new dataset "boston_all_log". On running both these methods, they arrive at the same models with 12 variables- crime, zn, chas, nox, rm, dis, rad, tax, ptratio, minor and lstat.

The model obtained has an R² value of 0.8854 and adjusted R² of 0.8805 which means that 87.77% of the variance of new_medv variable can be explained by the 12 independent variables and also 88.54% of the variance of new_medv variable is explained by them. The error RMSE value is 0.10332 which is a decrease from the previous model RMSE, although not much.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H₀ is that no independent variable has an effect on the dependent variable new_medv i.e., all $\beta_i = 0$. The alternate hypothesis H₁ is that there is at least one independent variable that has an effect on the dependent new_medv variable i.e., at least one $\beta_i \neq 0$. On running the goodness of fit test, we get F-value as 180.98 which has also increased from the previous model and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on new_medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables crime, zn, chas, rm, rad and minor have a positive effect on the dependent variable new_medv and the variables nox, dis, age, tax, ptratio, lstat have a negative effect on new_medv. The effect can be expanded into if the per capita crime rate of the town increase by 1%, then the median value of the house will increase by $(e^{0.05996} - 1) * 100\%$ i.e., 6.18%. This can be done for all the variables in the model.

Based on standardized beta estimates, rm is the independent variable with the highest influence on logmedv followed by dis, crime, nox, lstat, ptratio, tax, rad, age, minor, zn and chas. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix C, table 2)

The model equation is $new_medv = \log(medv) = 2.47906 + 0.05996crime + 0.00076665zn + 0.04569chas - 0.66316nox + 0.25313rm - 0.00123age - 0.03781dis + 0.01079rad - 0.00048448tax - 0.02516ptratio + 0.00074009minor - 0.01197lstat$, where chas = 1 if tract bounds river, 0 otherwise.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix C, table 2).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption is satisfied as confirmed by the scatterplots. The variables student residual plots are random around the origin. The student residuals vs predicted values plot has significant improvement with the transformation as already seen (Refer to Appendix C, fig. 18-30). The normal probability plot is a straight line. [4]

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. All the variables have variation inflation factor (VIF) lesser than 10. So, there is no issue of collinearity (Refer to Appendix C, table 2). There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 8, and 185 (Refer to Appendix C, fig. 17). The outliers are removed, and the model is run again.

On doing that, we get a model with R^2 value of 0.8938 and adjusted R^2 0.8893. The values have slightly increased from the previous model. The RMSE is 0.09977 which has decreased from before. For the same null and alternate hypothesis, this has a F-value of 195.72 and p-value is less than 0.05. So, the null hypothesis is rejected, and this model passes the goodness of fit test.

The signs of the effect of the variables using parameter estimates remains the same from the previous model. By increasing crime rate by 1%, the medv increases by $(e^{0.0621} - 1) * 100 = 6.41\%$. This can be calculated for all variables in the model.

Based on standardized beta estimates, rm is the independent variable with the highest influence on logmedv followed by dis, crime, lstat, nox, ptratio, age, tax, rad, minor, zn and chas. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix C, table 3)

The model equation is $\text{new_medv} = \log(\text{medv}) = 2.42099 + 0.0621\text{crime} + 0.0008005\text{zn} + 0.04938\text{chas} - 0.63361\text{nox} + 0.25741\text{rm} - 0.00143\text{age} - 0.03883\text{dis} + 0.00988\text{rad} - 0.00045523\text{tax} - 0.02333\text{ptratio} + 0.0007221\text{minor} - 0.01217\text{lstat}$, where chas = 1 if tract bounds rive, otherwise 0.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix C, table 3).

All four assumptions – Linearity, Constant Variance, Independency and Normality are satisfied which can be seen from the residual plots (Refer to Appendix C, fig. 32-45)

The model diagnostics have been checked. All VIF values are less than 10 so there is no collinearity issue. There is only one outlier, observation number 148, which can be left in the model.

This model is tested on test set and then the test model parameters are computed. The R^2 value for the test set is $(0.92171)^2 = 0.8495$. The error terms MAE has a value of 22.6199 and RMSE of 9.57828. The cross-validated R^2 is calculated by taking the absolute value of the difference between the squares of training R^2 and testing R^2 . The cross-validated R^2 value for this model is 0.0772 which is lesser than 0.3. Hence, it is a good model based on this parameter. (Refer to Appendix C, table 4)

Comparing the training and testing performance, the training set has higher R^2 and adjusted R^2 value but there is not a huge difference. Training does better with respect to RMSE as well. Overall, although training seems to be doing better than test, testing still performed well and hence this is a good model.

VI. QUADRATIC POLYNOMIAL REGRESSION

Based on the student residuals vs predicted plot with the U-shape, quadratic polynomial regression could be an option. This method is being exercised to see if this gets better results than the models with no transformation or with log transformation.

On running the full model using procedure glmselect, the resulting model had 91 independent variables and 1 dependent variable medv.

The R^2 value of the model is 0.9522 and Adj- R^2 is 0.9364 that means that 95.22% of the variance of the medv variable is explained by the 91 independent variables. The error terms

RMSE has value 2.15487, AIC has 1010.60656, AICC has 1074.65051 and SBC has a value of 1000.89985.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 is that no independent variable has an effect on the dependent variable medv i.e., all $\beta_i = 0$. Alternate hypothesis is that there is at least 1 variable that has an effect on the independent variable. The F-value for this model is 60.20 and the corresponding p-value is less than 0.05. So, the null hypothesis is rejected and hence this model satisfies the goodness of fit test (Refer to Appendix D, table 1).

The data is then split in to training and testing sets through holdout method with 75% training and 25% testing using procedure surveyselect with seed 456729. A new variable is created for training new_medv for just the training set in the new dataset "boston_all_pol".

Following this, two model selection process was run on the training dataset – stepwise and forward model selection. Two methods arrive at two different models.

Looking at the model obtained by stepwise method, model 1, it has the independent variables rm, dis*tax (dis_tax), rm*ptratio (rm_ptratio), lstat, rm*lstat (rm_lstat). Creating a new dataset with these variables along with the existing ones, the model is run again using the procedure regression to predict medv (Refer to Appendix D, table 2).

The model obtained has an R^2 value of 0.8635 and adjusted R^2 of 0.8610 which means that 86.35% of the variance of new_medv variable can be explained by the 5 independent variables and also 86.10% of the variance of new_medv variable is explained by them. The error RMSE value is 3.11001 which is an increase from the full model although not by much.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 and the alternate hypothesis H_1 is the same. On running the goodness of fit test, we get F-value as 341.57 which is a great increase from the full model and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on new_medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables rm and lstat have a positive effect on the dependent variable new_medv and the variables dis_tax, rm_ptratio, rm_lstat have a negative effect on new_medv. The effect can be expanded into if the average number of rooms increases by 1, then the median home value will increase by 12566.05\$. This can be done for all the variables in the model.

Based on standardized beta estimates, rm_lstat is the independent variable with the highest influence on medv followed by lstat, rm, dis_tax, and rm_ptratio. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix D, table 2)

The model equation is $new_medv = medv = -40.38857 + 12.56605rm - 0.00151dis*tax - 0.06595rm*ptratio + 2.33246lstat - 0.46491rm*lstat$.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix D, table 2).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption is satisfied as confirmed by the scatterplots. The variables student residual plots are random around the origin. The student residuals vs predicted values plot has significant improvement with the transformation as already seen (Refer to Appendix D, fig. 2-7). The normal probability plot seems to have a very slight S shape but that could be due to outliers.

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. The variables rm, dis_tax, rm_ptratio have VIF less than 10 so they are not collinear but the variables lstat and rm_lstat have VIF of 81.51468 and 69.53019 which are much greater than 10. So, these variables are collinear, but we can safely ignore this as one variable is actually the product of the other variable. So, this is to be expected and hence can be ignored [5]. There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 8,89,182,365 (Refer to Appendix D, fig. 1). The outliers are removed, and the model is run again.

The model obtained has an R^2 value of 0.8965 and adjusted R^2 of 0.8946 which means that 89.65% of the variance of new_medv variable can be explained by the 5 independent variables and also 89.46% of the variance of new_medv variable is explained by them. The error RMSE value has dropped to 2.67161.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 and the alternate hypothesis H_1 is the same. On running the goodness of fit test, we get F-value as 461.05 which is an increase from the previous model and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on new_medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables rm and lstat have a positive effect on the dependent variable new_medv and the variables dis_tax, rm_ptratio, rm_lstat have a negative effect on new_medv. The effect can be expanded into if the average number of rooms increases by 1, then the median home value will increase by $(12.74890 * 1000) \$ = 12748.90 \$$. This can be done for all the variables in the model.

Based on standardized beta estimates, rm_lstat is the independent variable with the highest influence on medv followed by lstat, rm, rm_ptratio and dis_tax. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix D, table 3)

The model equation is $new_medv = medv = -41.58632 + 12.74890rm - 0.00137dis*tax - 0.065953rm*ptratio + 2.37504lstat - 0.46988rm*lstat$.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix D, table 3).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption is satisfied as confirmed by the scatterplots. The variables student residual plots are random around the origin. The student residuals vs predicted values plot has significant improvement with the transformation as already seen (Refer to Appendix D, fig. 9-15). The normal probability plot is a straight line now.

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. The variables lstat and rm_lstat are still collinear and it can be ignored here as well [5]. There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 183, 184 and 339(Refer to Appendix D, fig. 8). The outliers are not removed.

This model is tested on test set and then the test model parameters are computed. The R^2 value for the test set is $(0.92611)^2 = 0.8577$. The error terms MAE has a value of 2.36930 and RMSE of 3.46342. The cross-validated R^2 value for this model is 0.068 which is lesser than 0.3. Hence, it is a good model based on this parameter. (Refer to Appendix D, table 4)

Comparing the training and testing performance, the training set has higher R^2 and adjusted R^2 value but there is not a huge difference. Testing has lower RMSE than training. Overall, looking at all parameters, it is good model.

Moving on to the second model that was obtained by forward selection method, it has 7 independent variables – rm, age*dis (age_dis), dis*tax (dis_tax), rm*pseudo_ratio (rm_pseudo_ratio), lstat, nox*lstat (nox_lstat) and rm*lstat (rm_lstat). Creating a new dataset like the previous model, regression is carried out again through procedure regression.

The model obtained has an R^2 value of 0.8723 and adjusted R^2 of 0.8690 which means that 87.23% of the variance of new_medv variable can be explained by the 7 independent variables and also 86.90% of the variance of new_medv variable is explained by them. The error RMSE value is 3.01872 which is an increase from the full model although not by much.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 and the alternate hypothesis H_1 is the same. On running the goodness of fit test, we get F-value as 261.62 and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on new_medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables rm and lstat have a positive effect on the dependent variable new_medv and the variables age_dis, dis_tax, rm_pseudo_ratio, nox_lstat, rm_lstat have a negative effect on new_medv. The effect can be expanded into if the average number of rooms

increases by 1, then the median home value will increase by $(13.17032 * 1000) = 13170.32\$$. This can be done for all the variables in the model.

Based on standardized beta estimates, lstat is the independent variable with the highest influence on medv followed by rm_lstat, rm, rm_ptratio, dis_tax, and age_dis. This is estimated using the absolute value of the standardized beta estimates. (Refer to Appendix D, table 4)

The model equation is $\text{new_medv} = \text{medv} = -42.33980 + 13.17032\text{rm} - 0.00482\text{age}*\text{dis} - 0.00165\text{dis}*\text{tax} - 0.8043\text{rm}*\text{ptratio} + 2.99927\text{lstat} - 0.59786\text{nox}*\text{lstat} - 0.50394\text{rm}*\text{lstat}$

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix D, table 4).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption is satisfied as confirmed by the scatterplots. The variables student residual plots are random around the origin. The student residuals vs predicted values plot has improvement as already seen (Refer to Appendix D, fig. 17-25). The normal probability plot seems to have a very slight S shape but that could be due to outliers.

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. The variables rm, age_dis, dis_tax, rm_ptratio have VIF less than 10 so they are not collinear but the variables lstat, nox_lstat and rm_lstat have VIF greater than 10. So, these variables are collinear, but we can safely ignore it for the same reason as before. [5] There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 8, 182 and 365 (Refer to Appendix D, fig. 16). The outliers are removed, and the model is run again.

The model obtained has an R^2 value of 0.9021 and adjusted R^2 of 0.8995 which means that 90.21% of the variance of new_medv variable can be explained by the 7 independent variables and also 89.95% of the variance of new_medv variable is explained by them. The error RMSE value has dropped to 2.60421.

The hypothesis to check if the model is good or not is framed and a goodness of fit test is run. The null hypothesis H_0 and the alternate hypothesis H_1 is the same. On running the goodness of fit test, we get F-value as 348.75 which is an increase from the previous model and p-value is less 0.05 which means that the null hypothesis is rejected and there is at least one independent variable that has an effect on new_medv.

Using the signs of the parameter estimates or beta coefficients, it can be said that the independent variables rm and lstat have a positive effect on the dependent variable new_medv and the variables dis_tax, rm_ptratio, rm_lstat have a negative effect on new_medv. The effect can be expanded into if the average number of rooms increases by 1,

then the median home value will increase by $(13.17032 * 1000) \$ = 13170.32 \$$. This can be done for all the variables in the model.

Based on standardized beta estimates, the order of influence of the independent variables on the dependent variable is same as the previous model. (Refer to Appendix D, table 5)

The model equation is $new_medv = medv = -43.37581 + 13.33592rm - 0.00549age*dis - 0.00144dis*tax - 0.08550rm*pseudo + 3.03364lstat - 0.61369nox*lstat - 0.50391rm*lstat$.

Significance test is also run on each independent variable. From the p-values for these variables, it can be seen that all variables in the model are significant as their p-values are less than 0.05 in this model (Refer to Appendix D, table 5).

The four assumptions of the regression model are checked using student residual plots and normality plot. The linearity assumption is satisfied as confirmed by the scatterplots. The variables student residual plots are random around the origin. The student residuals vs predicted values plot has significant improvement with the transformation as already seen (Refer to Appendix D, fig.27-35). The normal probability plot is a straight line now.

Using the VIF, r, influence add-ons in the model, the model diagnostics can be assessed. The variables lstat, nox_lstat and rm_lstat are still collinear and it can be ignored here as well. [5]There are outliers and influential points for this model. Observations that are both outliers and influential points are observation numbers 88 and 185(Refer to Appendix D, fig. 26). The outliers are not removed.

This model is tested on test set and then the test model parameters are computed. The R^2 value for the test set is $(0.93033)^2 = 0.8655$. The error terms MAE has a value of 2.22168 and RMSE of 3.36385. The cross-validated R^2 value for this model is 0.064 which is lesser than 0.3. Hence, it is a good model based on this parameter. (Refer to Appendix D, table 6)

Comparing the training and testing performance, the training set has higher R^2 and adjusted R^2 value but there is not a huge difference. Training has lower RMSE than testing. Overall, looking at all parameters, training seems to be doing better than testing but it still is a good model as there isn't any huge disparities.

Comparing the stepwise model 1 and forward selection model 2, looking at the training test parameters for both models, model 2 has higher R^2 and adj R^2 than 1, RMSE is also lower for model 2. F-value is higher for model 1. Multicollinearity is ignored in both. Considering all this, training for model 2 is better than 1. Comparing the test model parameters of both model, R^2 is higher for model 2. MAE, RMSE are lower for model 2. Both satisfy the cross-validated R^2 test. Looking at all this, model 2 test is better than 1. So, overall, in quadratic polynomial regression, model 2 is better.

In the next section, the log model and the polynomial model is compared and a conclusion is decided.

VII. COMPARING THE TWO FITTED MODELS

Comparing the training models, the polynomial model has the higher R² and Adj R² than the log transformation model although just 0.01 difference. The log transformation model has lower RMSE than the polynomial one. Polynomial model has the higher F-value and both of them satisfy the goodness of fit test. The polynomial model has multicollinearity issues but the other one doesn't.

Comparing the testing model parameters, polynomial model has R² 0.02 higher than the log transformation model. MAE and RMSE models have lower values for polynomial model. Both the models satisfy the cross-validated R² condition.

In spite of polynomial model working better than log model, log model also has good parameters and is a good model. Log transformation model is the simpler one out of the two. If both models work well, then choosing the simpler one is the good option. Hence, the final fitted model is the log transformed model(Refer to Appendix C, table 1).

VIII. THE FINAL MODEL

The final model is $\log(\text{medv}) = 2.50502 + 0.04650\text{crime} + 0.00070974\text{zn} + 0.03864\text{chas} - 0.52470\text{nox} + 0.24922\text{rm} - 0.00094060\text{age} - 0.03668\text{dis} + 0.01153\text{rad} - 0.00052892\text{tax} - 0.02459\text{ptratio} + 0.00060319\text{minor} - 0.01494\text{lstat}$, where chas = 1 if tract bounds rive, otherwise 0. (Run on the full dataset boston_log) (Refer to Appendix E, table 1)

All assumptions are satisfied. The model diagnostics was checked and there are some outliers which can be ignored. The effect of the variables remains the same. But the variable chas is insignificant as the p-value is higher than 0.05, it is 0.0723. Removing that variable and running the model again, we get the model equation

$\log(\text{medv}) = 2.50664 + 0.04657\text{crime} + 0.00071733\text{zn} - 0.52204\text{nox} + 0.25034\text{rm} - 0.00092748\text{age} - 0.03722\text{dis} + 0.01231\text{rad} - 0.00054065\text{tax} - 0.02506\text{ptratio} + 0.00060768\text{minor} - 0.01482\text{lstat}$ (Refer to Appendix E, table 2)

All assumptions – linearity, constant variance, independency and normality are satisfied. The model diagnostics also look good. There are some outliers and influential points which we can ignore. There is no issue of multicollinearity. The other parameters such as R² etc. also look good.

Interpreting the variables, increase in crime rate by 1% will lead to $(e^{0.04657} - 1) * 100 = 4.77\%$ increase in median home value. Increase in zn variable by 1 will lead to $(e^{0.00070974} - 1) * 100 = 0.07\%$ increase in medv. When nitric oxide concentration increases by 1, then the medv will decrease by $(e^{0.52204} - 1) * 100 = 68.55\%$. When the average number of rooms increases by 1, the median home value will increase by $(e^{0.25034} - 1) * 100 = 28.45\%$. When the age proportion increases by 1, the medv variable decreases by $(e^{0.00092748} - 1) * 100 = 0.093\%$. When the weighted distance to the employment center increases by 1, then the medv will decrease by $(e^{0.03722} - 1) * 100 = 3.79\%$. When the index of accessibility to highways increases by 1, then the medv will increase by $(e^{0.01231} - 1) * 100 = 1.24\%$. When the tax

variable increases by 1 i.e., 10000\$, the medv variable decreases by $(e^{0.00054065} - 1) * 100 = 0.0054\%$. When the pupil-teacher ratio increases by 1, then the medv variable will decrease by $(e^{0.02506} - 1) * 100 = 2.54\%$. When the minor variable increases by 1, the medv variable increases by $(e^{0.00060768} - 1) * 100 = 0.06\%$. When the lower status population increases by 1%, the medv variable decreases by $(e^{0.014182} - 1) * 100 = 1.43\%$.

Two predictions were made using this model. When the crime rate is 0.5%, proportion for zoned residential land is 15, nitric oxides concentration(parts per 10 million) is 0.5, average number of rooms per dwelling is 3, proportion of homes built prior to 1940 is 50, weighted distance to closest employment centers is 4, with index for accessibility to highways as 4, paying 400(per 10000\$) tax, with pupil-teacher ratio as 15, with minor as 300 and percentage of lower status homes as 25% , the median value for home is $(e^{2.1043}) = 8.20136$ i.e., 8201.36\$ is the median value of the home with 95% confidence interval equal to $(e^{2.0220}, e^{2.1865})$ i.e., (7.5534, 8.9039) and prediction interval equal to $(e^{1.8716}, e^{2.3370})$ i.e., (6.4987,10.3501) which means that the predicted value for medv can increase from 6498.7\$ to 10350.1\$. (Refer to Appendix E, fig. 1)

Another prediction is when the crime rate is 18.5%, proportion for zoned residential land is 25, nitric oxides concentration(parts per 10 million) is 0.4, average number of rooms per dwelling is 4, proportion of homes built prior to 1940 is 70, weighted distance to closest employment centers is 5, with index for accessibility to highways as 10, paying 350(per 10000\$) tax, with pupil-teacher ratio as 20, with minor as 250 and percentage of lower status homes as 25% , the median value for home is $(e^{3.1416}) = 23.14086$ i.e., 23140.86\$ is the median value of the home95% confidence interval equal to $(e^{2.8739}, e^{3.4094})$ i.e., (17.7059, 30.24709) and prediction interval equal to $(e^{2.7966}, e^{3.4867})$ i.e., (16.38883, 32.67793) which means that the predicted value for medv can increase from 16388.83\$ to 32677.93\$. (Refer to Appendix E, fig. 1)

IX. CONCLUSION

The initial model did not satisfy certain assumptions. So, log transformation and polynomial regression was adopted for this dataset. After various steps, the log transformed model was chosen as the final model.

Based on the model, there seems to be some inconsistency as increase in crime would not lead to increase in the value of a home. Collecting more data could solve this issue. Also, the data is outdated. There are variables like nox, crime, dis, tax, rad etc. that changes over time and as cities and states develop, a lot of these parameters change as well. For example, there are probably more employment centers in Boston now than in 1978. Tax policies could have changed, more highways would've been constructed. Crime rates have also changed since 1978. All these factors make the model not useful for the present.

The important avenue with respect to this dataset and model for the future is to keep collecting and updating data as this is more subjective to time and the model would have to

be adopted to the changes, if not, it becomes useless. This shows that the current result isn't enough to make decisions or give recommendations for the present.

All these topics and suggestion were easier to identify and make because there was lots of background research involved. As already mentioned, the articles helped put things into perspective on what's important and what's not. The scientific articles gave some detailed insights on the statistical side of this process.

REFERENCES

- [1] D. L. R. David Harrison Jr., "Hedonic housing prices and the demand for clean air," *Journal of Environmental Economics and Management*, vol. 5, no. 1, 1978.
- [2] J. Gomaz, "8 critical factors that influence a home's value," 2019.
- [3] M. M. J. E. M. a. P. M. Tranmer, "Multiple Linear Regression (2nd Edition)," *Cathie Marsh Institute Working Paper 2020-01*, 2020.
- [4] H. W. N. L. T. C. H. H. Y. L. X. M. T. Changyong Feng, "Log-transformation and its implications for data analysis," *Shanghai Arch Psychiatry*, vol. 26, no. 2, pp. 105 - 109, 2014.
- [5] P. Allison, "When Can You Safely Ignore Multicollinearity?," 2012.

APPENDIX

APPENDIX A – Data Exploratory Stage

Fig 1.

Fitted model assumptions and Diagnostics check

Obs	crime	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	minor	Istat	medv
1	0.00632	18.0	2.31	0	0.538	6.575	65.2	4.0900	1	296	15.3	396.90	4.98	24.0
2	0.02731	0.0	7.07	0	0.469	6.421	78.9	4.9671	2	242	17.8	396.90	9.14	21.6
3	0.02729	0.0	7.07	0	0.469	7.185	61.1	4.9671	2	242	17.8	392.83	4.03	34.7
4	0.03237	0.0	2.18	0	0.458	6.998	45.8	6.0622	3	222	18.7	394.63	2.94	33.4
5	0.06905	0.0	2.18	0	0.458	7.147	54.2	6.0622	3	222	18.7	396.90	5.33	36.2
6	0.02985	0.0	2.18	0	0.458	6.430	58.7	6.0622	3	222	18.7	394.12	5.21	28.7
7	0.08829	12.5	7.87	0	0.524	6.012	66.6	5.5605	5	311	15.2	395.60	12.43	22.9
8	0.14455	12.5	7.87	0	0.524	6.172	96.1	5.9505	5	311	15.2	396.90	19.15	27.1
9	0.21124	12.5	7.87	0	0.524	5.631	100.0	6.0821	5	311	15.2	386.63	29.93	16.5
10	0.17004	12.5	7.87	0	0.524	6.004	85.9	6.5921	5	311	15.2	386.71	17.10	18.9

Fig. 2

Descriptives

The MEANS Procedure

Variable	Mean	Minimum	Maximum	Std Dev	Std Error	25th Pctl	50th Pctl	75th Pctl	Lower 95% CL for Mean	Upper 95% CL for Mean	Mode
crime	0.7182351	0.0063200	18.4982000	1.7627831	0.0910297	0.0612900	0.1405200	0.5370000	0.5392408	0.8972293	0.0150100
zn	15.3333333	0	100.0000000	25.9505048	1.3400783	0	0	22.0000000	12.6983009	17.9683657	0
indus	8.8374400	0.4600000	25.6500000	6.3134986	0.3260277	4.0500000	6.9100000	10.5900000	8.1963629	9.4785171	19.5800000
chas	0.0933333	0	1.0000000	0.2912876	0.0150420	0	0	0	0.0637558	0.1229109	0
nox	0.5194659	0.3850000	0.8710000	0.1096108	0.0056603	0.4379000	0.4930000	0.5470000	0.5083359	0.5305958	0.5380000
rm	6.3685067	3.5610000	8.7800000	0.7220313	0.0372855	5.9270000	6.2450000	6.7260000	6.2951911	6.4418222	6.1270000
age	61.9800000	2.9000000	100.0000000	28.7008645	1.4821063	36.6000000	65.1000000	89.8000000	59.0656941	64.8943059	100.0000000
dis	4.3744125	1.1296000	12.1265000	2.1365919	0.1103331	2.5961000	4.0123000	5.8700000	4.1574615	4.5913636	3.4952000
rad	5.4586667	1.0000000	24.0000000	4.5653629	0.2357543	4.0000000	5.0000000	5.0000000	4.9950965	5.9222368	5.0000000
tax	328.6986667	187.0000000	666.0000000	102.2730296	5.2813565	270.0000000	307.0000000	398.0000000	318.3137917	339.0835416	307.0000000
ptratio	17.8581333	12.6000000	22.0000000	2.2173916	0.1145056	16.4000000	18.0000000	19.6000000	17.6329778	18.0832888	14.7000000
minor	379.4128533	70.8000000	396.9000000	41.5068023	2.1434021	380.3400000	392.2000000	396.0600000	375.1982236	383.6274831	396.9000000
Istat	10.4724533	1.7300000	37.9700000	6.0966282	0.3148285	6.0500000	9.3800000	13.2800000	9.8533974	11.0915092	6.3600000
medv	25.1944000	11.8000000	50.0000000	8.7709472	0.4529298	19.5000000	22.8000000	28.7000000	24.3037919	26.0850081	50.0000000

Fig. 3

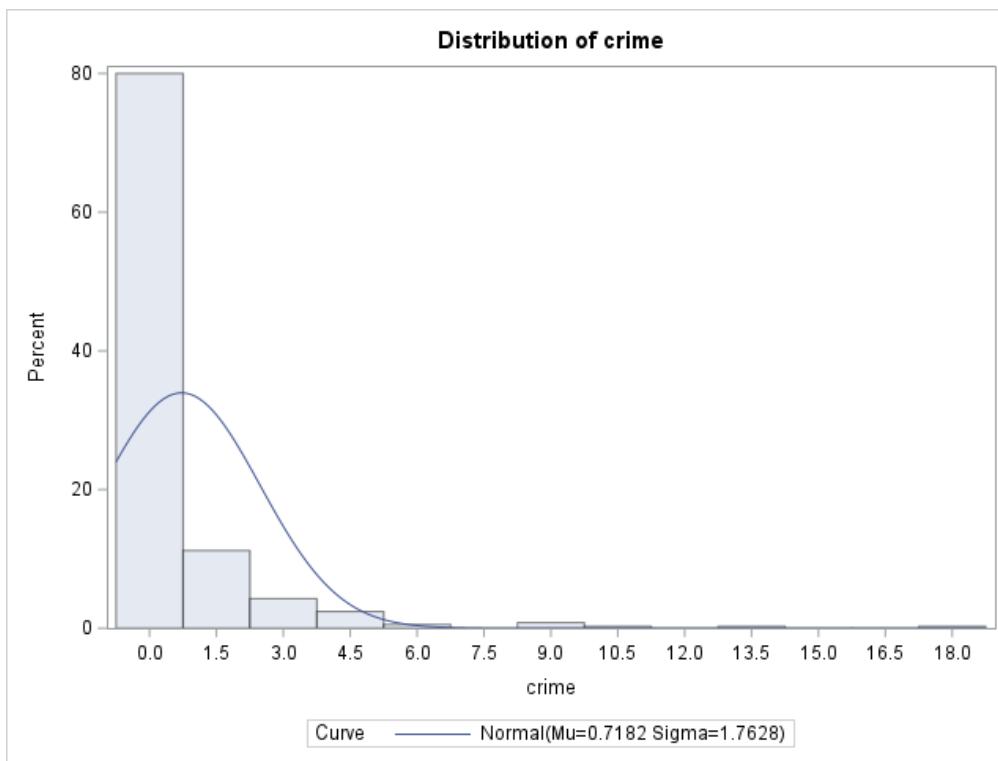


Fig. 4

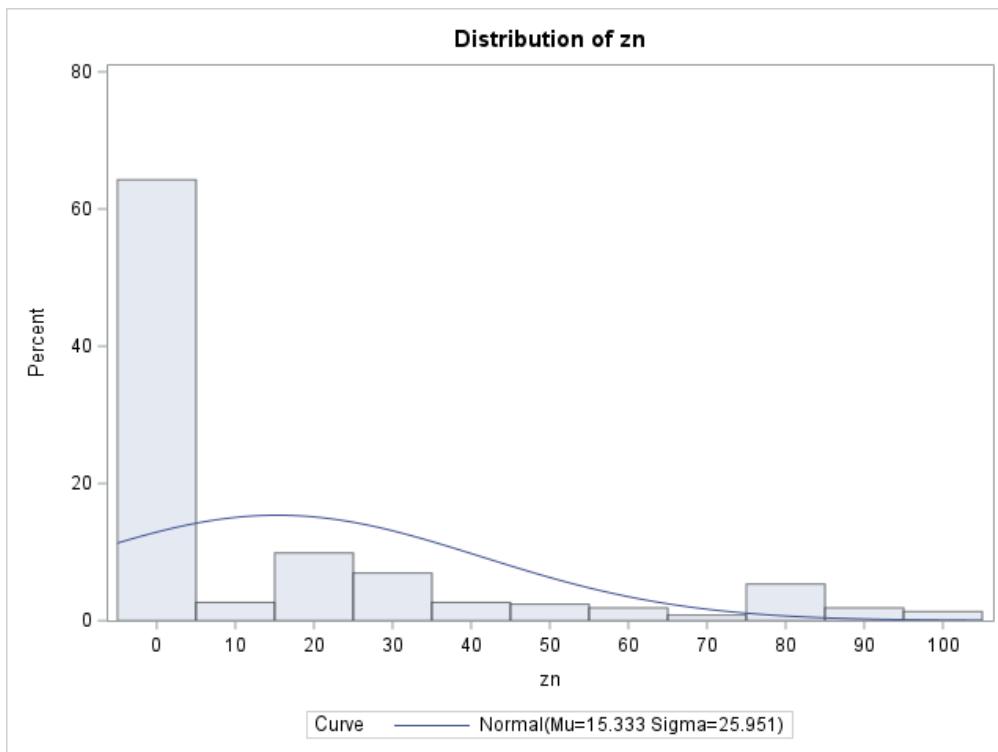


Fig. 5

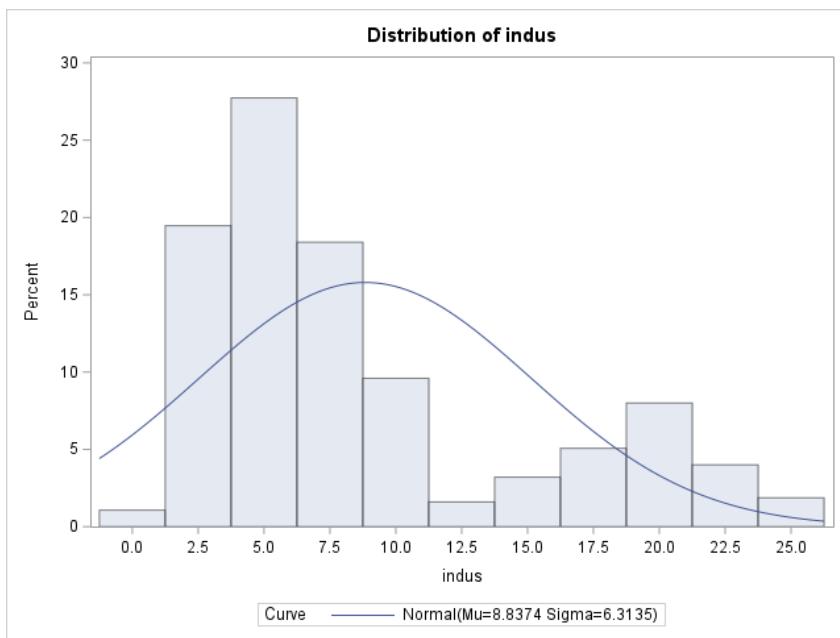


Fig. 6

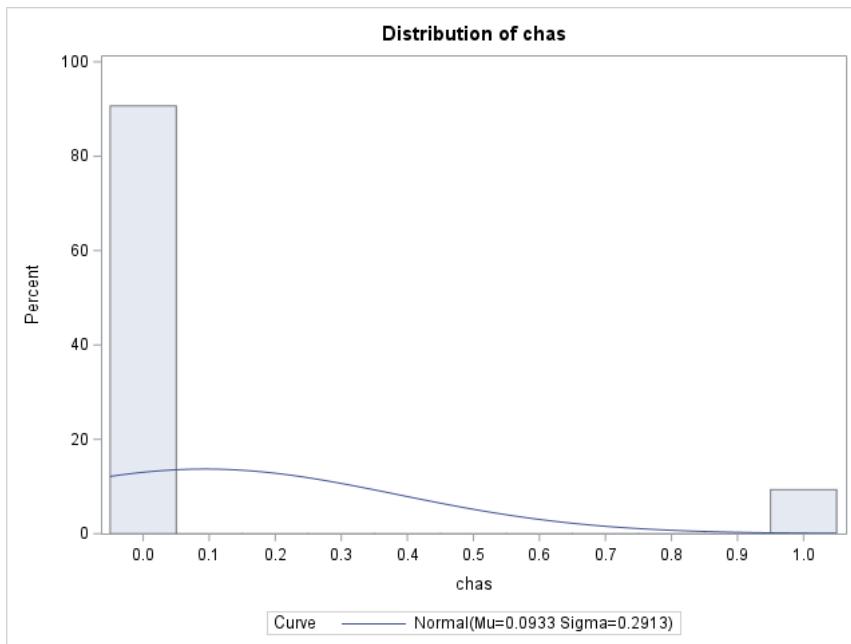


Fig. 7

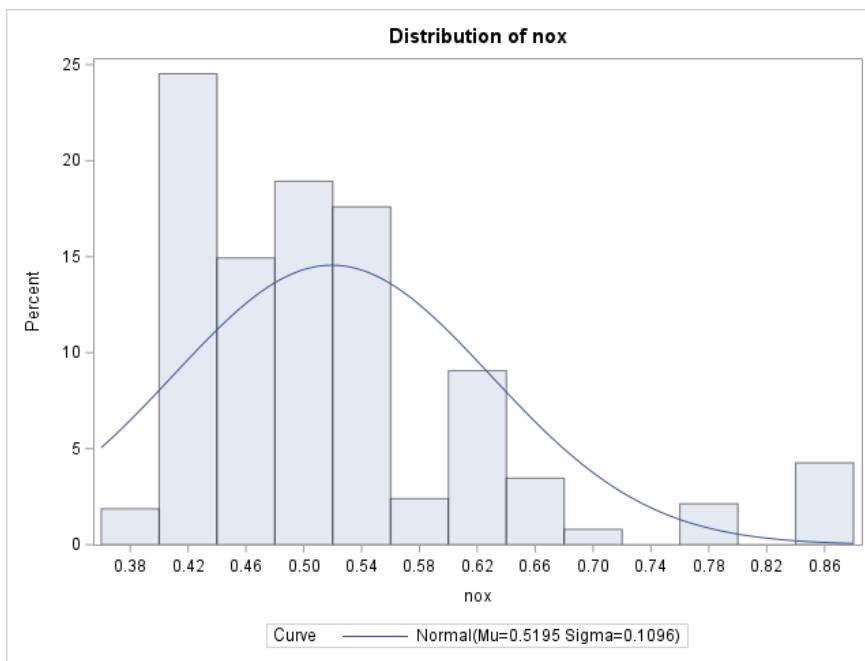


Fig. 8

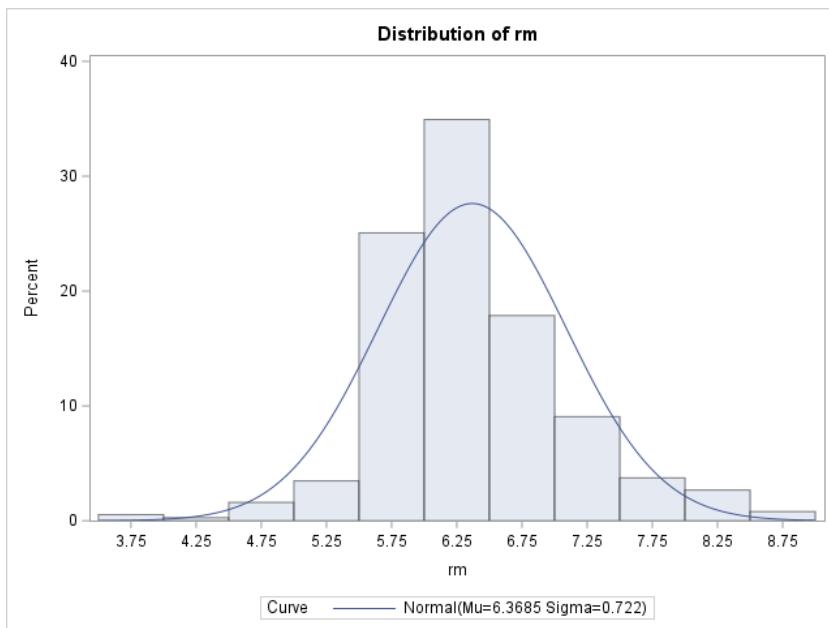


Fig. 9

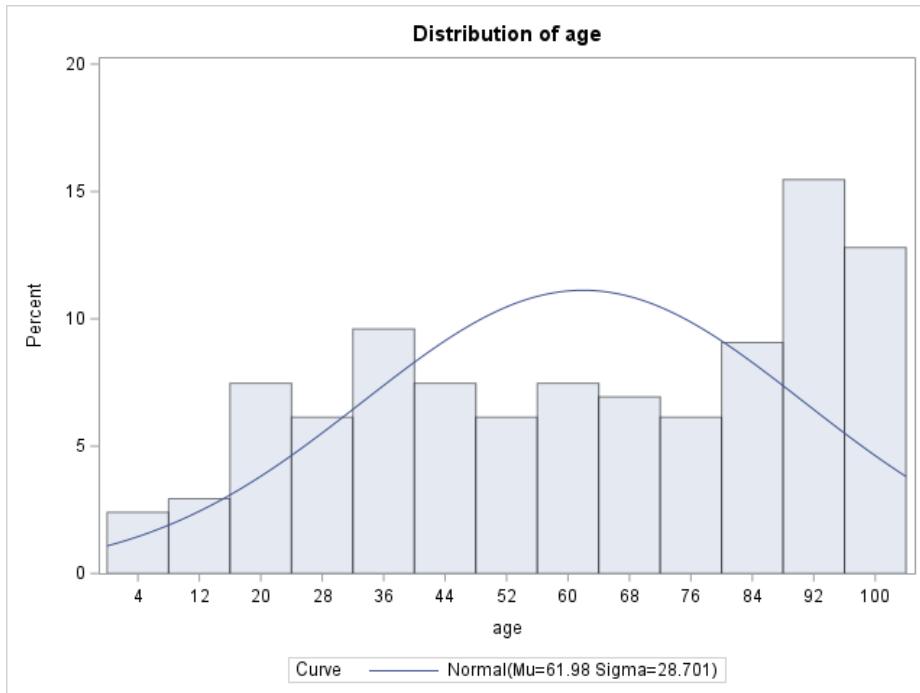


Fig. 10

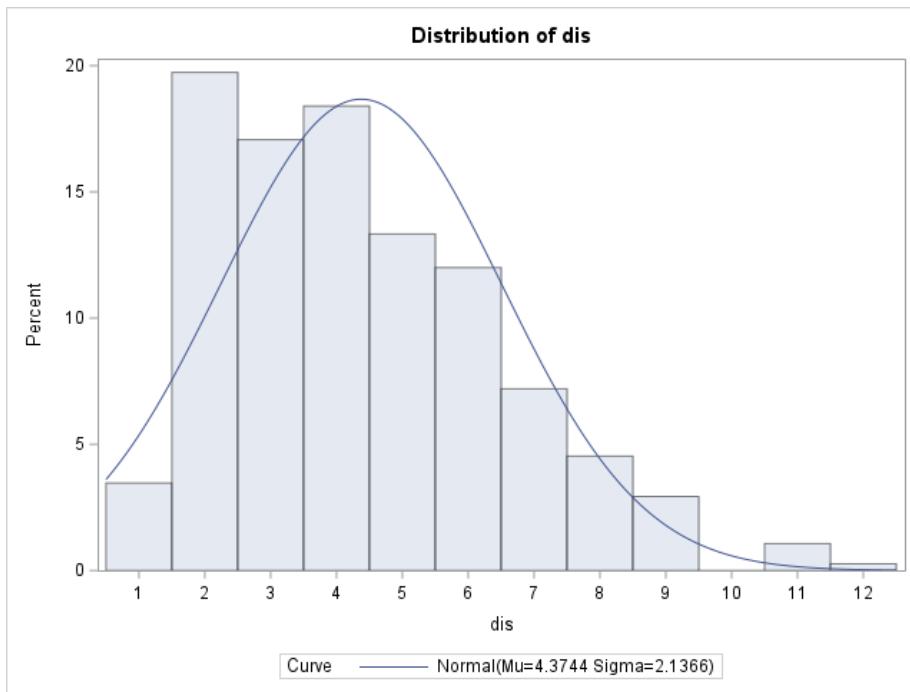


Fig. 11

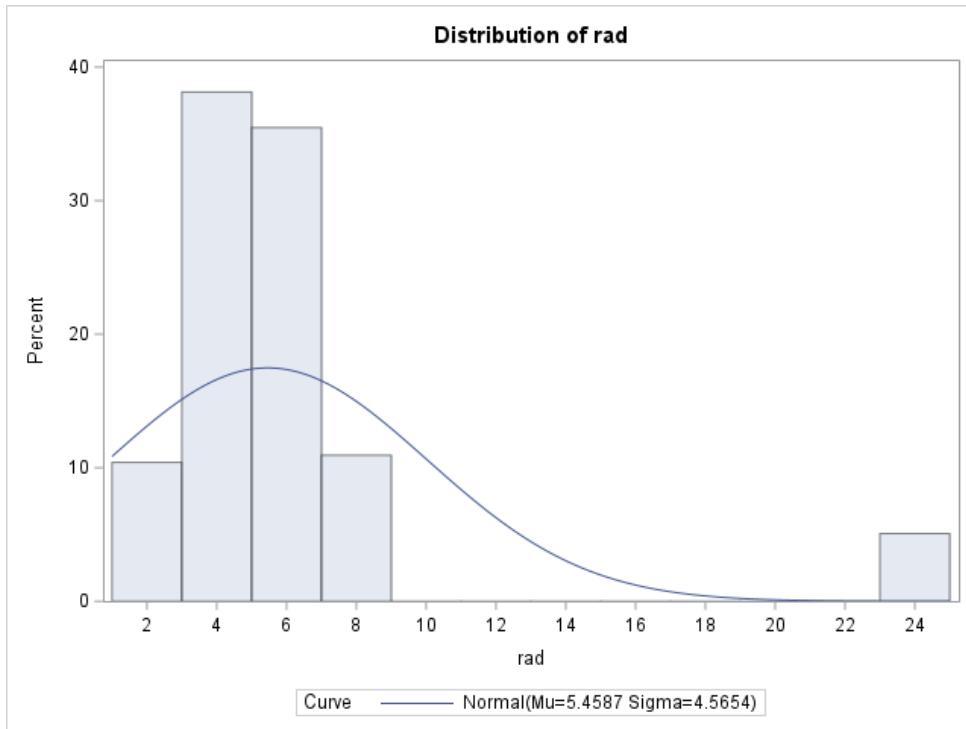


Fig. 12

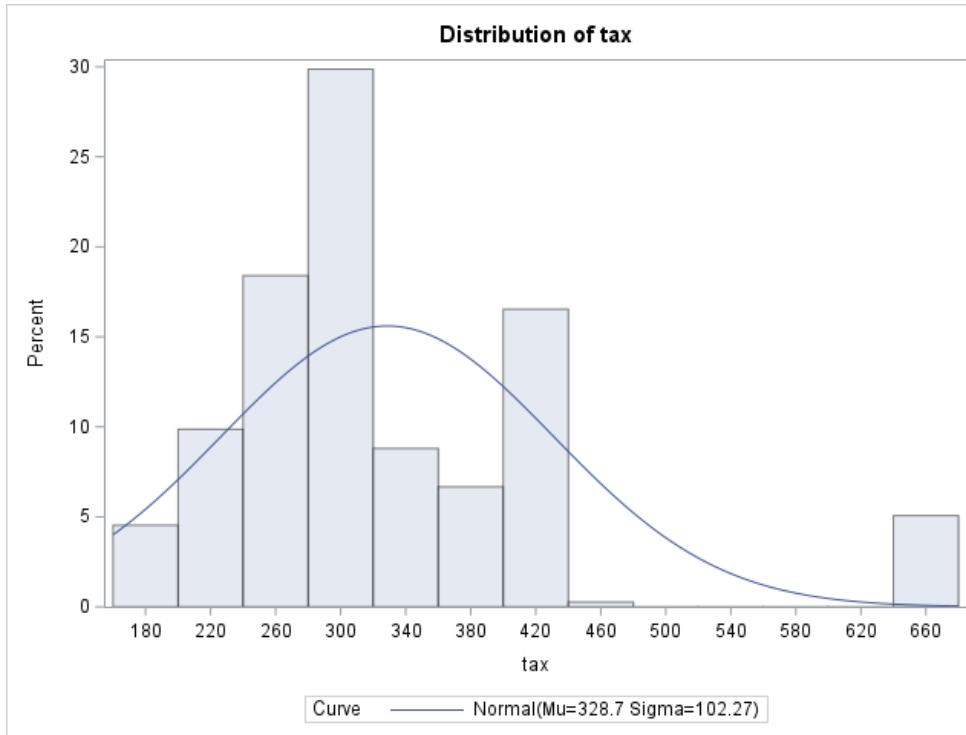


Fig. 13

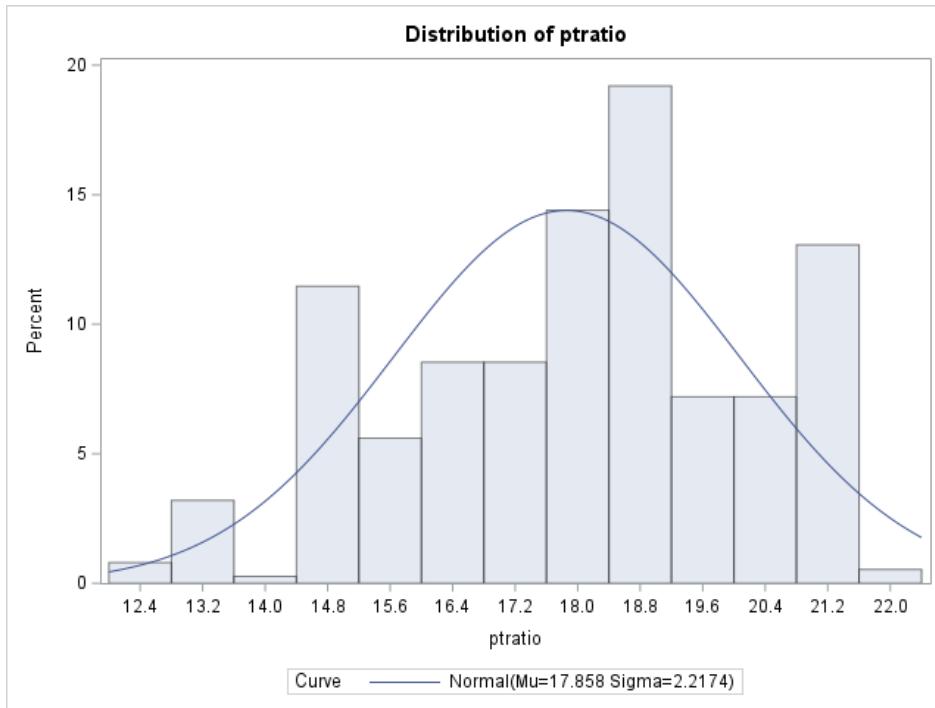


Fig. 14

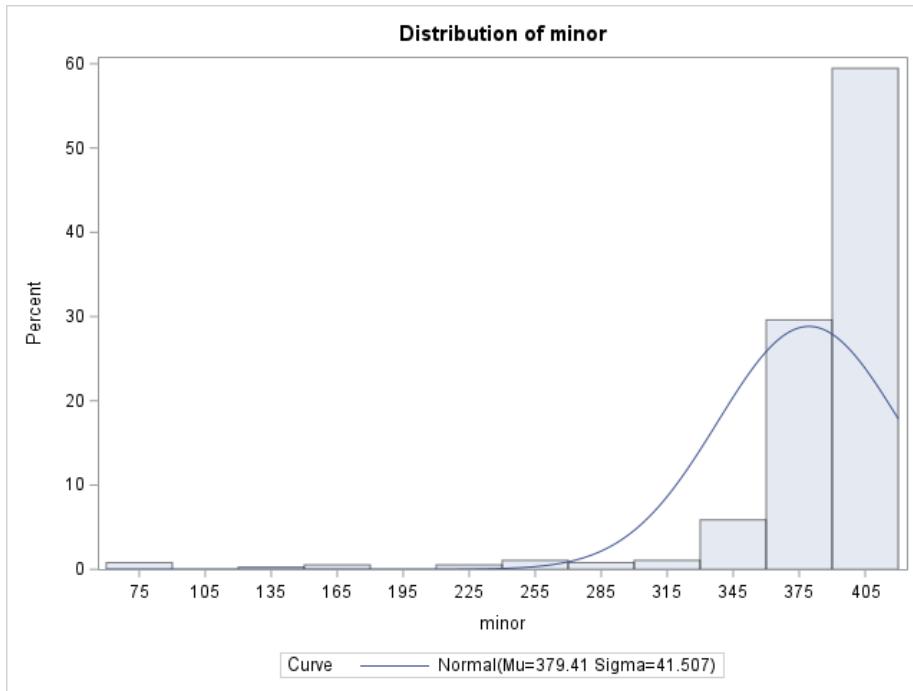


Fig. 15

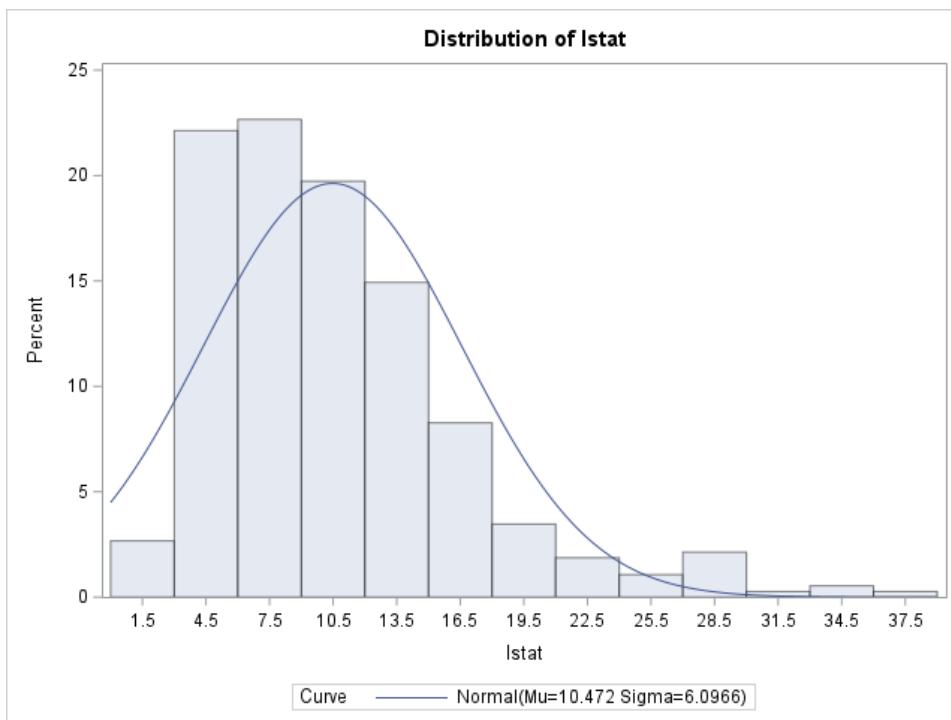


Fig. 16

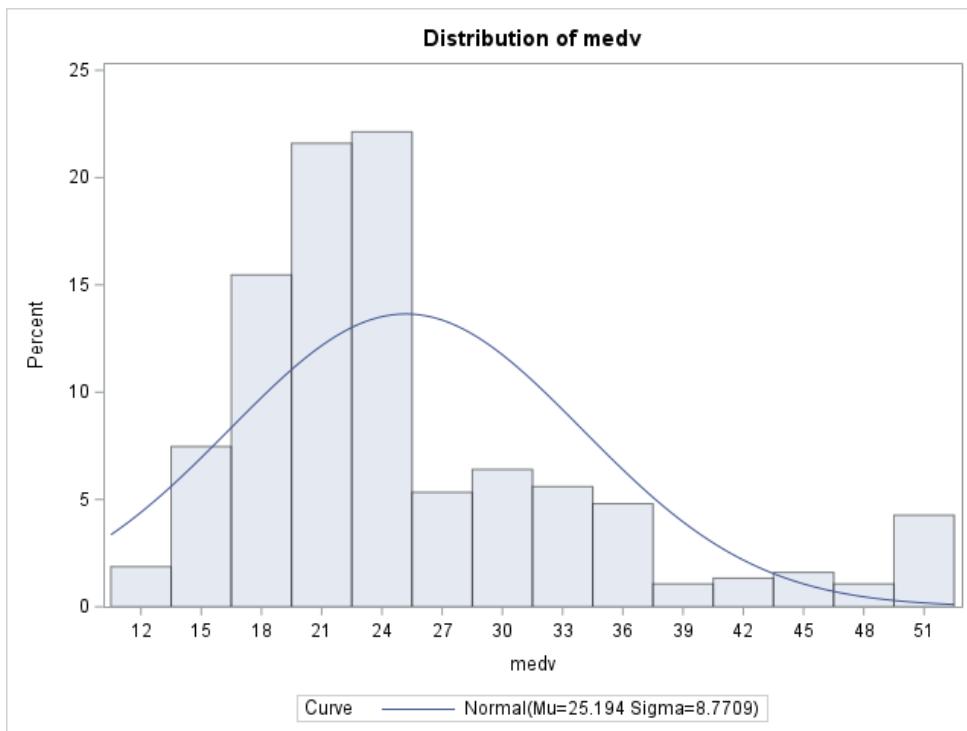


Fig. 17

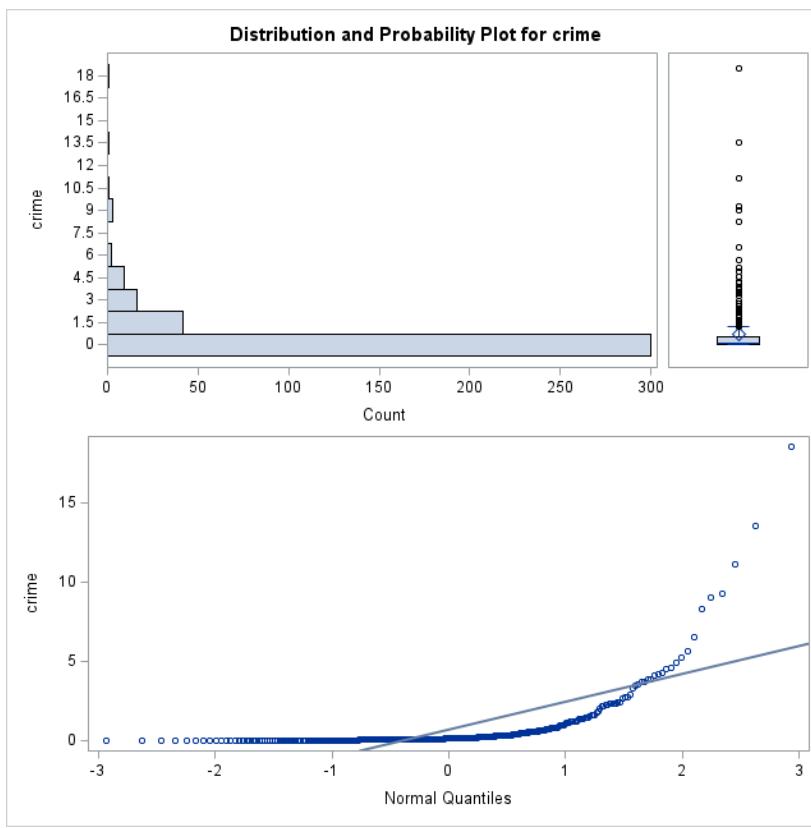


Fig. 18

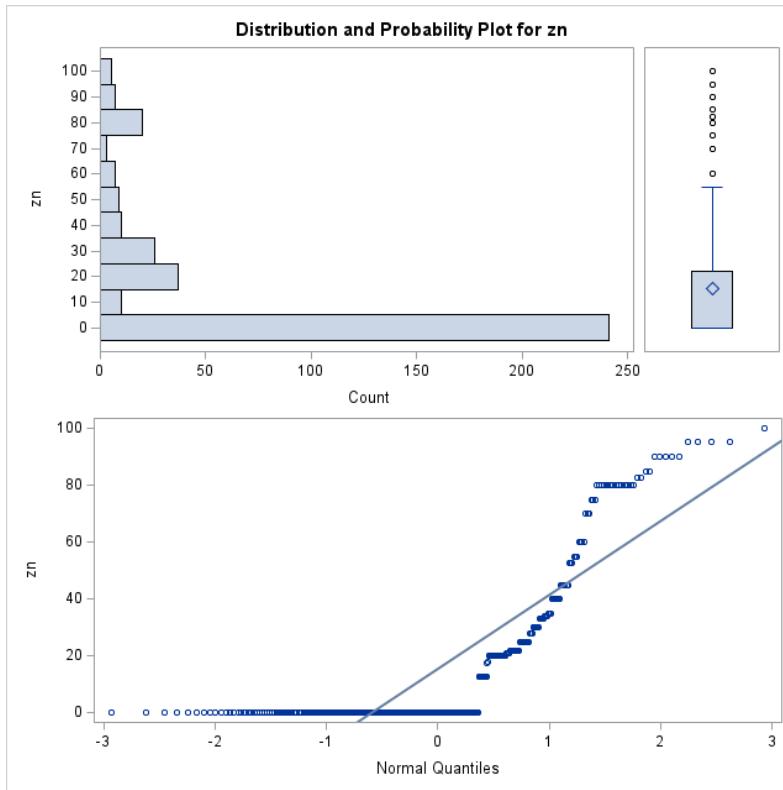


Fig. 19

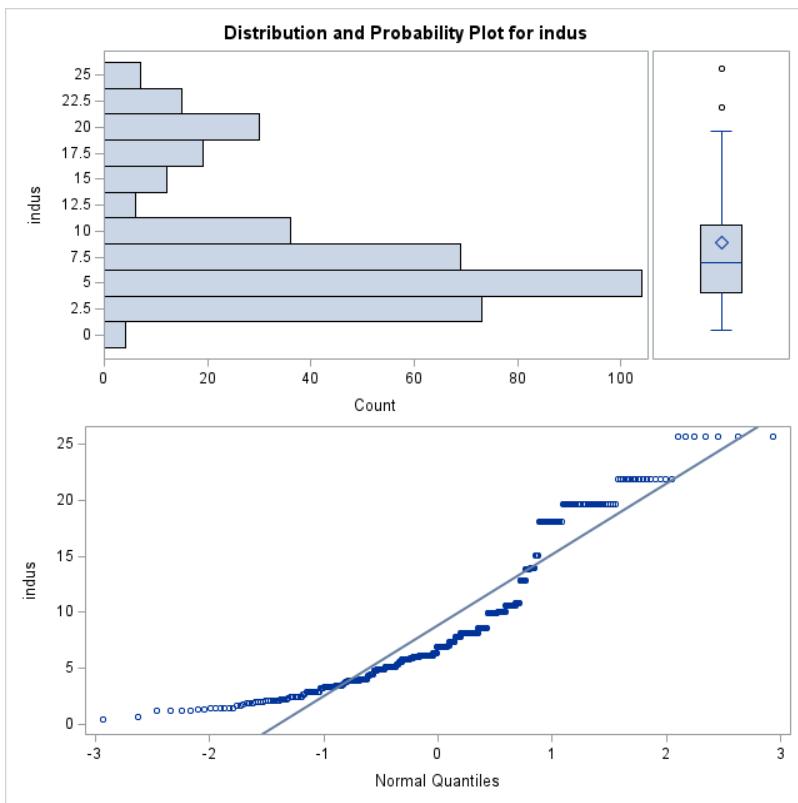


Fig. 20

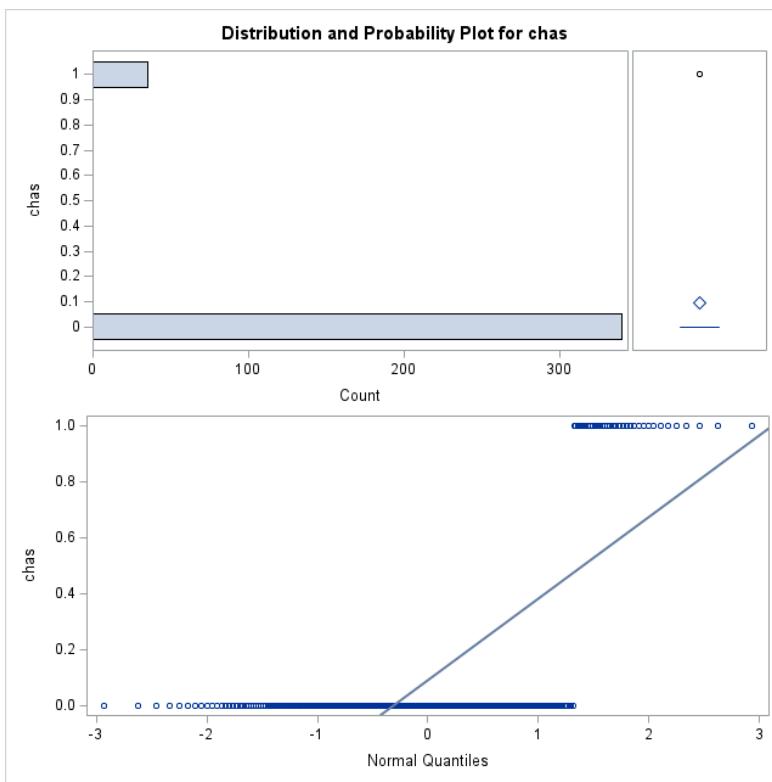


Fig. 21

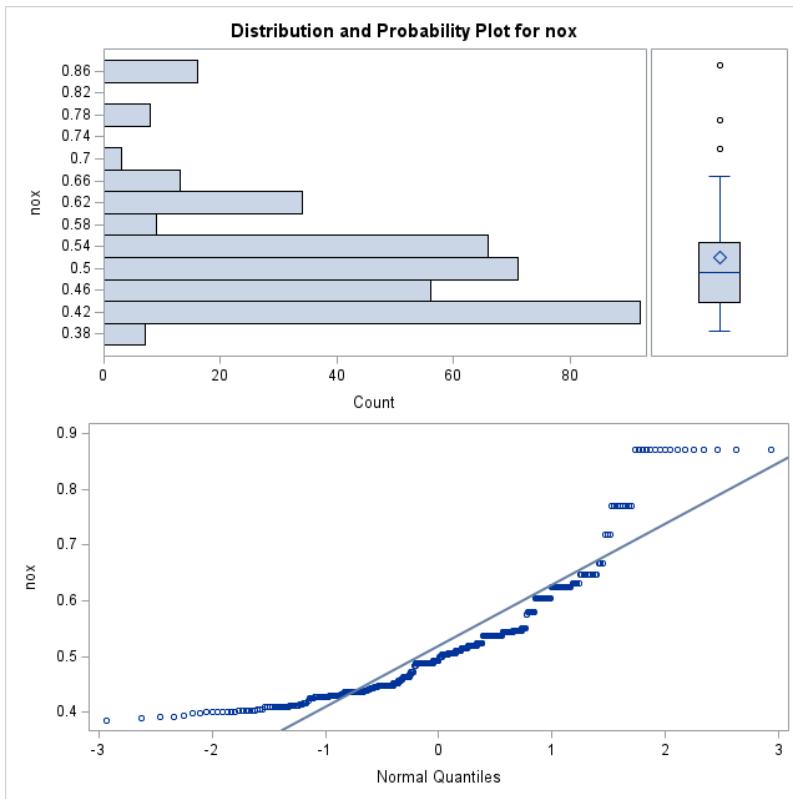


Fig. 22

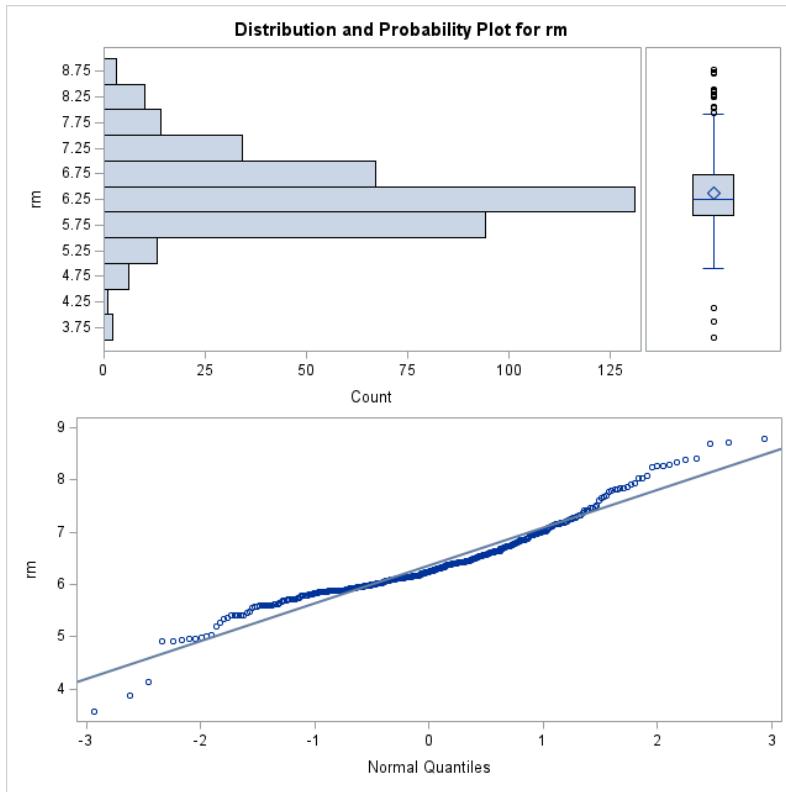


Fig. 23

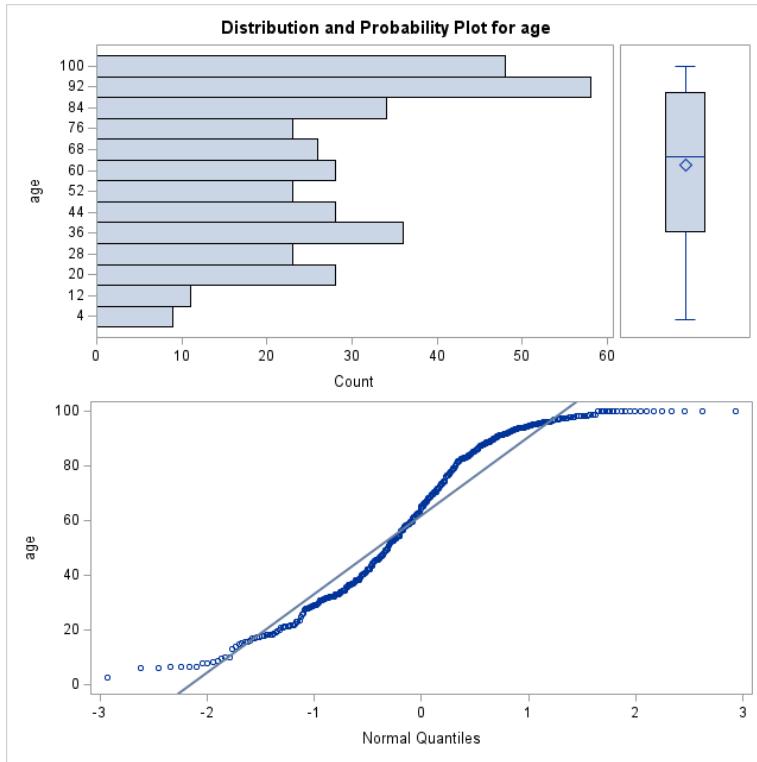


Fig. 24

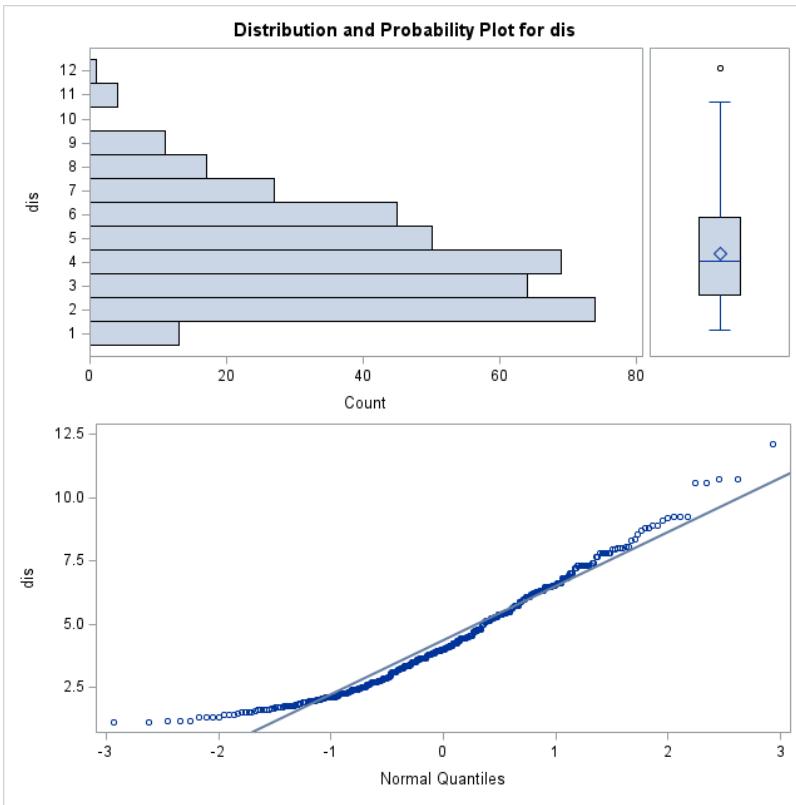


Fig. 25

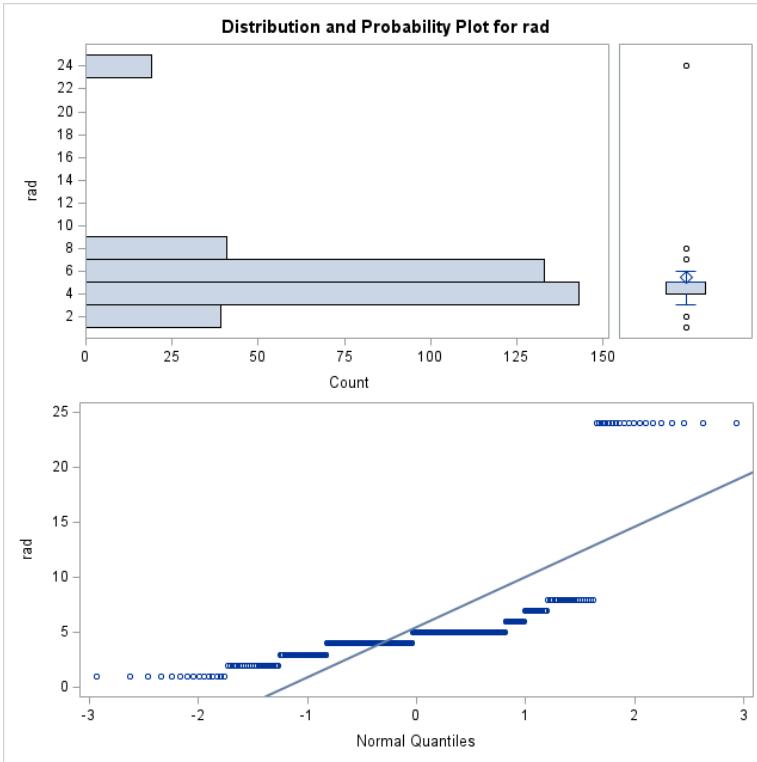


Fig. 26

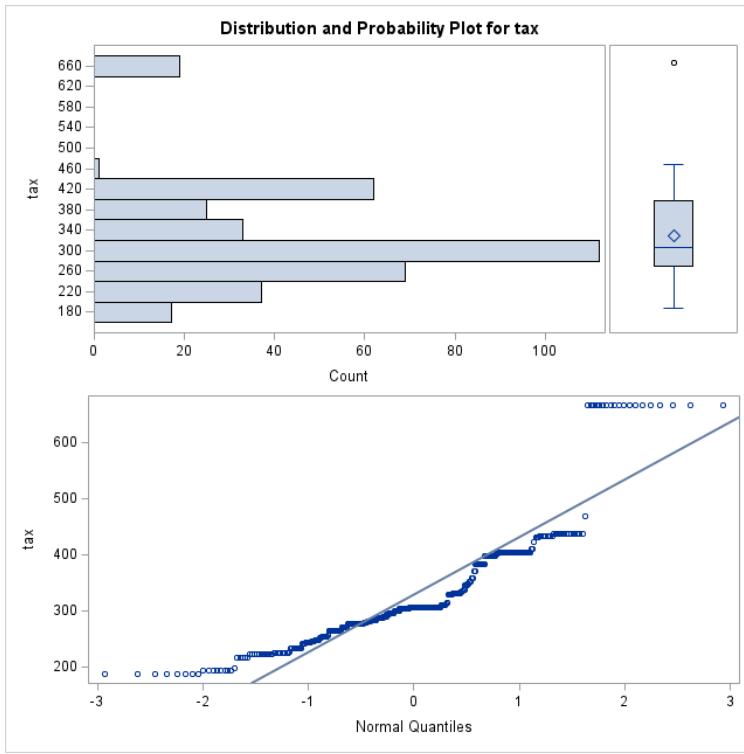


Fig. 27

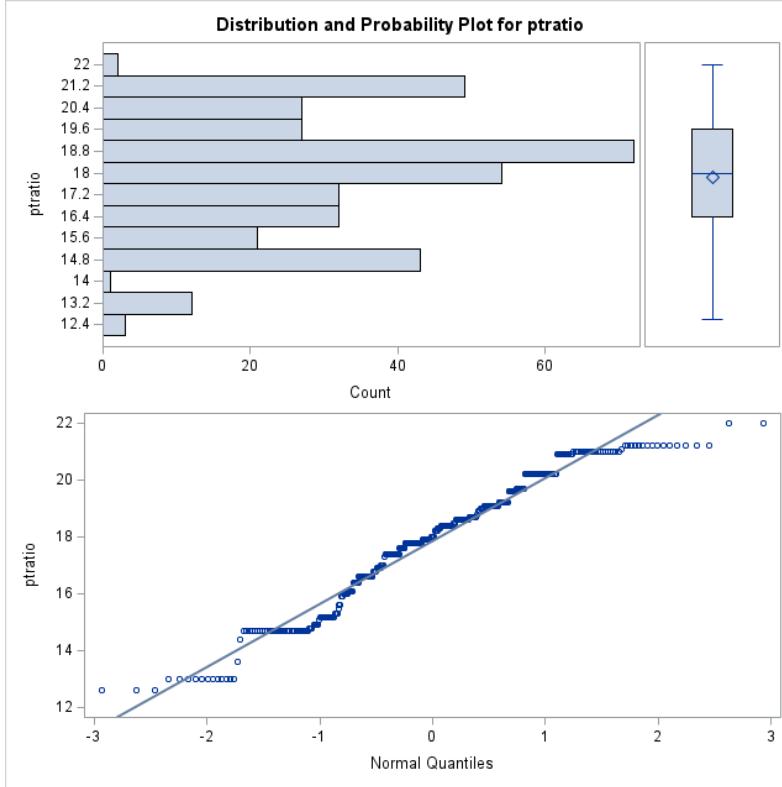


Fig. 28

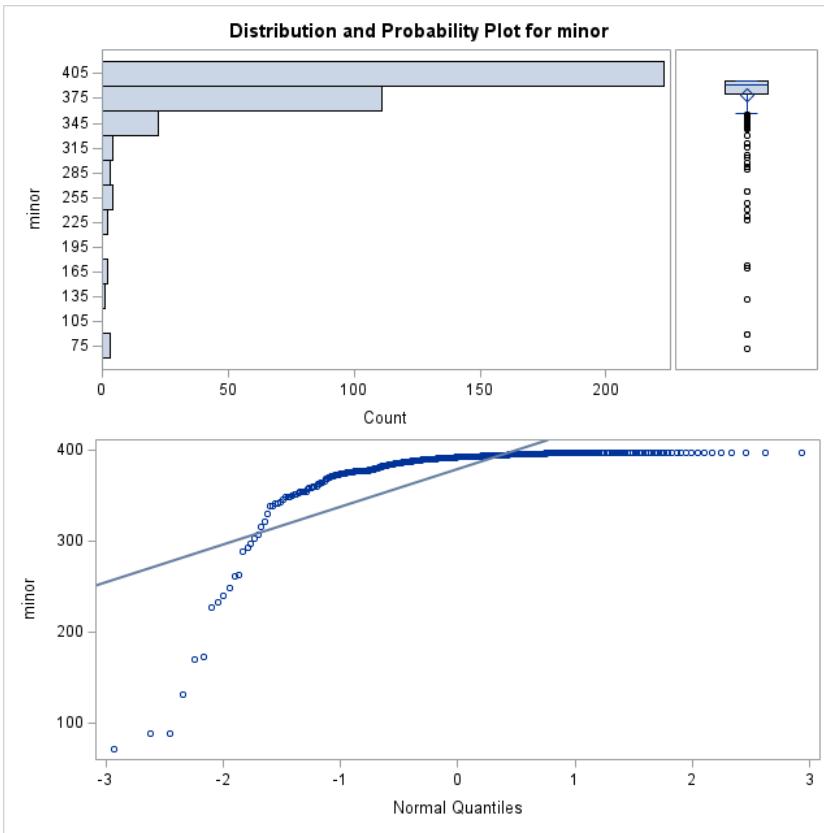


Fig. 29

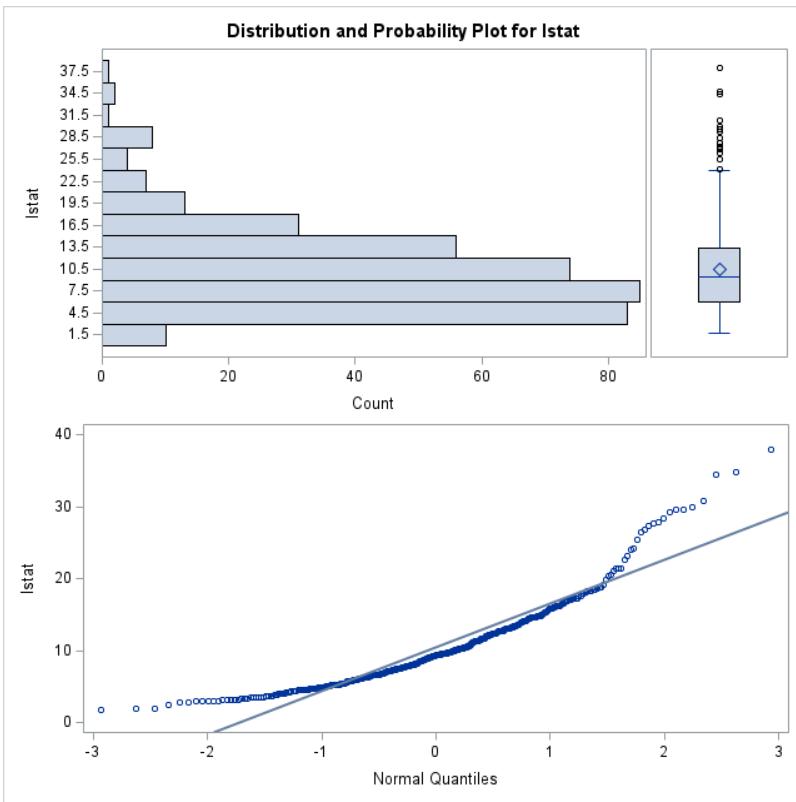


Fig. 30

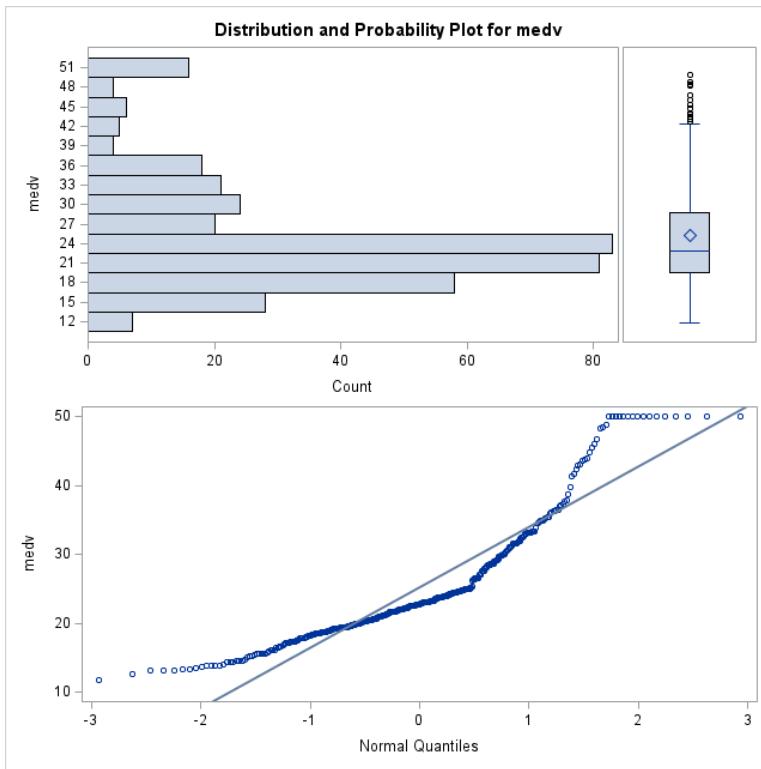


Fig. 31

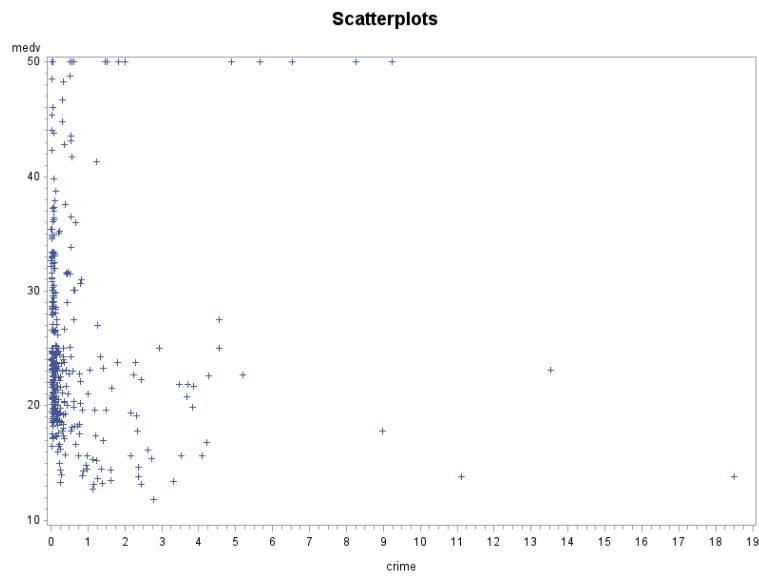


Fig. 32

Scatterplots

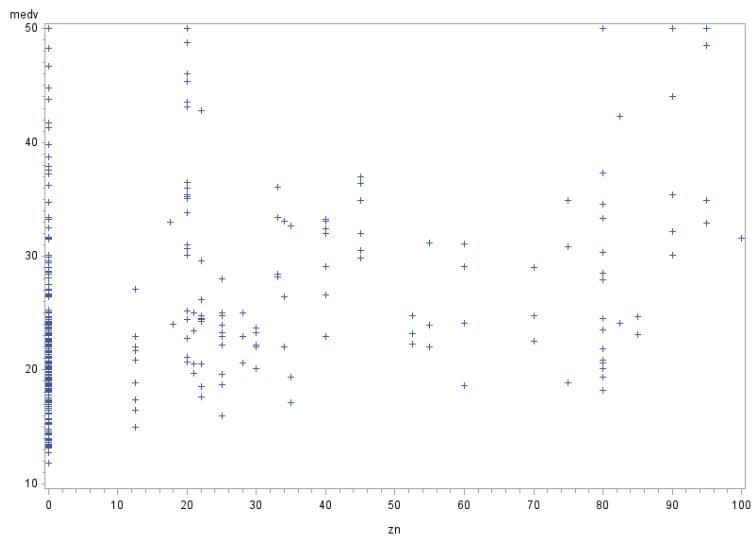


Fig. 33

Scatterplots

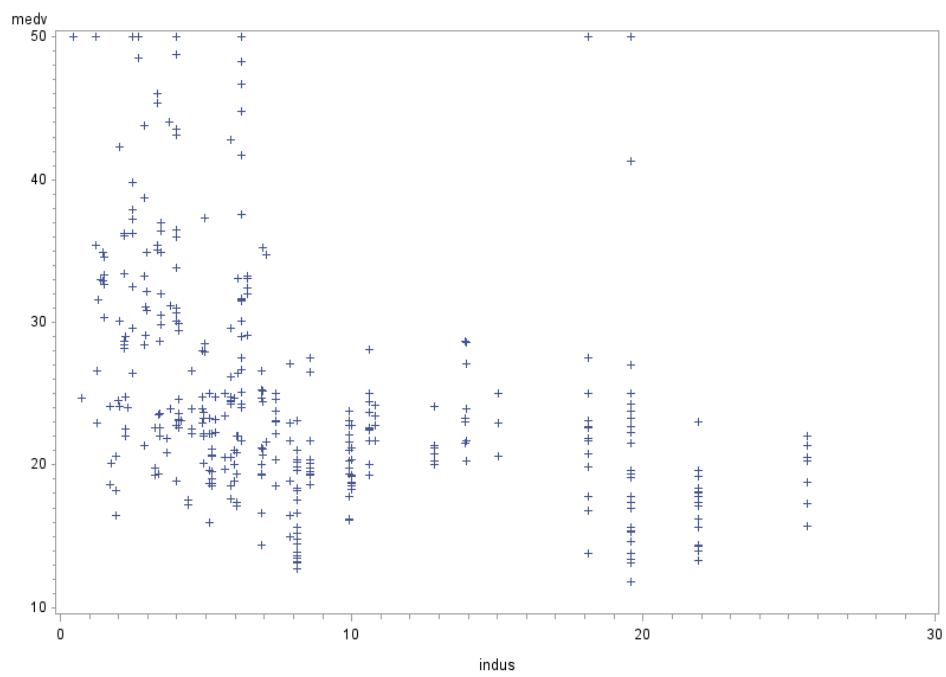


Fig. 34

Scatterplots

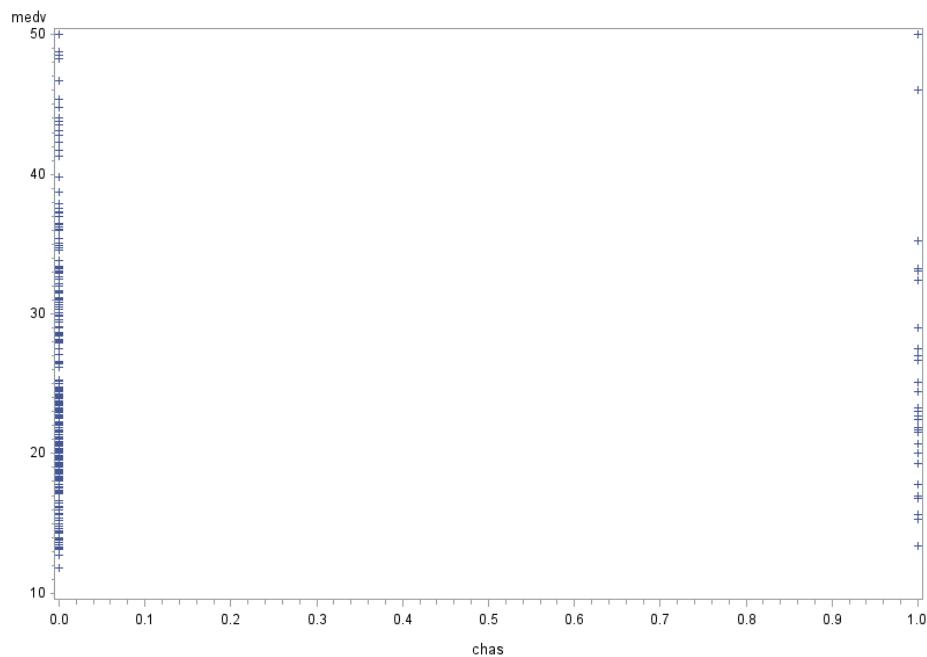


Fig. 35

Scatterplots

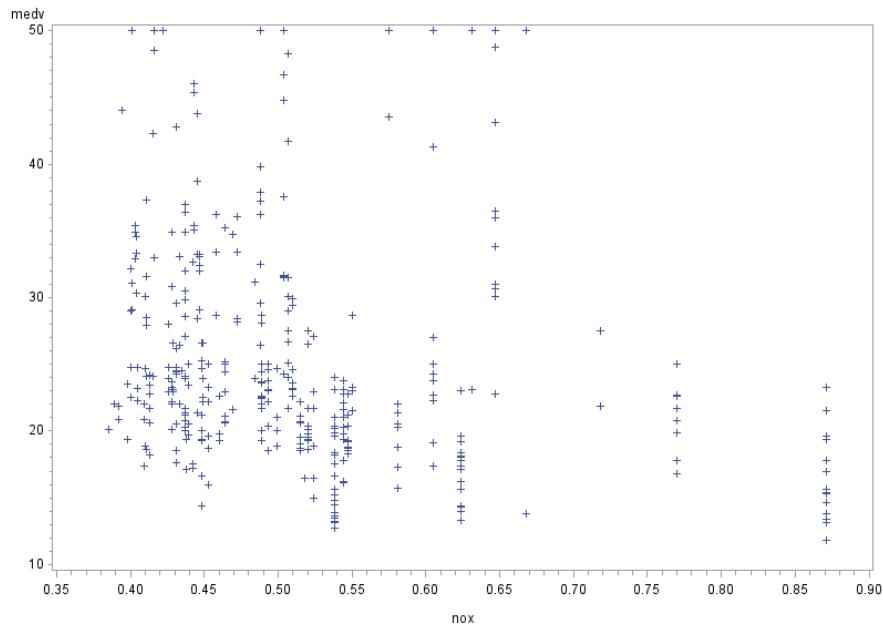


Fig. 36

Scatterplots

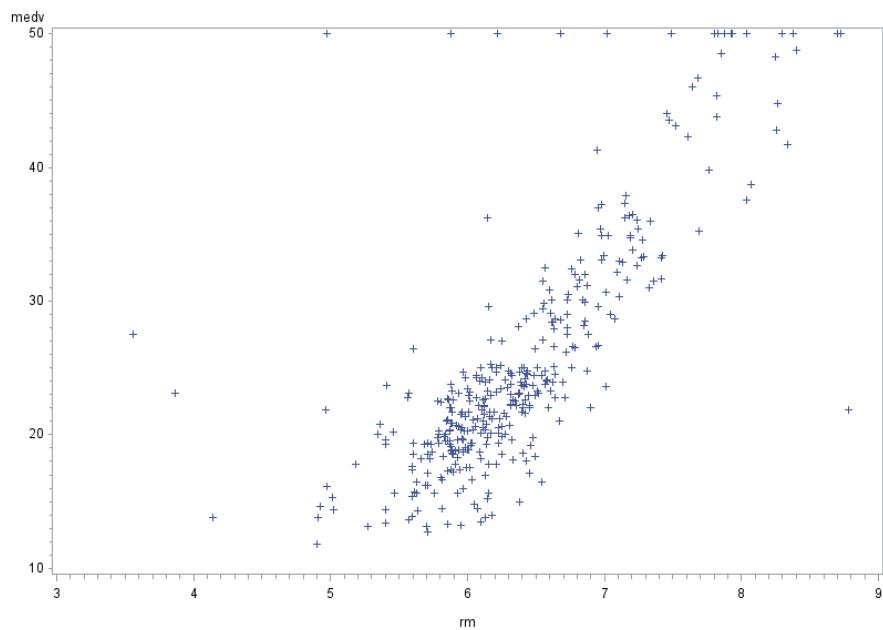


Fig. 37

Scatterplots

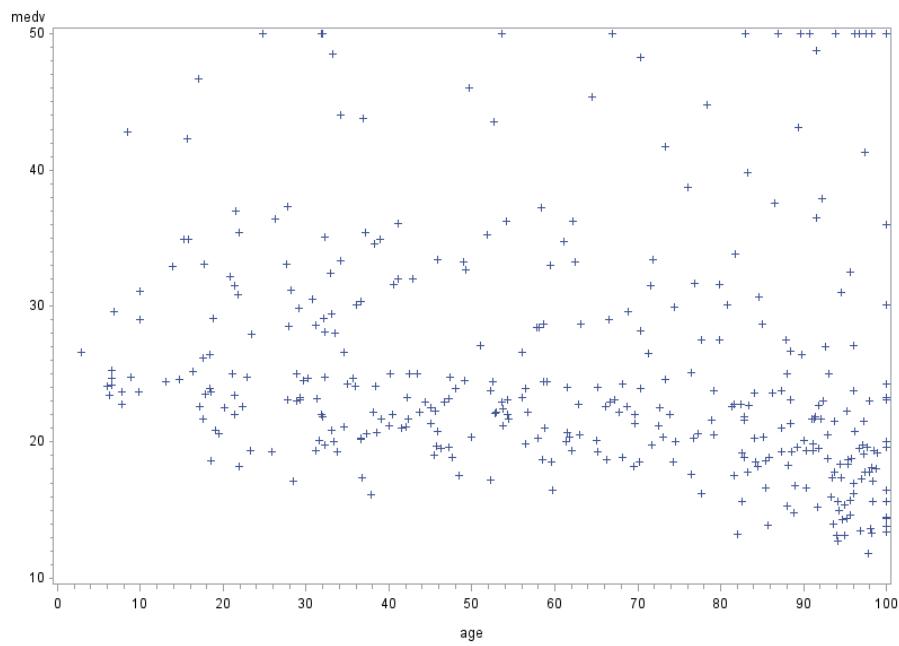


Fig. 38

Scatterplots

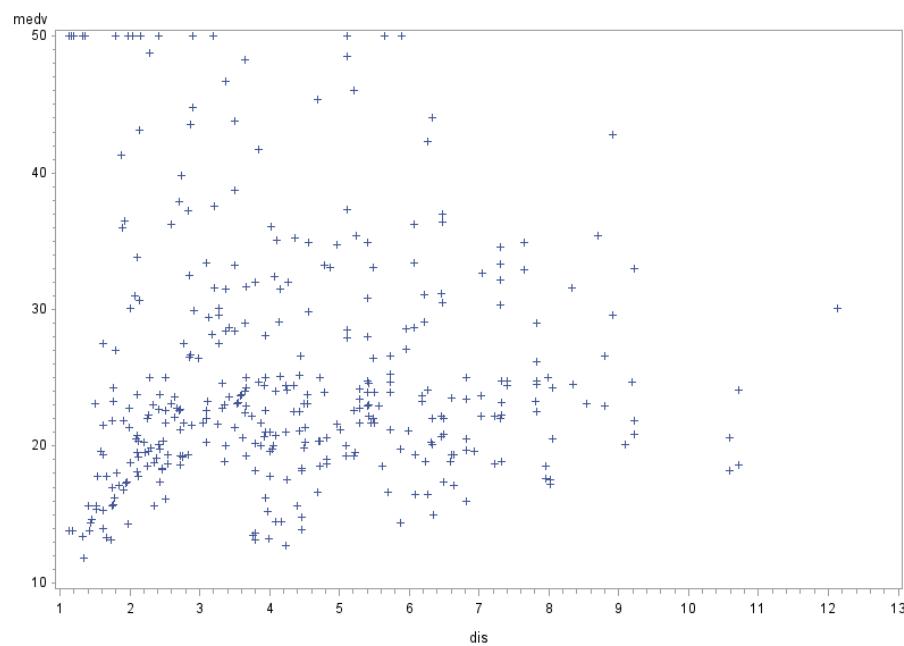


Fig. 39

Scatterplots

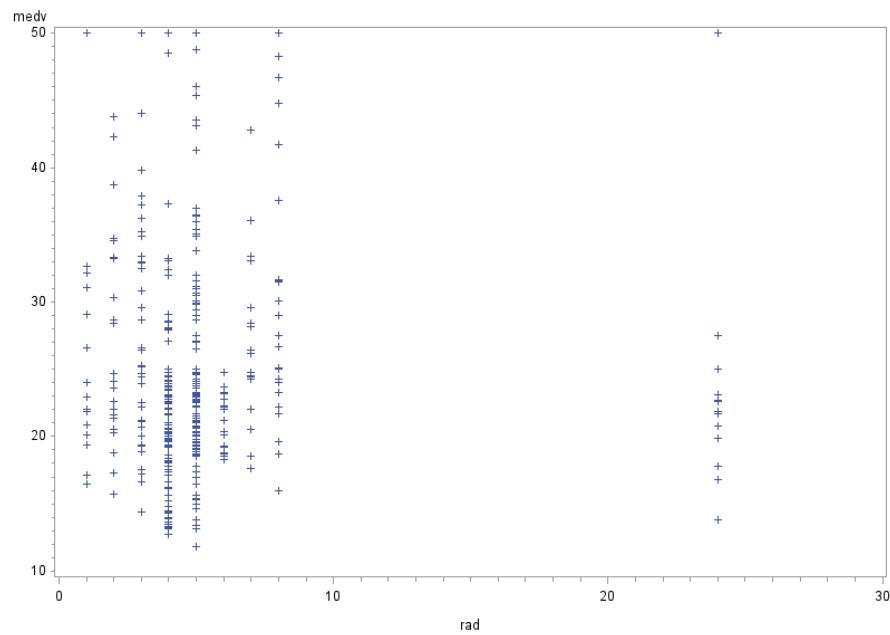


Fig. 40

Scatterplots

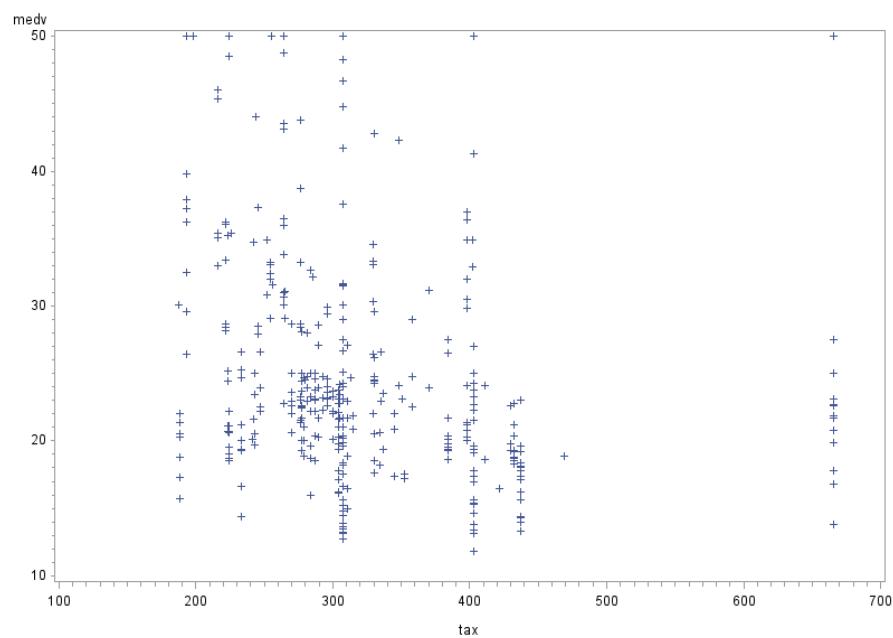


Fig. 41

Scatterplots

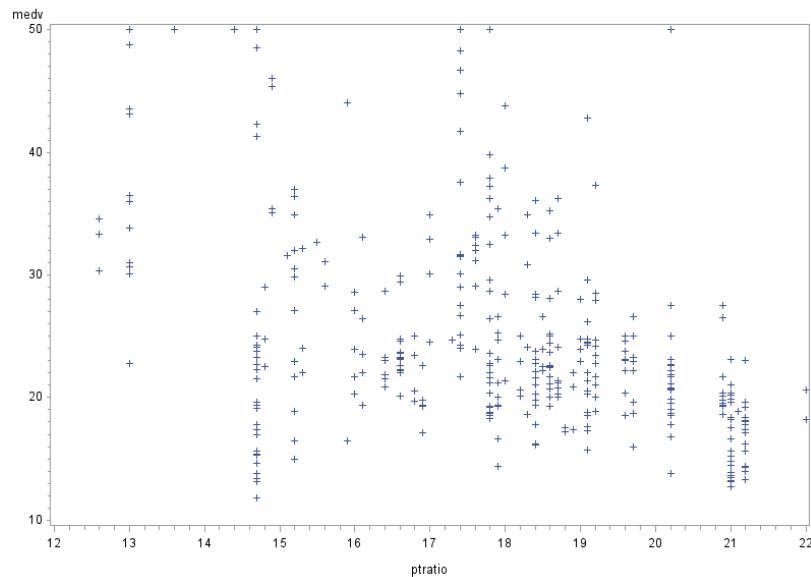


Fig. 42

Scatterplots

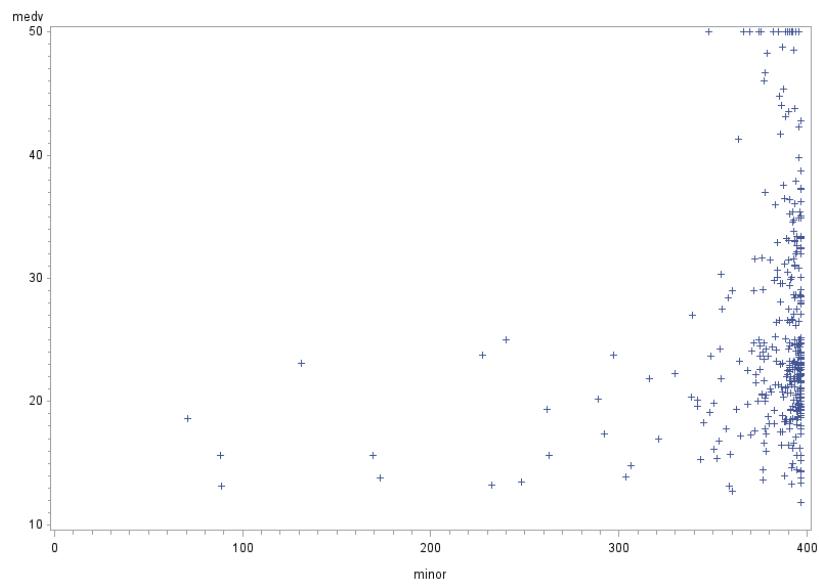


Fig. 43

Scatterplots

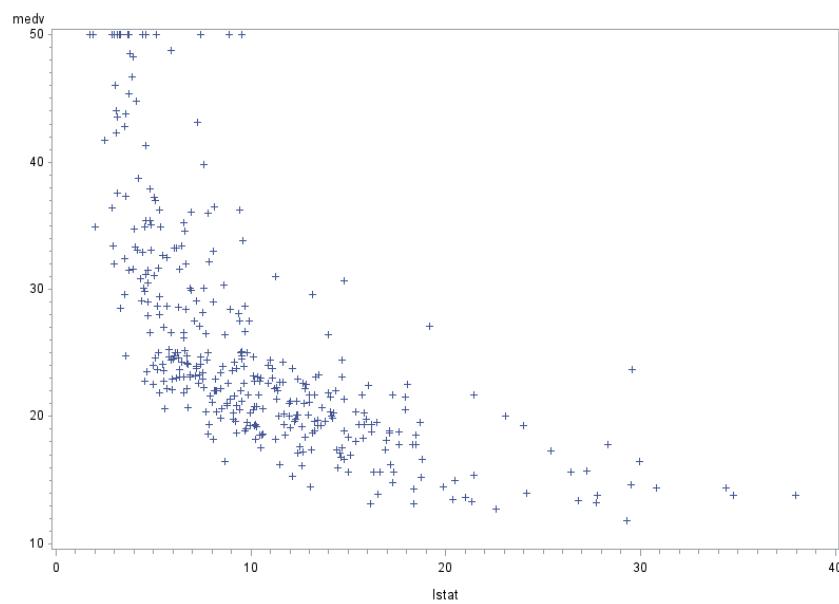


Fig. 44

Pearson Correlation Coefficients, N = 375 Prob > r under H0: Rho=0														
	crime	zn	indus	chas	nox	rm	age	dis	rad	tax	ptratio	minor	Istat	medv
crime	1.00000 <.0001	-0.21399 <.0001	0.45386 <.0001	0.20657 <.0001	0.53360 <.0001	-0.31044 <.0001	0.37829 <.0001	-0.39877 <.0001	0.75936 <.0001	0.69373 <.0001	0.11588 0.0248	-0.30261 <.0001	0.32048 <.0001	-0.01481 0.7749
zn	-0.21399 <.0001	1.00000 <.0001	-0.48816 0.0681	-0.09433 <.0001	-0.47234 <.0001	0.29961 <.0001	-0.54172 <.0001	0.63397 <.0001	-0.19114 0.0002	-0.16831 0.0011	-0.30445 <.0001	0.14393 0.0052	-0.37513 <.0001	0.27888 <.0001
indus	0.45386 <.0001	-0.48816 <.0001	1.00000 0.0001	0.19753 <.0001	0.71166 <.0001	-0.38523 <.0001	0.56979 <.0001	-0.63166 <.0001	0.32548 <.0001	0.51616 <.0001	0.16211 0.0016	-0.31283 <.0001	0.48576 <.0001	-0.28576 <.0001
chas	0.20657 <.0001	-0.09433 0.0681	0.19753 0.0001	1.00000 <.0001	0.21678 0.1939	0.06723 0.0007	0.17373 0.0001	-0.20220 <.0001	0.27133 <.0001	0.18081 0.0004	-0.05313 0.3048	-0.04966 0.3376	0.04054 0.4338	0.11888 0.0213
nox	0.53360 <.0001	-0.47234 <.0001	0.71166 <.0001	0.21678 <.0001	1.00000 <.0001	-0.27901 <.0001	0.68604 <.0001	-0.72091 <.0001	0.41944 <.0001	0.53015 <.0001	-0.07488 0.1479	-0.39915 <.0001	0.45703 <.0001	-0.19826 0.0001
rm	-0.31044 <.0001	0.29961 <.0001	-0.38523 0.1939	0.06723 <.0001	-0.27901 0.0153	1.00000 0.0036	-0.20323 0.0001	0.12517 0.0153	-0.15020 0.0036	-0.26861 <.0001	-0.34207 0.0001	0.20477 0.0001	-0.64347 0.0001	0.77073 <.0001
age	0.37829 <.0001	-0.54172 <.0001	0.56979 <.0001	0.17373 0.0007	0.68604 <.0001	-0.20323 <.0001	1.00000 0.0001	-0.70018 <.0001	0.27608 <.0001	0.34772 0.0001	0.09062 0.0797	-0.23857 0.0001	0.53111 0.0001	-0.19727 0.0001
dis	-0.39877 <.0001	0.63397 <.0001	-0.63166 <.0001	-0.20220 0.0001	-0.72091 0.0153	0.12517 0.0153	-0.70018 0.0001	1.00000 0.0001	-0.29710 0.0001	-0.33734 0.0001	-0.01774 0.7321	0.23309 0.0001	-0.34094 0.0001	-0.00473 0.9273
rad	0.75936 <.0001	-0.19114 0.0002	0.32548 <.0001	0.27133 0.0001	0.41944 0.0036	-0.15020 0.0001	0.27608 0.0001	-0.29710 0.0001	1.00000 0.0001	0.76603 0.0001	0.22047 0.0001	-0.12609 0.0145	0.07382 0.1536	0.10109 0.0504
tax	0.69373 <.0001	-0.16831 0.0011	0.51616 0.0001	0.18081 0.0004	0.53015 0.0001	-0.26861 0.0001	0.34772 0.0001	-0.33734 0.0001	0.76603 0.0001	1.00000 0.0003	0.18418 0.0003	-0.24978 0.0001	0.20910 0.0001	-0.11132 0.0311
ptratio	0.11588 0.0248	-0.30445 <.0001	0.16211 0.0016	-0.05313 0.3048	-0.07488 0.1479	-0.34207 0.0001	0.09062 0.0797	-0.01774 0.7321	0.22047 0.0001	0.18418 0.0003	1.00000 0.0003	0.07069 0.1719	0.20737 0.0001	-0.38236 0.0001
minor	-0.30261 <.0001	0.14393 0.0052	-0.31283 <.0001	-0.04966 0.3376	-0.39915 0.0001	0.20477 0.0001	-0.23857 0.0001	0.23309 0.0001	-0.12609 0.0145	-0.24978 0.0001	0.07069 0.1719	1.00000 0.0001	-0.20106 0.0007	0.17426 0.0007
Istat	0.32048 <.0001	-0.37513 <.0001	0.48576 0.0001	0.04054 0.4338	0.45703 0.0001	-0.64347 0.0001	0.53111 0.0001	-0.34094 0.0001	0.07382 0.1536	0.20910 0.0001	0.20737 0.0001	-0.20106 0.0001	1.00000 0.0001	-0.65111 0.0001
medv	-0.01481 0.7749	0.27888 <.0001	-0.28576 <.0001	0.11888 0.0213	-0.19826 0.0001	0.77073 0.0001	-0.19727 0.0001	-0.00473 0.9273	0.10109 0.0504	-0.11132 0.0311	-0.38236 0.0001	0.17426 0.0007	-0.65111 0.0001	1.00000

APPENDIX B – Full Model

Table 1

Full model - Model assumptions and diagnostics

The REG Procedure

Model: MODEL1

Dependent Variable: medv

Number of Observations Read	375
Number of Observations Used	375

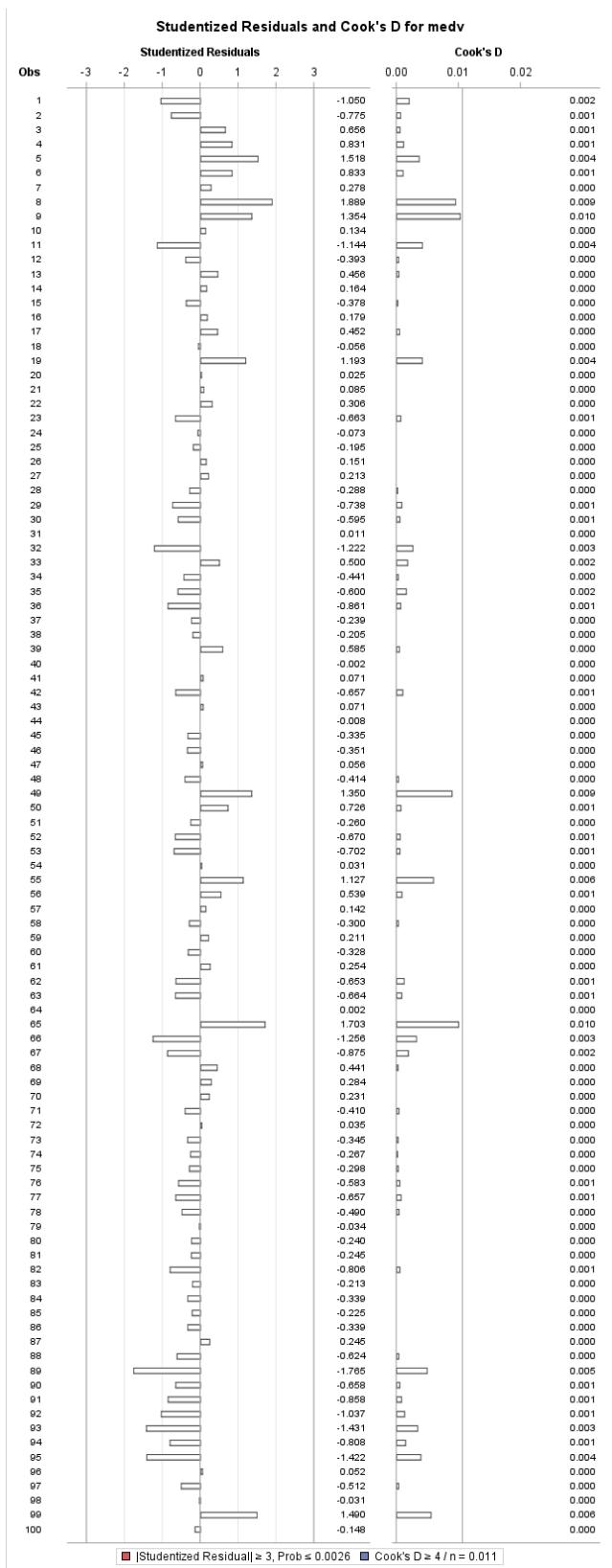
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	13	22446	1726.60036	98.53	<.0001
Error	361	6325.83351	17.52308		

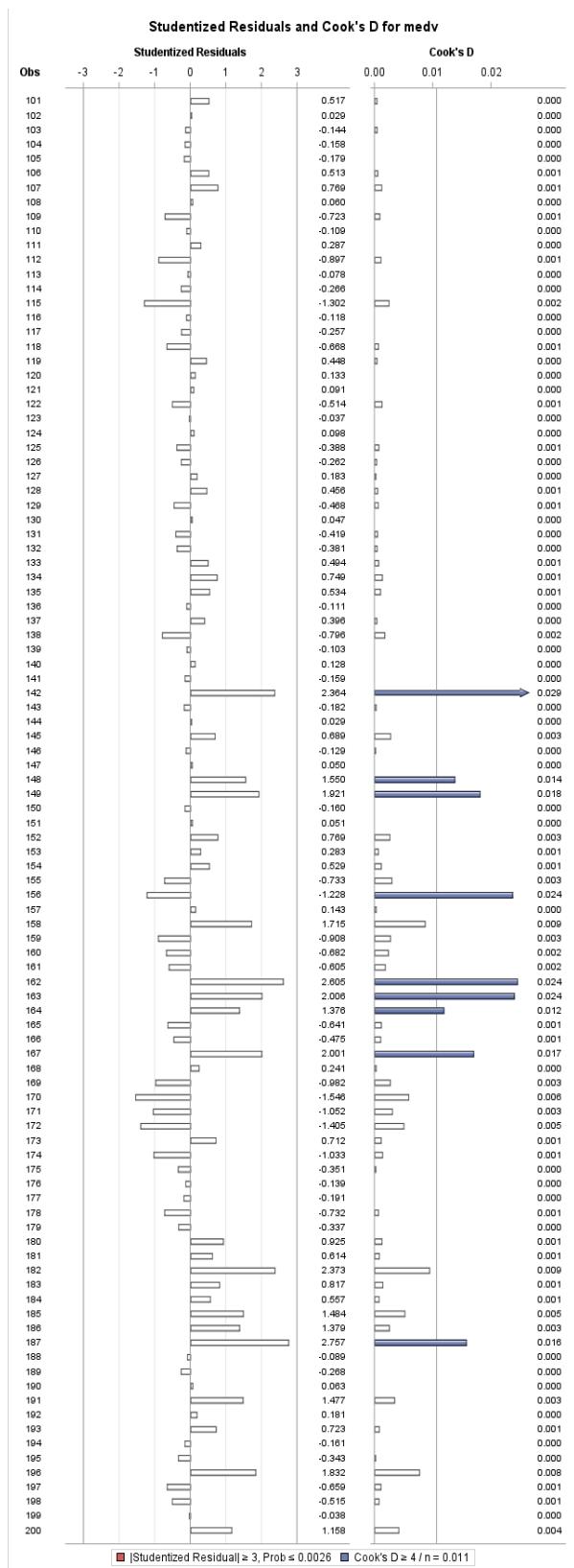
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Corrected Total	374	28772			

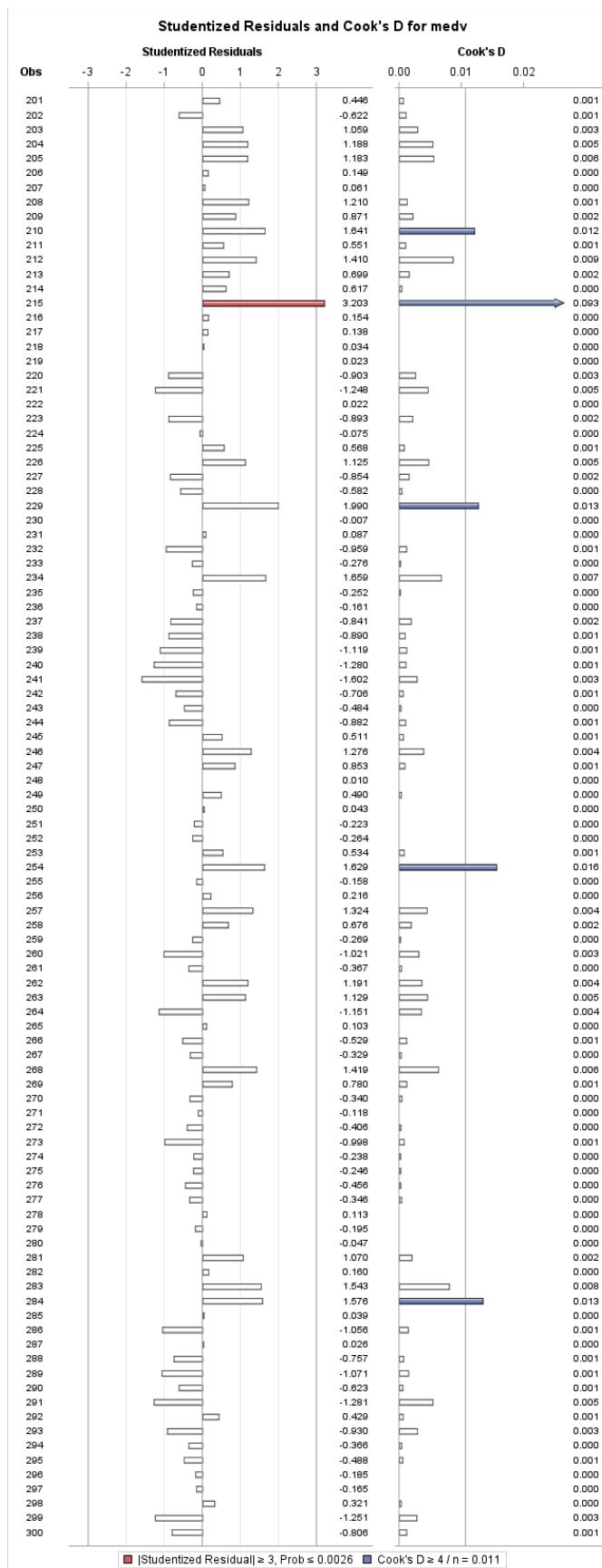
Root MSE	4.18606	R-Square	0.7801
Dependent Mean	25.19440	Adj R-Sq	0.7722
Coeff Var	16.61503		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	12.63324	5.55329	2.27	0.0235	0	0
crime	1	1.28883	0.22380	5.76	<.0001	0.25903	3.32187
zn	1	0.03579	0.01241	2.88	0.0042	0.10588	2.21198
indus	1	0.07362	0.05727	1.29	0.1994	0.05300	2.79005
chas	1	0.12124	0.79815	0.15	0.8794	0.00403	1.15365
nox	1	-17.35683	3.98974	-4.35	<.0001	-0.21691	4.08185
rm	1	6.41932	0.44839	14.32	<.0001	0.52844	2.23708
age	1	-0.00650	0.01269	-0.51	0.6088	-0.02127	2.83139
dis	1	-1.18892	0.18716	-6.35	<.0001	-0.28962	3.41311
rad	1	0.39406	0.09801	4.02	<.0001	0.20511	4.27283
tax	1	-0.01565	0.00389	-4.02	<.0001	-0.18253	3.37728
ptratio	1	-0.71652	0.12293	-5.83	<.0001	-0.18114	1.58574
minor	1	0.01167	0.00591	1.97	0.0491	0.05521	1.28351
lstat	1	-0.42714	0.05824	-7.33	<.0001	-0.29690	2.69046

Fig. 1







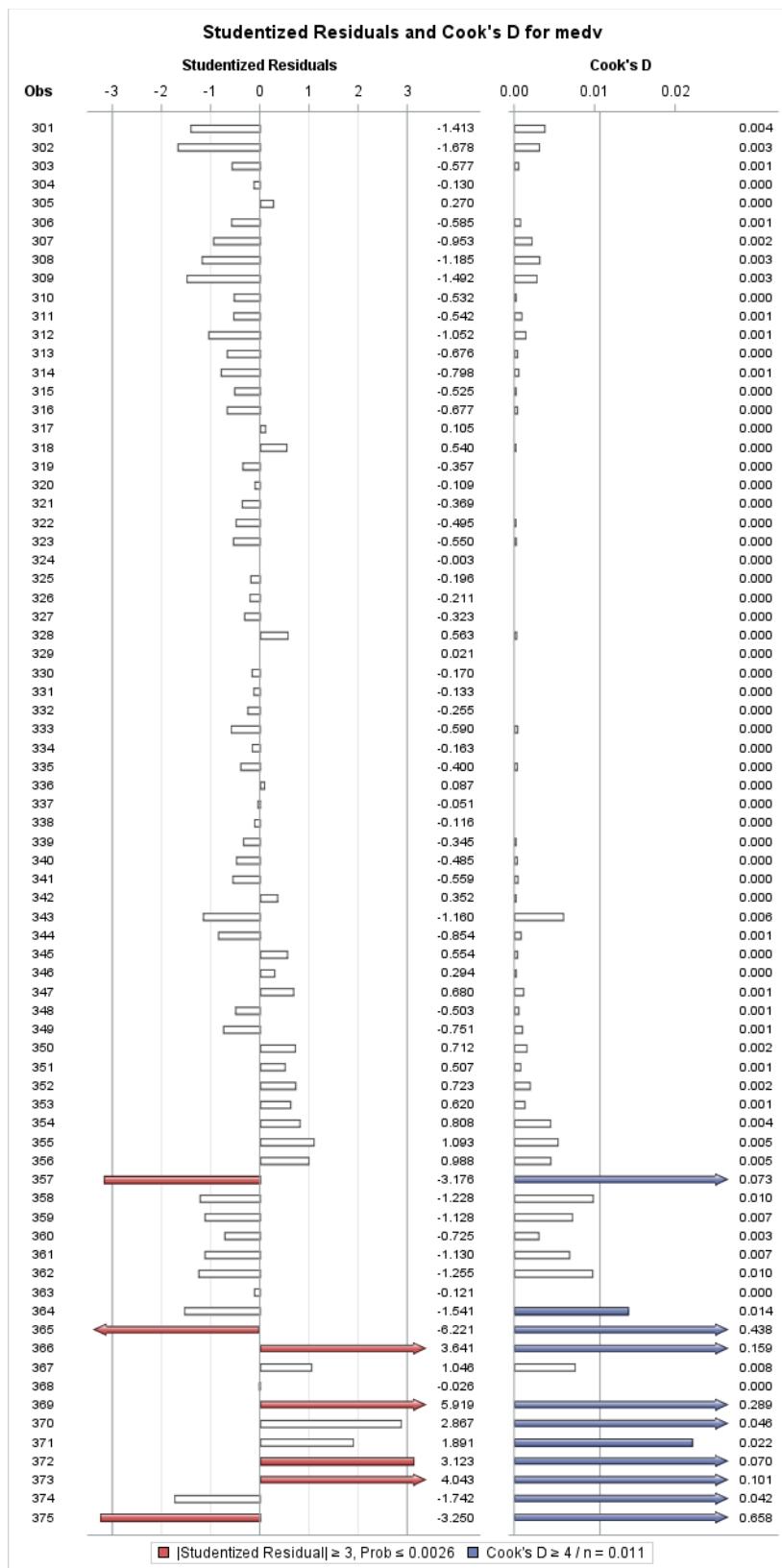


Fig. 2

Full model - Model assumptions and diagnostics

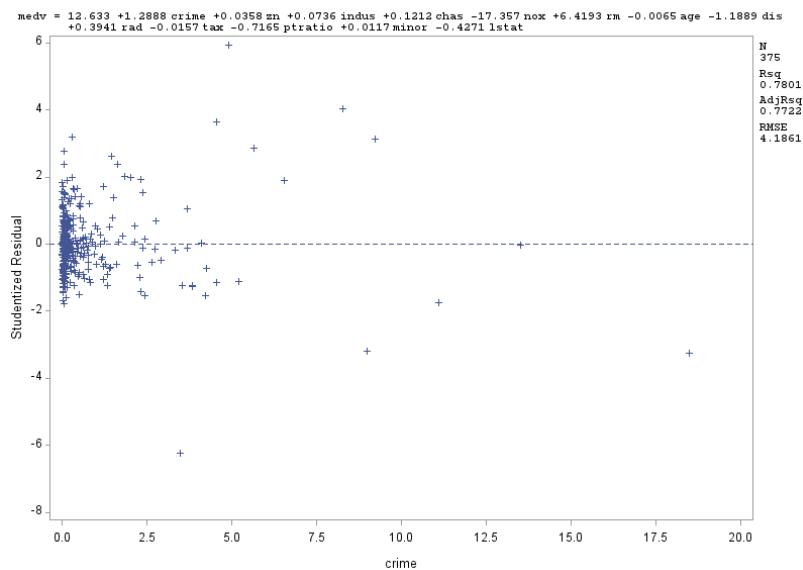


Fig. 3

Full model - Model assumptions and diagnostics

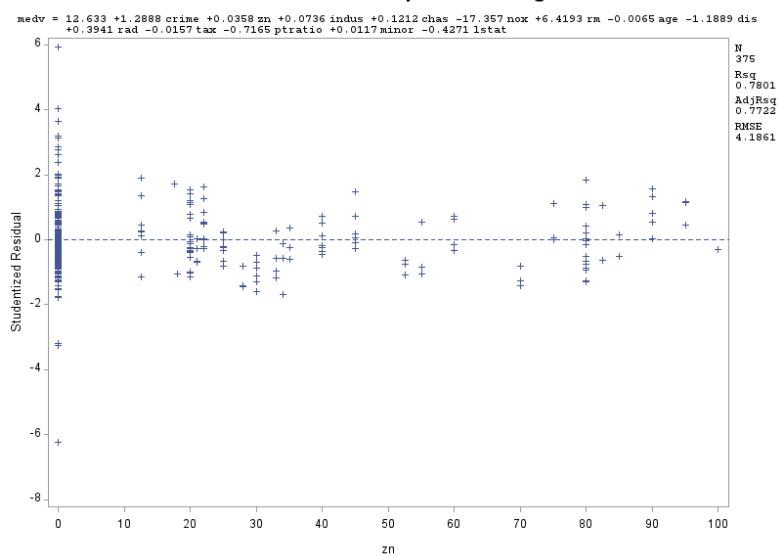


Fig. 4

Full model - Model assumptions and diagnostics

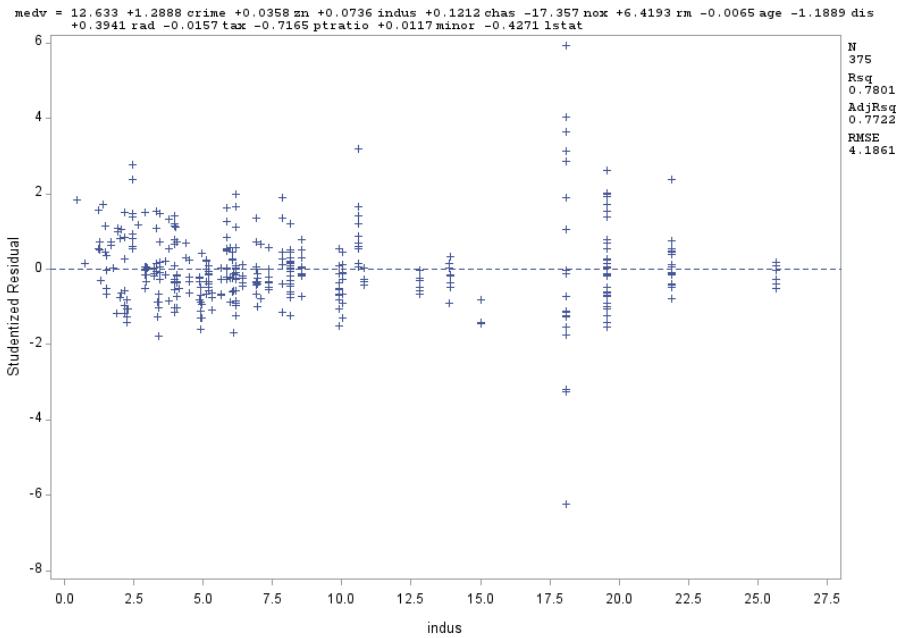


Fig. 5

Full model - Model assumptions and diagnostics

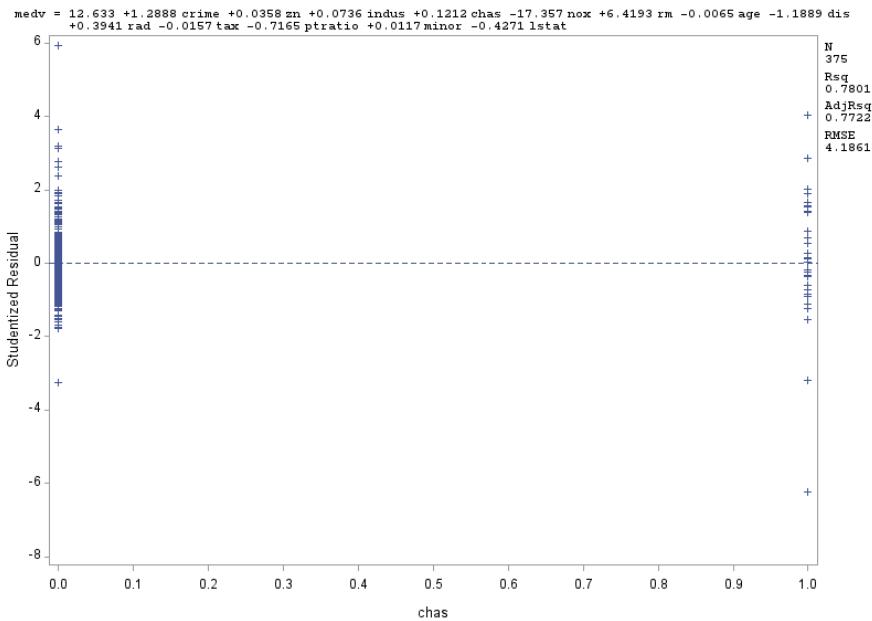


Fig. 6

Full model - Model assumptions and diagnostics

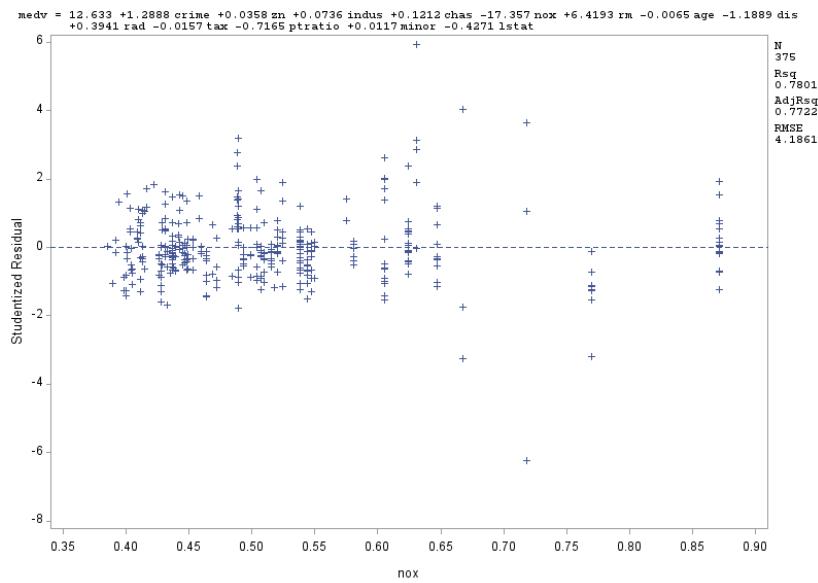


Fig. 7

Full model - Model assumptions and diagnostics

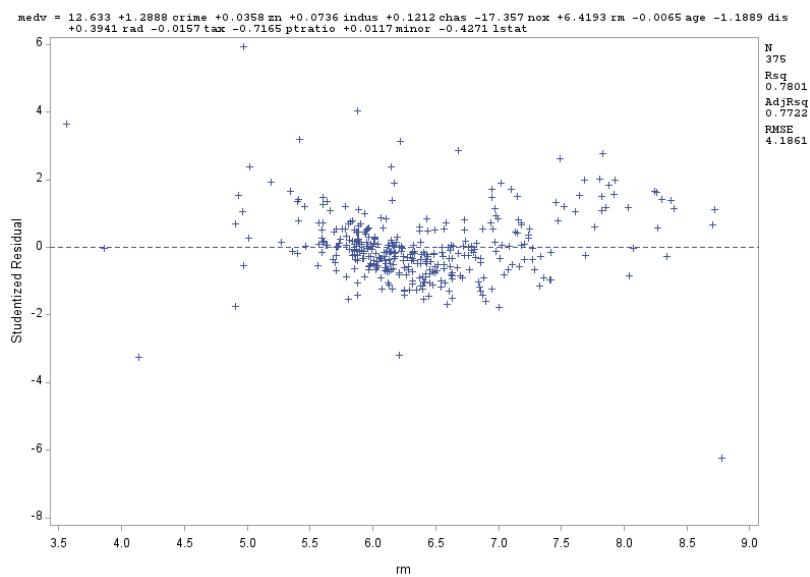


Fig. 8

Full model - Model assumptions and diagnostics

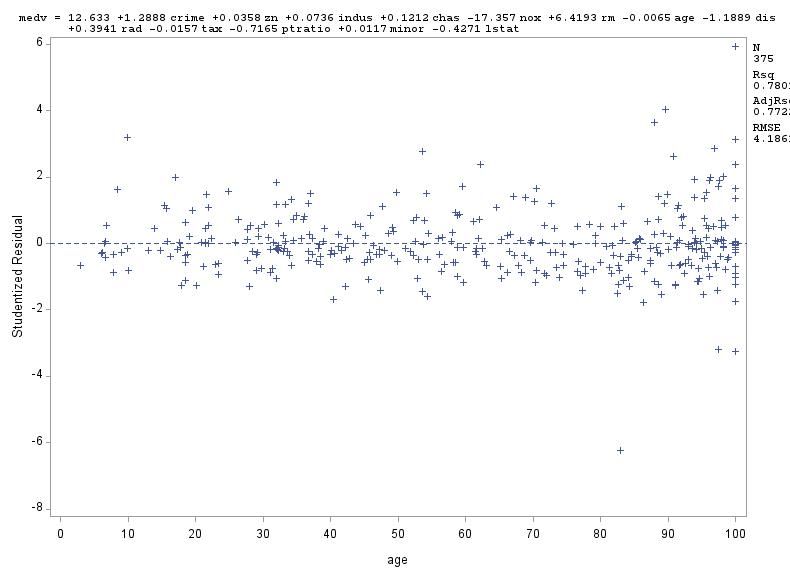


Fig. 9

Full model - Model assumptions and diagnostics

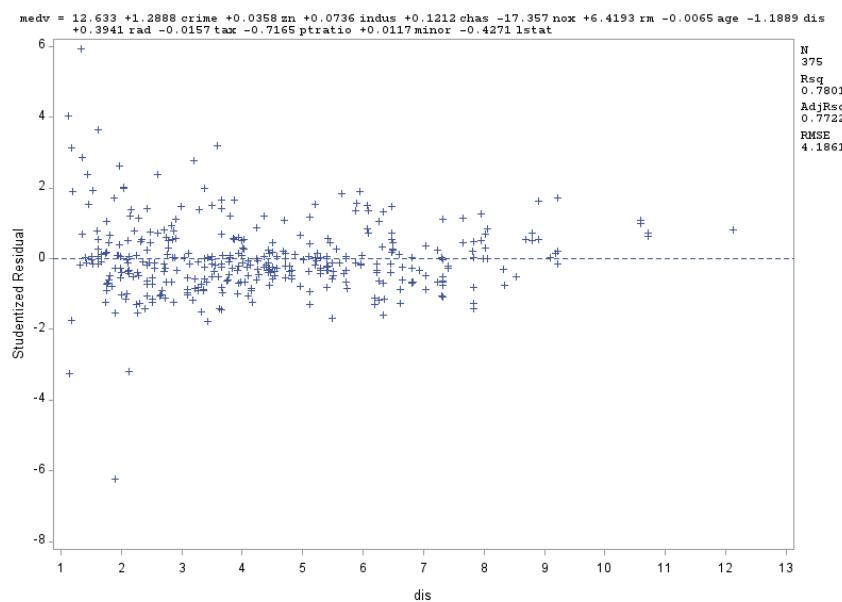


Fig. 10

Full model - Model assumptions and diagnostics

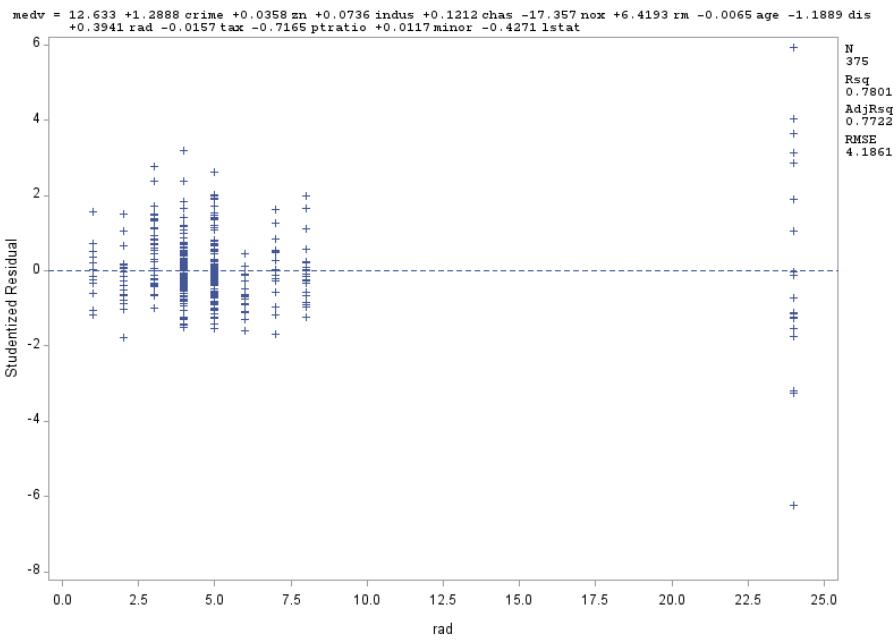


Fig. 11

Full model - Model assumptions and diagnostics

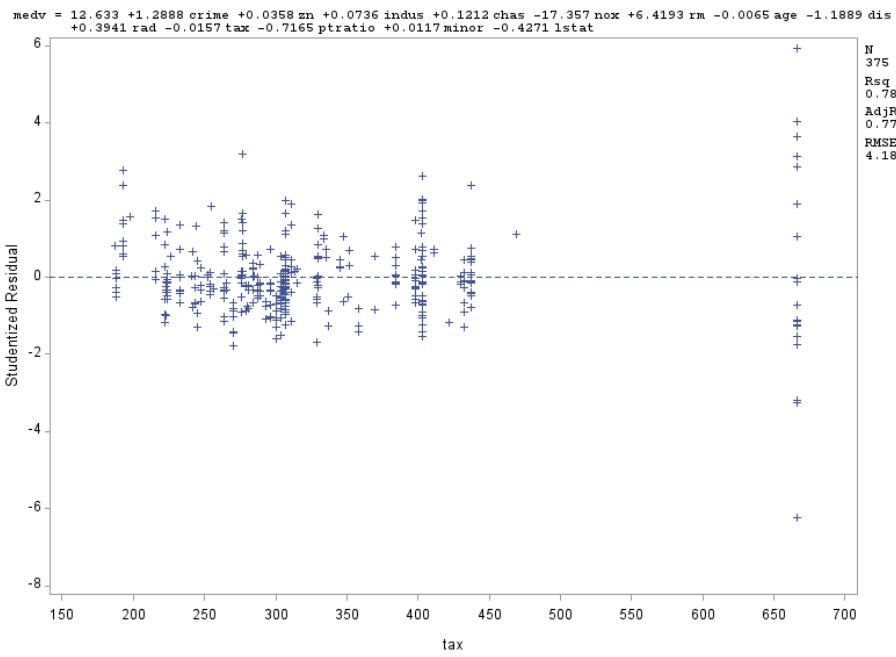


Fig. 12

Full model - Model assumptions and diagnostics

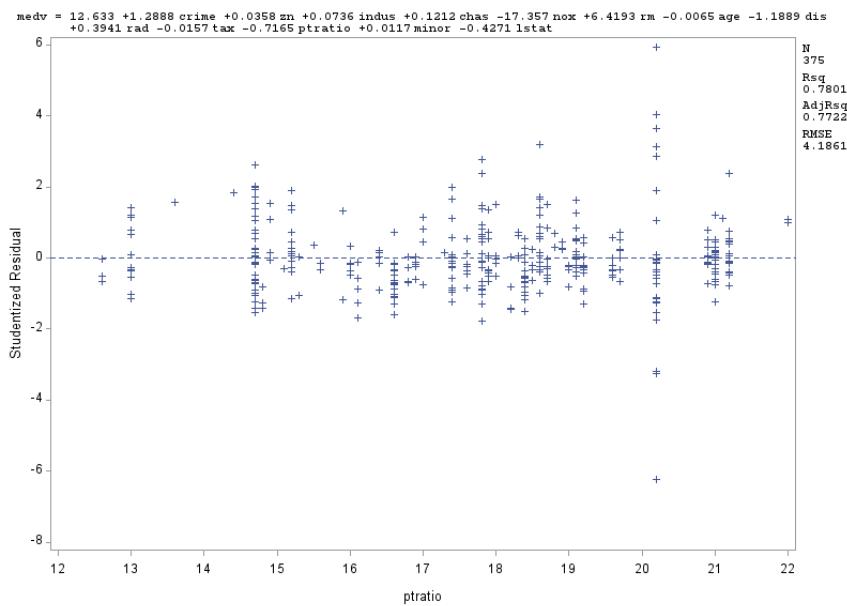


Fig. 13

Full model - Model assumptions and diagnostics

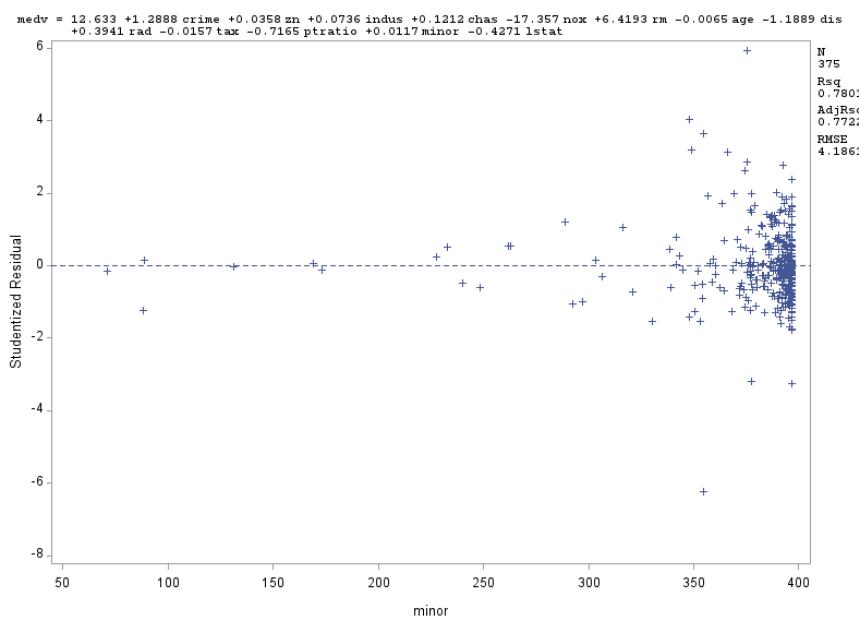


Fig. 14

Full model - Model assumptions and diagnostics

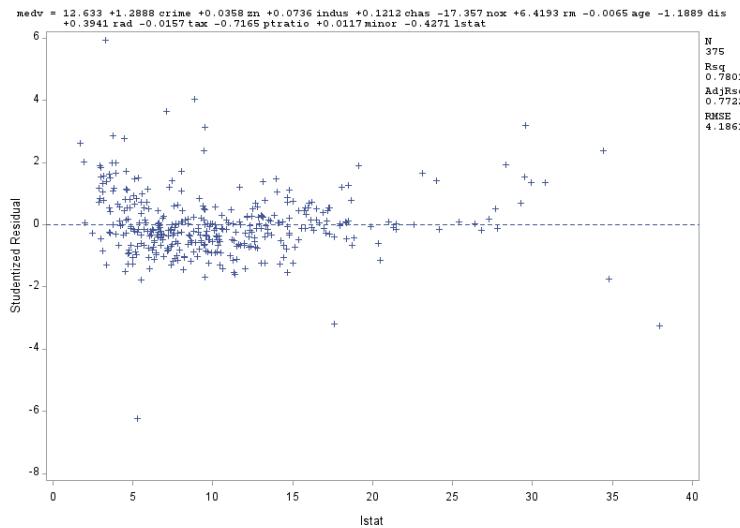


Fig. 15

Full model - Model assumptions and diagnostics

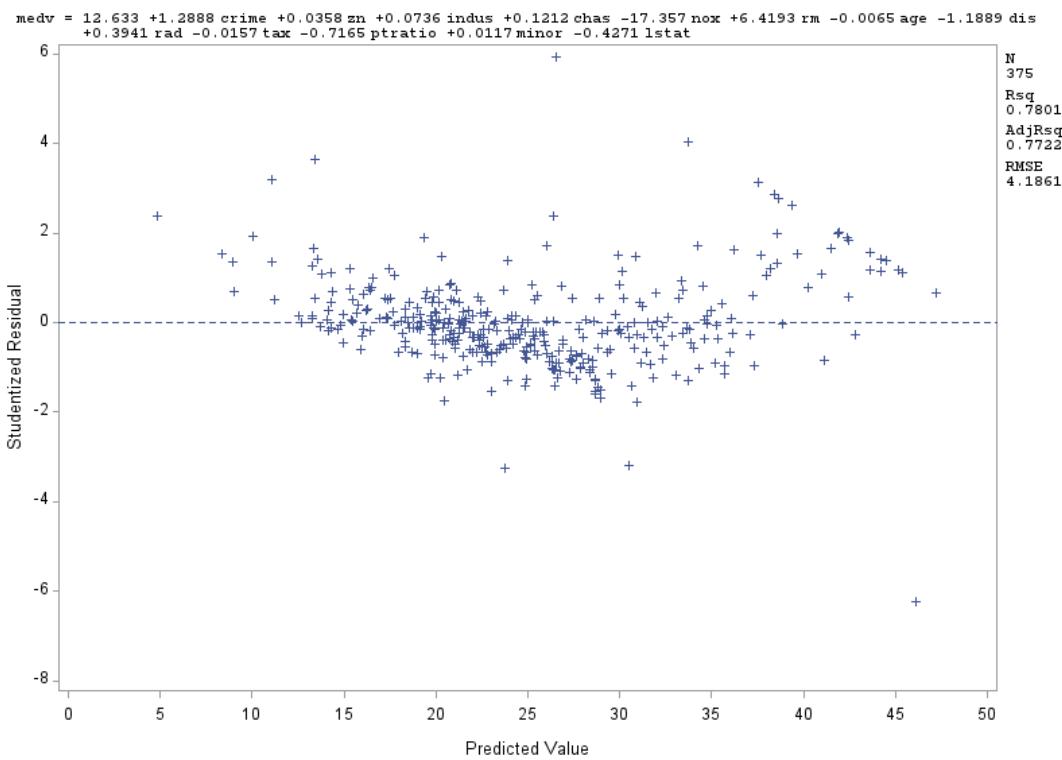


Fig. 16

Full model - Model assumptions and diagnostics

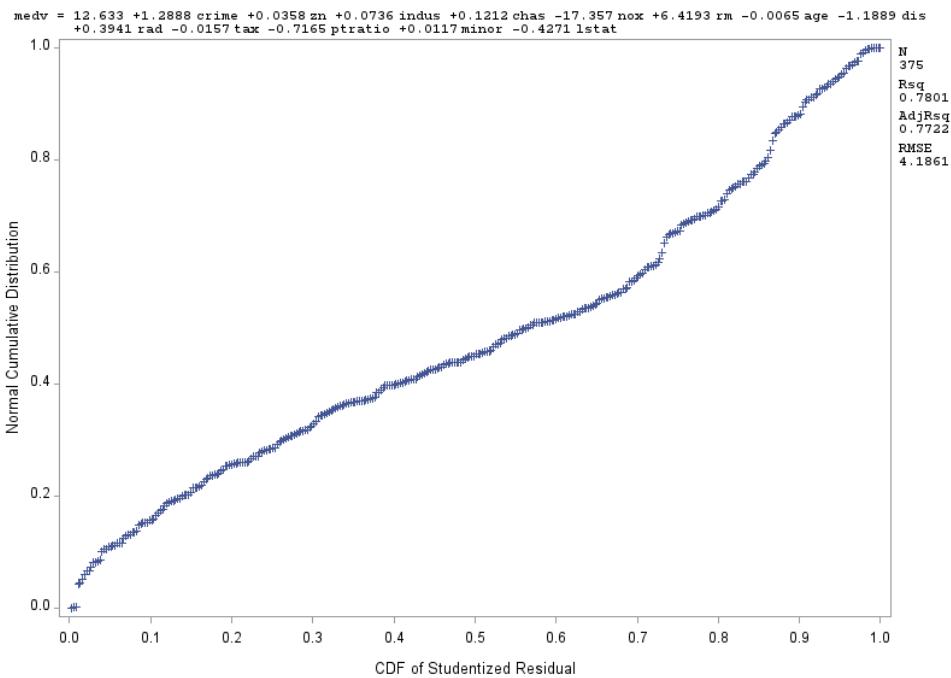


Table 2

Checking Model assumptions and diagnostics after removing outliers

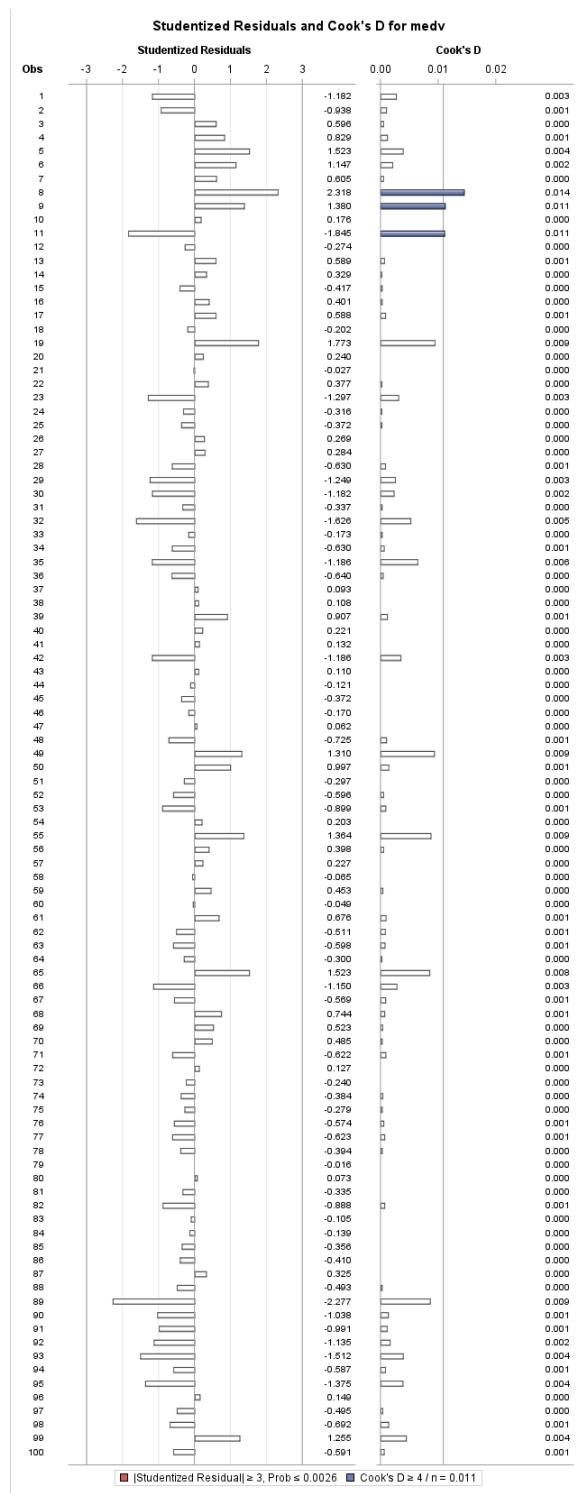
The REG Procedure
Model: MODEL1
Dependent Variable: medv
Number of Observations Read
367
Number of Observations Used
367

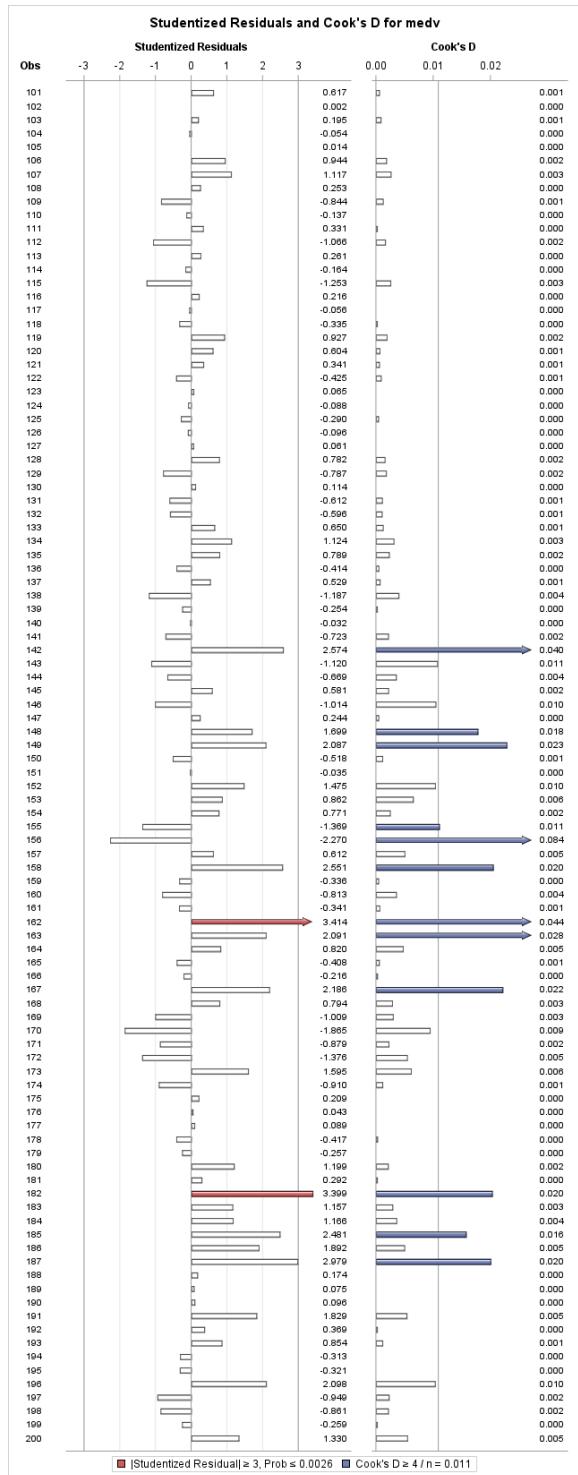
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	13	22972	1767.08614	166.65	<.0001
Error	353	3742.95757	10.60328		
Corrected Total	366	26715			

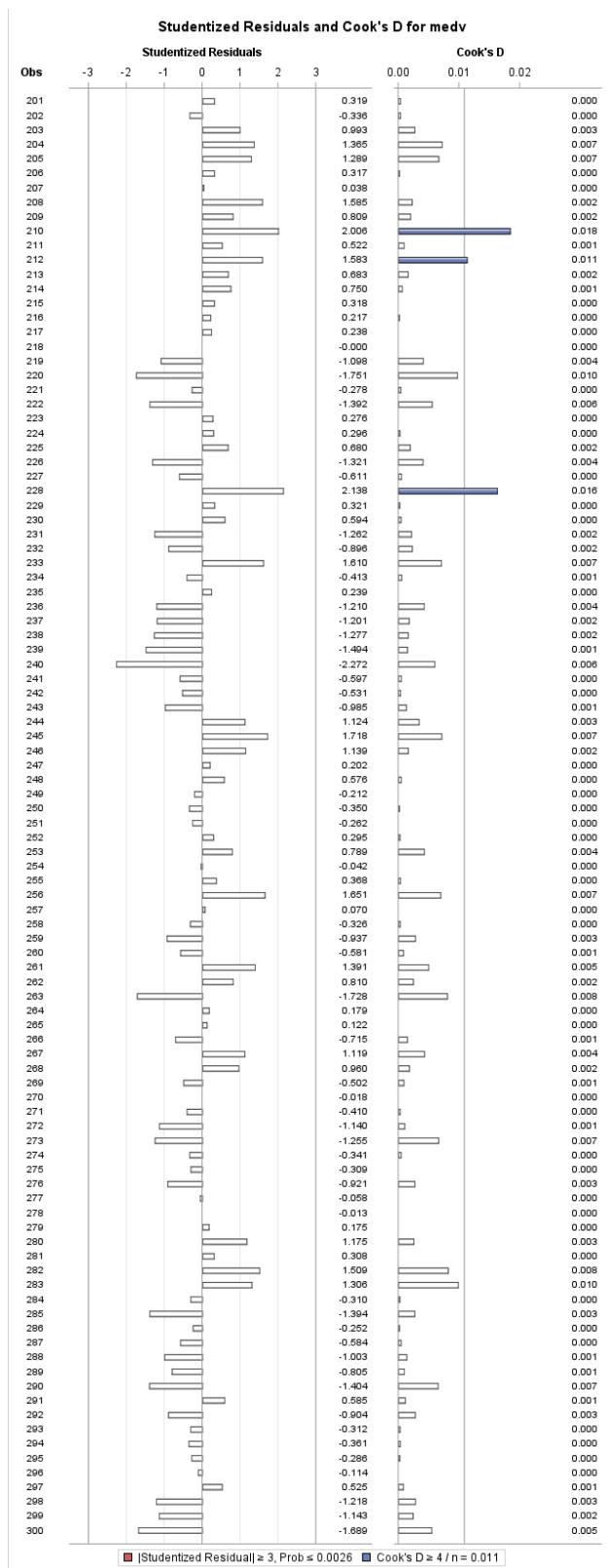
Root MSE	3.25627	R-Square	0.8599
Dependent Mean	25.04959	Adj R-Sq	0.8547
Coeff Var	12.99928		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	-8.01046	4.57775	-1.75	0.0810	0	0
crime	1	1.72577	0.24493	7.05	<.0001	0.25880	3.39914
zn	1	0.02307	0.00972	2.37	0.0181	0.07057	2.22582
indus	1	0.05365	0.04475	1.20	0.2314	0.03924	2.69961
chas	1	1.08971	0.63738	1.71	0.0882	0.03603	1.11915
nox	1	-12.53887	3.17136	-3.95	<.0001	-0.15868	4.05816
rm	1	8.51415	0.39053	21.80	<.0001	0.68525	2.48915
age	1	-0.02453	0.01041	-2.36	0.0190	-0.08189	3.04143
dis	1	-0.90588	0.14755	-6.14	<.0001	-0.22487	3.38009
rad	1	0.21037	0.08159	2.58	0.0103	0.09407	3.35379
tax	1	-0.01464	0.00303	-4.82	<.0001	-0.15767	2.69156
ptratio	1	-0.57630	0.09643	-5.98	<.0001	-0.14956	1.57765
minor	1	0.01385	0.00477	2.90	0.0039	0.06783	1.37604
lstat	1	-0.24693	0.05024	-4.91	<.0001	-0.16980	3.00784

Fig. 17







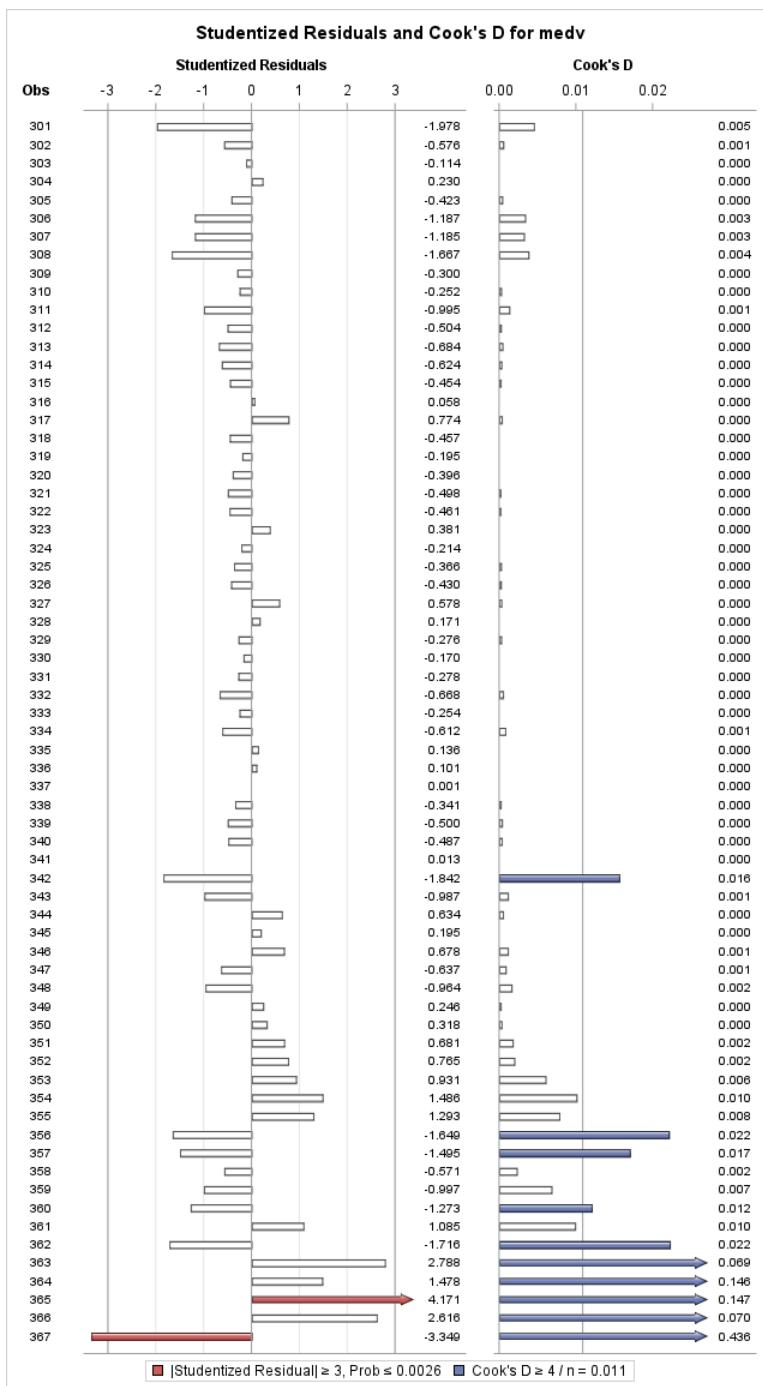


Fig. 18

Checking Model assumptions and diagnostics after removing outliers

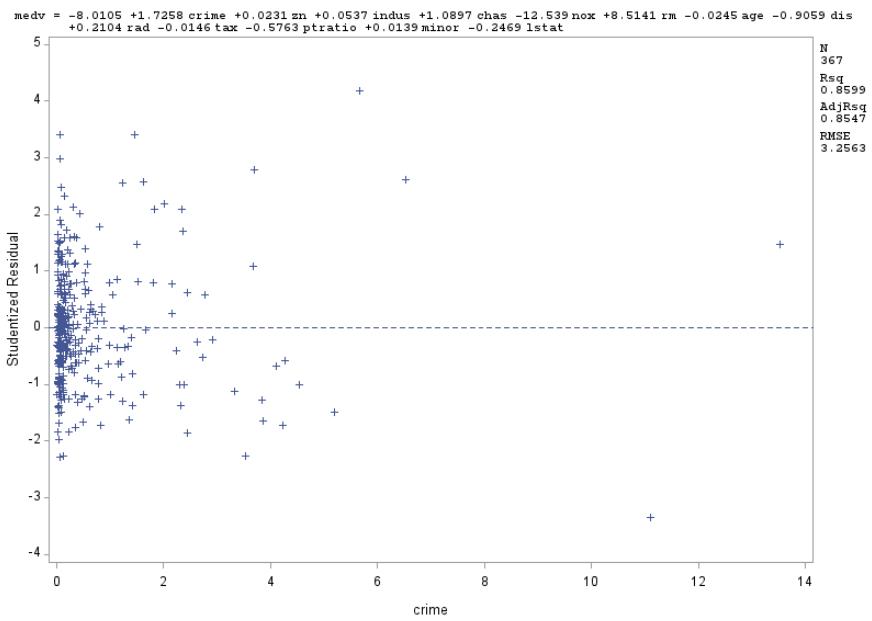


Fig. 19

Checking Model assumptions and diagnostics after removing outliers

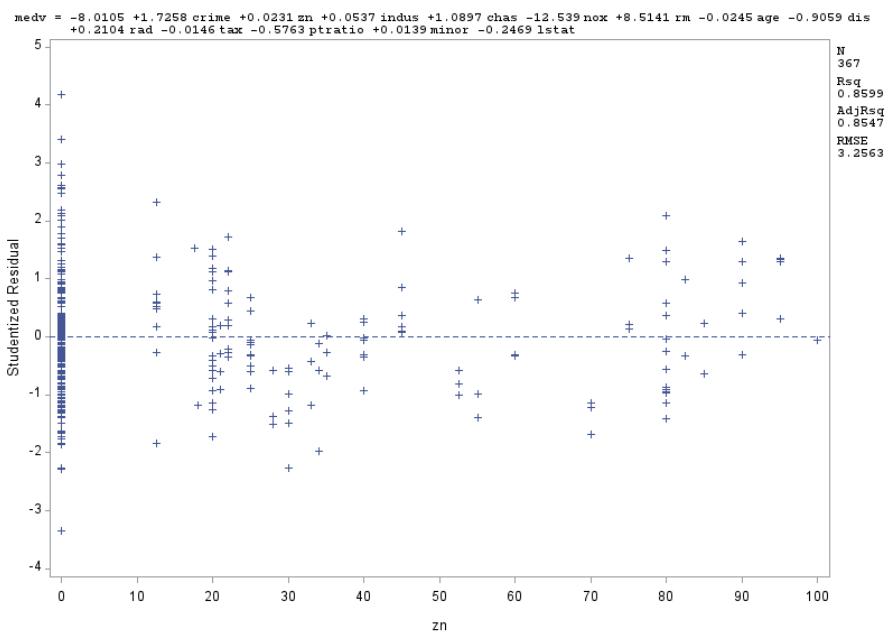


Fig. 20

Checking Model assumptions and diagnostics after removing outliers

```
medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat
```

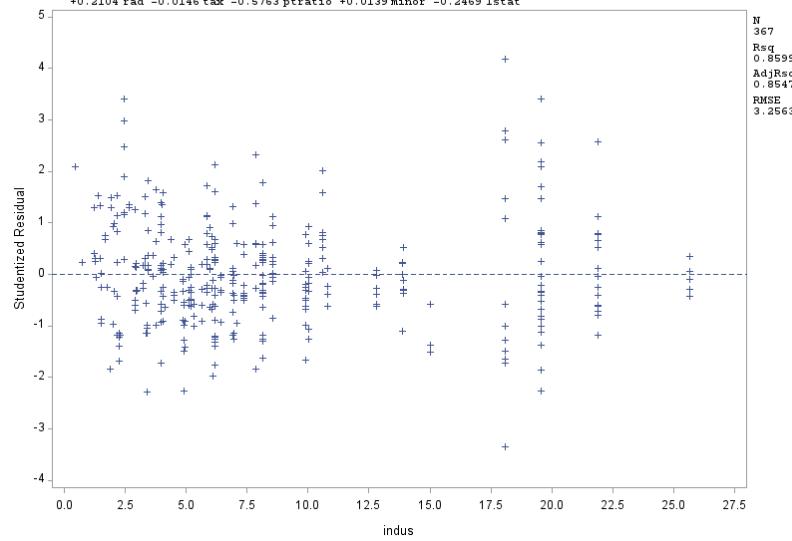


Fig. 21

Checking Model assumptions and diagnostics after removing outliers

```
medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat
```

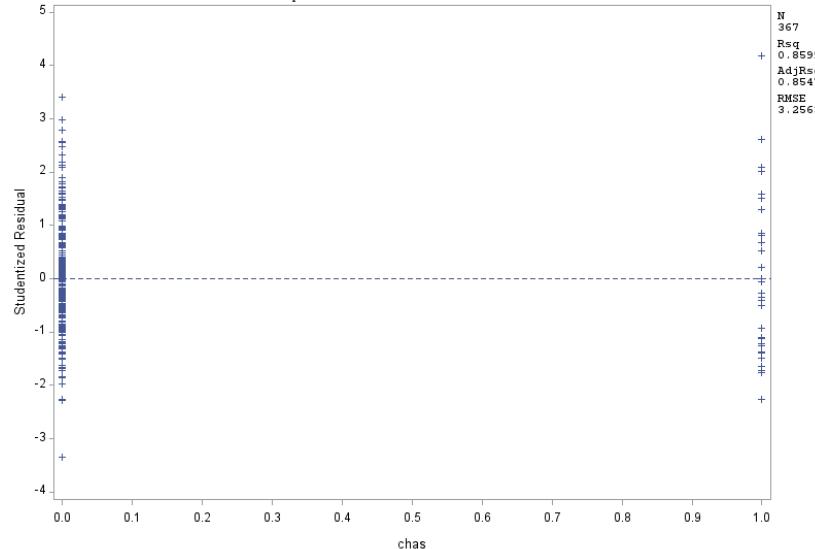


Fig. 22

Checking Model assumptions and diagnostics after removing outliers

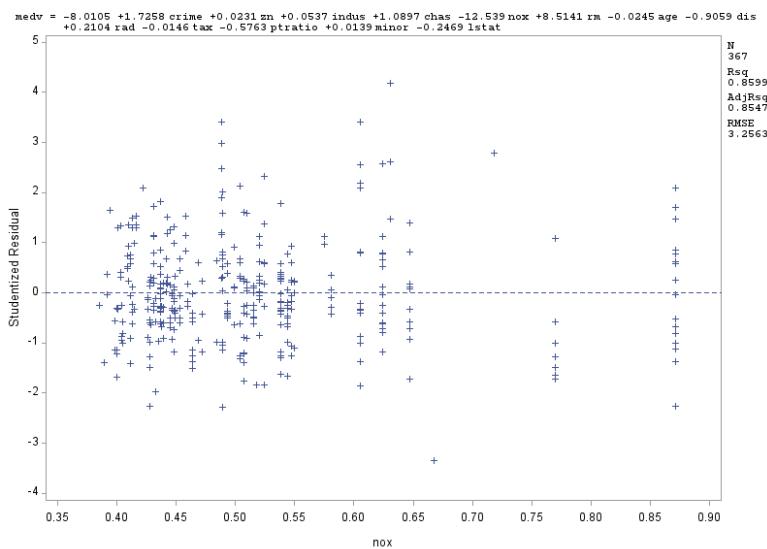


Fig. 23

Checking Model assumptions and diagnostics after removing outliers

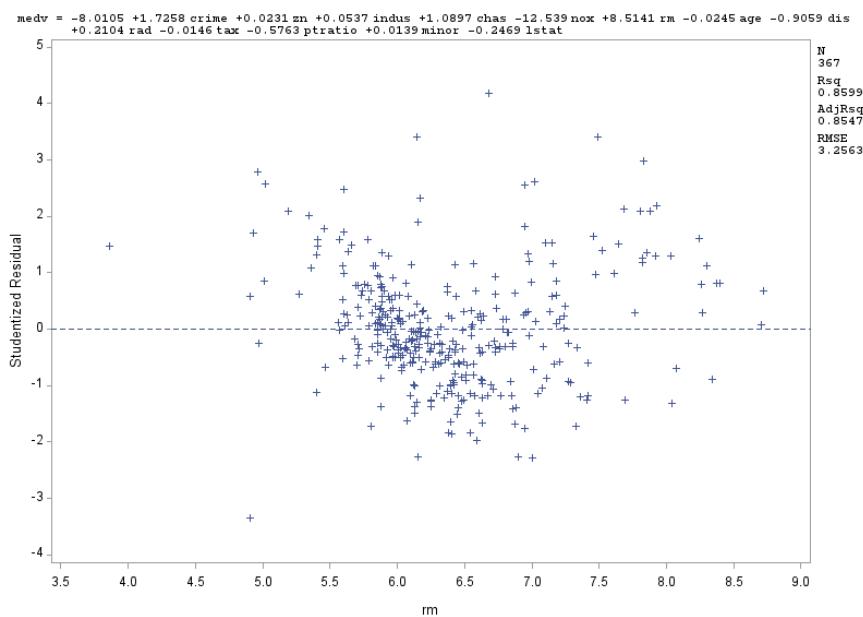


Fig. 24

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

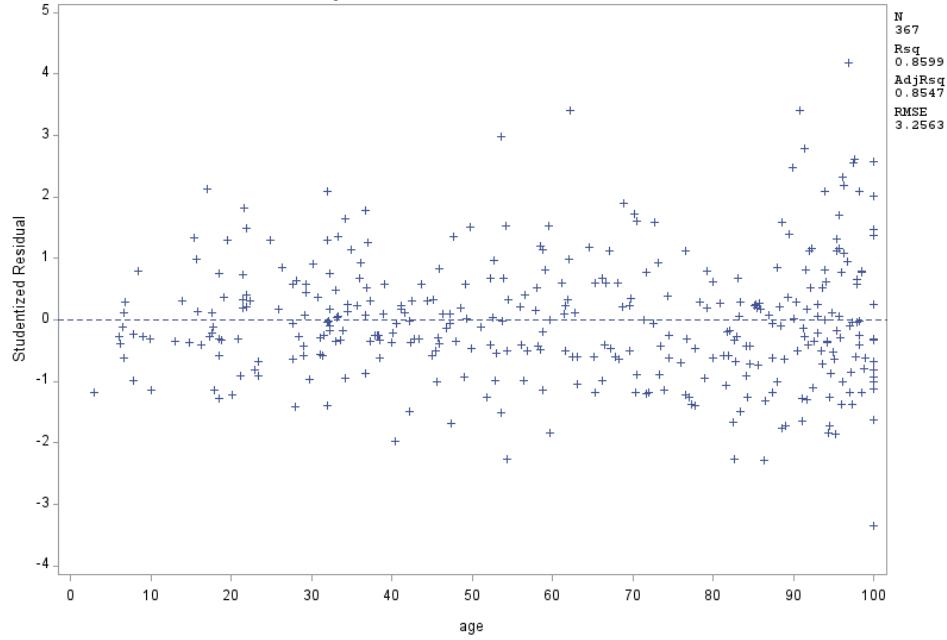


Fig. 25

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

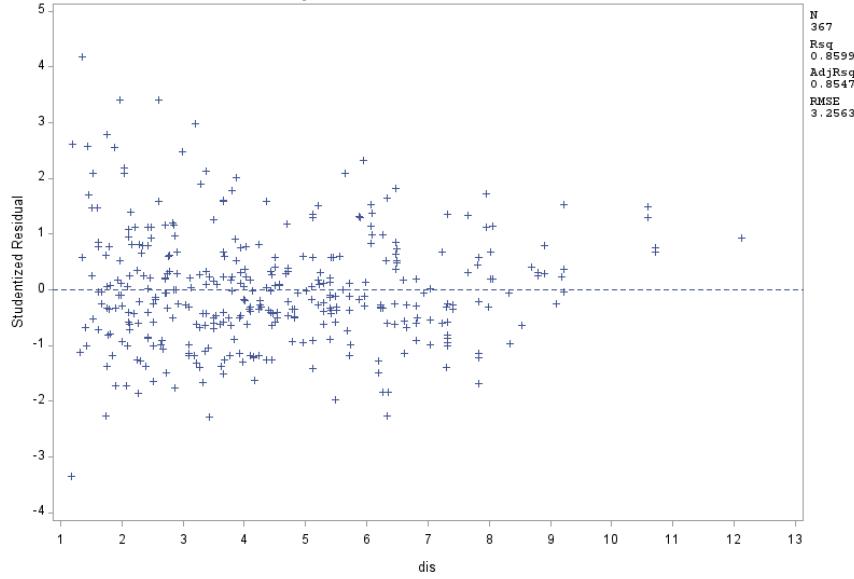


Fig. 26

Checking Model assumptions and diagnostics after removing outliers

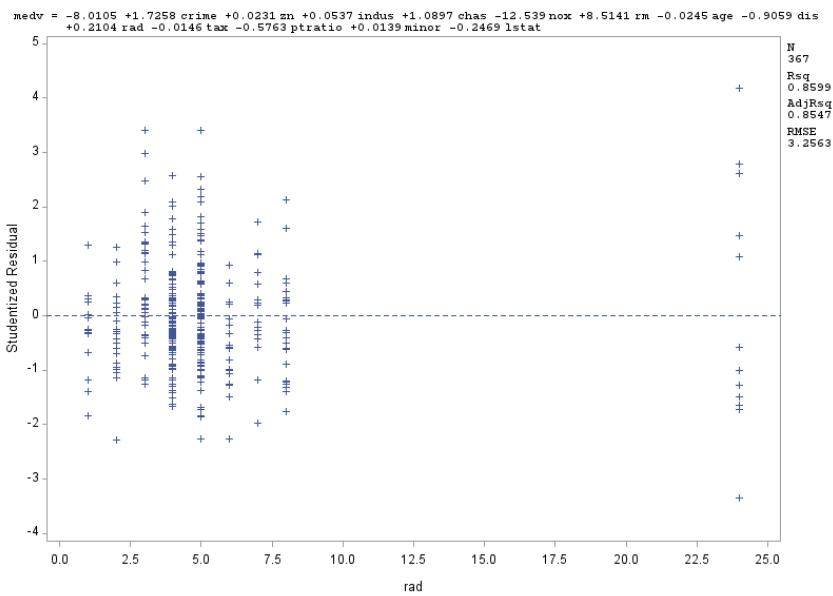


Fig. 27

Checking Model assumptions and diagnostics after removing outliers

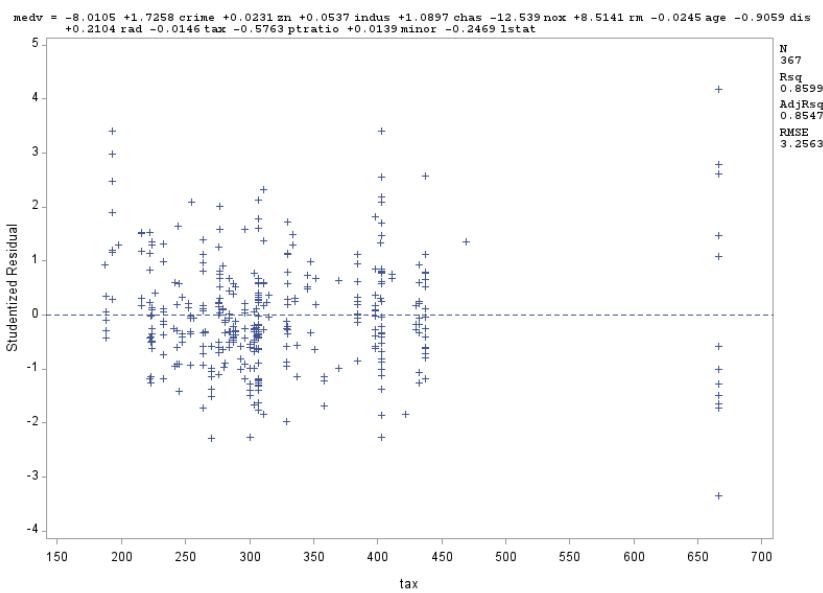


Fig. 28

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

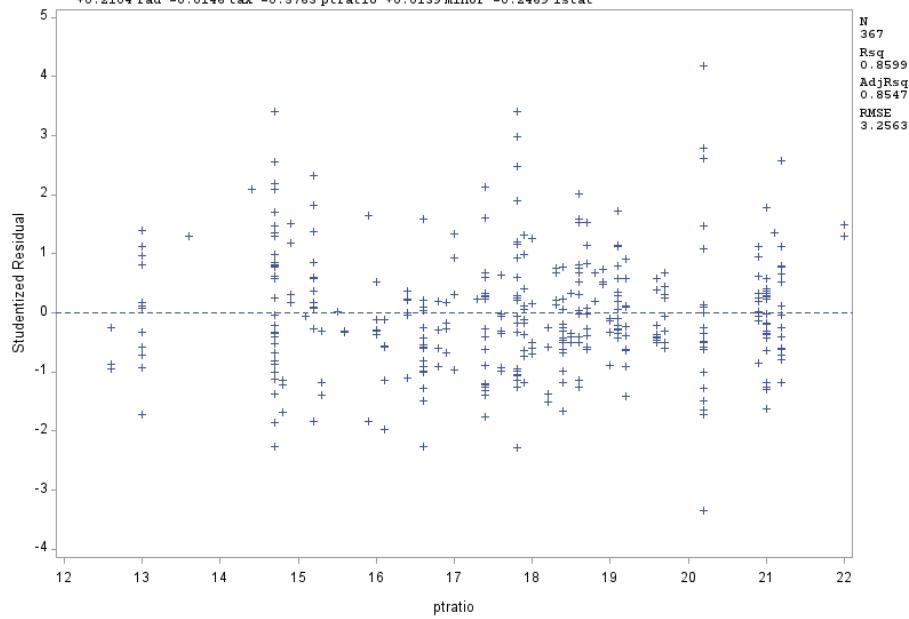


Fig. 29

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

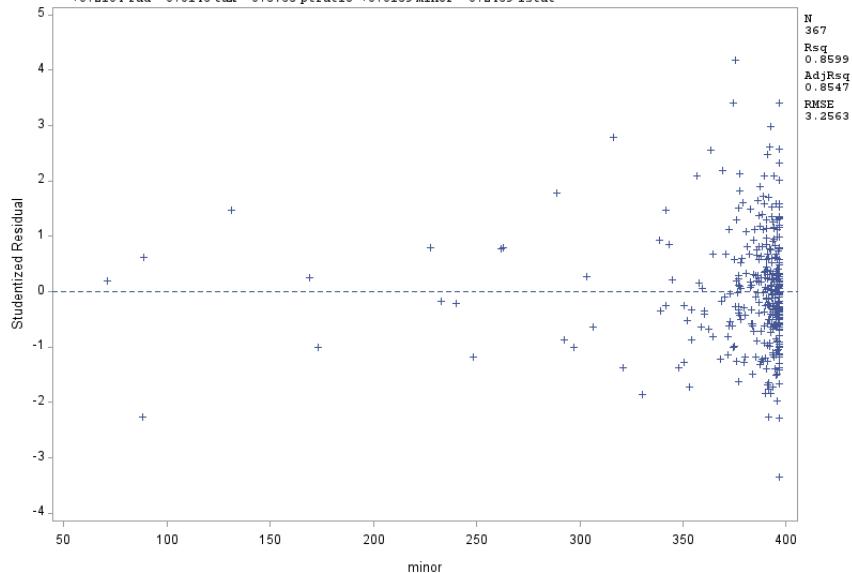


Fig. 30

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

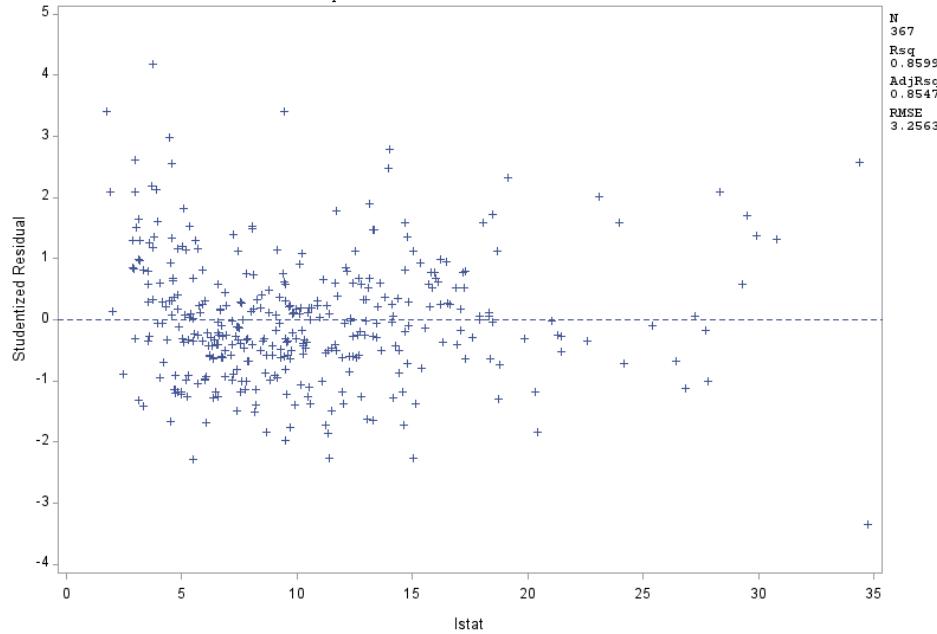


Fig. 31

Checking Model assumptions and diagnostics after removing outliers

medv = -8.0105 +1.7258 crime +0.0231 zn +0.0537 indus +1.0897 chas -12.539 nox +8.5141 rm -0.0245 age -0.9059 dis
+0.2104 rad -0.0146 tax -0.5763 ptratio +0.0139 minor -0.2469 lstat

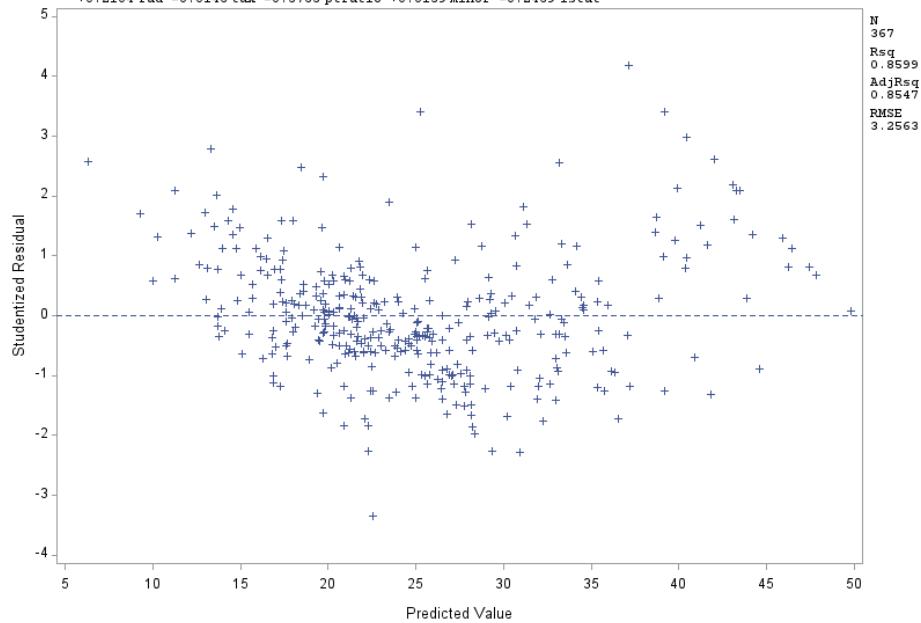
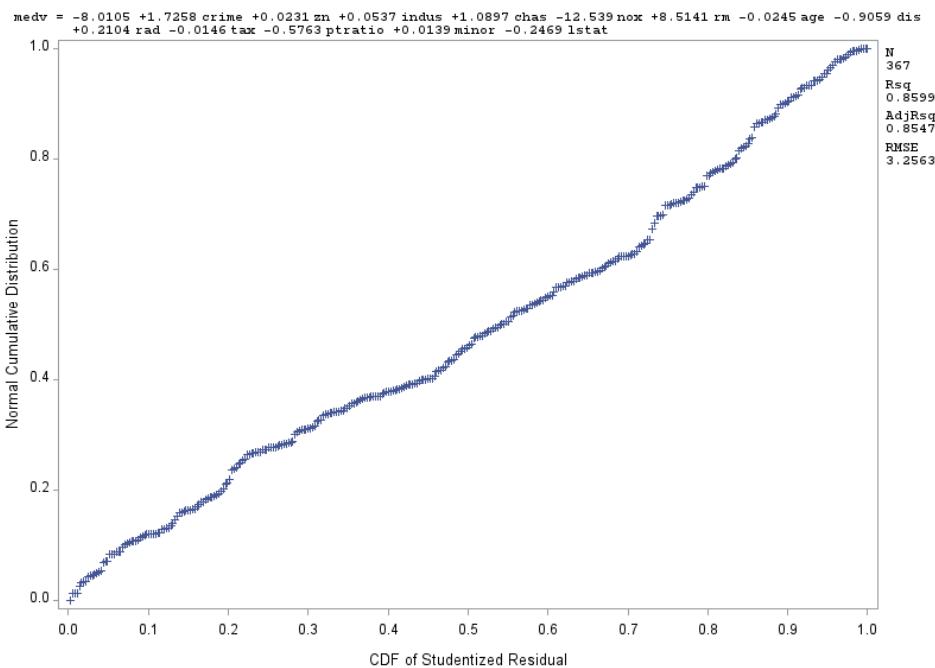


Fig. 32

Checking Model assumptions and diagnostics after removing outliers



APPENDIX C – Log transformation

Table 1

Checking Model assumptions and diagnostics after transformation

The REG Procedure
Model: MODEL1
Dependent Variable: logmedv
Number of Observations Read 367
Number of Observations Used 367

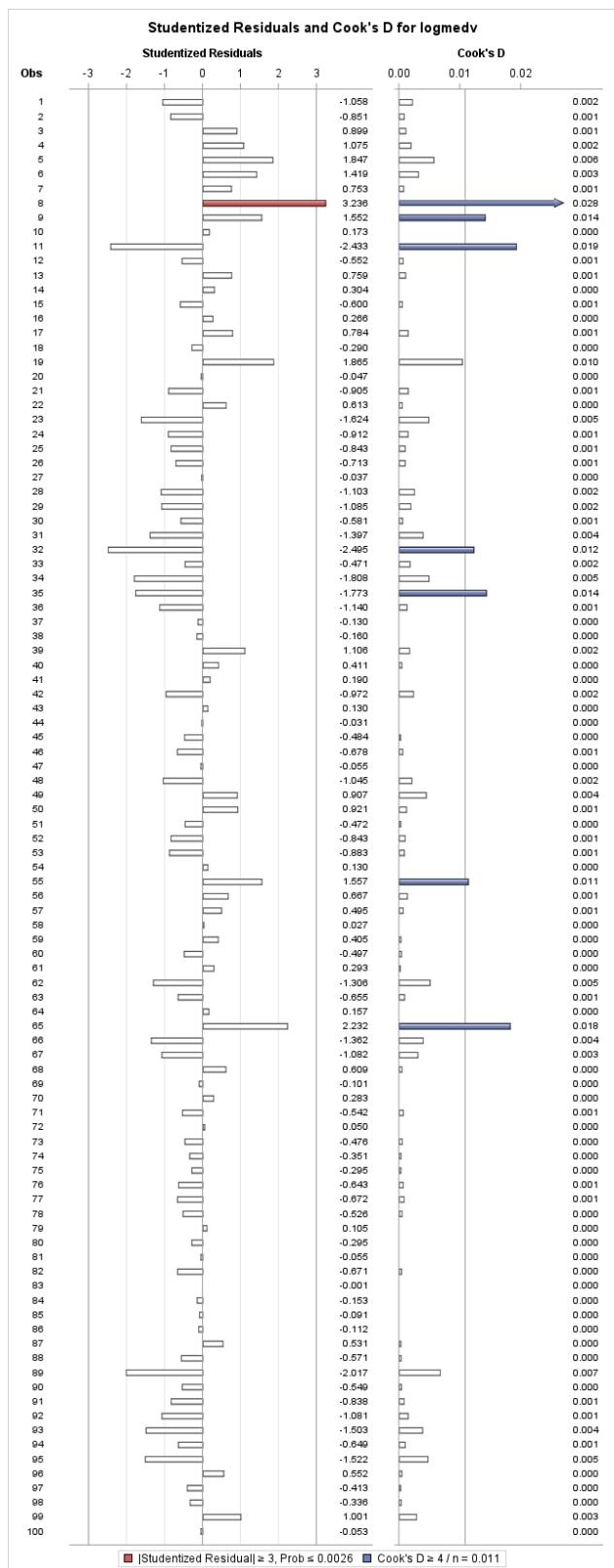
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	13	30.73803	2.36446	194.96	<.0001
Error	353	4.28116	0.01213		
Corrected Total	366	35.01918			

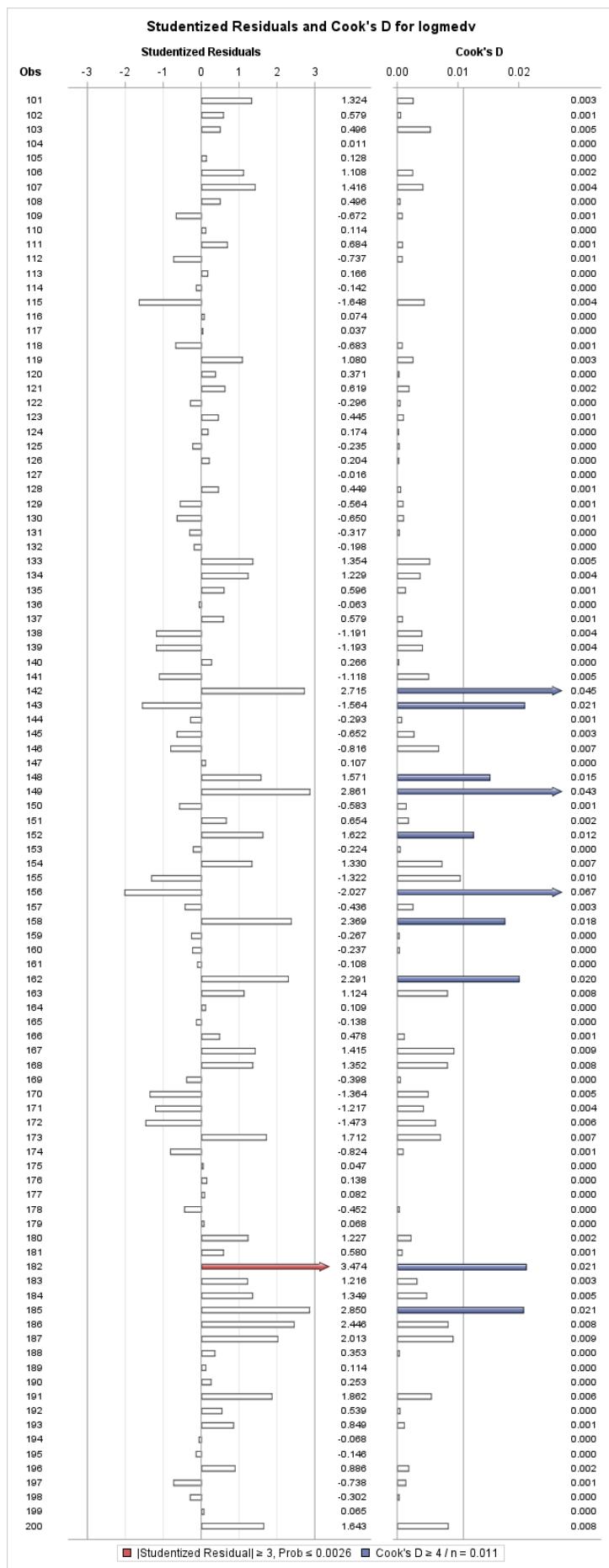
Root MSE	0.11013	R-Square	0.8777
Dependent Mean	3.17054	Adj R-Sq	0.8732

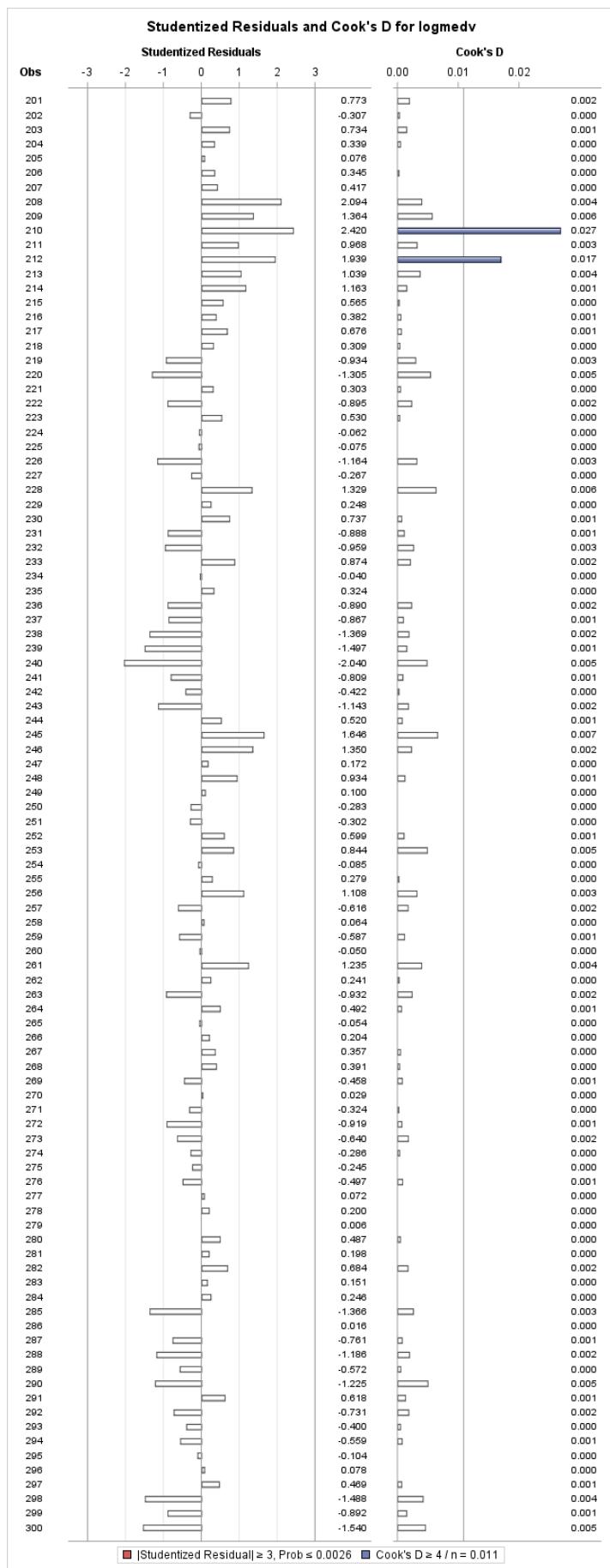
Coeff Var	3.47344		
-----------	---------	--	--

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	2.51620	0.15482	16.25	<.0001	0	0
crime	1	0.04567	0.00828	5.51	<.0001	0.18916	3.39914
zn	1	0.00074642	0.00032864	2.27	0.0237	0.06306	2.22582
indus	1	0.00227	0.00151	1.50	0.1353	0.04577	2.69961
chas	1	0.03477	0.02156	1.61	0.1077	0.03175	1.11915
nox	1	-0.57152	0.10726	-5.33	<.0001	-0.19976	4.05816
rm	1	0.25083	0.01321	18.99	<.0001	0.55760	2.48915
age	1	-0.00092026	0.00035193	-2.61	0.0093	-0.08486	3.04143
dis	1	-0.03512	0.00499	-7.04	<.0001	-0.24076	3.38009
rad	1	0.01240	0.00276	4.49	<.0001	0.15311	3.35379
tax	1	-0.00057326	0.00010264	-5.58	<.0001	-0.17051	2.69156
ptratio	1	-0.02534	0.00326	-7.77	<.0001	-0.18166	1.57765
minor	1	0.00060553	0.00016140	3.75	0.0002	0.08190	1.37604
lstat	1	-0.01516	0.00170	-8.92	<.0001	-0.28803	3.00784

Fig. 1







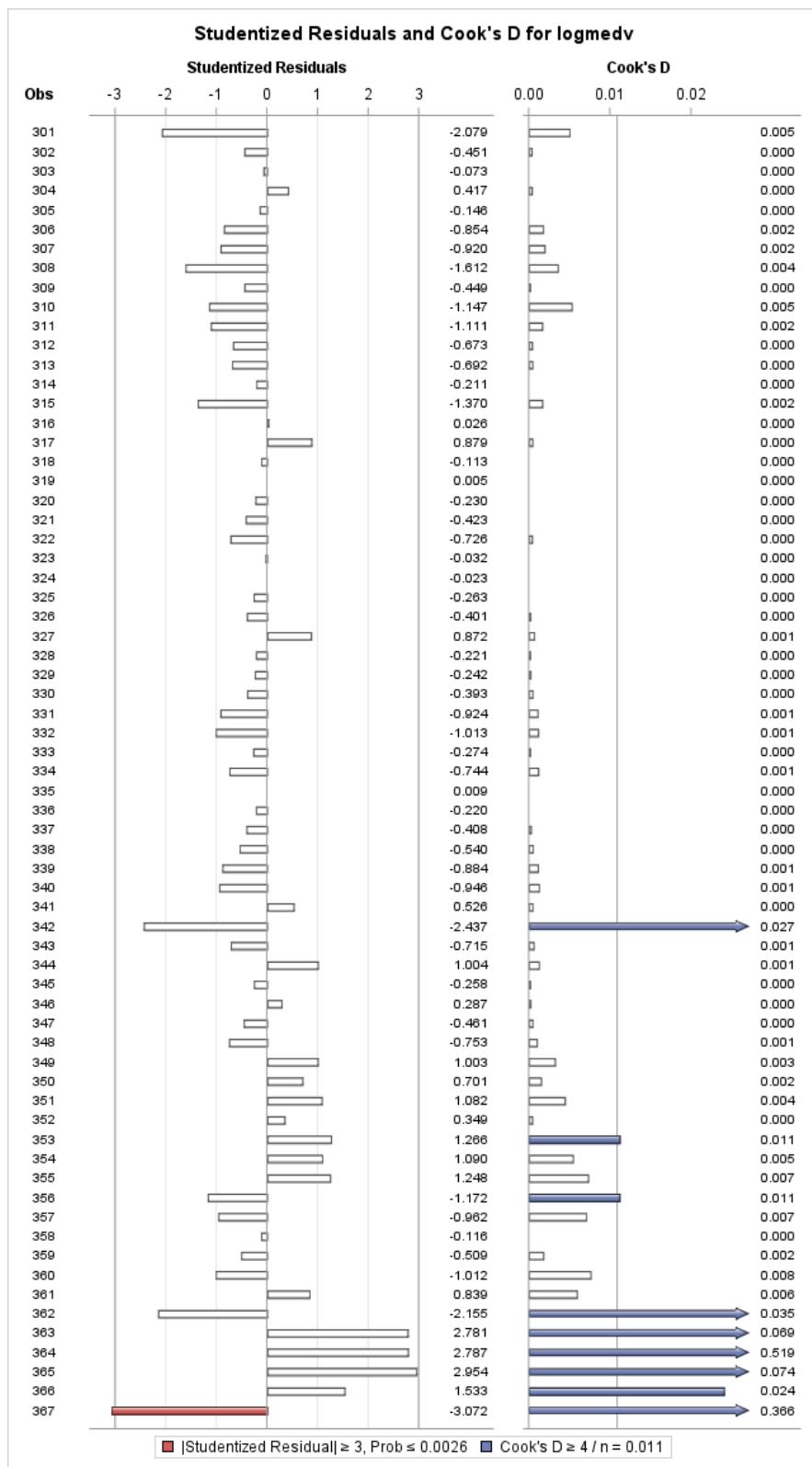


Fig. 2

Checking Model assumptions and diagnostics after transformation

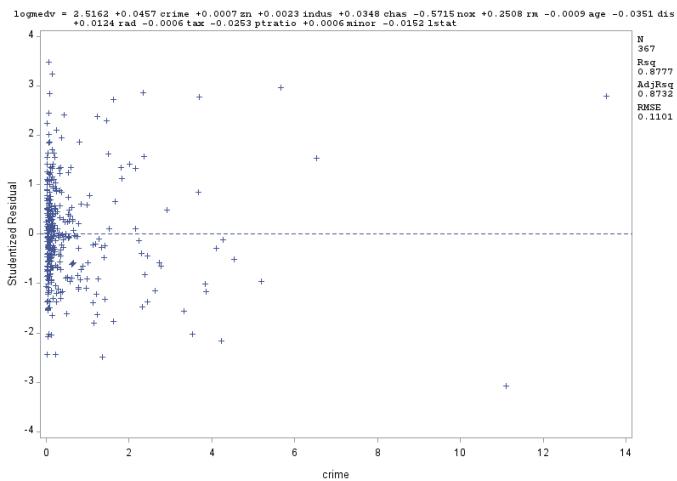


Fig. 3

Checking Model assumptions and diagnostics after transformation

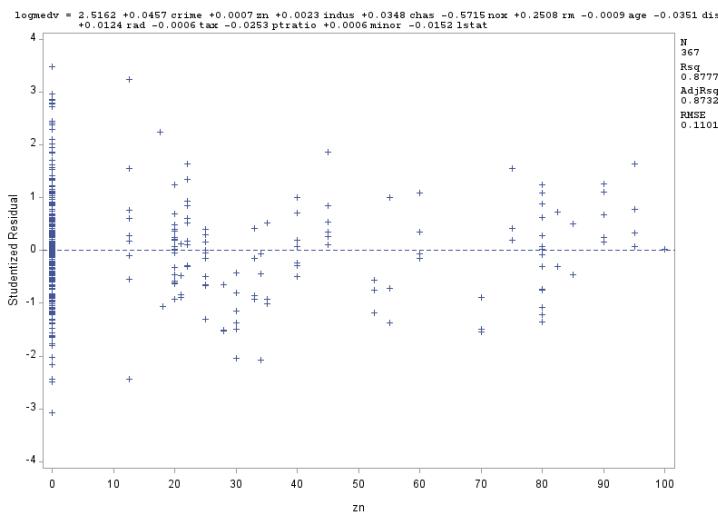


Fig. 4

Checking Model assumptions and diagnostics after transformation

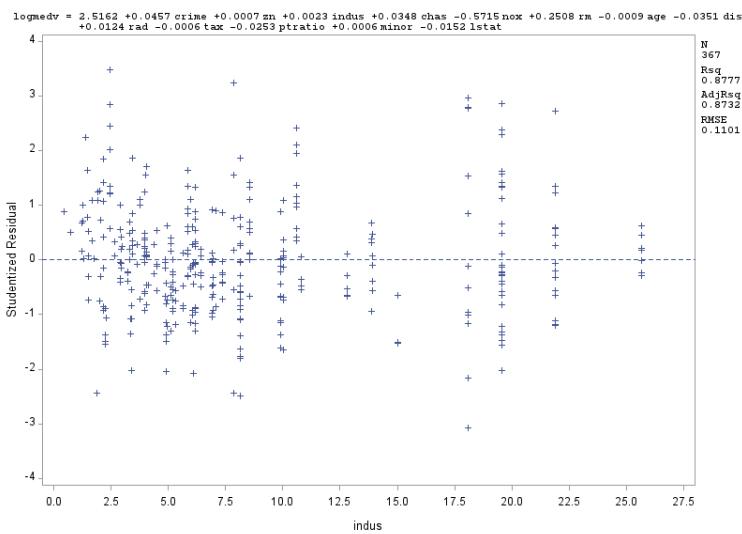


Fig. 5

Checking Model assumptions and diagnostics after transformation

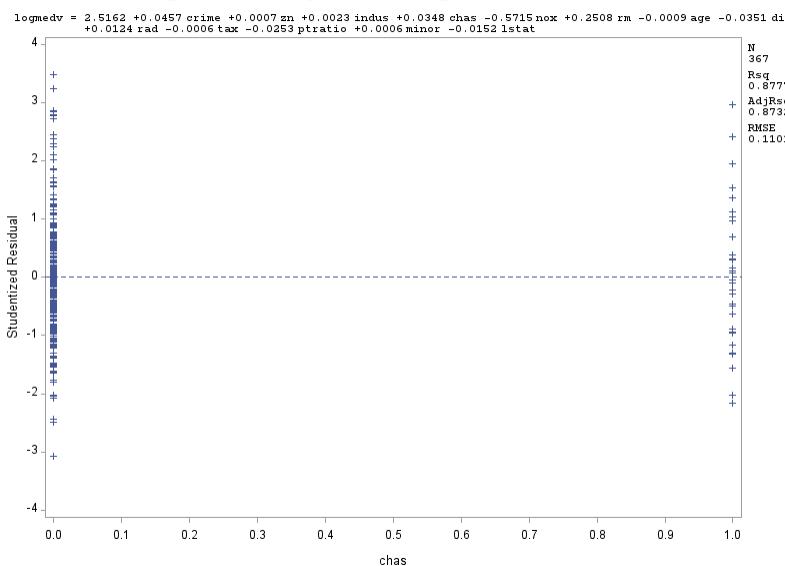


Fig. 6

Checking Model assumptions and diagnostics after transformation

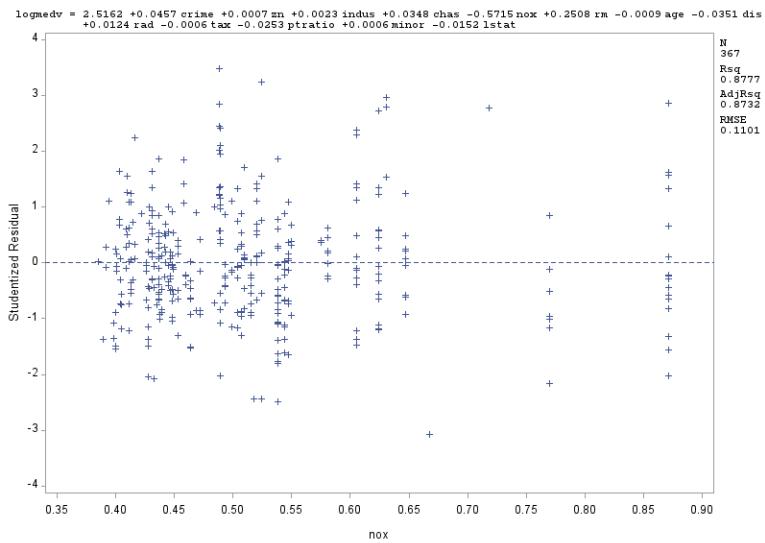


Fig. 7

Checking Model assumptions and diagnostics after transformation

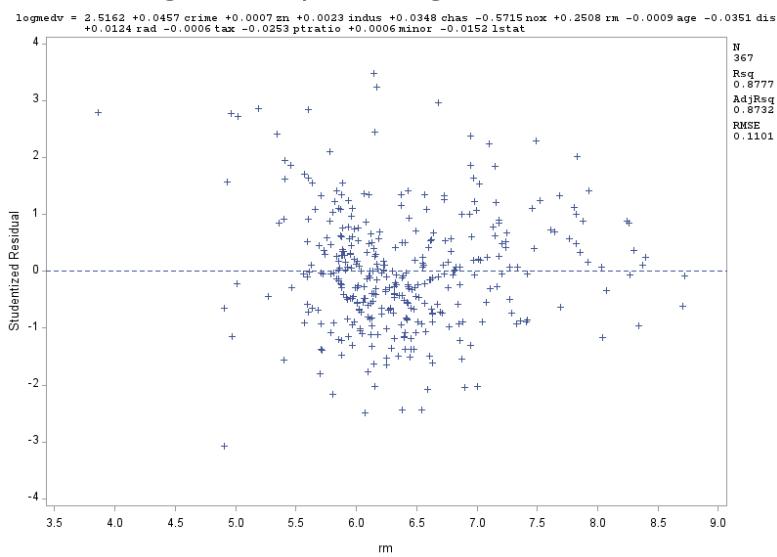


Fig. 8

Checking Model assumptions and diagnostics after transformation

```
logmedv = 2.5162 +0.0457 crime +0.0007 zn +0.0023 indus +0.0348 chas -0.5715 nox +0.2508 rm -0.0009 age -0.0351 dis
+0.0124 rad -0.0006 tax -0.0253 ptratio +0.0006 minor -0.0152 lstat
```

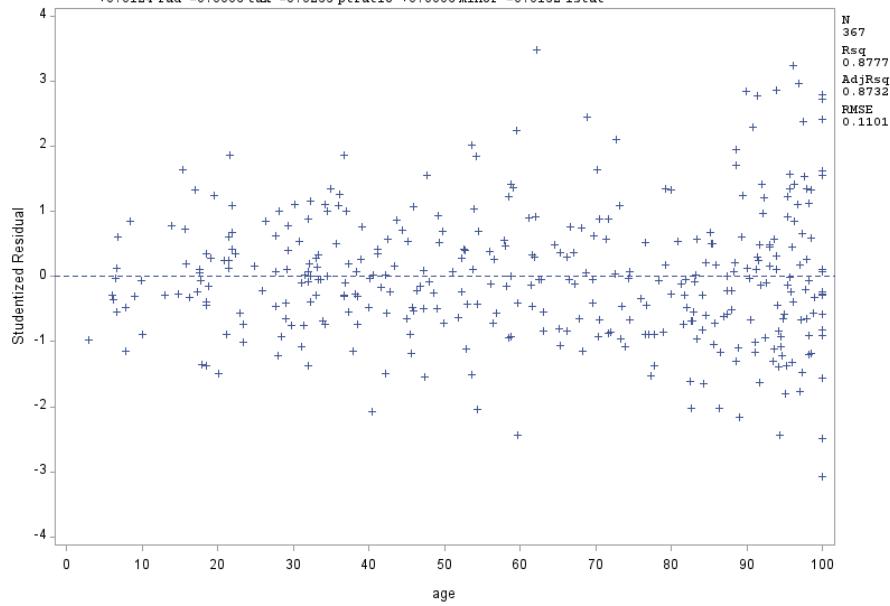


Fig. 9

Checking Model assumptions and diagnostics after transformation

```
logmedv = 2.5162 +0.0457 crime +0.0007 zn +0.0023 indus +0.0348 chas -0.5715 nox +0.2508 rm -0.0009 age -0.0351 dis
+0.0124 rad -0.0006 tax -0.0253 ptratio +0.0006 minor -0.0152 lstat
```

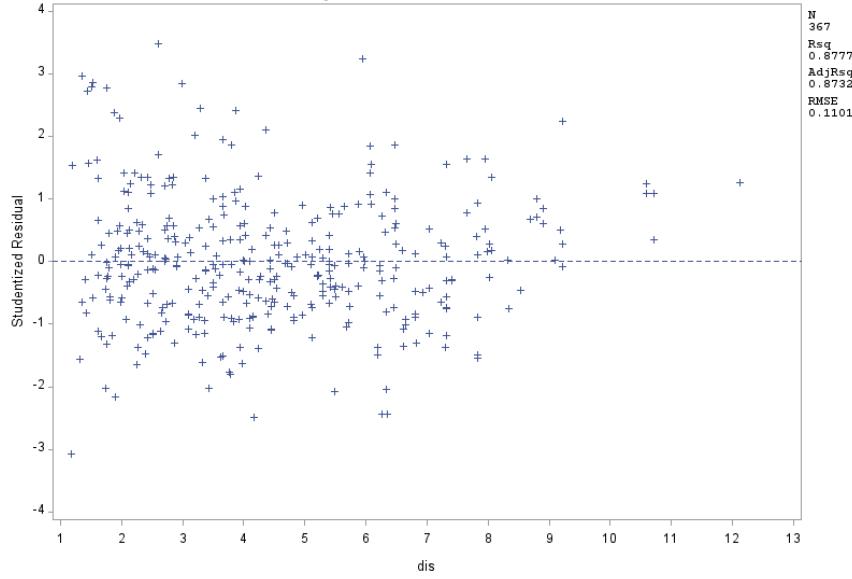


Fig. 10

Checking Model assumptions and diagnostics after transformation

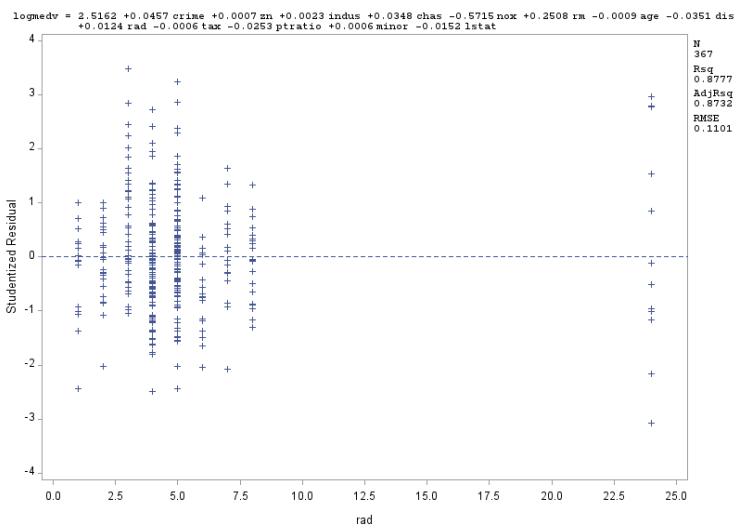


Fig. 11

Checking Model assumptions and diagnostics after transformation

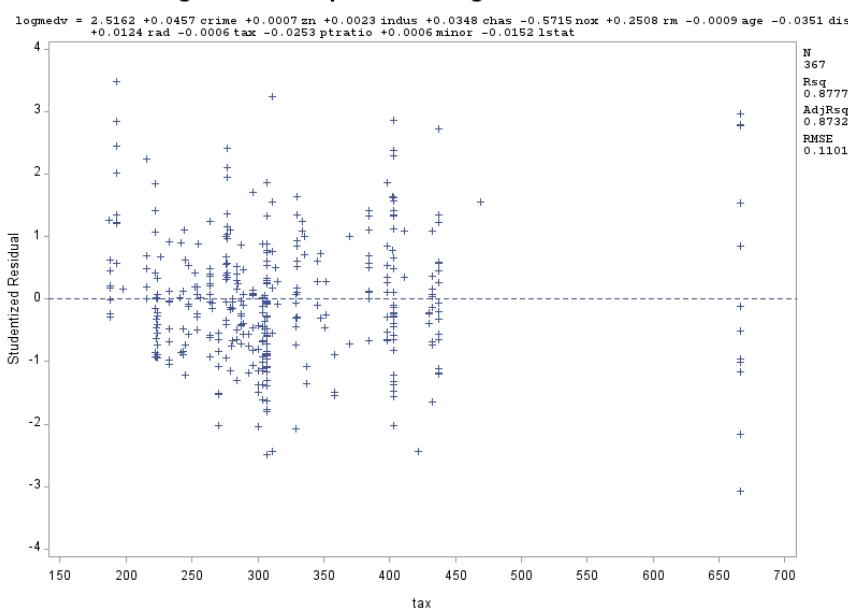


Fig. 12

Checking Model assumptions and diagnostics after transformation

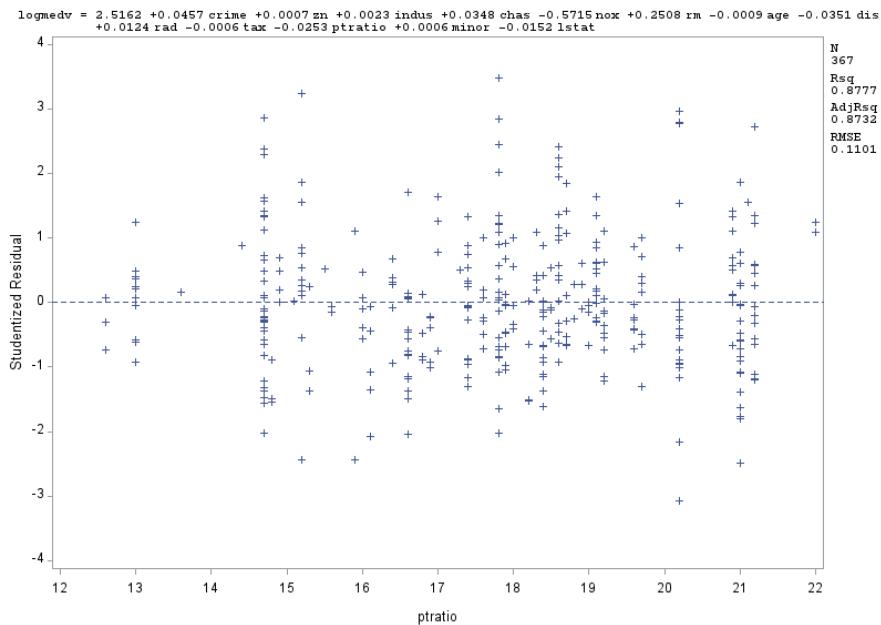


Fig. 13

Checking Model assumptions and diagnostics after transformation

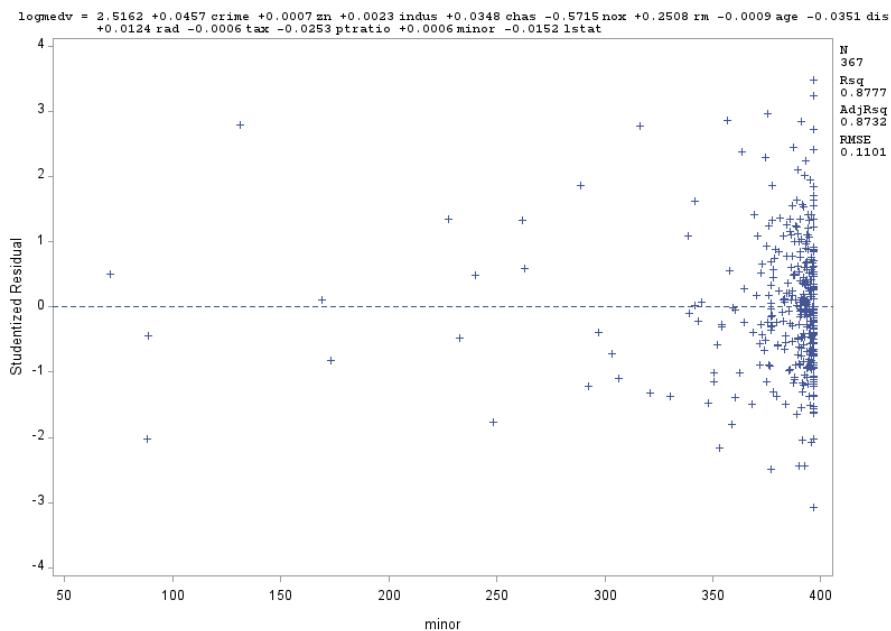


Fig. 14

Checking Model assumptions and diagnostics after transformation

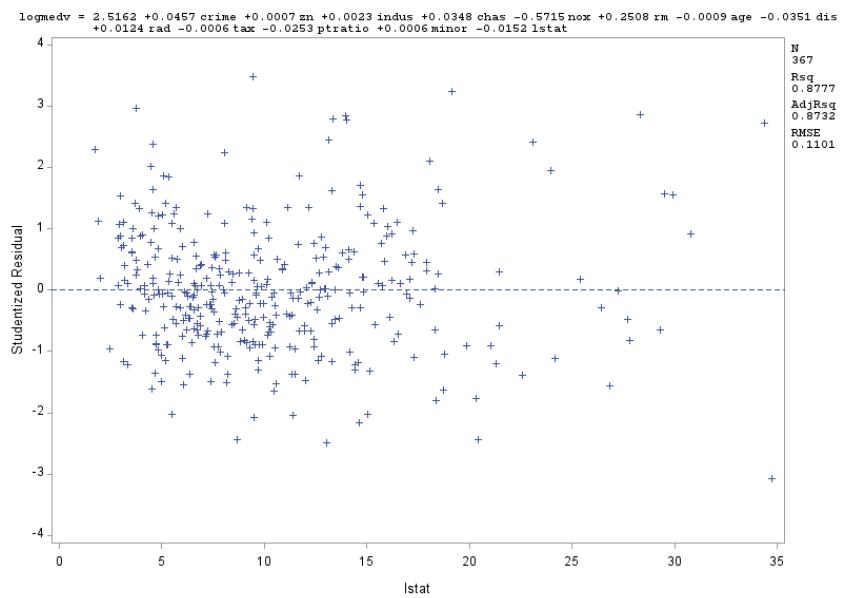


Fig. 15

Checking Model assumptions and diagnostics after transformation

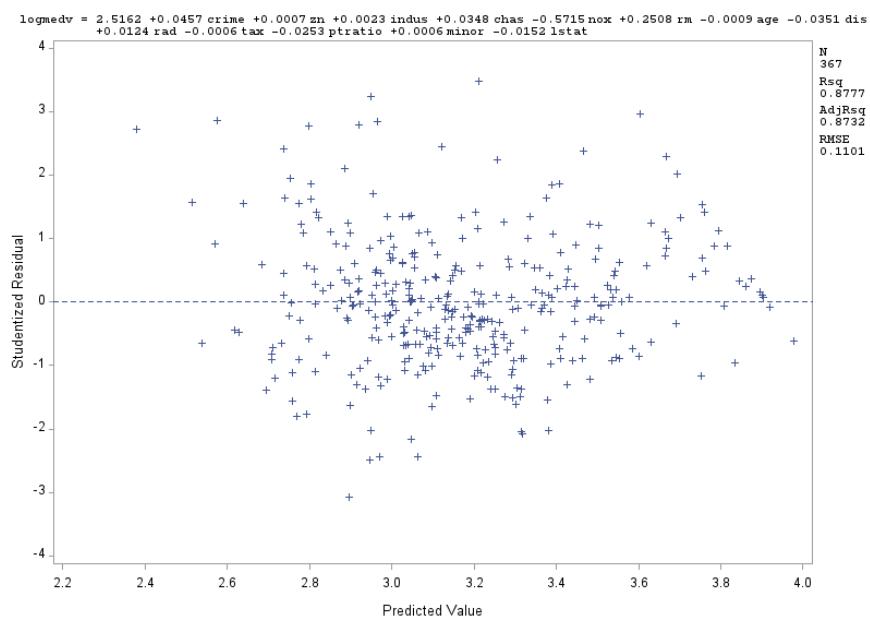


Fig. 16

Checking Model assumptions and diagnostics after transformation

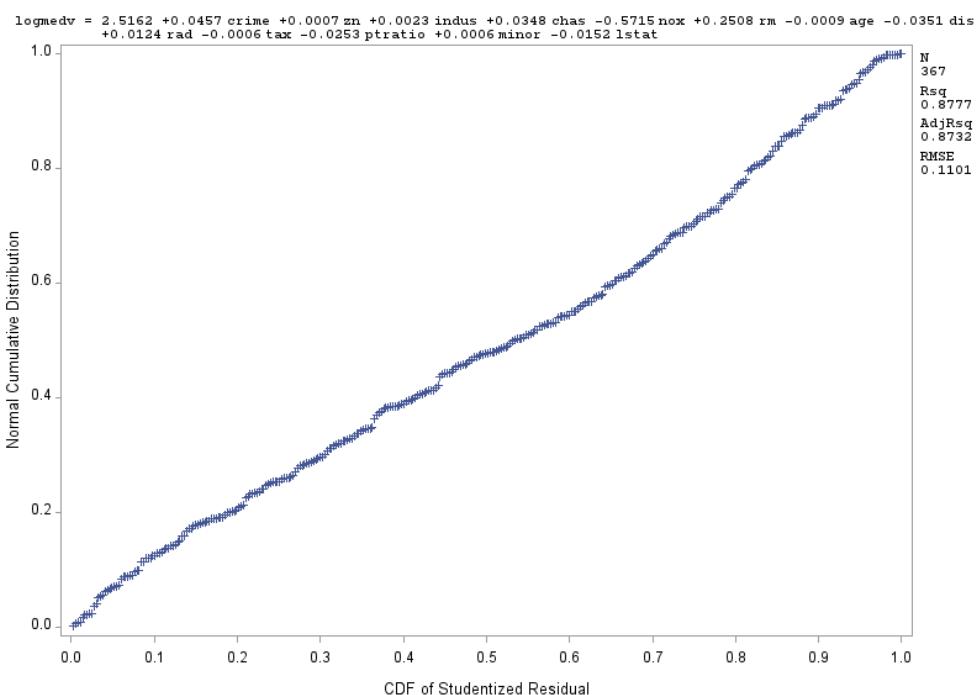


Table 2

Model 1

The REG Procedure
Model: MODEL1
Dependent Variable: new_medv

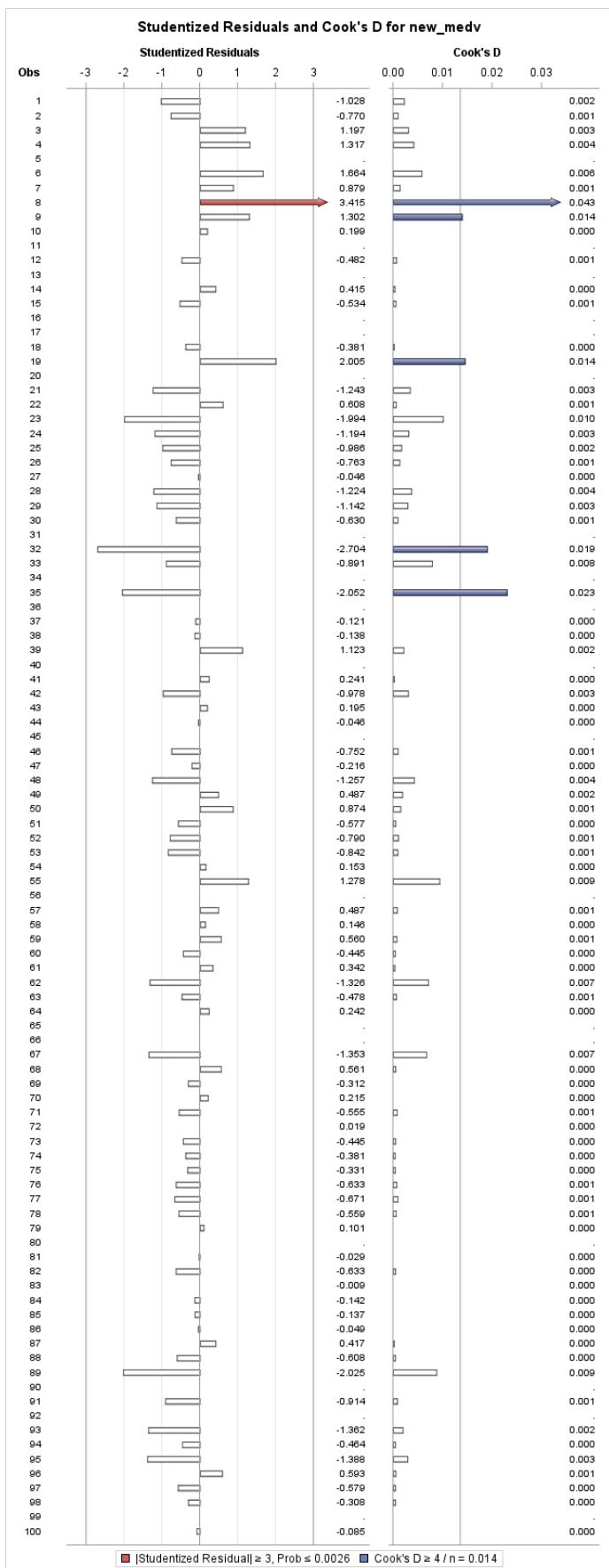
Number of Observations Read	367
Number of Observations Used	294
Number of Observations with Missing Values	73

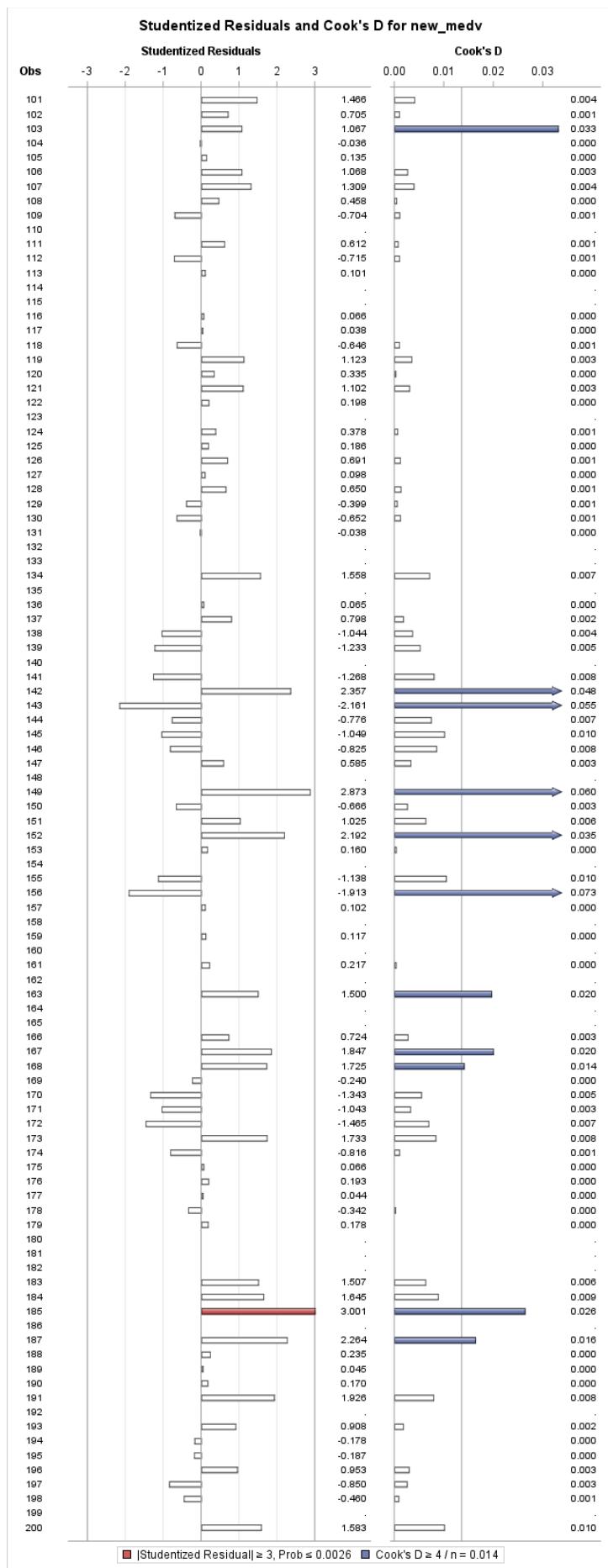
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	12	23.18421	1.93202	180.98	<.0001
Error	281	2.99983	0.01068		
Corrected Total	293	26.18404			

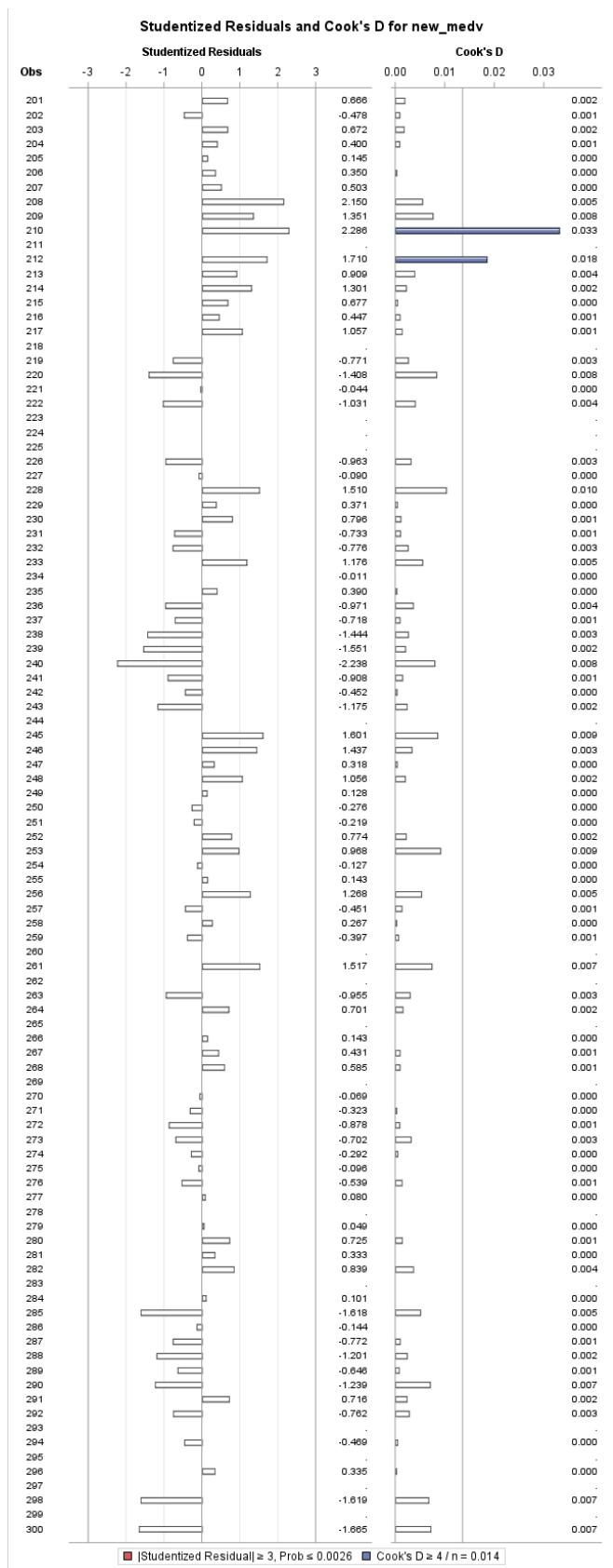
Root MSE	0.10332	R-Square	0.8854
Dependent Mean	3.16658	Adj R-Sq	0.8805
Coeff Var	3.26291		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	2.47906	0.16384	15.13	<.0001	0	0
crime	1	0.05996	0.00955	6.28	<.0001	0.24271	3.66359
zn	1	0.00076665	0.00035056	2.19	0.0296	0.06692	2.29668
chas	1	0.04569	0.02259	2.02	0.0440	0.04347	1.13248
nox	1	-0.66316	0.10826	-6.13	<.0001	-0.23820	3.70877
rm	1	0.25313	0.01476	17.15	<.0001	0.55179	2.53851
age	1	-0.00123	0.00035849	-3.42	0.0007	-0.11843	2.94335
dis	1	-0.03781	0.00520	-7.27	<.0001	-0.26690	3.30945
rad	1	0.01079	0.00304	3.55	0.0005	0.12819	3.20278
tax	1	-0.00048448	0.00010412	-4.65	<.0001	-0.14233	2.29480
ptratio	1	-0.02516	0.00341	-7.37	<.0001	-0.18400	1.52954
minor	1	0.00074009	0.00016305	4.54	<.0001	0.11146	1.47913
lstat	1	-0.01197	0.00182	-6.57	<.0001	-0.23457	3.12529

Fig. 17







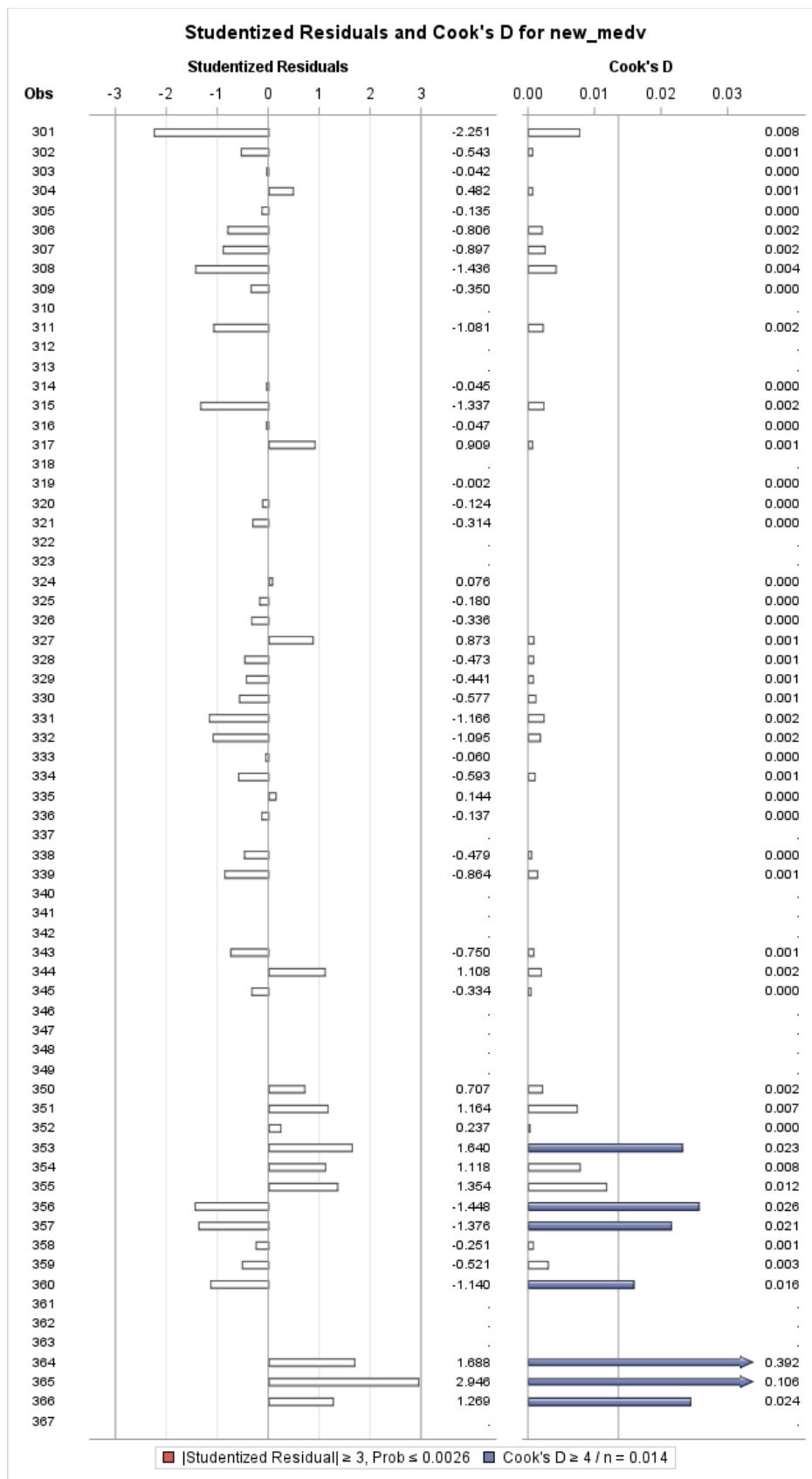


Fig. 18

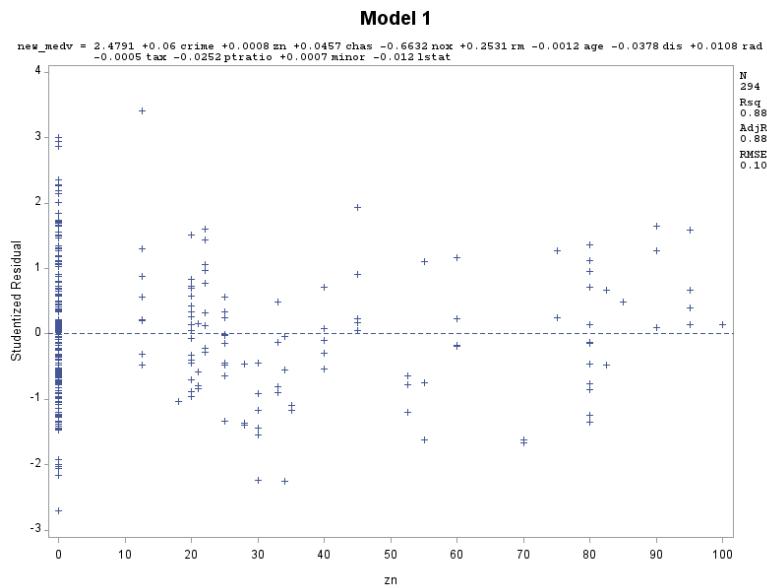


Fig. 19

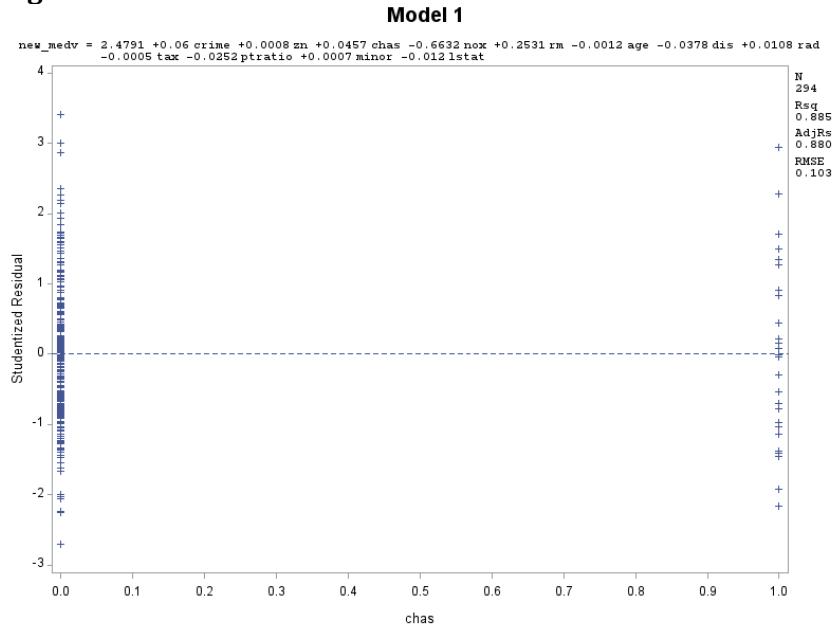
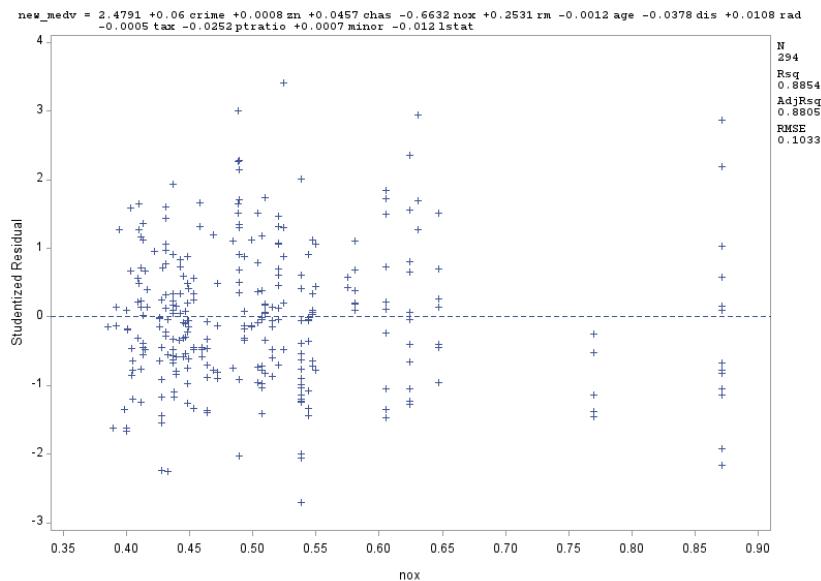
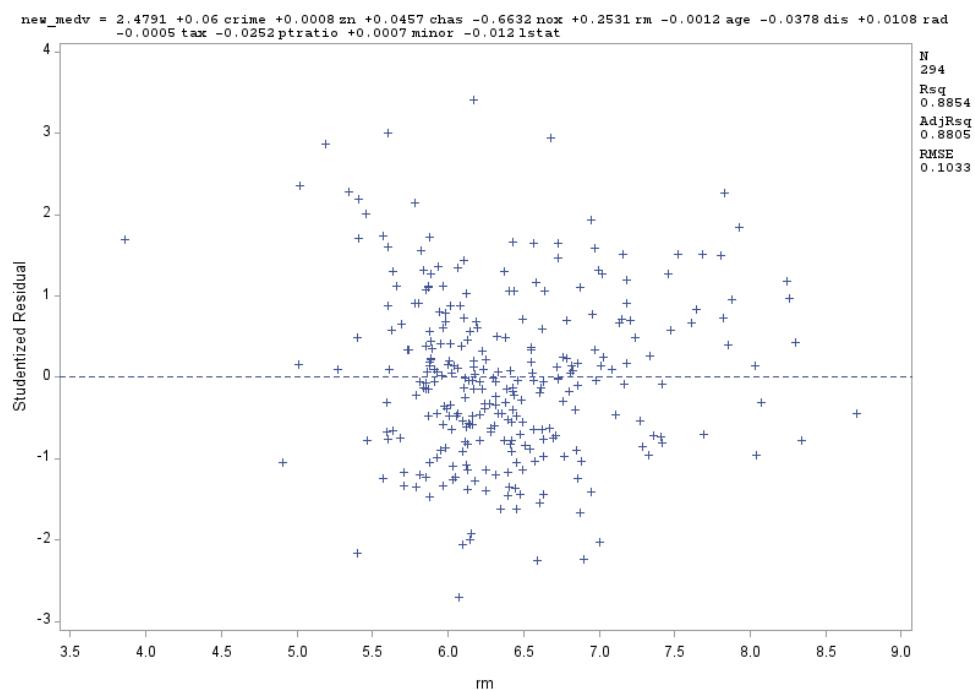


Fig. 20

Model 1**Fig. 21****Model 1****Fig. 22**

Model 1

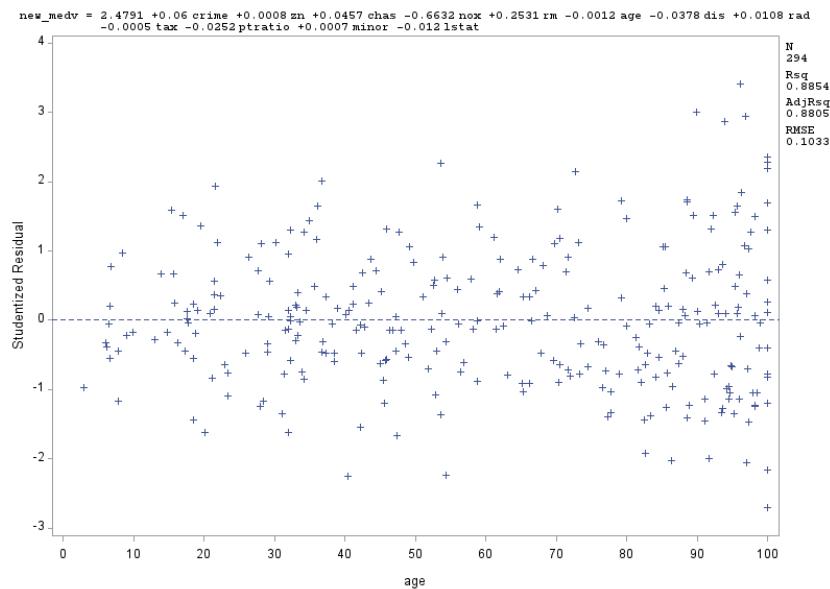


Fig. 23

Model 1

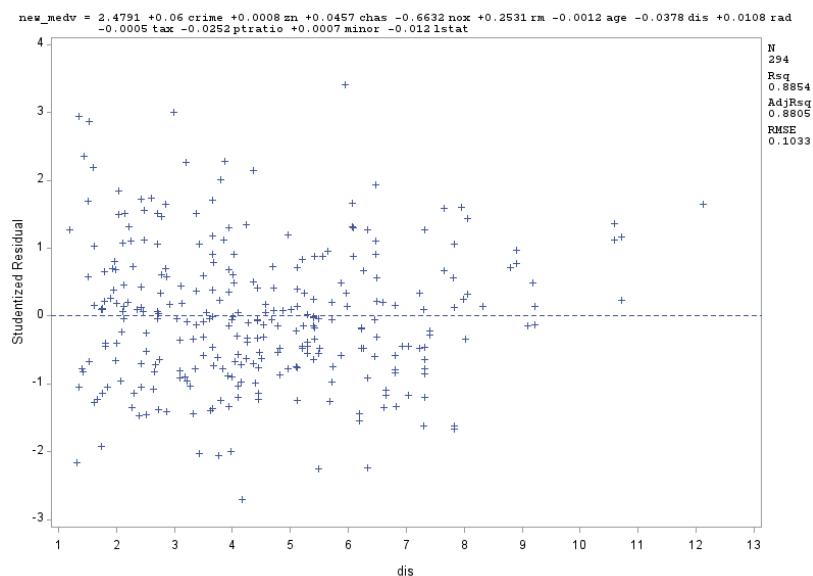


Fig. 24

Model 1

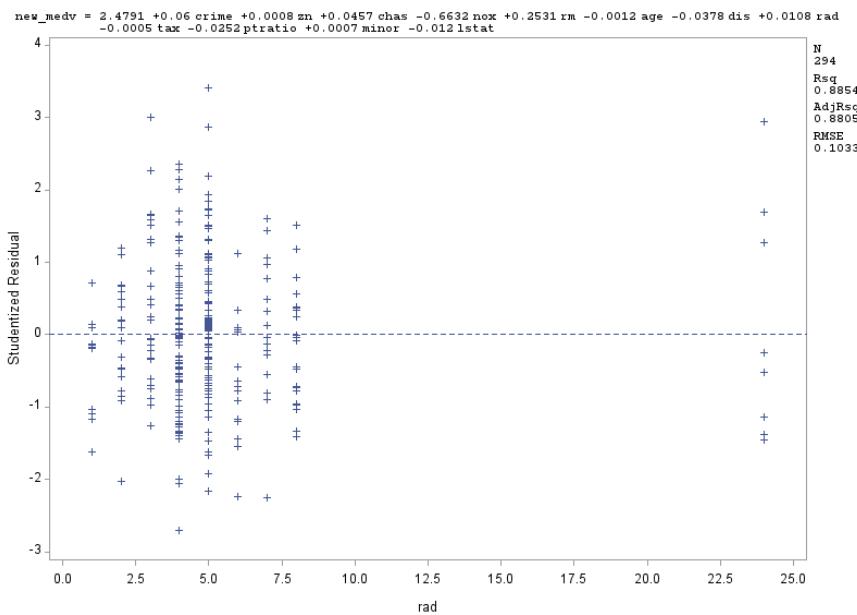


Fig. 25

Model 1

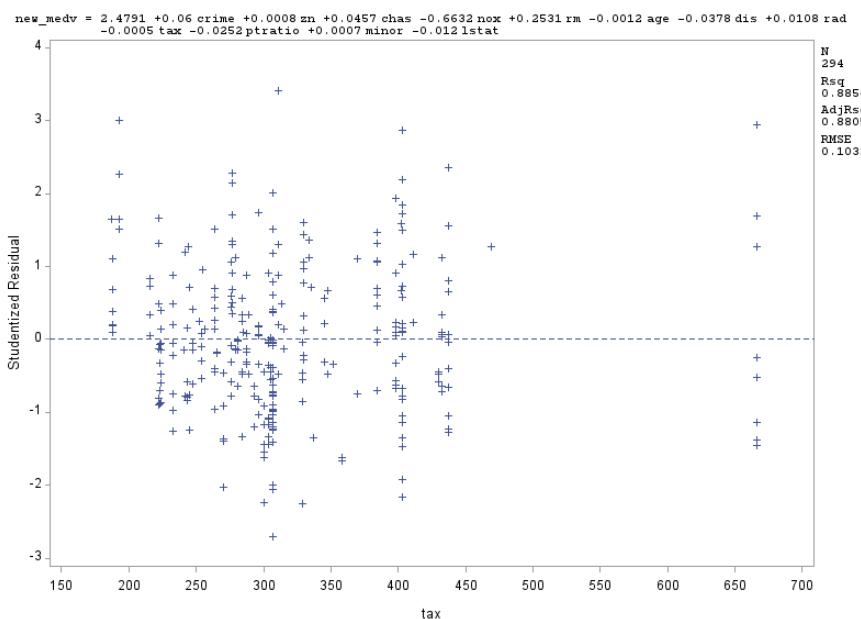


Fig. 26

Model 1

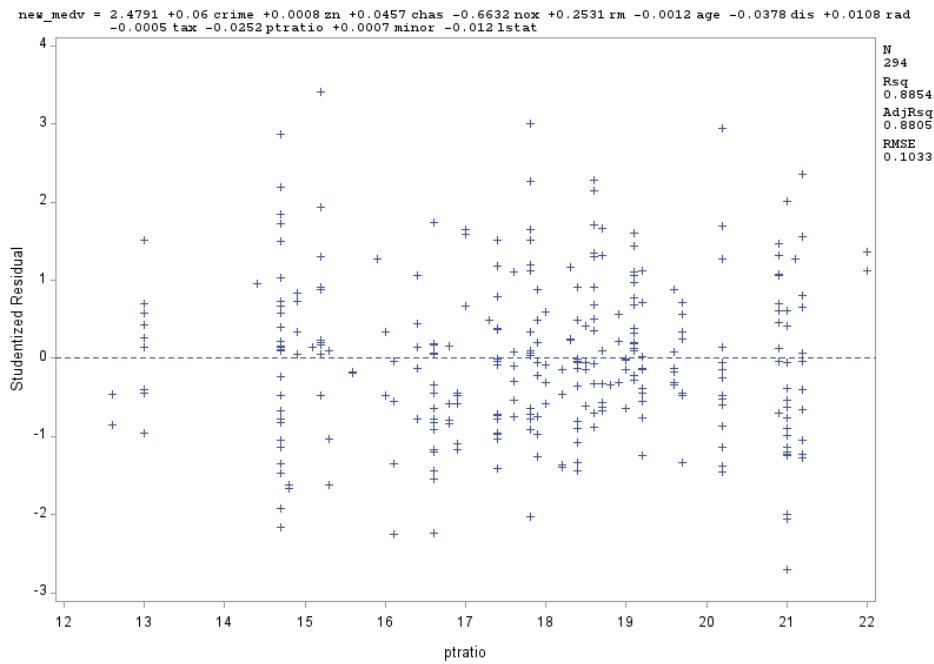


Fig. 27

Model 1

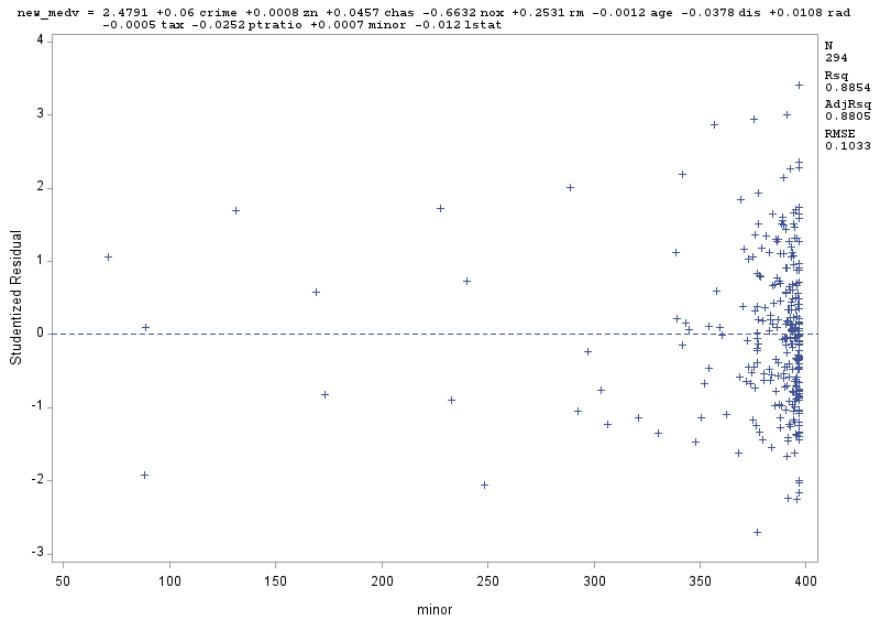
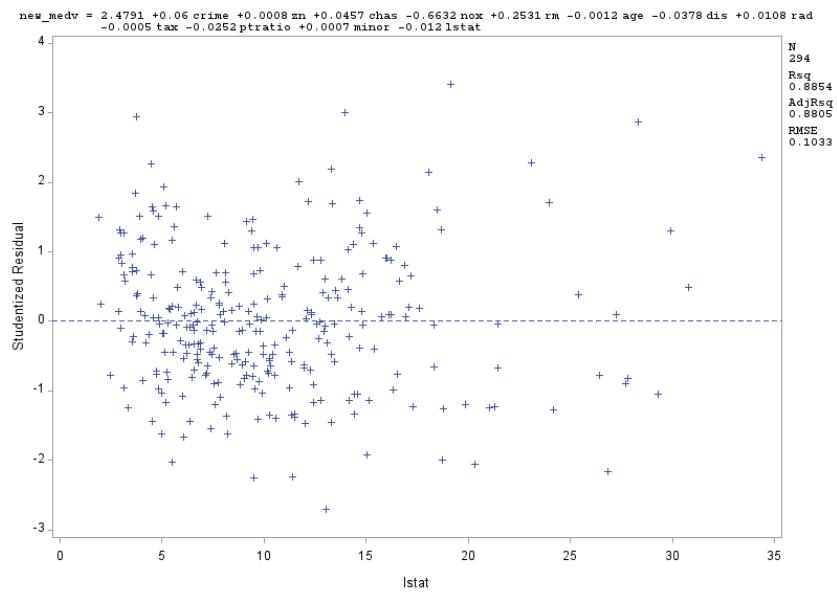
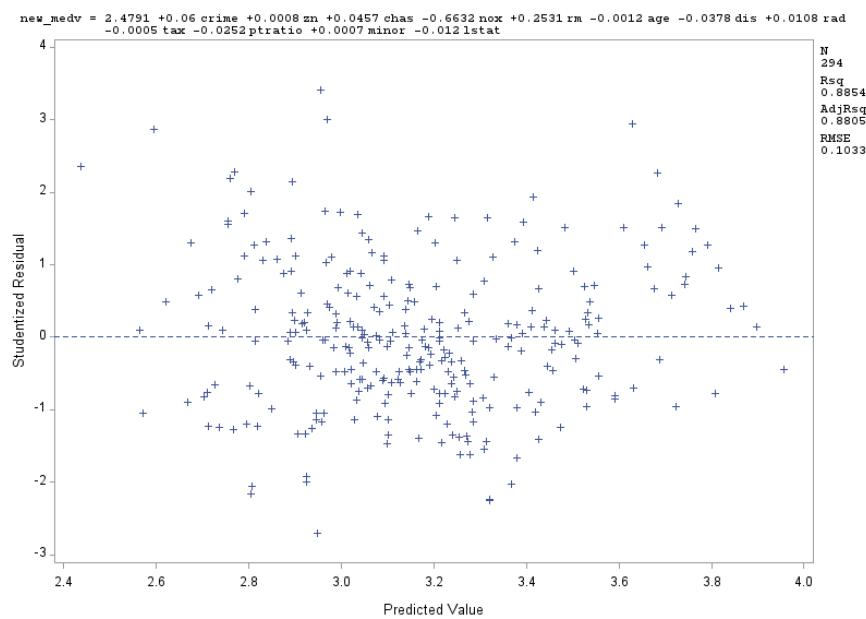


Fig. 28

Model 1**Fig. 29****Model 1****Fig. 30**

Model 1

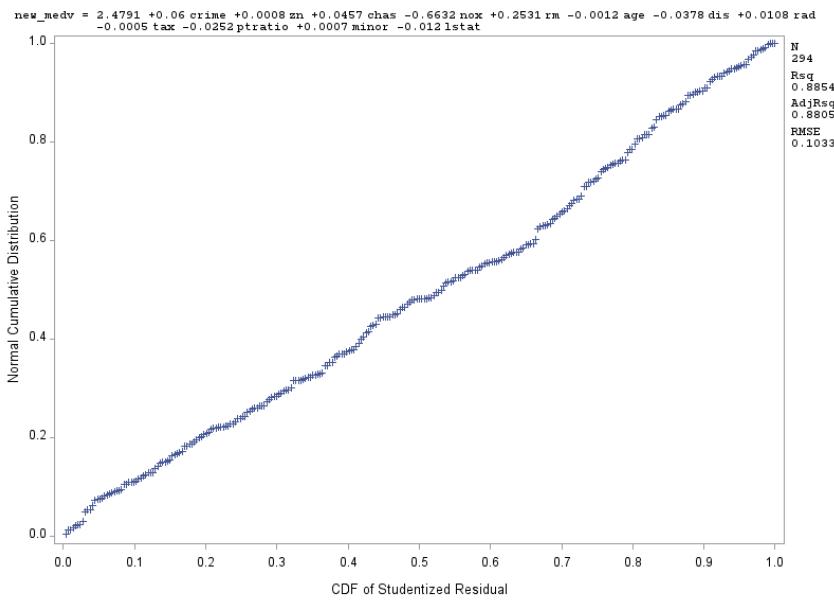


Table 3

Model 1 after removing outliers

The REG Procedure

Model: MODEL1

Dependent Variable: new_medv

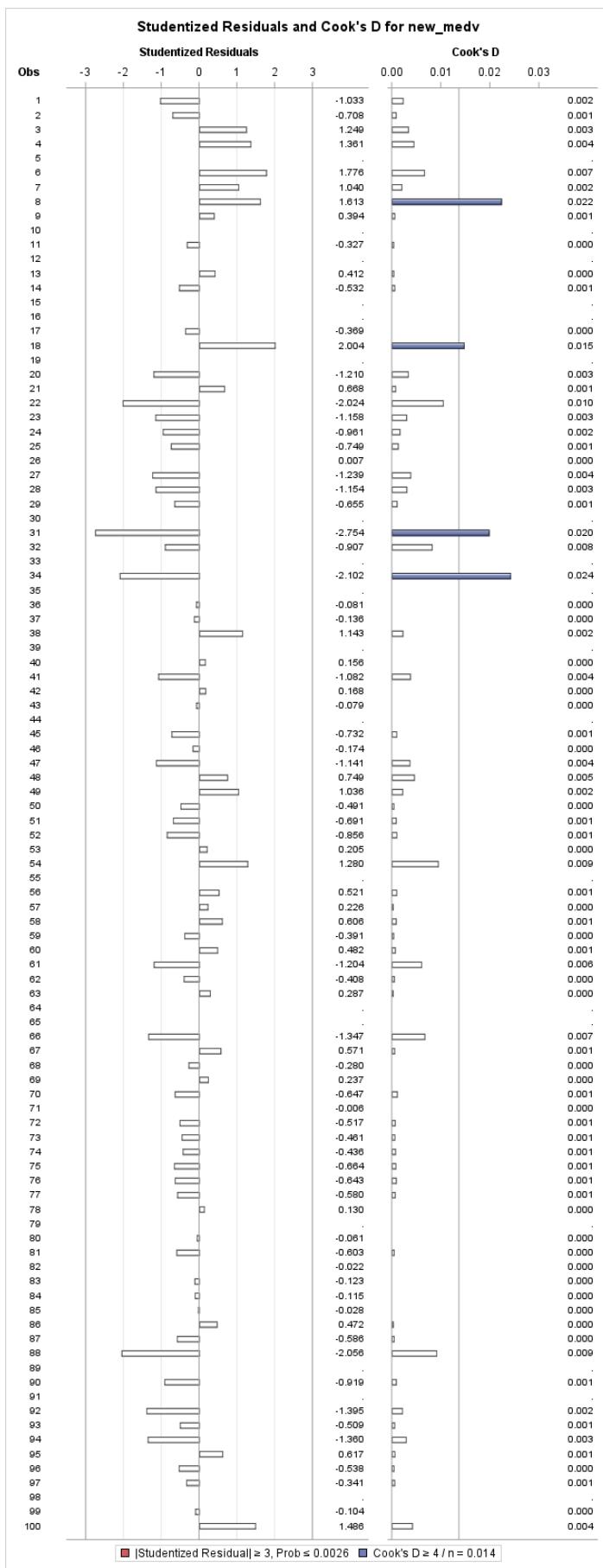
Number of Observations Read	365
Number of Observations Used	292
Number of Observations with Missing Values	73

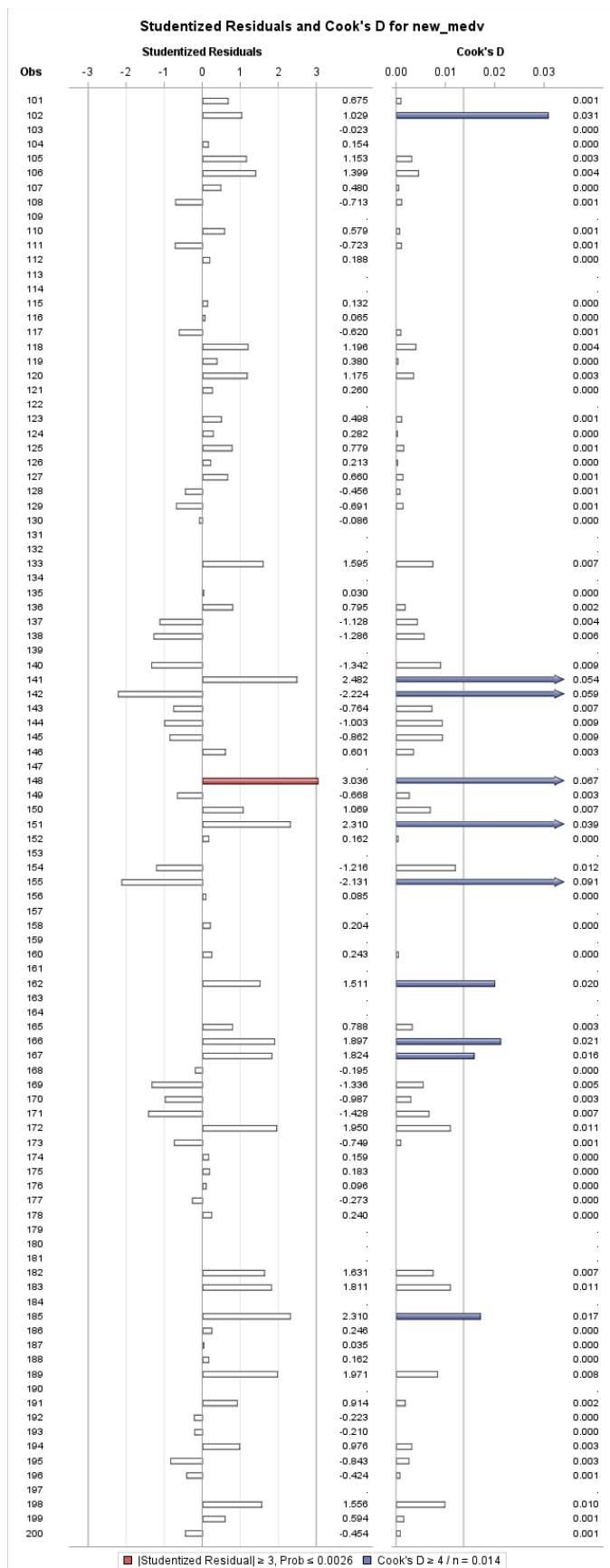
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	12	23.37772	1.94814	195.72	<.0001
Error	279	2.77704	0.00995		
Corrected Total	291	26.15476			

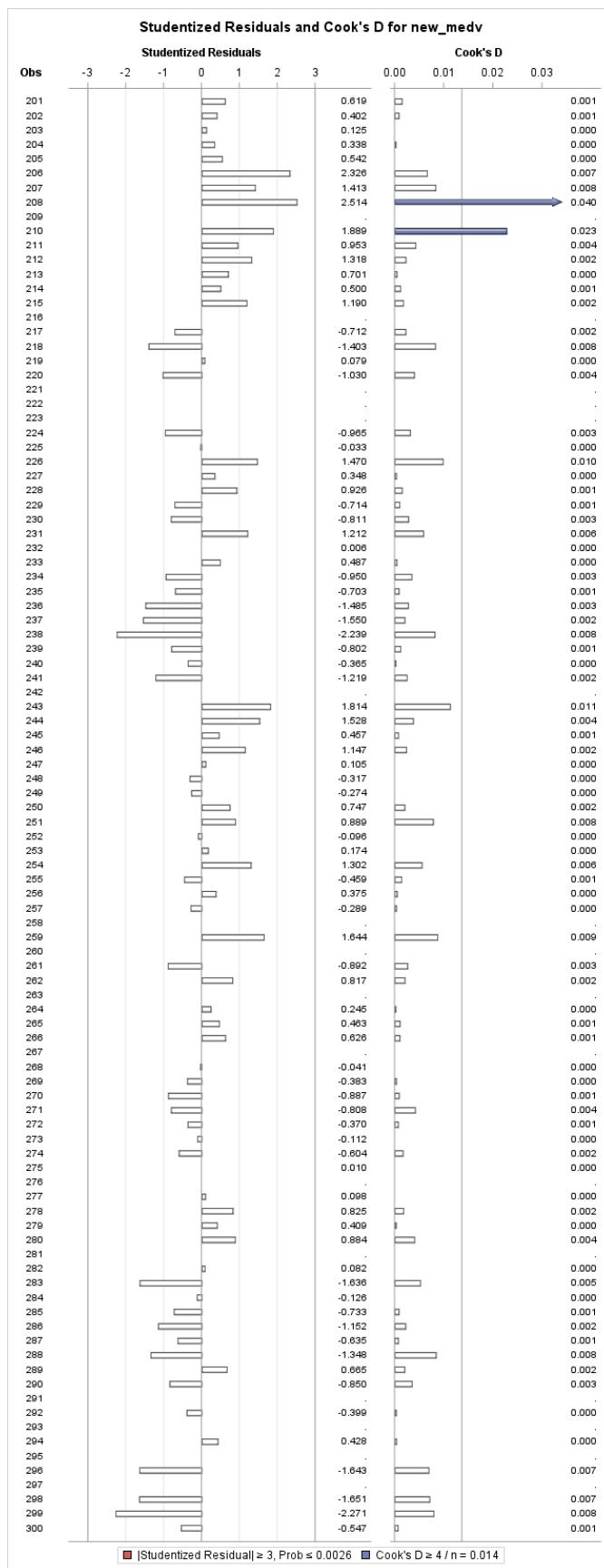
Root MSE	0.09977	R-Square	0.8938
Dependent Mean	3.16576	Adj R-Sq	0.8893
Coeff Var	3.15146		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	2.42099	0.15932	15.20	<.0001	0	0
crime	1	0.06120	0.00922	6.64	<.0001	0.24779	3.66380
zn	1	0.00080050	0.00033922	2.36	0.0190	0.06986	2.30320
chas	1	0.04938	0.02182	2.26	0.0244	0.04699	1.13319
nox	1	-0.63361	0.10475	-6.05	<.0001	-0.22768	3.72284
rm	1	0.25741	0.01436	17.93	<.0001	0.55996	2.56263
age	1	-0.00143	0.00034877	-4.09	<.0001	-0.13738	2.96147
dis	1	-0.03883	0.00506	-7.67	<.0001	-0.27379	3.34696
rad	1	0.00988	0.00294	3.36	0.0009	0.11737	3.21295
tax	1	-0.00045523	0.00010106	-4.50	<.0001	-0.13333	2.30206
ptratio	1	-0.02333	0.00332	-7.02	<.0001	-0.17029	1.54708
minor	1	0.00072210	0.00015749	4.58	<.0001	0.10877	1.47888
lstat	1	-0.01217	0.00177	-6.89	<.0001	-0.23753	3.12320

Fig. 31







Studentized Residuals and Cook's D for new_medv

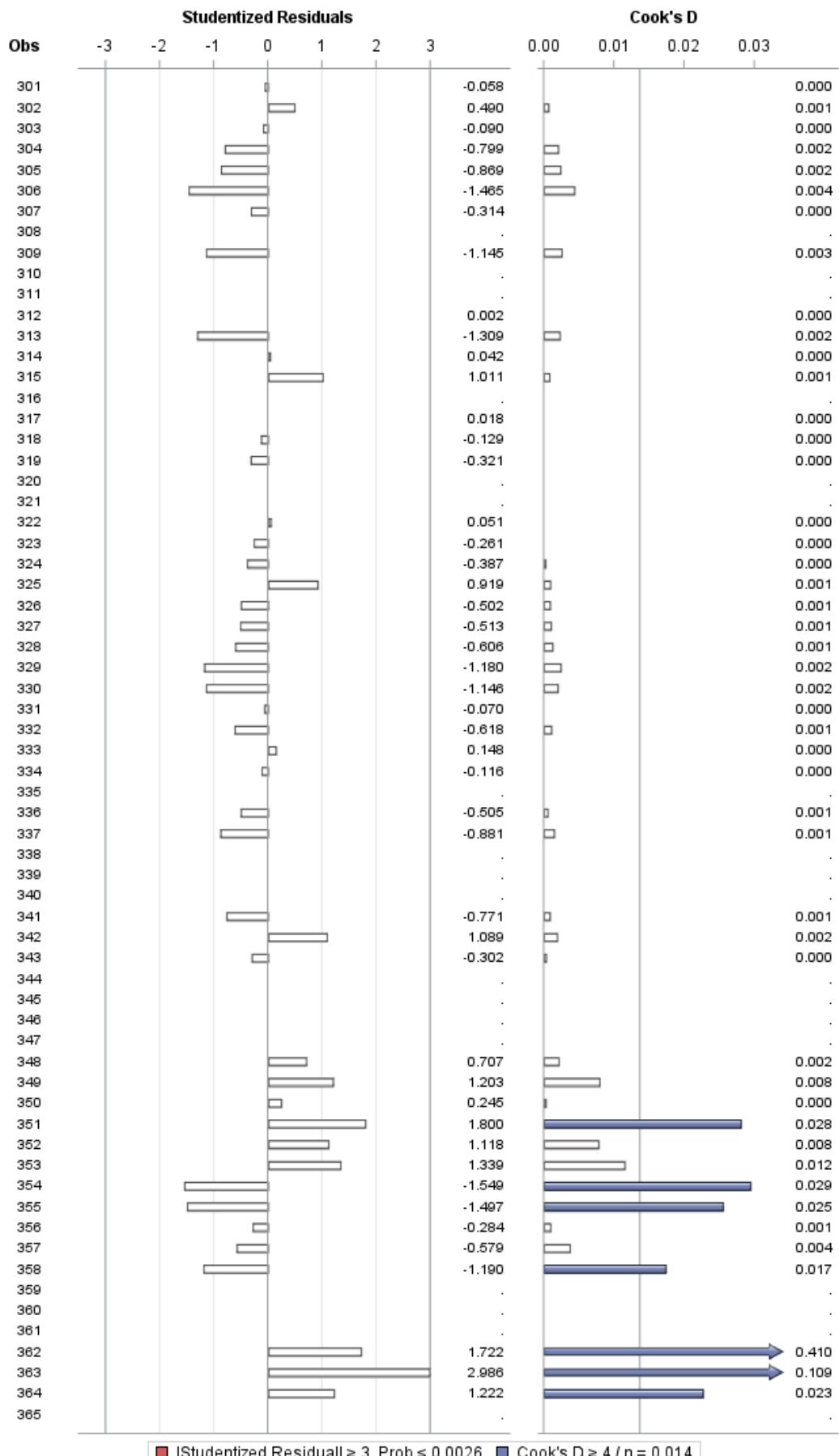


Fig. 32

Model 1 after removing outliers

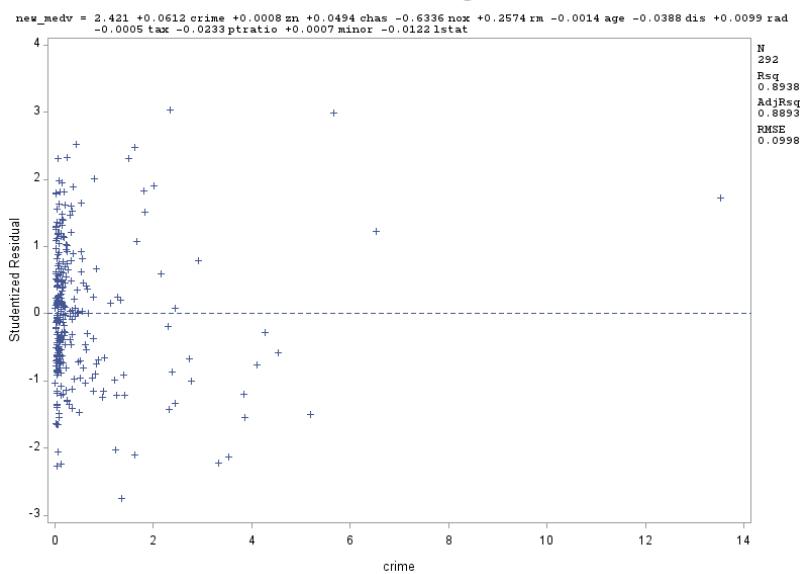


Fig. 33

Model 1 after removing outliers

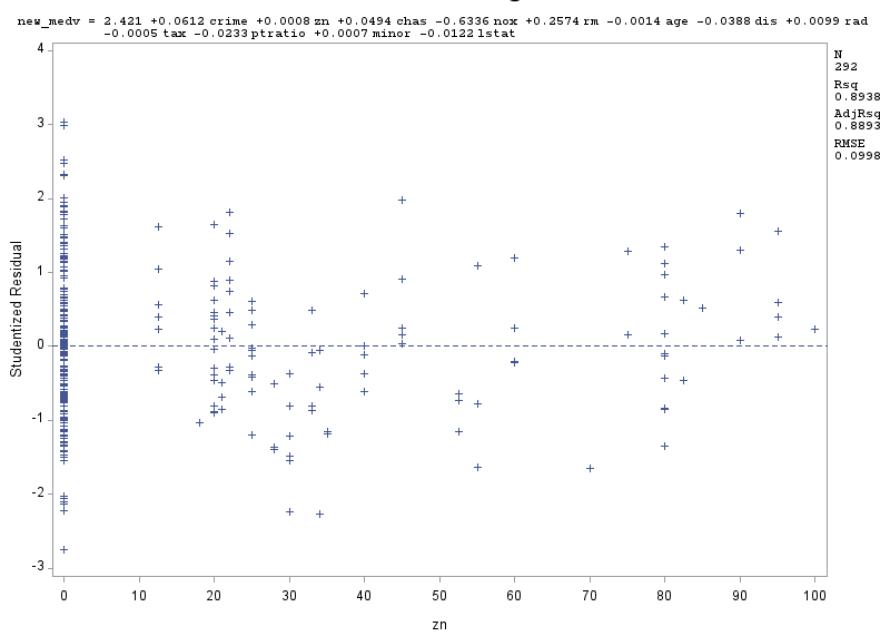


Fig. 34

Model 1 after removing outliers

new_medv = 2.421 +0.0612 crime +0.0008 zn +0.0494 chas -0.6336 nox +0.2574 rm -0.0014 age -0.0388 dis +0.0099 rad
 -0.0005 tax -0.0233 ptratio +0.0007 minor -0.0122 lstat

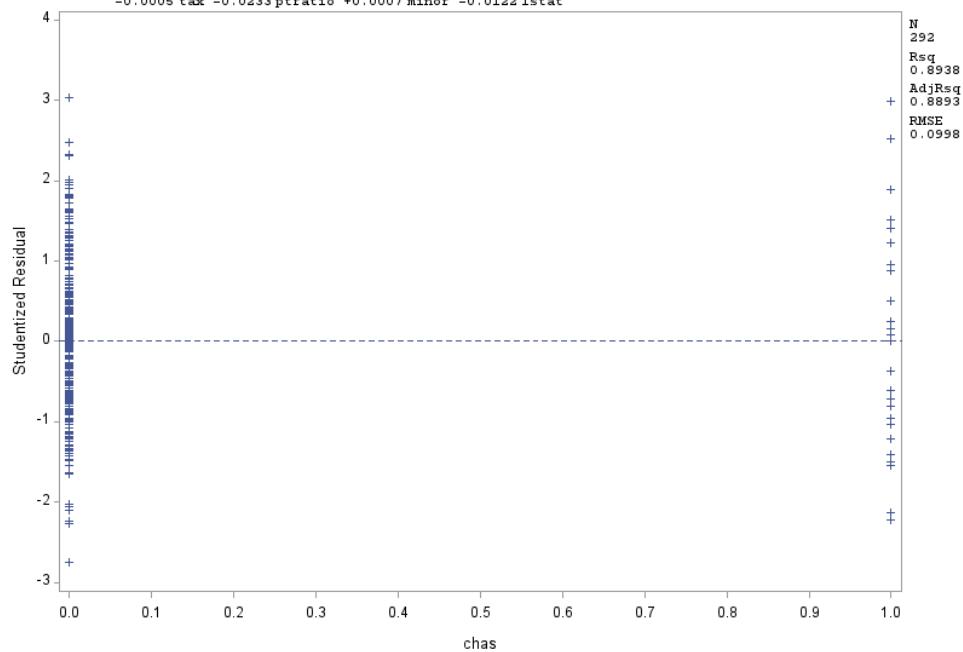


Fig. 35

Model 1 after removing outliers

new_medv = 2.421 +0.0612 crime +0.0008 zn +0.0494 chas -0.6336 nox +0.2574 rm -0.0014 age -0.0388 dis +0.0099 rad
 -0.0005 tax -0.0233 ptratio +0.0007 minor -0.0122 lstat

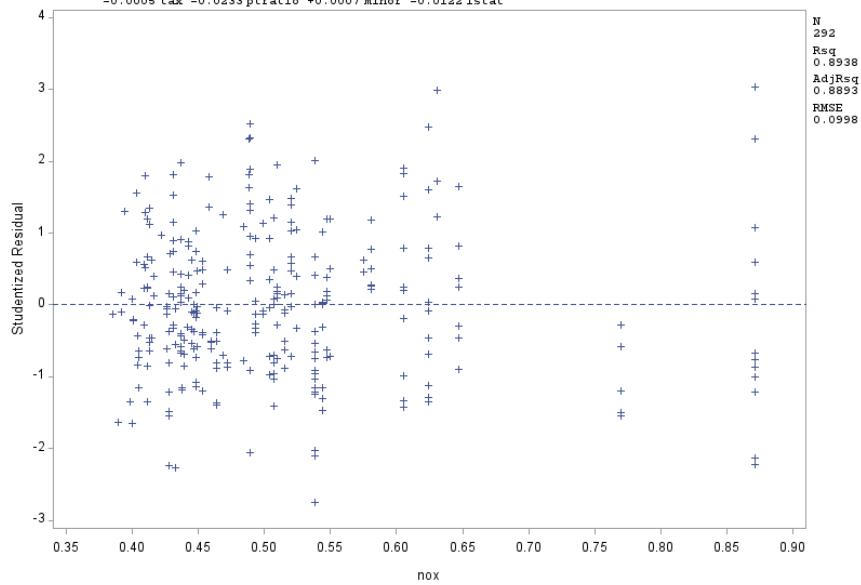


Fig. 36

Model 1 after removing outliers

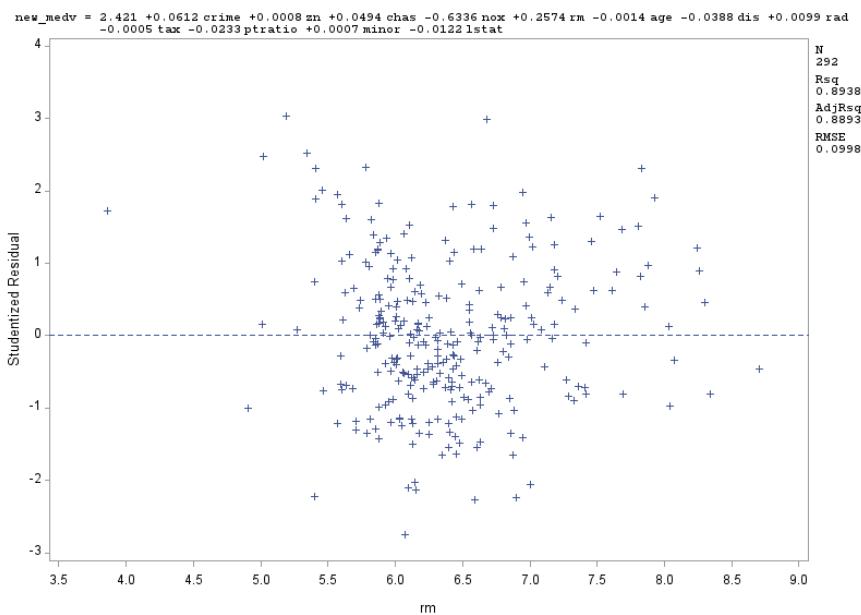


Fig. 37

Model 1 after removing outliers

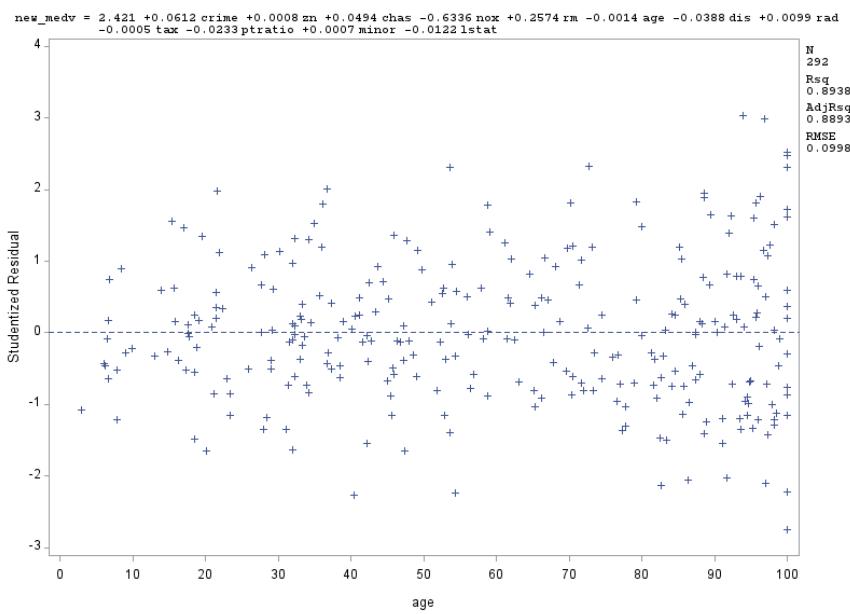


Fig. 38

Model 1 after removing outliers

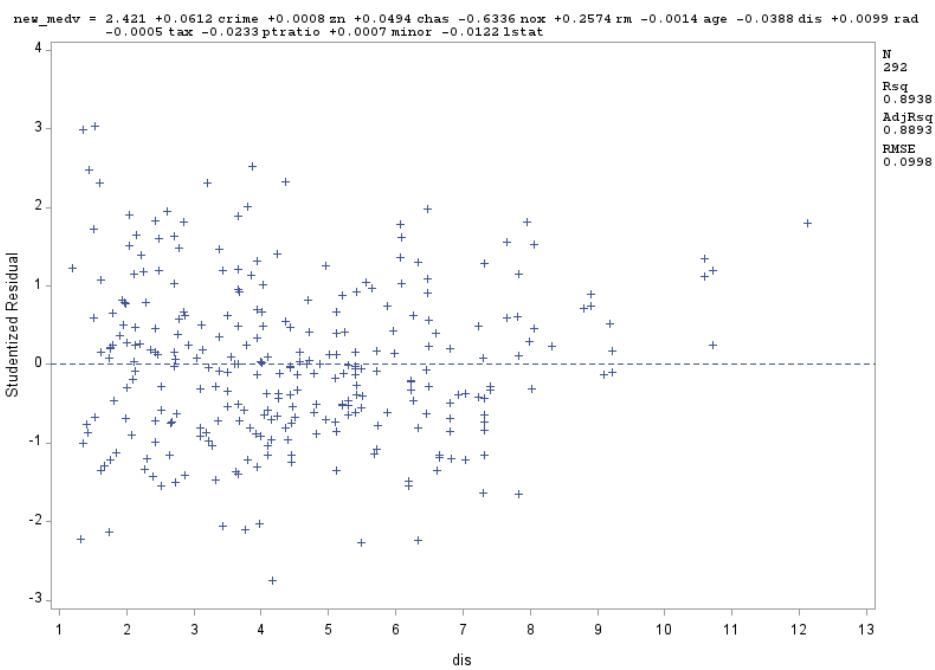


Fig. 39

Model 1 after removing outliers

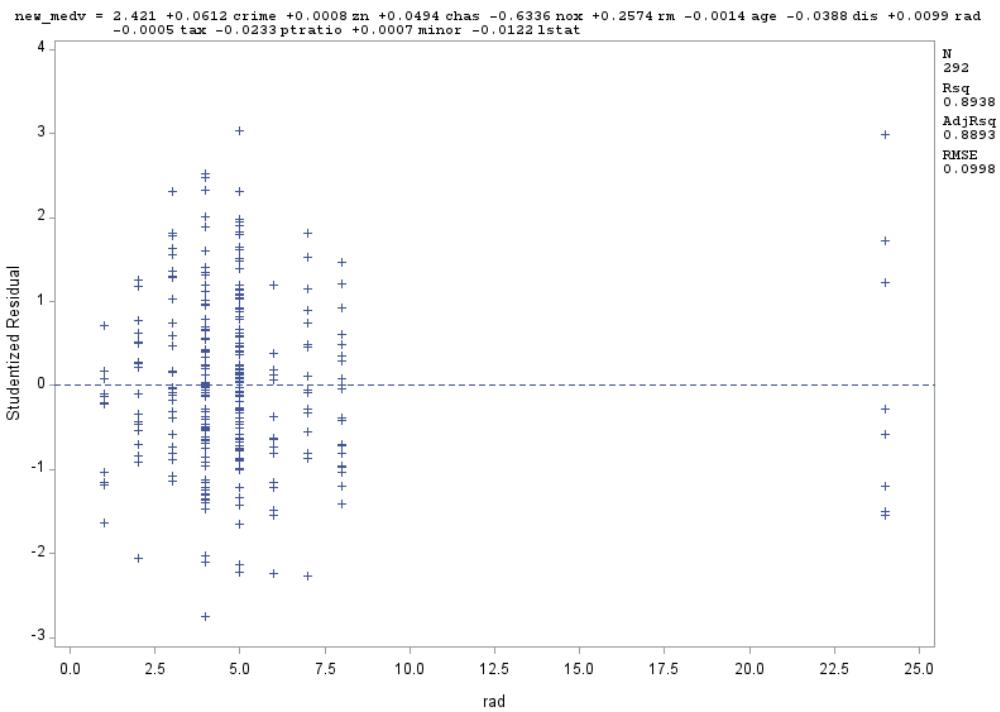


Fig. 40

Model 1 after removing outliers

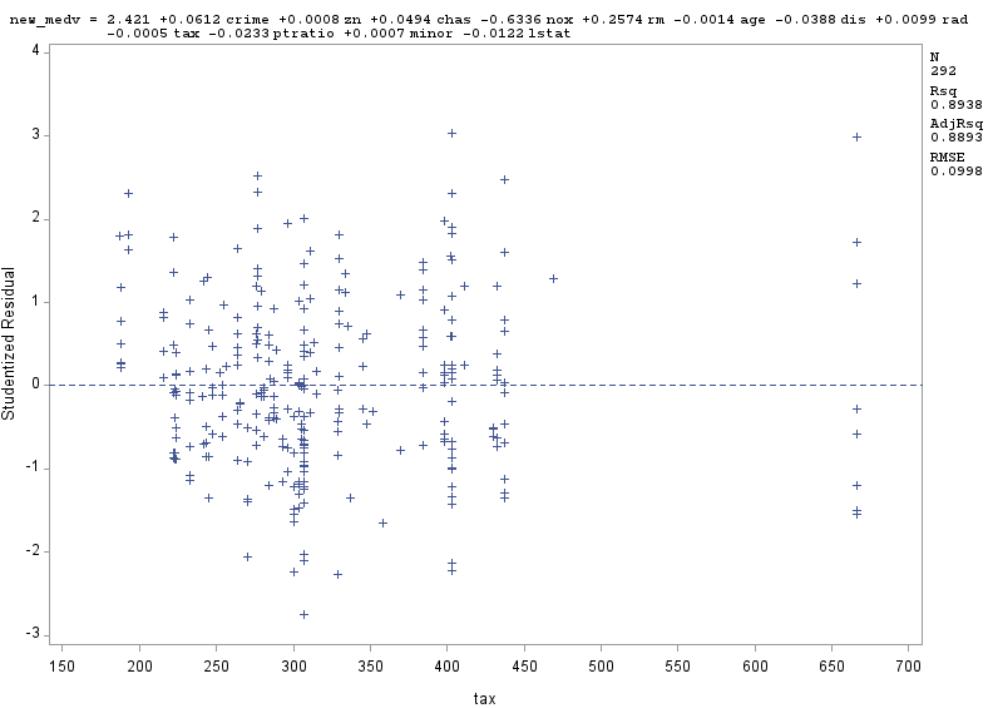


Fig. 41

Model 1 after removing outliers

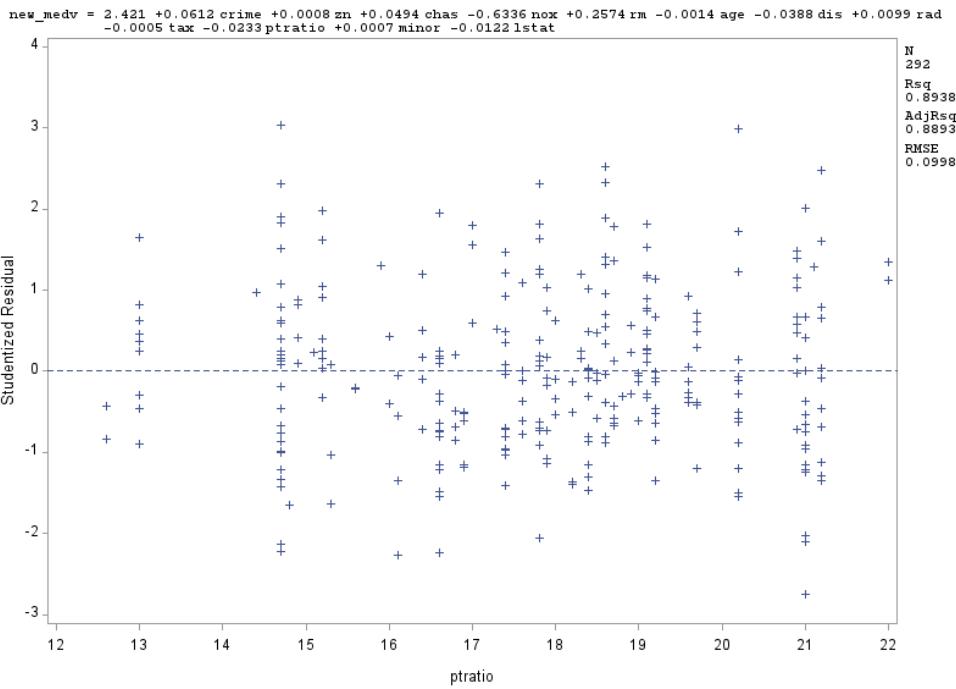


Fig. 42

Model 1 after removing outliers

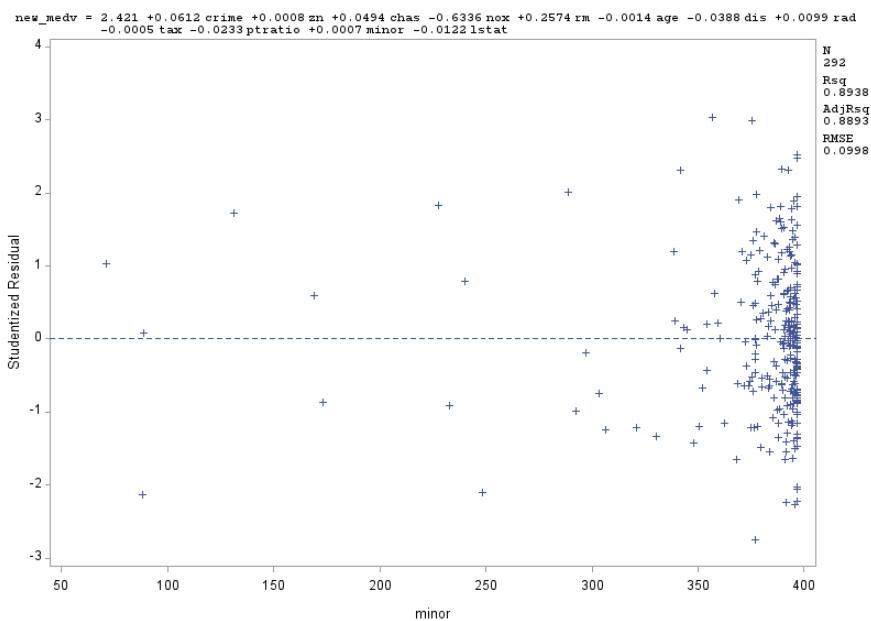


Fig. 43

Model 1 after removing outliers

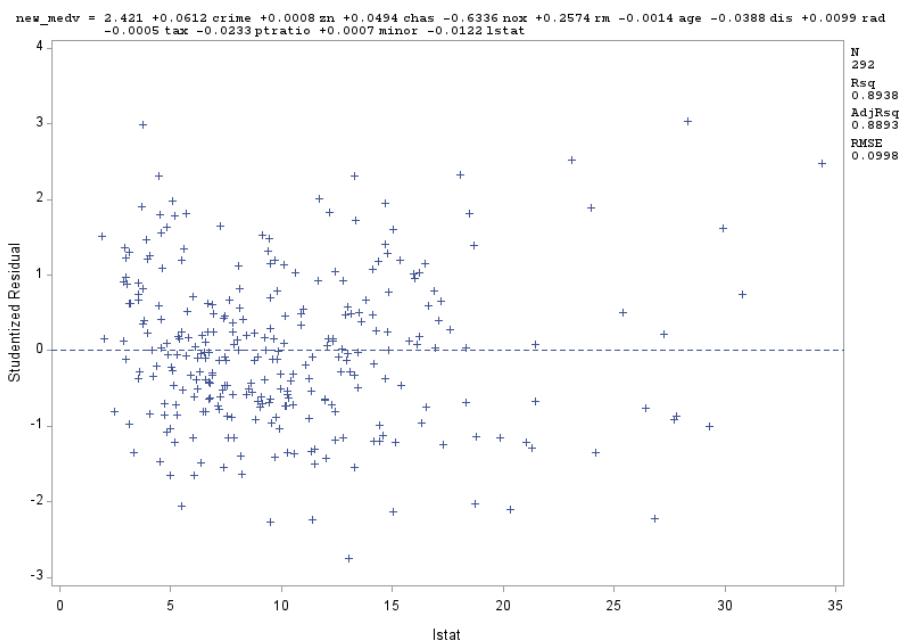


Fig. 44

Model 1 after removing outliers

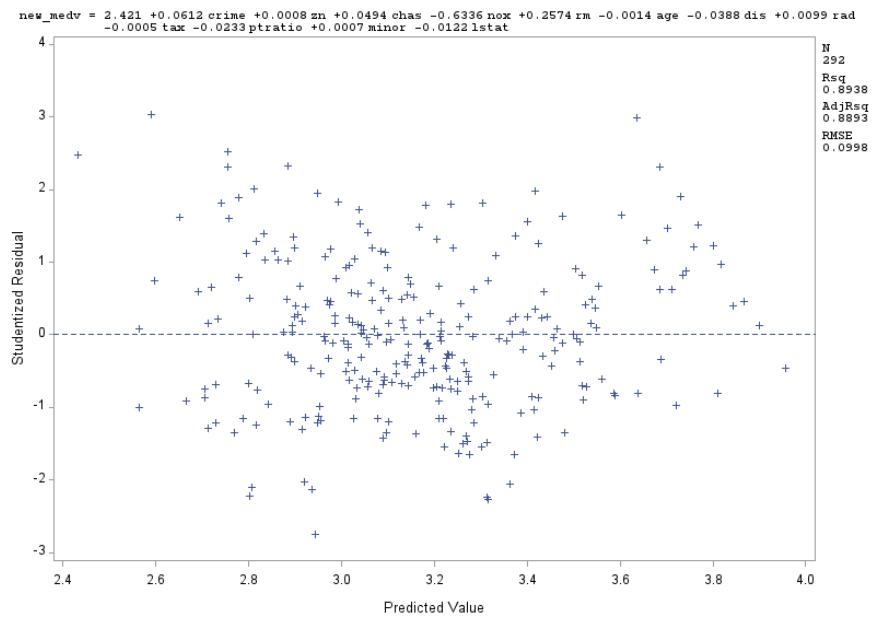


Fig. 45

Model 1 after removing outliers

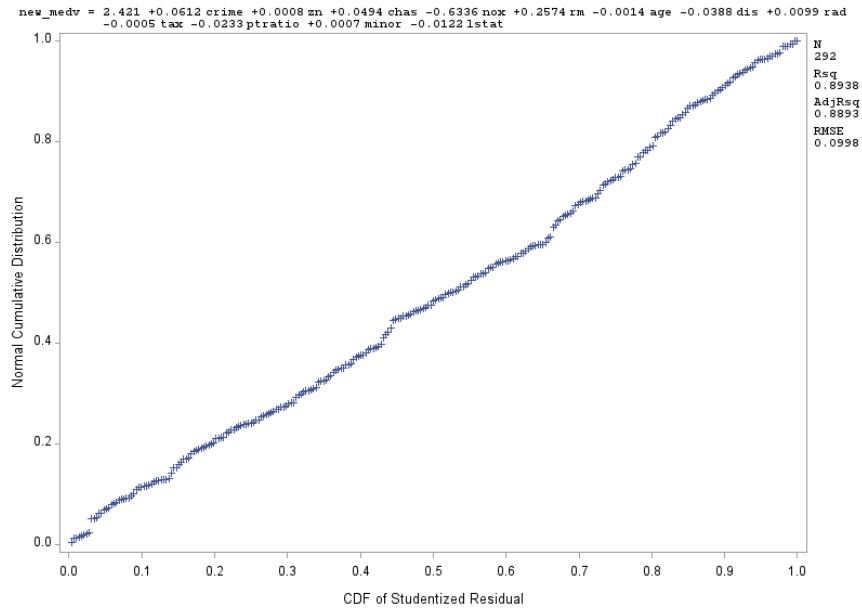


Table 4

Performance test parameters - Model 1

Obs	_TYPE_	_FREQ_	rmse	mae
1	0	73	9.57828	22.6199

Performance test parameters - Model 1

The CORR Procedure

2 Variables:	medv yhat
---------------------	-----------

Simple Statistics							
Variable	N	Mean	Std Dev	Sum	Minimum	Maximum	Label
medv	73	25.79041	9.85777	1883	12.70000	50.00000	
yhat	73	3.17053	0.30402	231.44860	2.53534	3.91443	Predicted Value of new_medv

Pearson Correlation Coefficients, N = 73 Prob > r under H0: Rho=0		
	medv	yhat
medv	1.00000	0.92171 <.0001
yhat Predicted Value of new_medv	0.92171 <.0001	1.00000

APPENDIX D – Quadratic Polynomial Regression

Table 1

Quadratic Polynomial Regression

The GLMSELECT Procedure
Least Squares Model (No Selection)

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	91	25438	279.53963	60.20	<.0001
Error	275	1276.97106	4.64353		
Corrected Total	366	26715			

Root MSE	2.15489
Dependent Mean	25.04959
R-Square	0.9522
Adj R-Sq	0.9364
AIC	1010.60656
AICC	1074.65051
SBC	1000.89985

Table 2

Model 1- Checking Model assumptions and diagnostics

The REG Procedure

Model: MODEL1

Dependent Variable: new_medv

Number of Observations Read	367
Number of Observations Used	276
Number of Observations with Missing Values	91

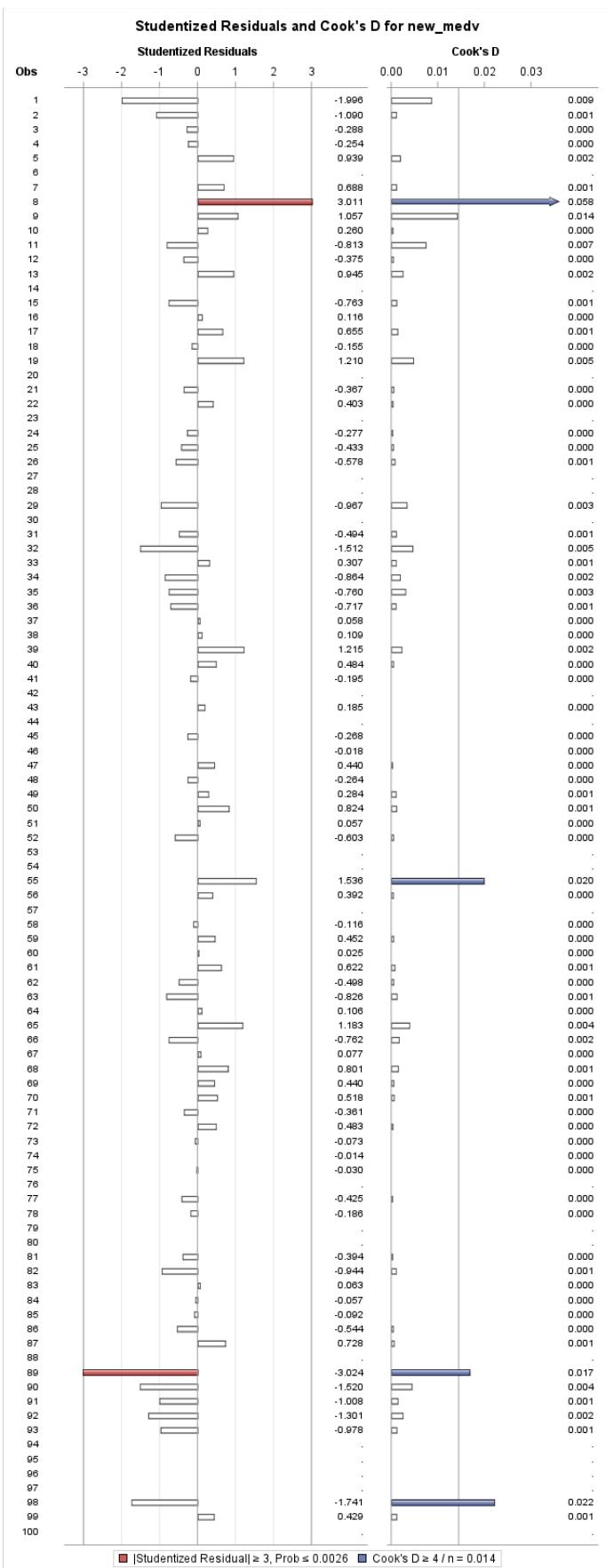
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	16519	3303.77401	341.57	<.0001
Error	270	2611.49196	9.67219		
Corrected Total	275	19130			

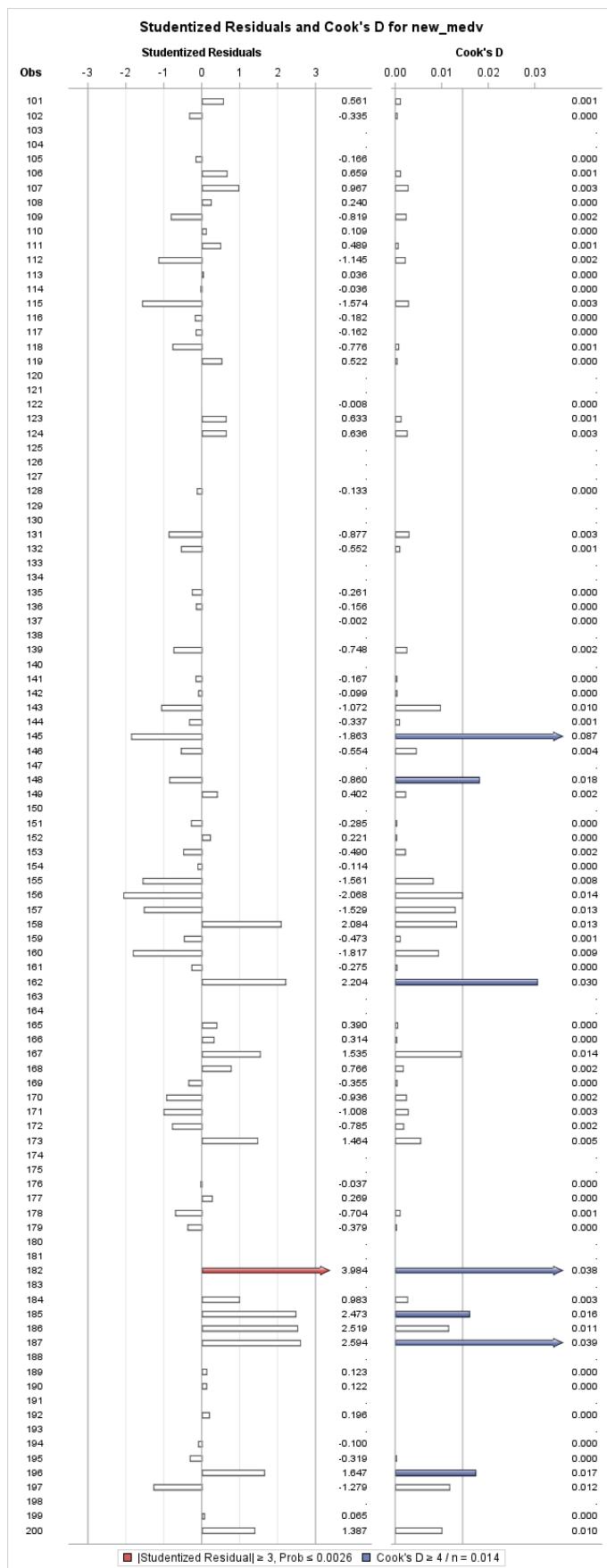
Root MSE	3.11001	R-Square	0.8635
Dependent Mean	24.78551	Adj R-Sq	0.8610
Coeff Var	12.54772		

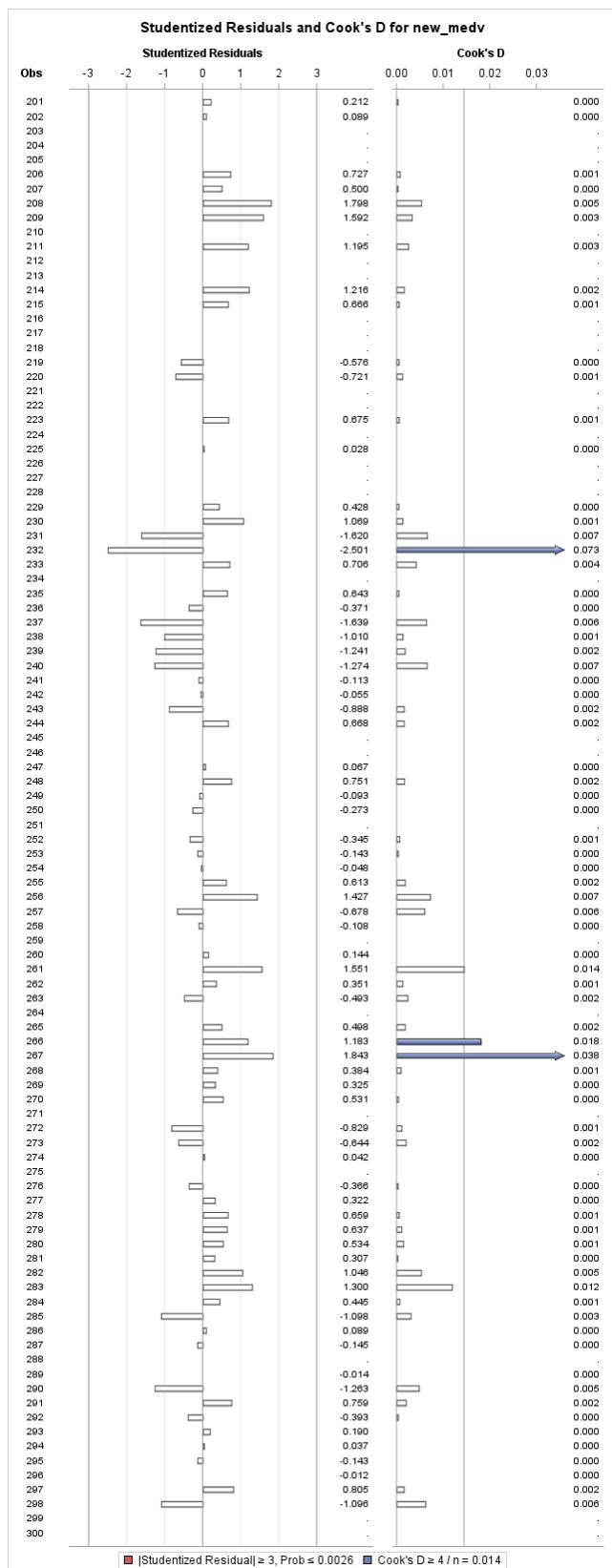
Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	-40.38857	3.41815	-11.82	<.0001	0	0
rm	1	12.56605	0.48493	25.91	<.0001	1.00384	2.96823

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
dis_tax	1	-0.00151	0.00029123	-5.18	<.0001	-0.12445	1.14352
rm_ptratio	1	-0.06595	0.01397	-4.72	<.0001	-0.12075	1.29385
lstat	1	2.33246	0.28251	8.26	<.0001	1.67609	81.51468
rm_lstat	1	-0.46491	0.04915	-9.46	<.0001	-1.77338	69.53019

Fig. 1







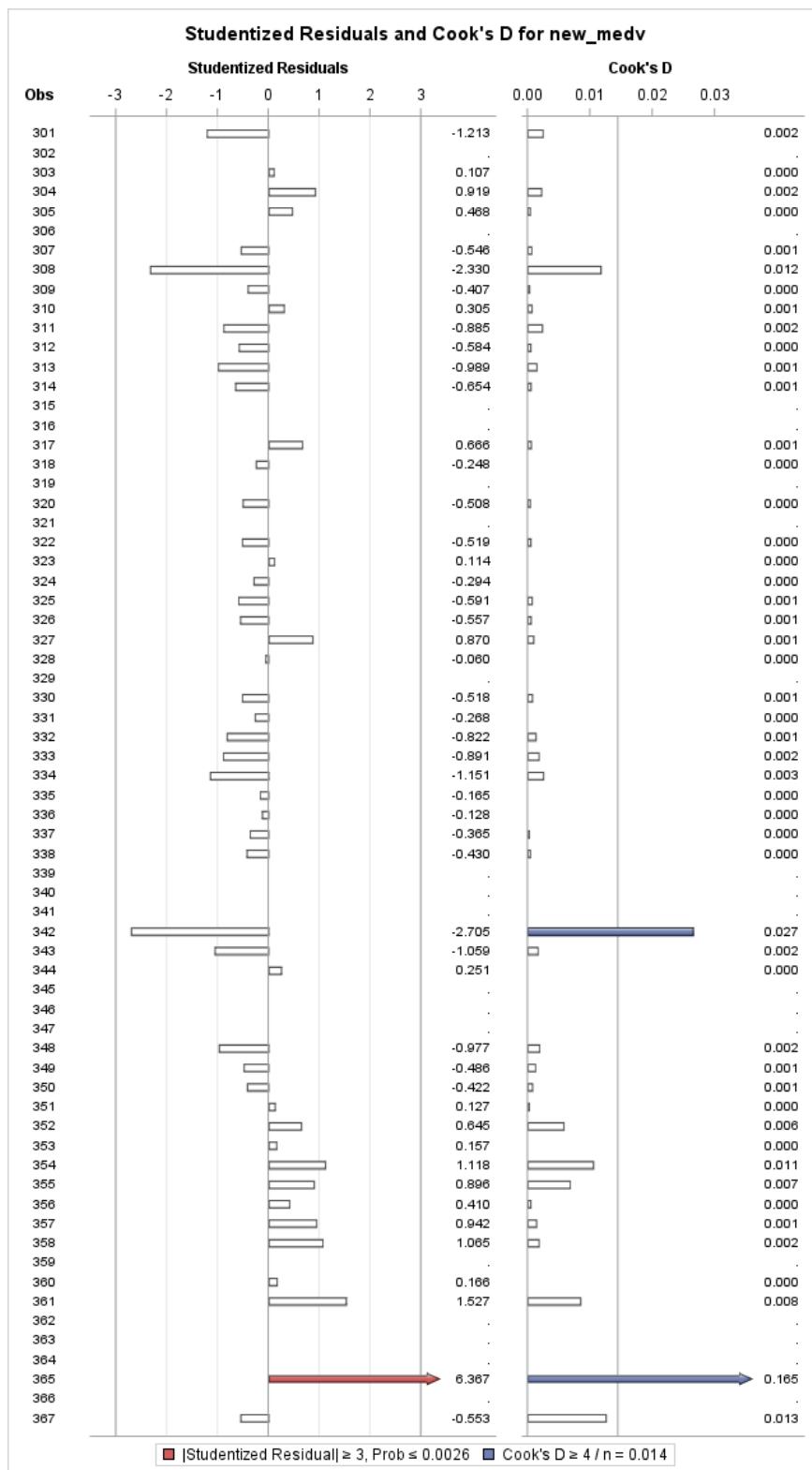


Fig. 2

Model 1- Checking Model assumptions and diagnostics

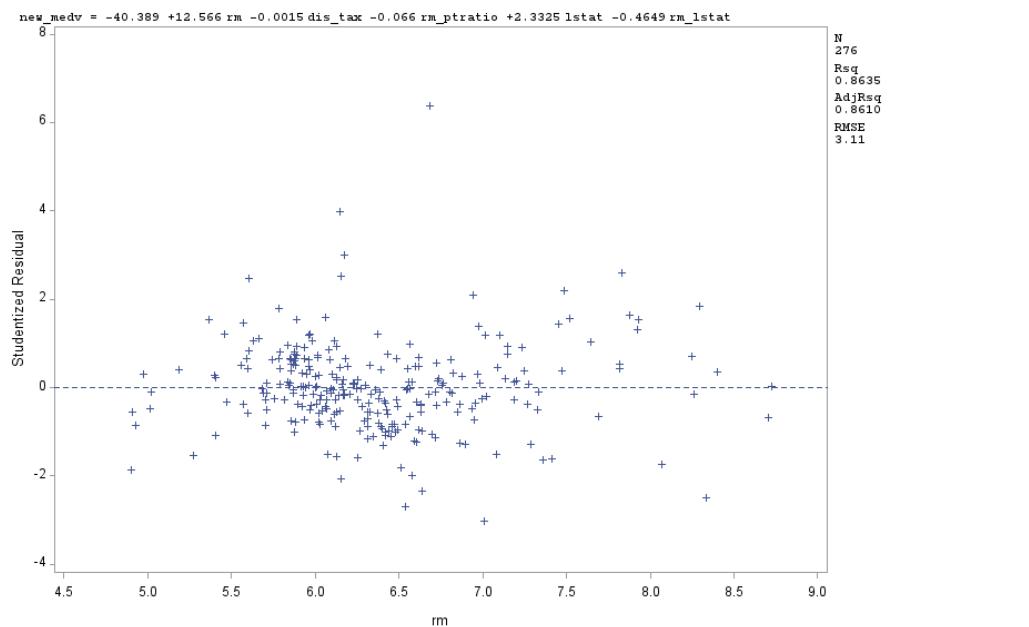


Fig. 3

Model 1- Checking Model assumptions and diagnostics

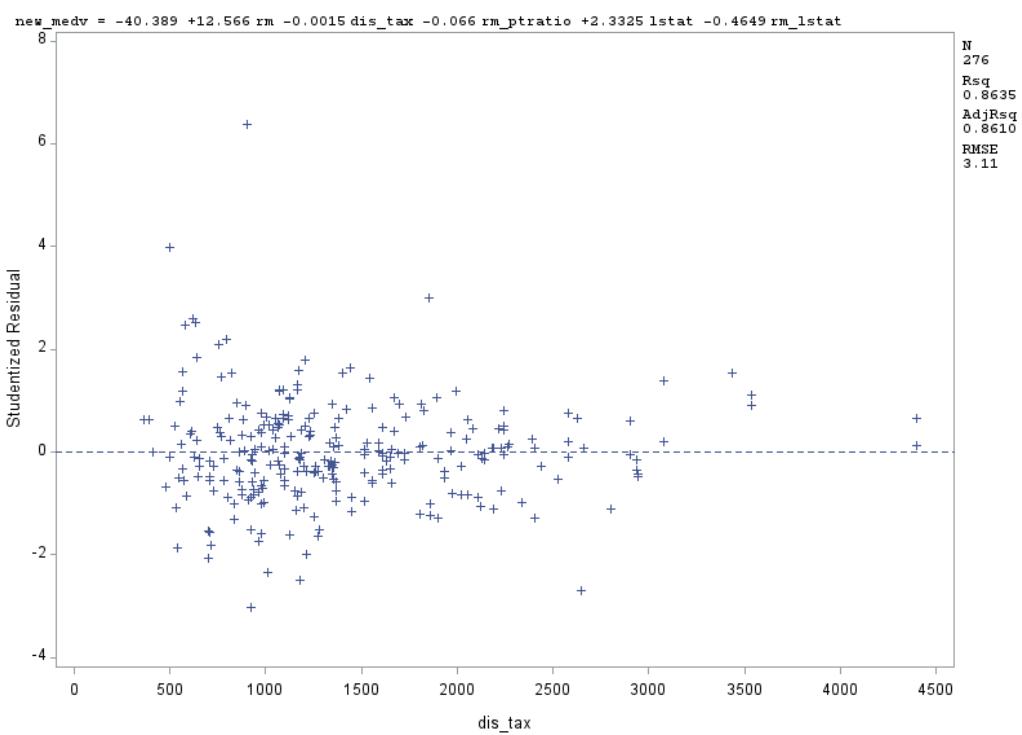


Fig. 4

Model 1- Checking Model assumptions and diagnostics

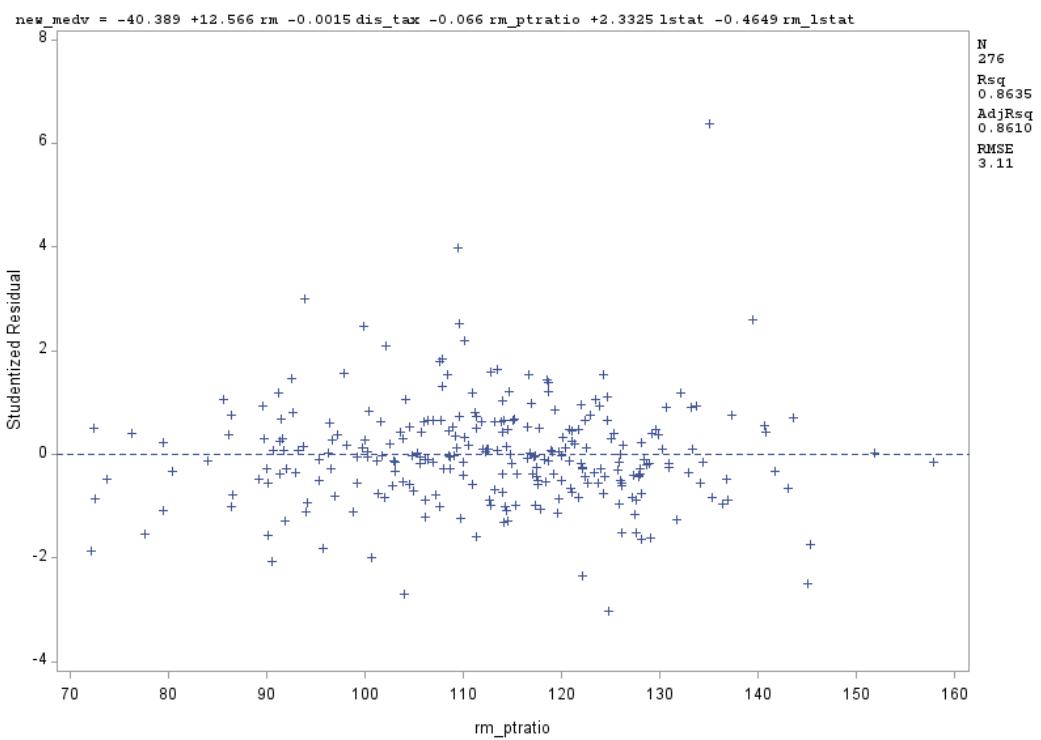


Fig. 5

Model 1- Checking Model assumptions and diagnostics

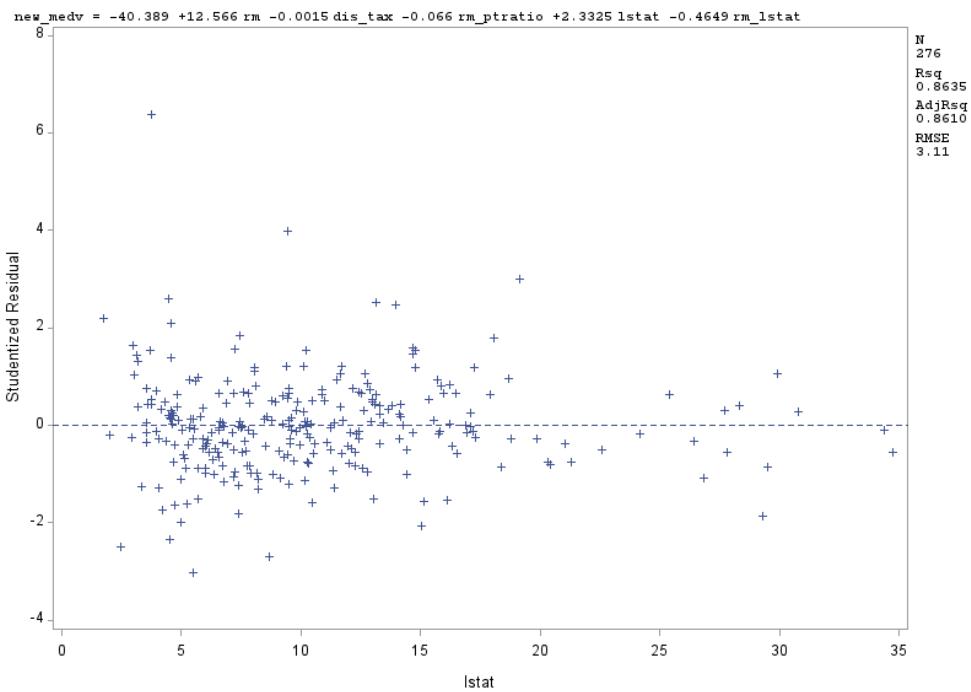


Fig. 6

Model 1- Checking Model assumptions and diagnostics

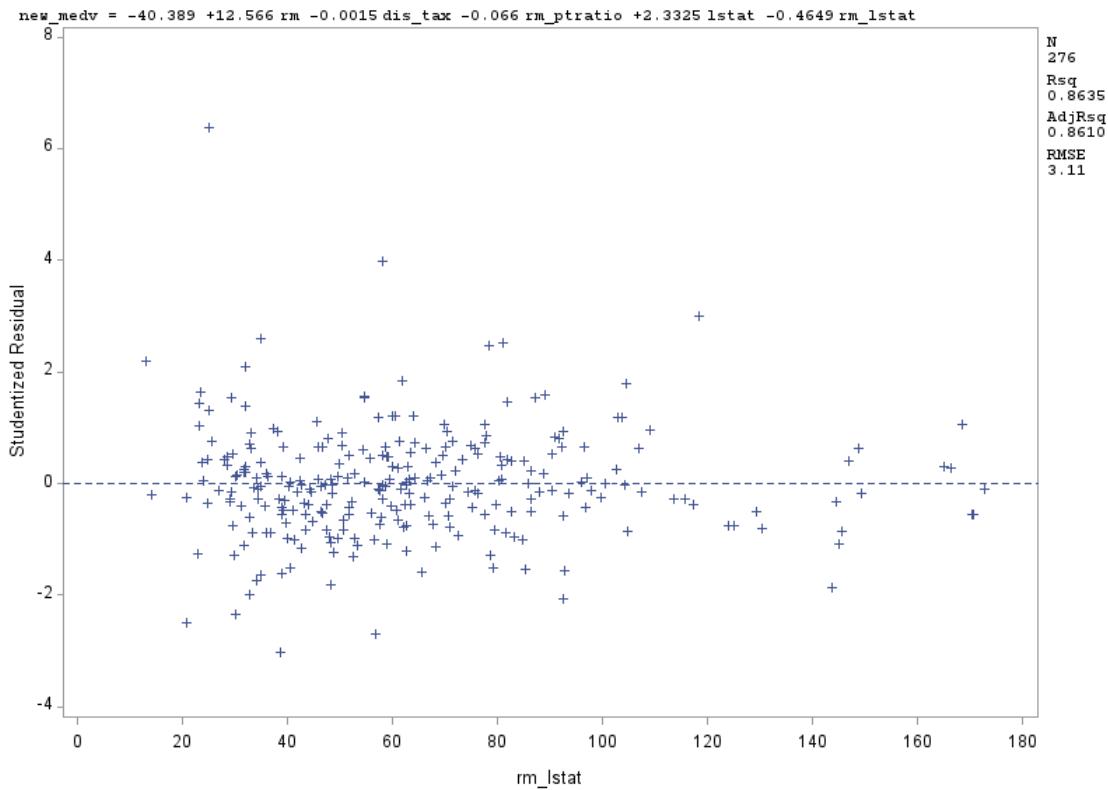


Fig. 7

Model 1- Checking Model assumptions and diagnostics

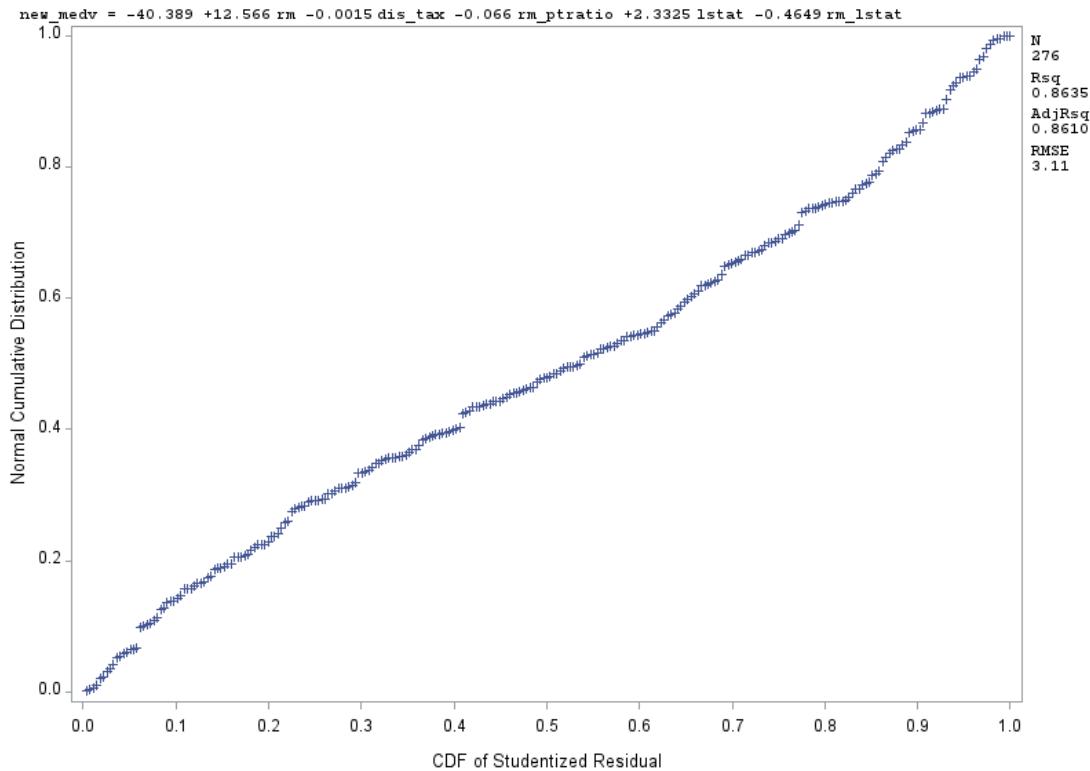


Table 3

The REG Procedure

Model: MODEL1

Dependent Variable: new_medv

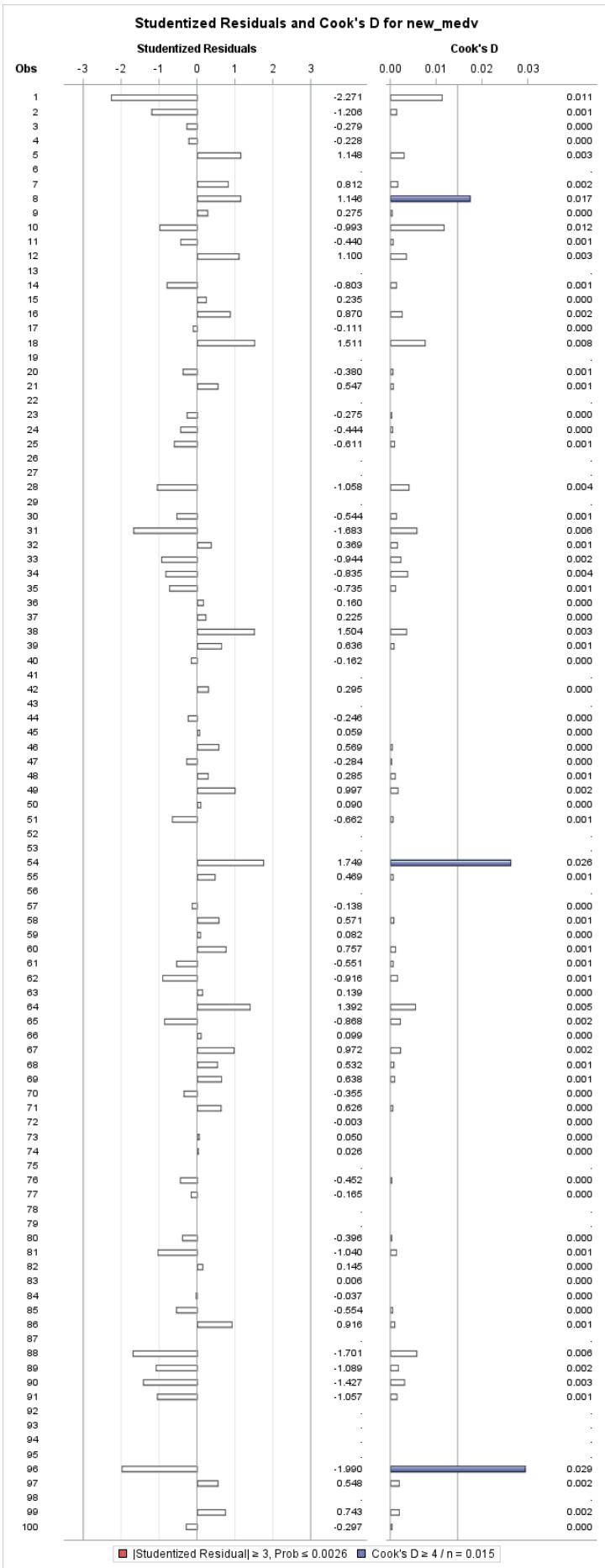
Number of Observations Read	363
Number of Observations Used	272
Number of Observations with Missing Values	91

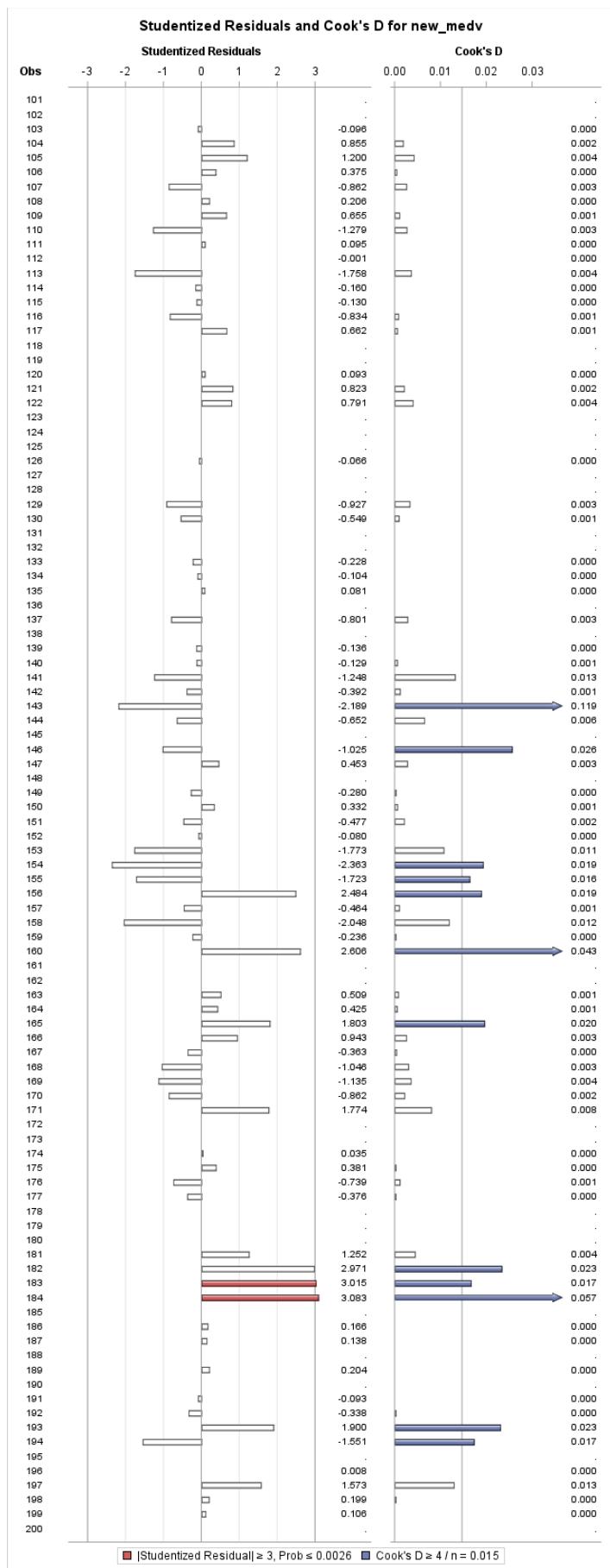
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	5	16454	3290.74414	461.05	<.0001
Error	266	1898.57631	7.13750		
Corrected Total	271	18352			

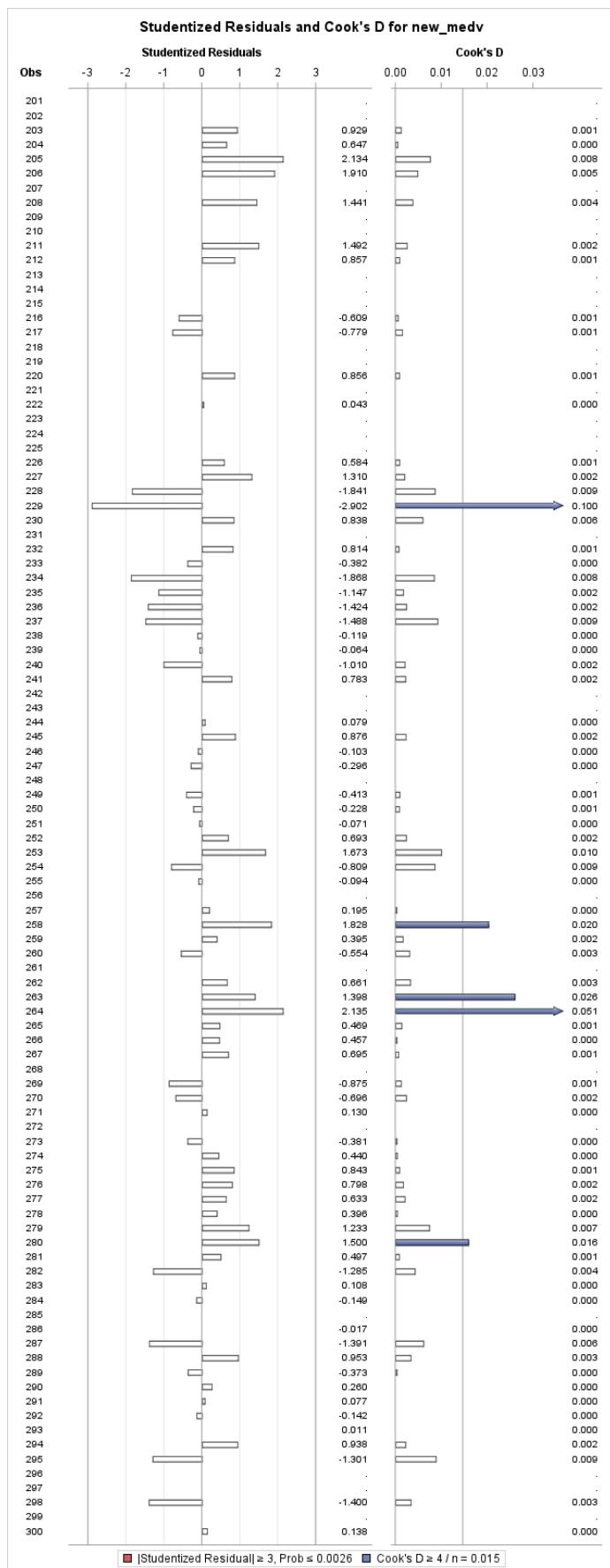
Root MSE	2.67161	R-Square	0.8965
Dependent Mean	24.64669	Adj R-Sq	0.8946
Coeff Var	10.83963		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	-41.58632	2.94417	-14.12	<.0001	0	0
rm	1	12.74890	0.41748	30.54	<.0001	1.03728	2.96662
dis_tax	1	-0.00137	0.00025347	-5.41	<.0001	-0.11484	1.16055
rm_ptratio	1	-0.06953	0.01211	-5.74	<.0001	-0.12894	1.29691
lstat	1	2.37504	0.24389	9.74	<.0001	1.72954	81.10798
rm_lstat	1	-0.46988	0.04254	-11.05	<.0001	-1.81295	69.26354

Fig. 8







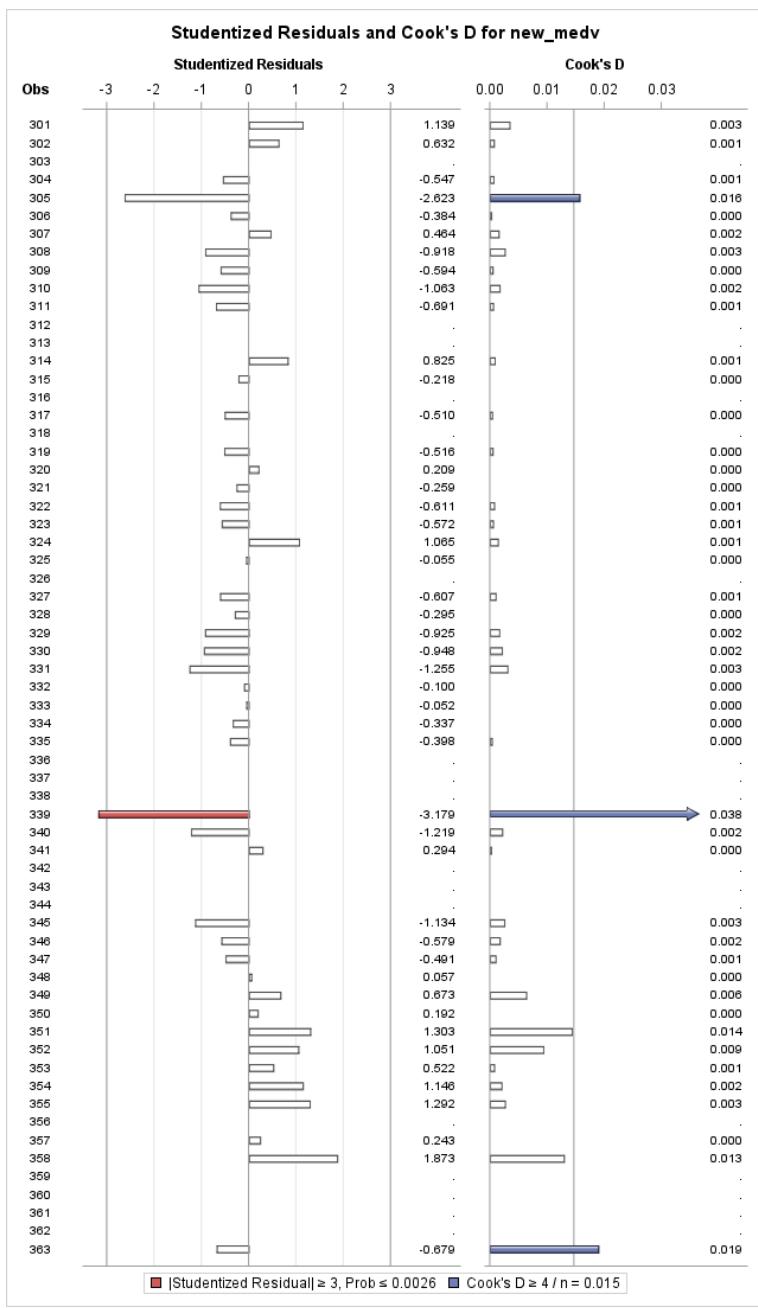


Fig. 9

Model 1 after removing outliers

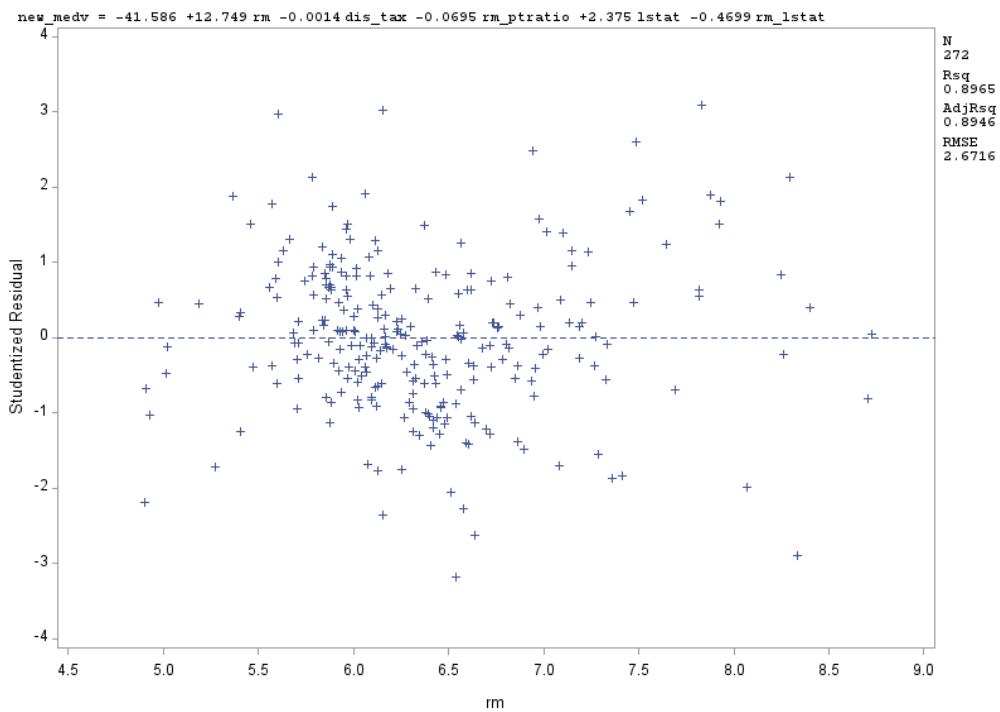


Fig. 10

Model 1 after removing outliers

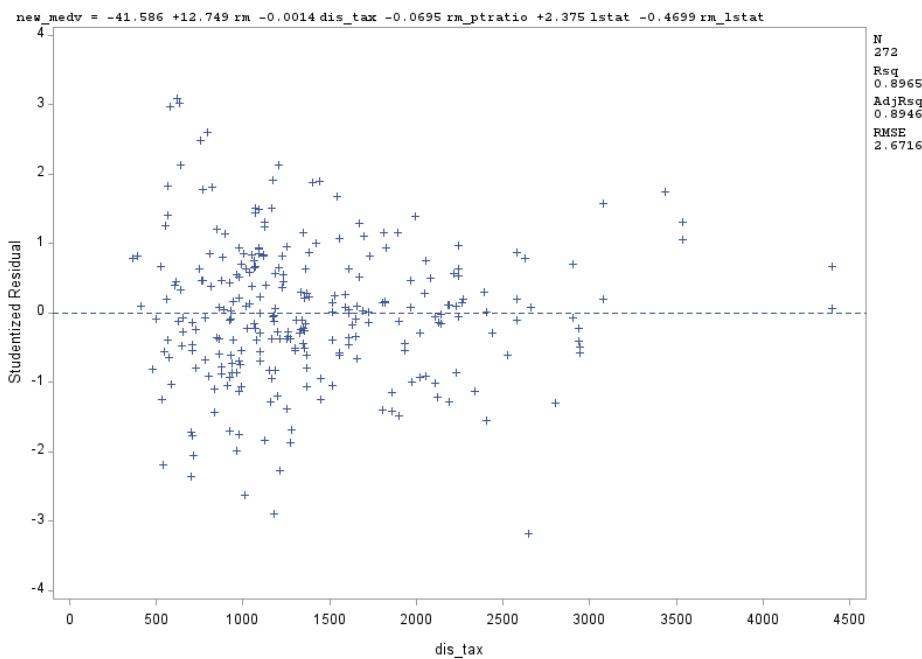


Fig. 11

Model 1 after removing outliers

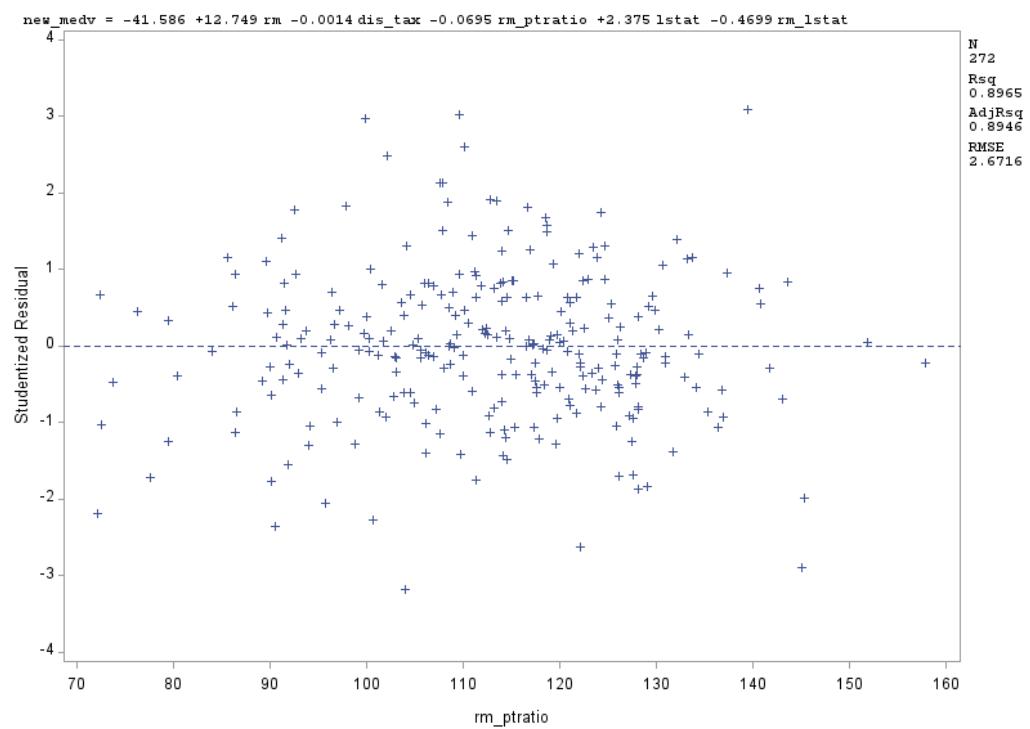


Fig. 12

Model 1 after removing outliers

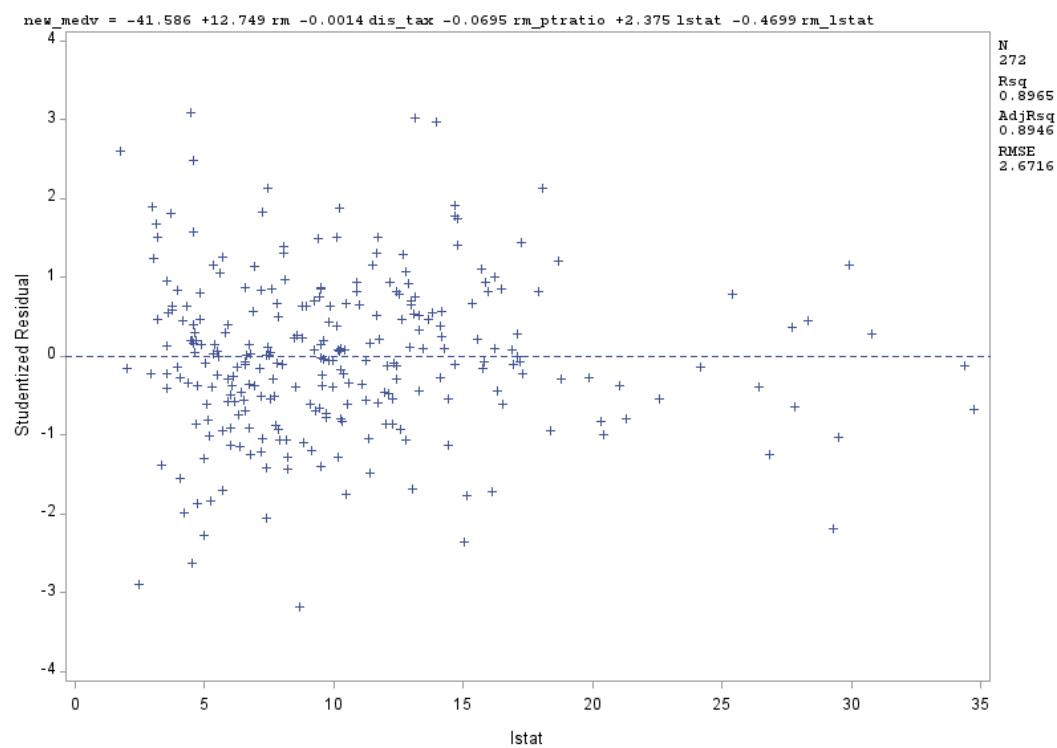


Fig. 13

Model 1 after removing outliers

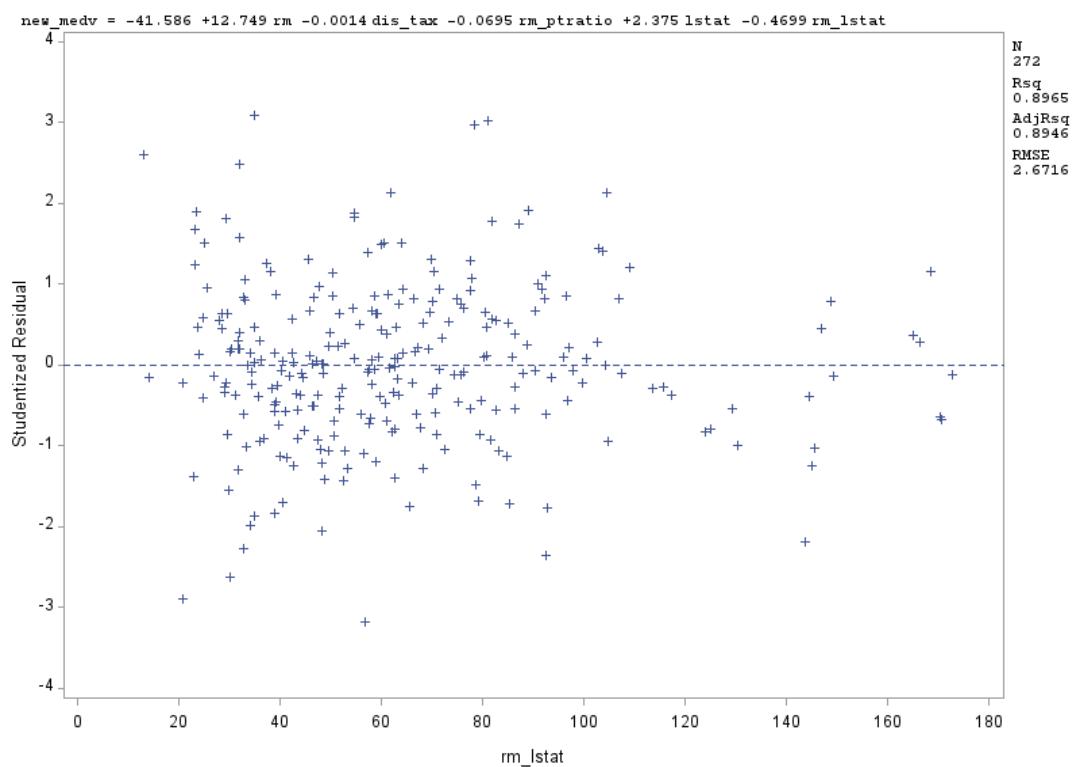


Fig. 14

Model 1 after removing outliers

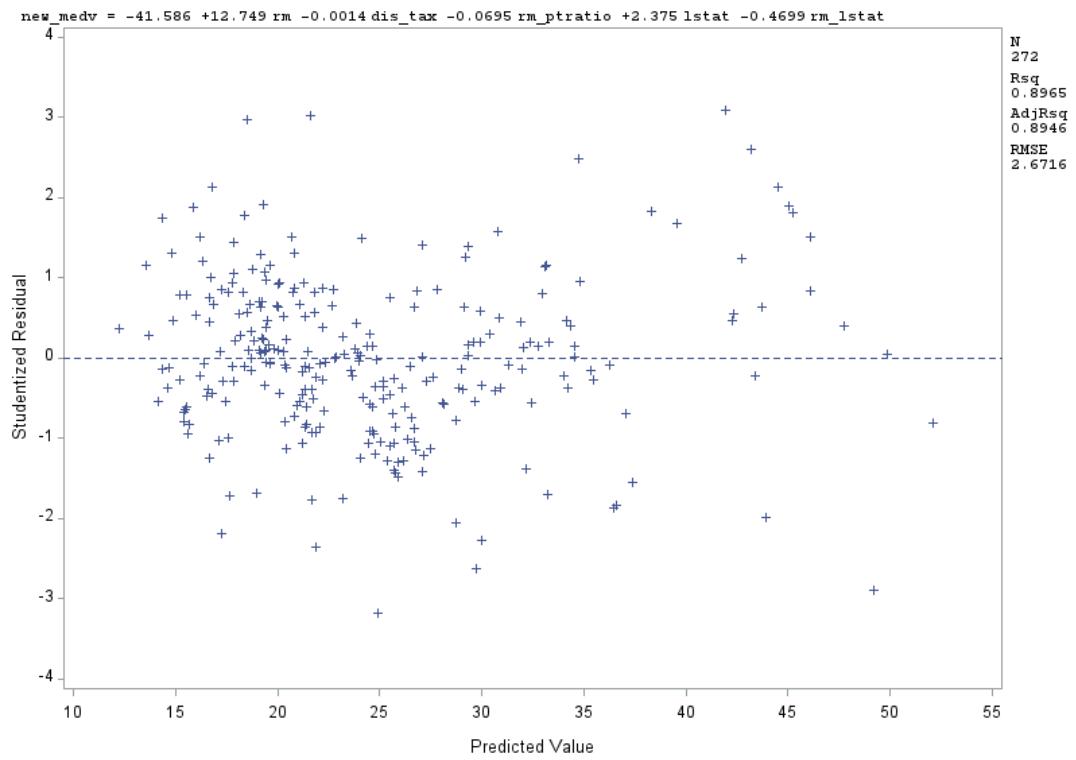


Fig. 15

Model 1 after removing outliers

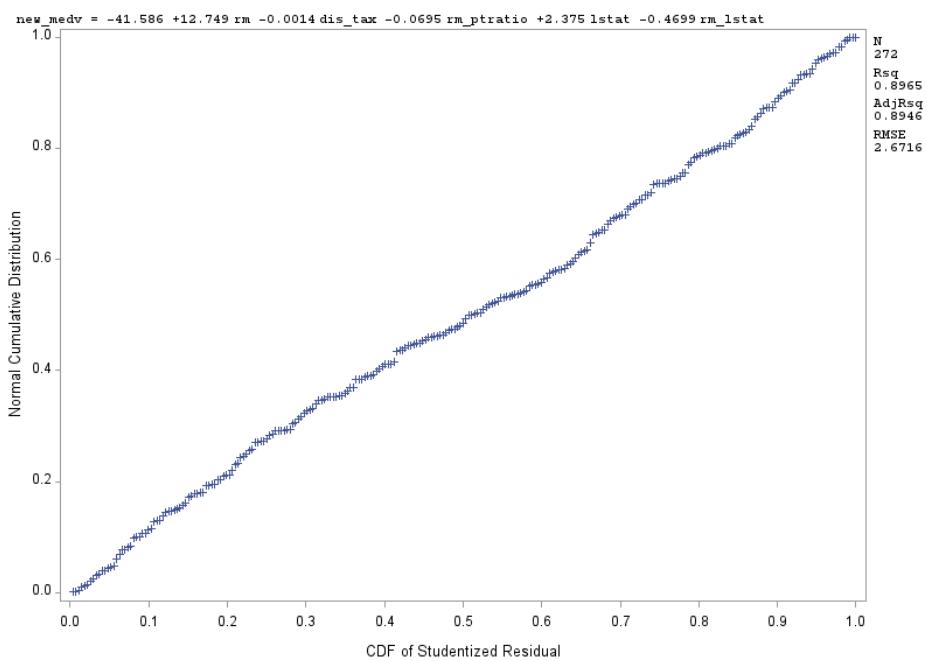


Table 4

Performance test parameters - Model 1

Obs	_TYPE_	_FREQ_	rmse	mae
1	0	91	3.46342	2.36930

Performance test parameters - Model 1

The CORR Procedure

2 Variables: medv yhat

Simple Statistics							
Variable	N	Mean	Std Dev	Sum	Minimum	Maximum	Label
medv	91	25.85055	9.13302	2352	14.30000	50.00000	
yhat	91	25.34226	8.81136	2306	8.32210	50.25216	Predicted Value of new_medv

Pearson Correlation Coefficients, N = 91		
Prob > r under H0: Rho=0		
	medv	yhat
medv	1.00000	0.92611 <.0001
yhat Predicted Value of new_medv	0.92611 <.0001	1.00000

Table 5

Model 2- Checking Model assumptions and diagnostics

The REG Procedure

Model: MODEL1

Dependent Variable: new_medv

Number of Observations Read	367
Number of Observations Used	276
Number of Observations with Missing Values	91

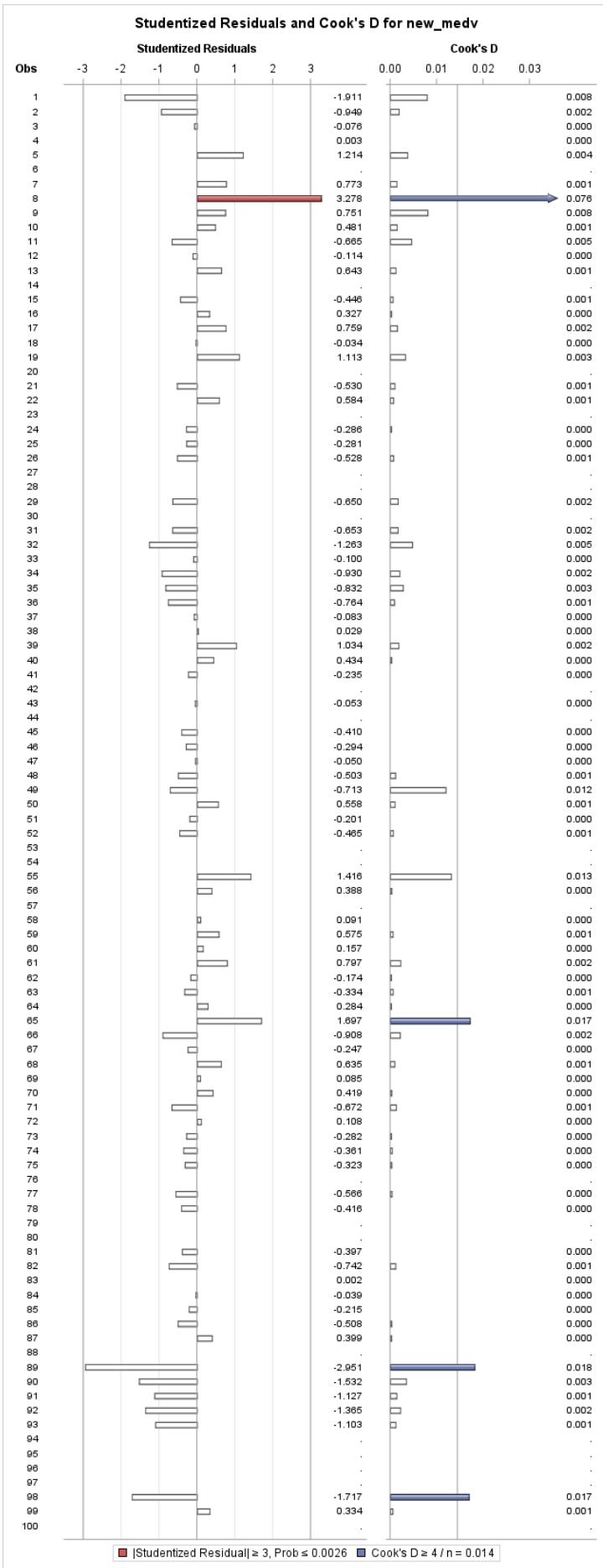
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	7	16688	2384.02414	261.62	<.0001
Error	268	2442.19304	9.11266		
Corrected Total	275	19130			

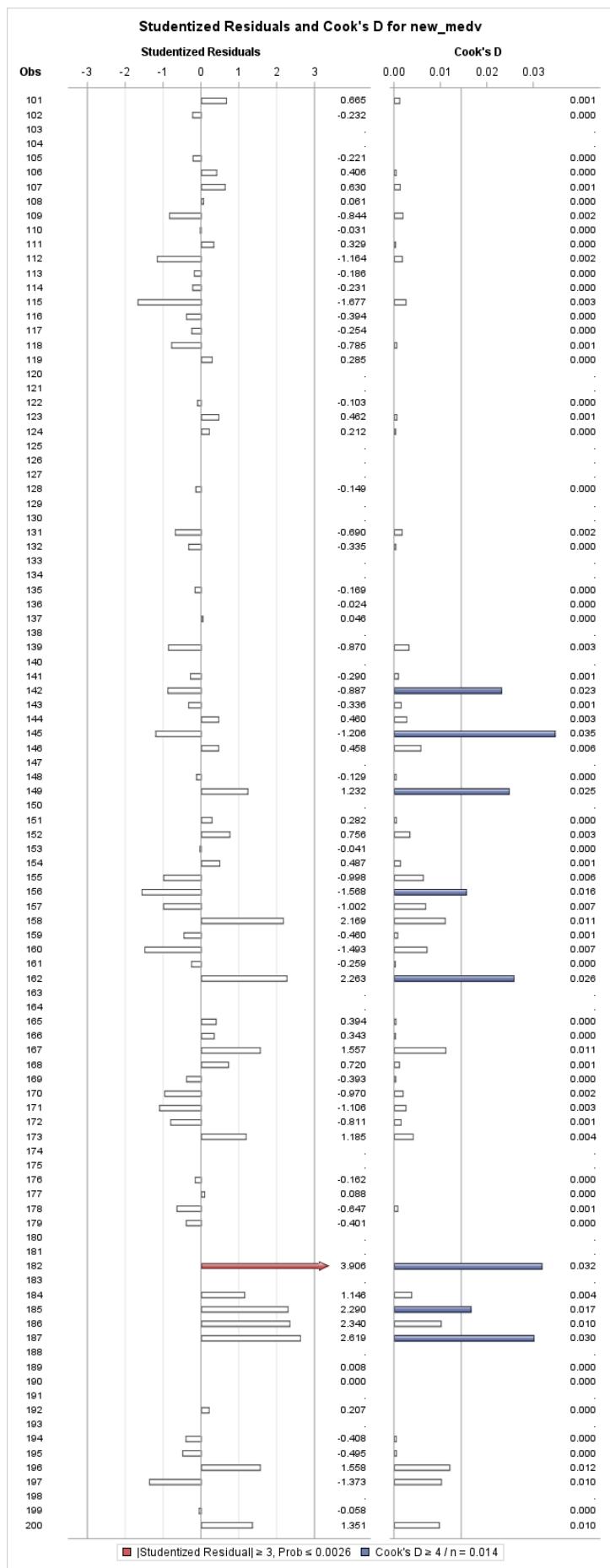
Root MSE	3.01872	R-Square	0.8723
Dependent Mean	24.78551	Adj R-Sq	0.8690
Coeff Var	12.17937		

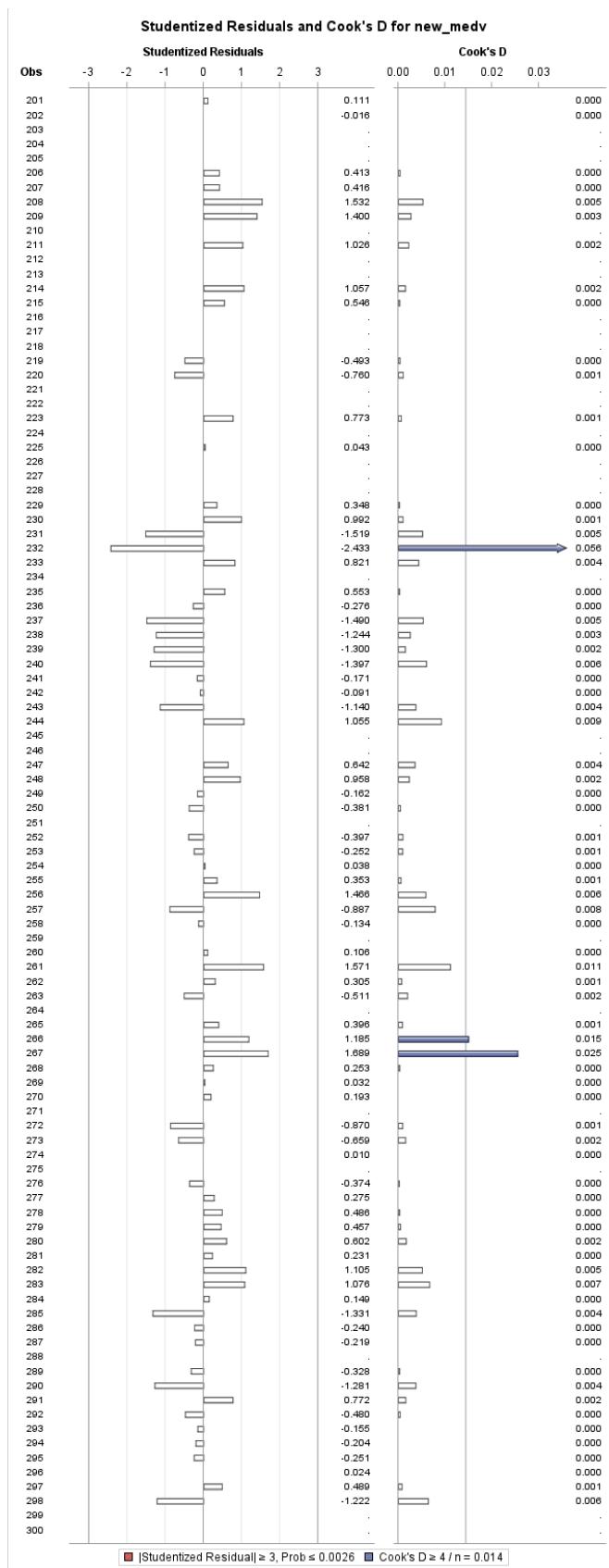
Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	-42.33980	3.35135	-12.63	<.0001	0	0
rm	1	13.17032	0.49332	26.70	<.0001	1.05212	3.26034

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
age_dis	1	-0.00482	0.00198	-2.43	0.0159	-0.06428	1.47185
dis_tax	1	-0.00165	0.00030821	-5.34	<.0001	-0.13590	1.35933
rm_ptratio	1	-0.08043	0.01408	-5.71	<.0001	-0.14726	1.39503
lstat	1	2.99927	0.32916	9.11	<.0001	2.15526	117.45268
nox_lstat	1	-0.59786	0.14380	-4.16	<.0001	-0.31722	12.22035
rm_lstat	1	-0.50394	0.05085	-9.91	<.0001	-1.92226	78.97871

Fig. 16







Studentized Residuals and Cook's D for new_medv

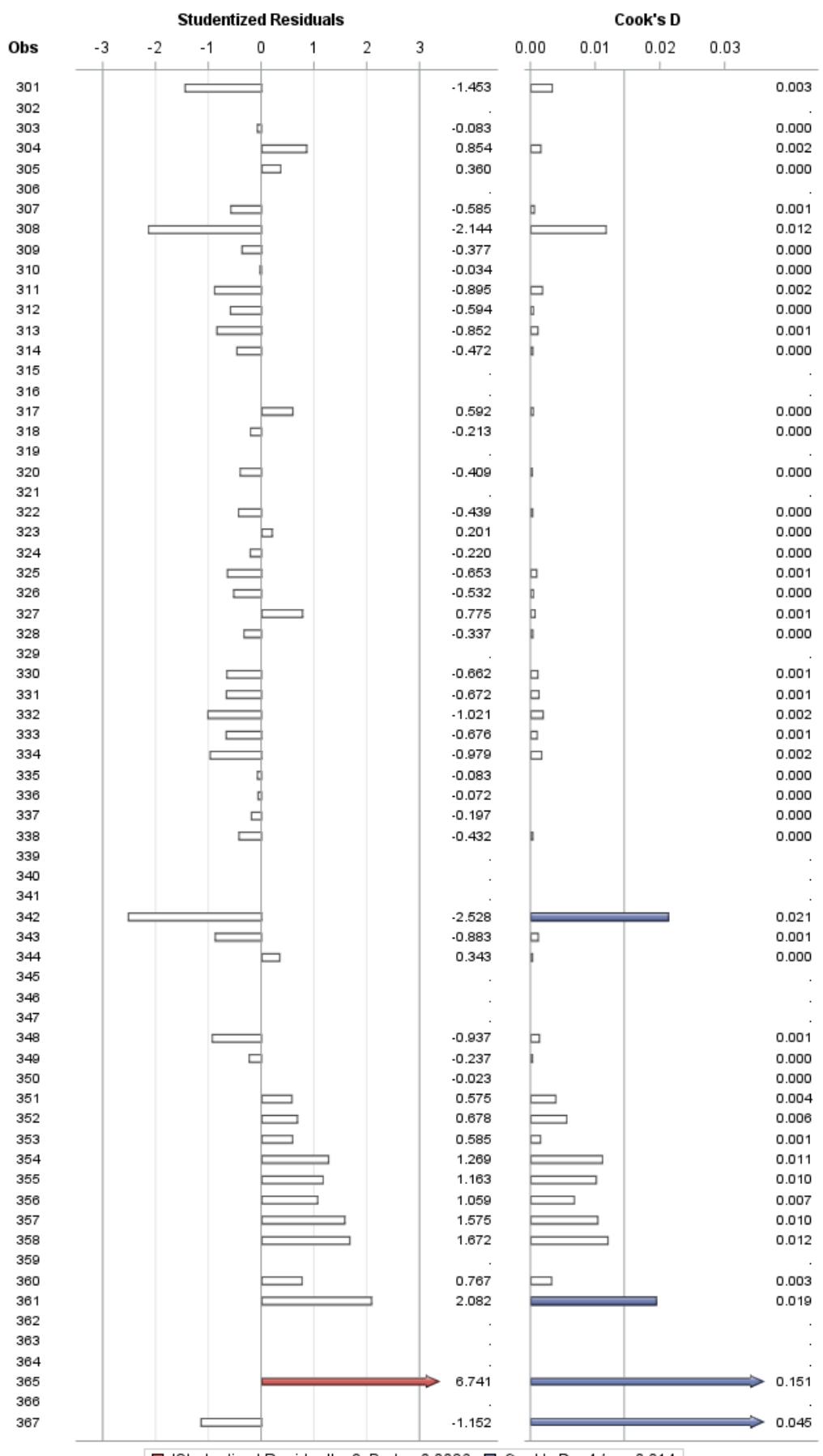


Fig. 17
Model 2- Checking Model assumptions and diagnostics

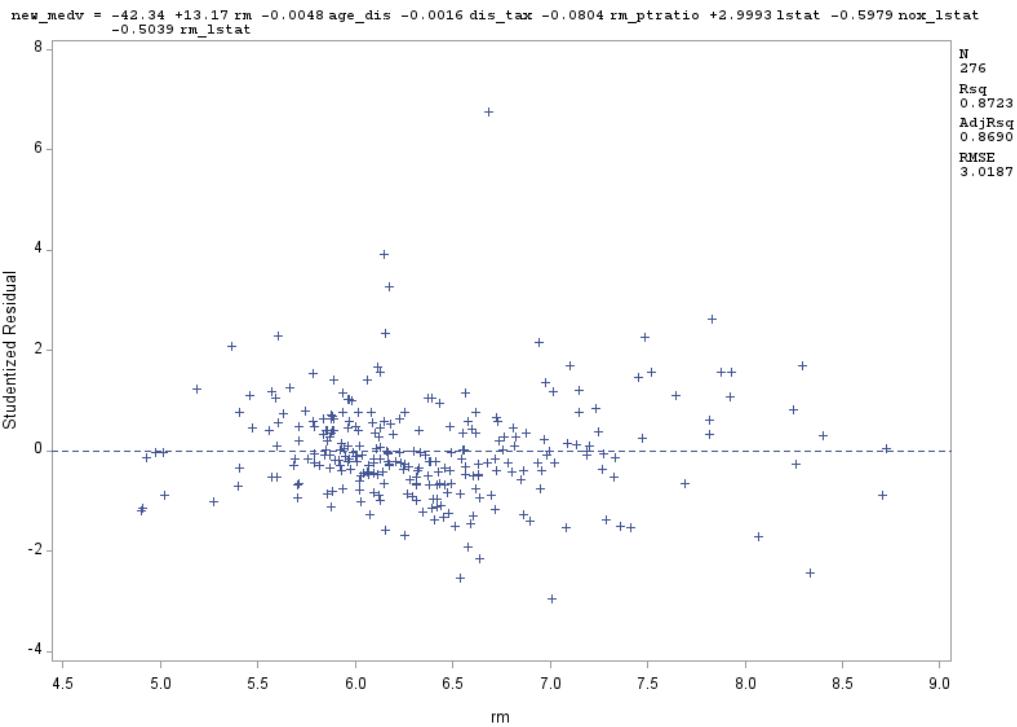


Fig. 18
Model 2- Checking Model assumptions and diagnostics

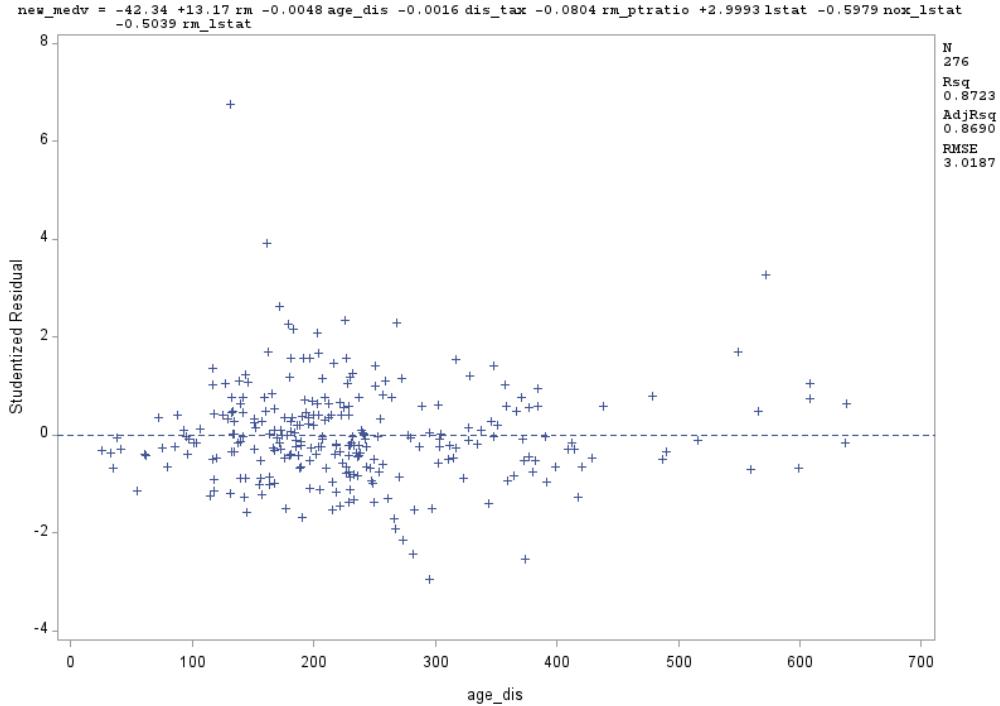


Fig. 19

Model 2- Checking Model assumptions and diagnostics

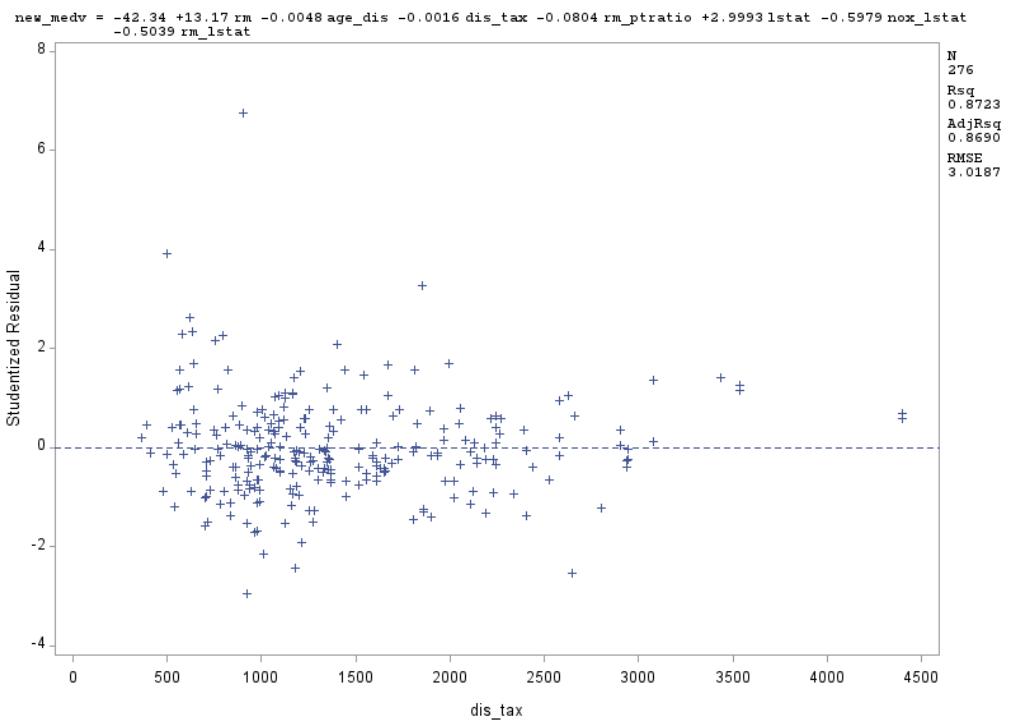


Fig. 20

Model 2- Checking Model assumptions and diagnostics

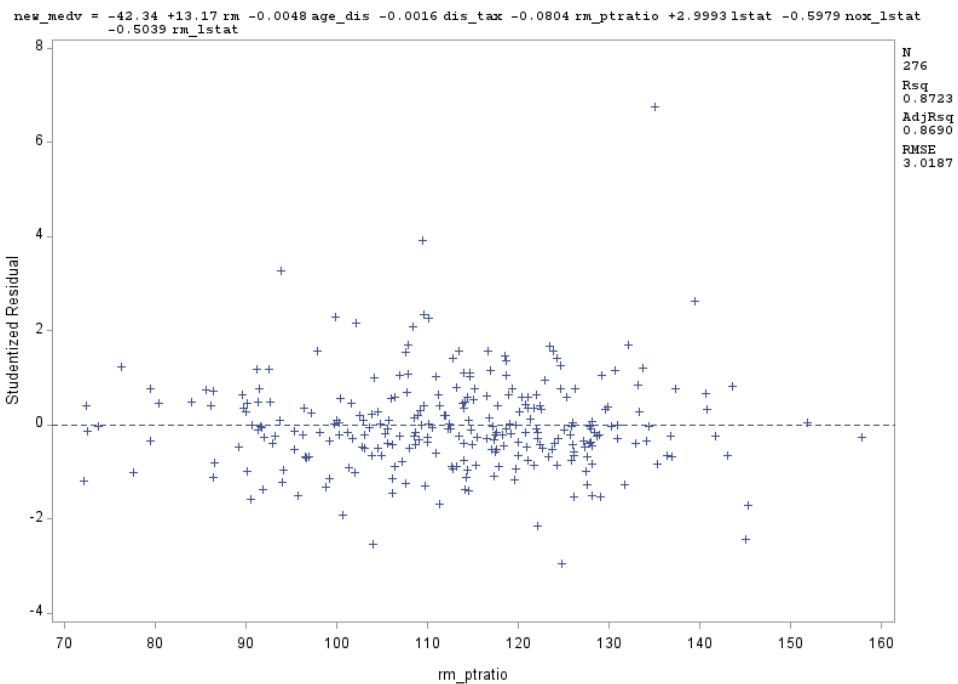


Fig. 21

Model 2- Checking Model assumptions and diagnostics

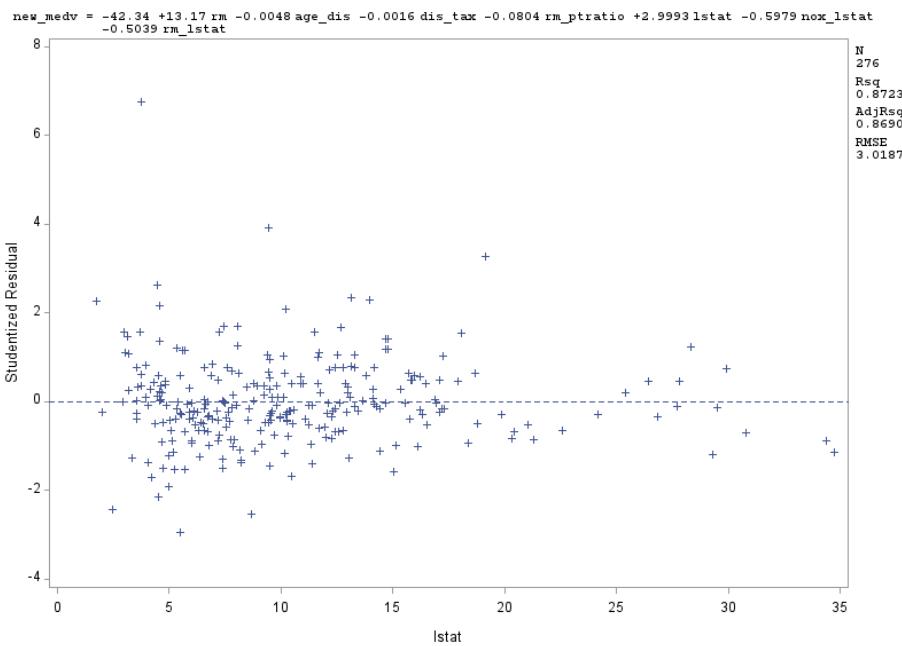


Fig. 22

Model 2- Checking Model assumptions and diagnostics

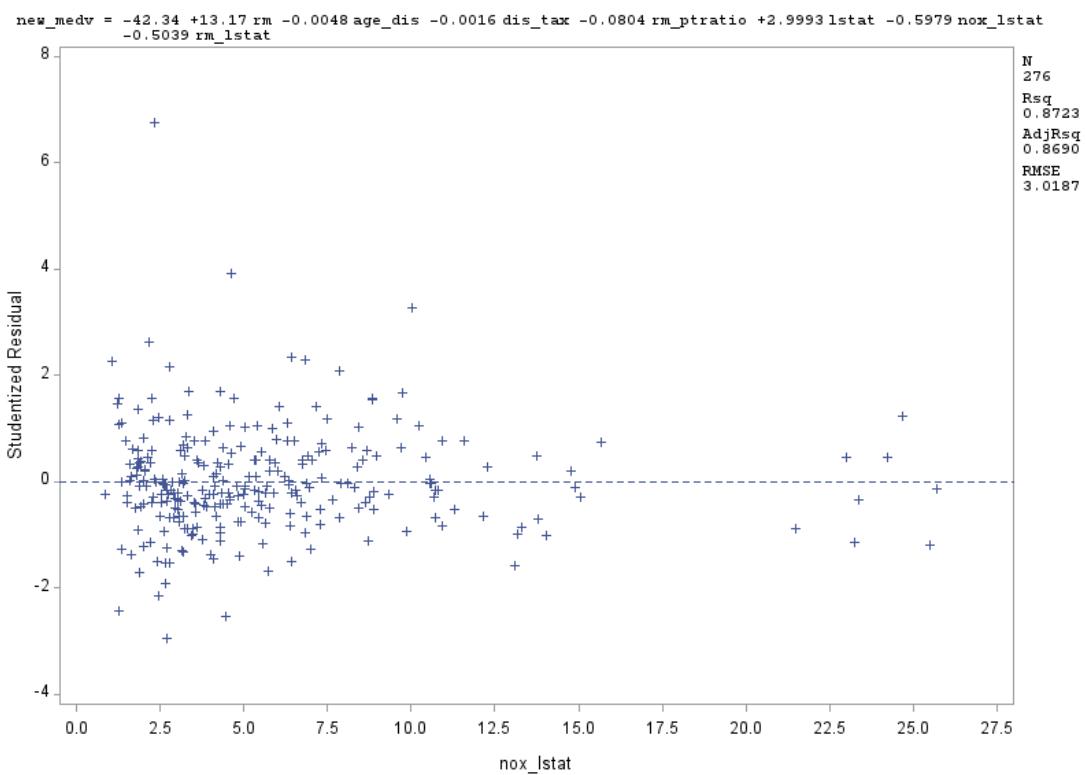


Fig. 23

Model 2- Checking Model assumptions and diagnostics

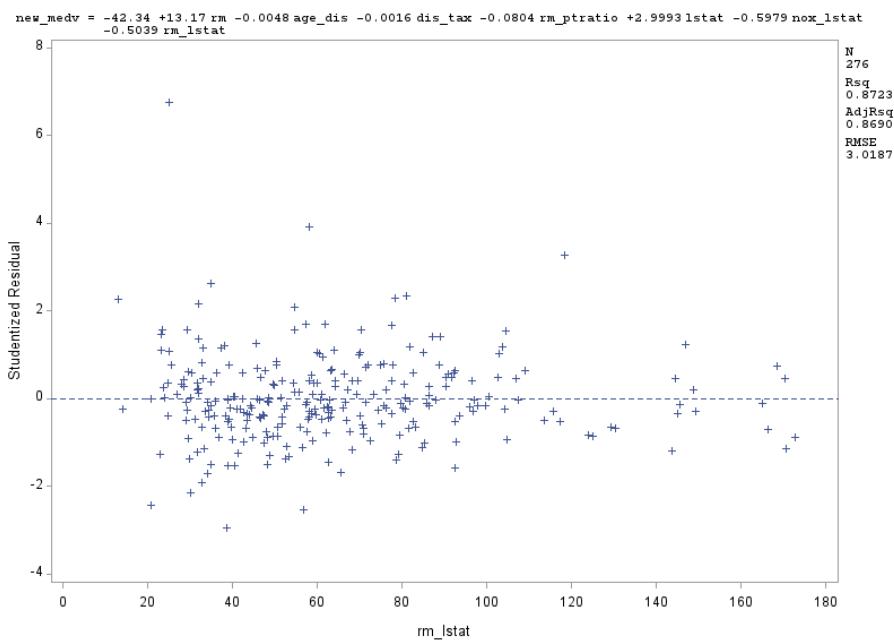


Fig. 24

Model 2- Checking Model assumptions and diagnostics

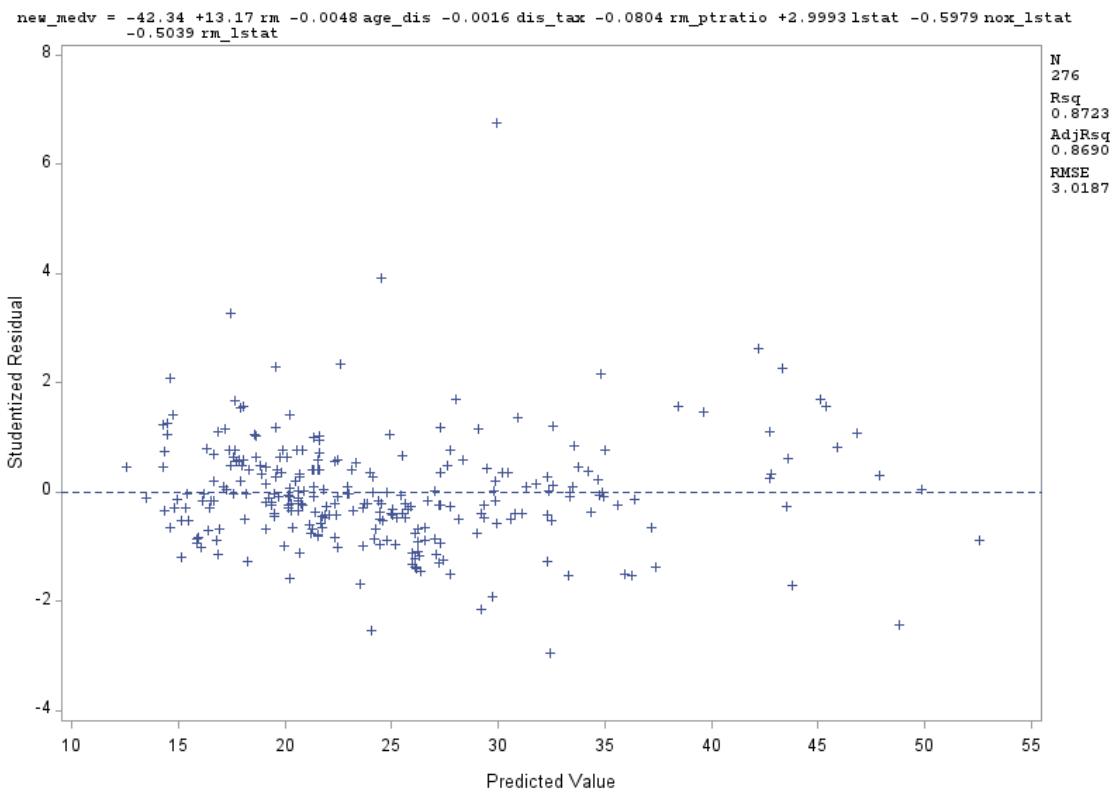


Fig. 25

Model 2- Checking Model assumptions and diagnostics

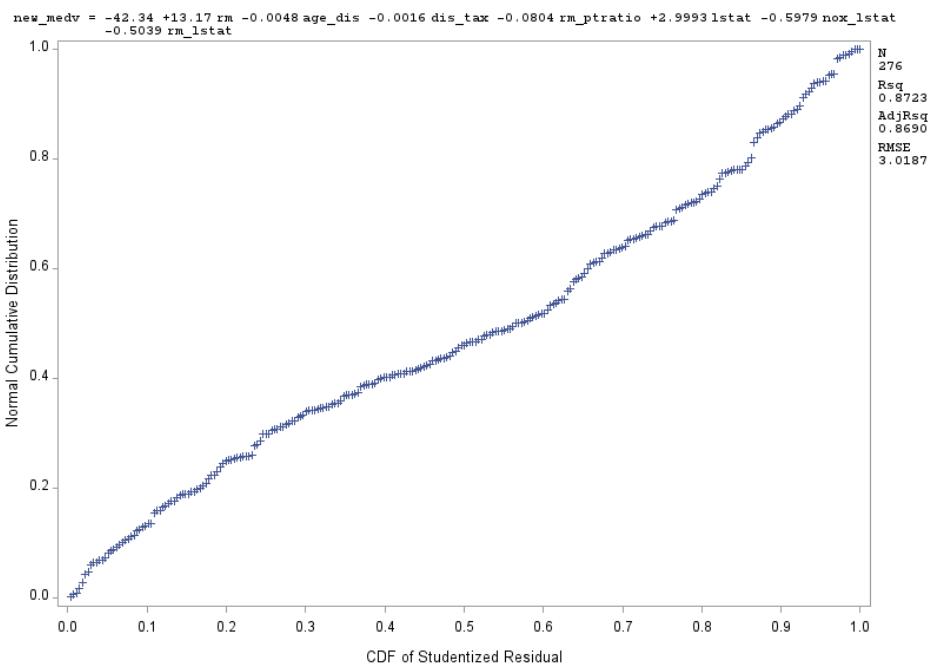


Table 5

Model 2 after removing outliers

The REG Procedure

Model: MODEL1

Dependent Variable: new_medv

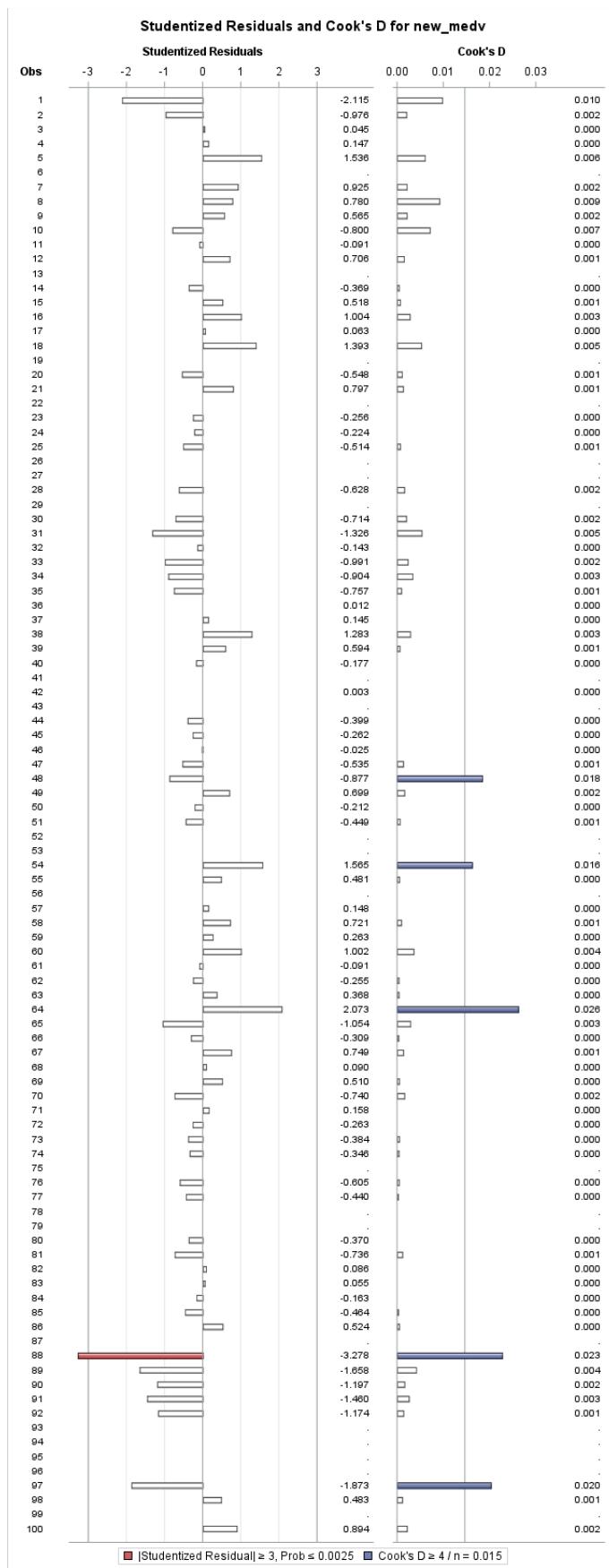
Number of Observations Read	364
Number of Observations Used	273
Number of Observations with Missing Values	91

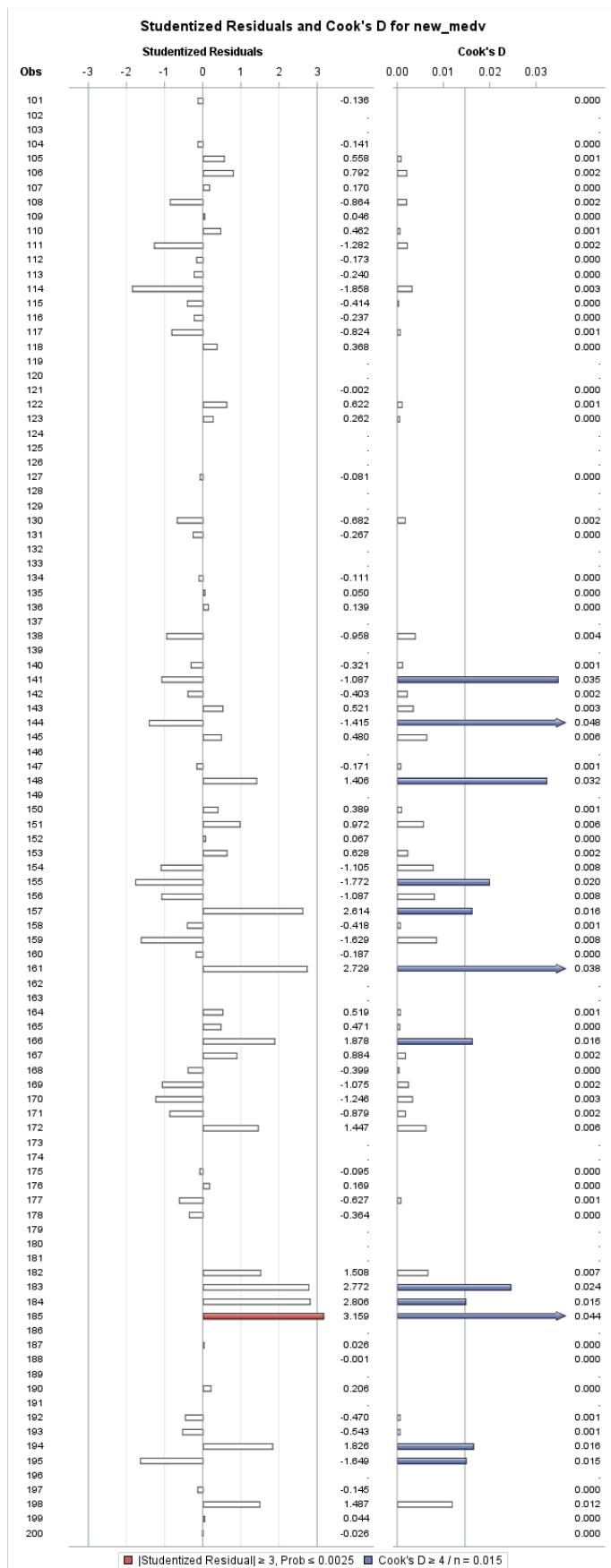
Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	7	16556	2365.16956	348.75	<.0001
Error	265	1797.20164	6.78189		
Corrected Total	272	18353			

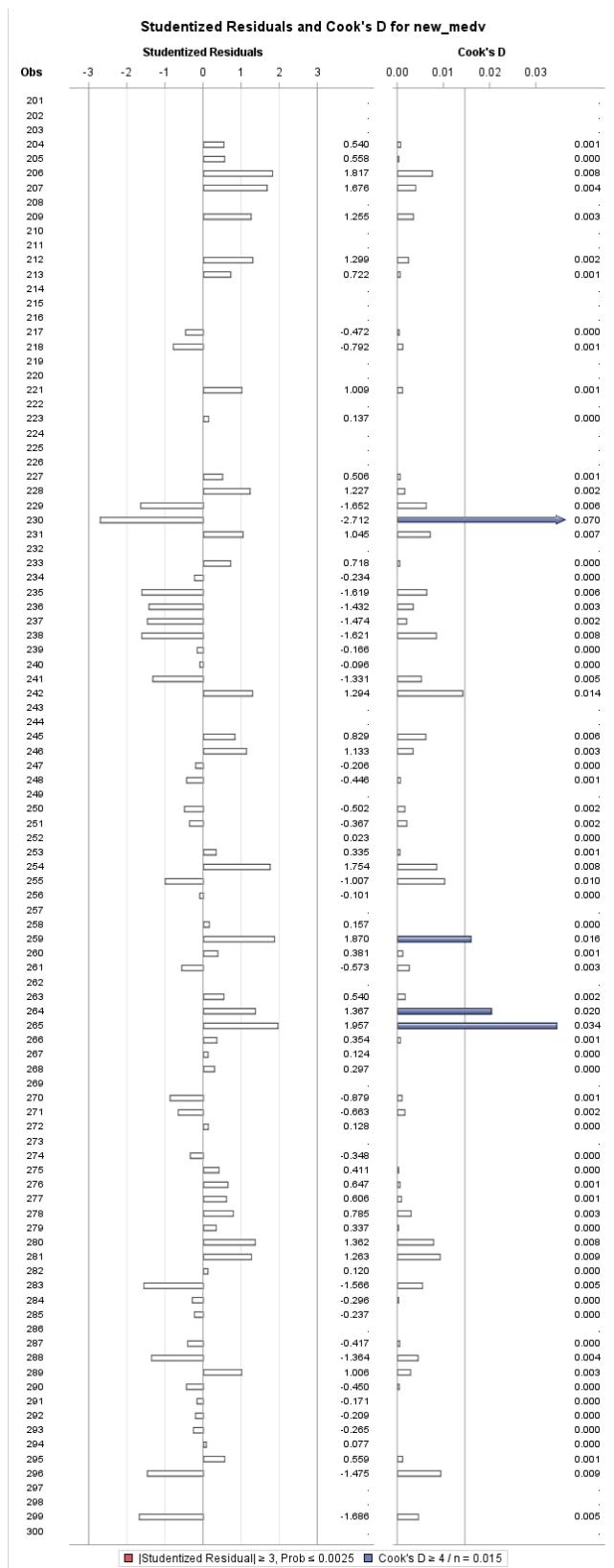
Root MSE	2.60421	R-Square	0.9021
Dependent Mean	24.64286	Adj R-Sq	0.8995
Coeff Var	10.56780		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	-43.37581	2.89730	-14.97	<.0001	0	0
rm	1	13.33592	0.42628	31.28	<.0001	1.08682	3.26602
age_dis	1	-0.00549	0.00173	-3.18	0.0016	-0.07338	1.43868
dis_tax	1	-0.00144	0.00026773	-5.38	<.0001	-0.12082	1.36506
rm_ptratio	1	-0.08550	0.01227	-6.97	<.0001	-0.15874	1.40408
lstat	1	3.03364	0.28433	10.67	<.0001	2.21184	116.30657
nox_lstat	1	-0.61369	0.12431	-4.94	<.0001	-0.33147	12.20018
rm_lstat	1	-0.50391	0.04400	-11.45	<.0001	-1.94646	78.17365

Fig. 26







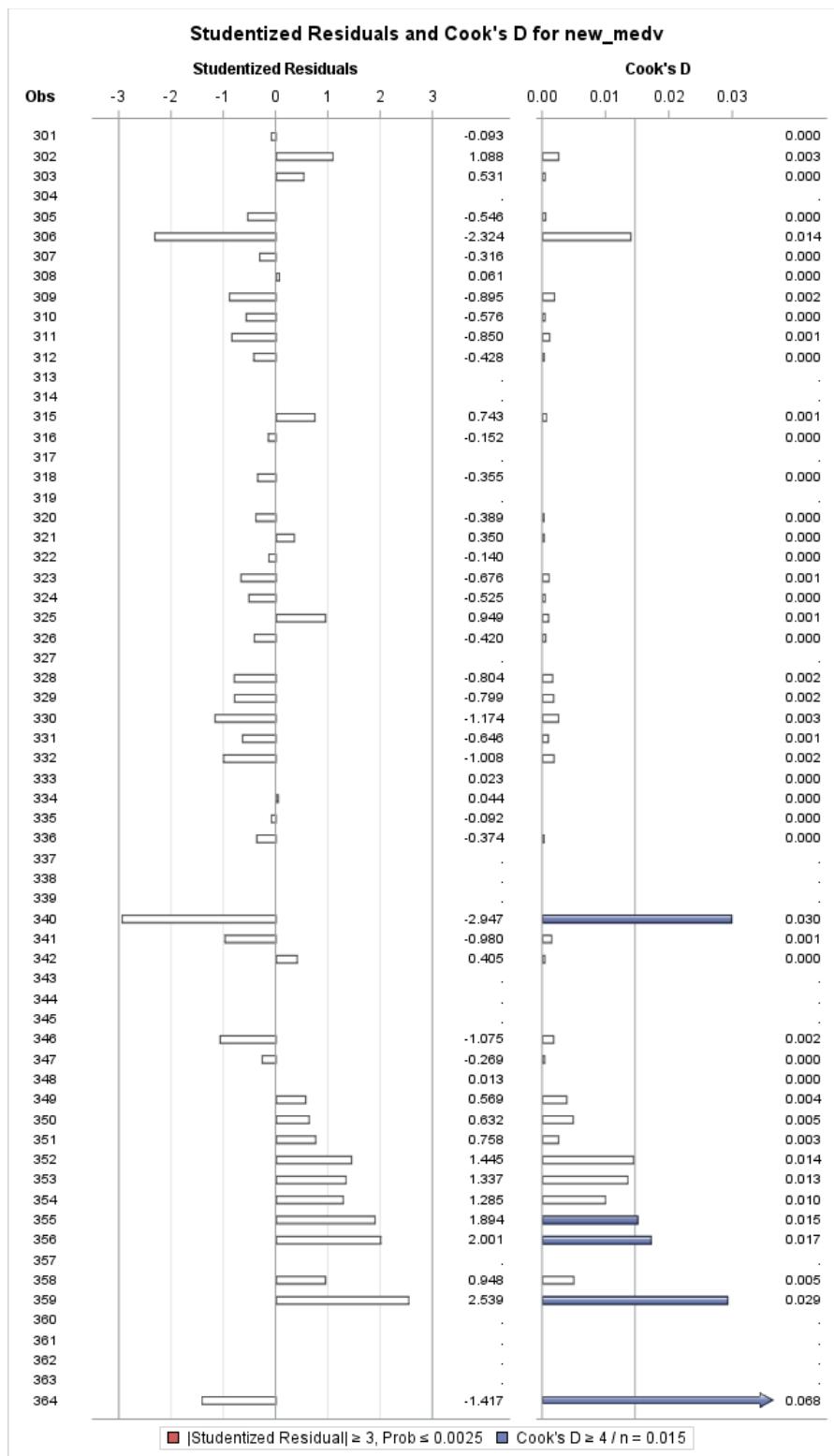


Fig. 27

Model 2 after removing outliers

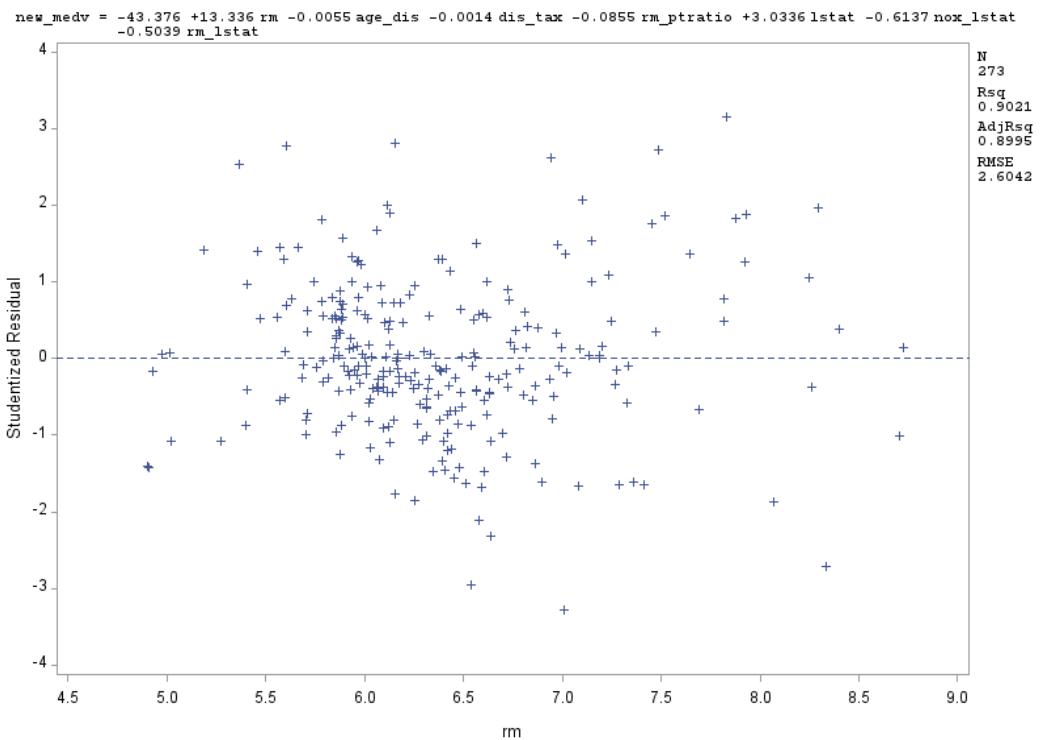


Fig. 28

Model 2 after removing outliers

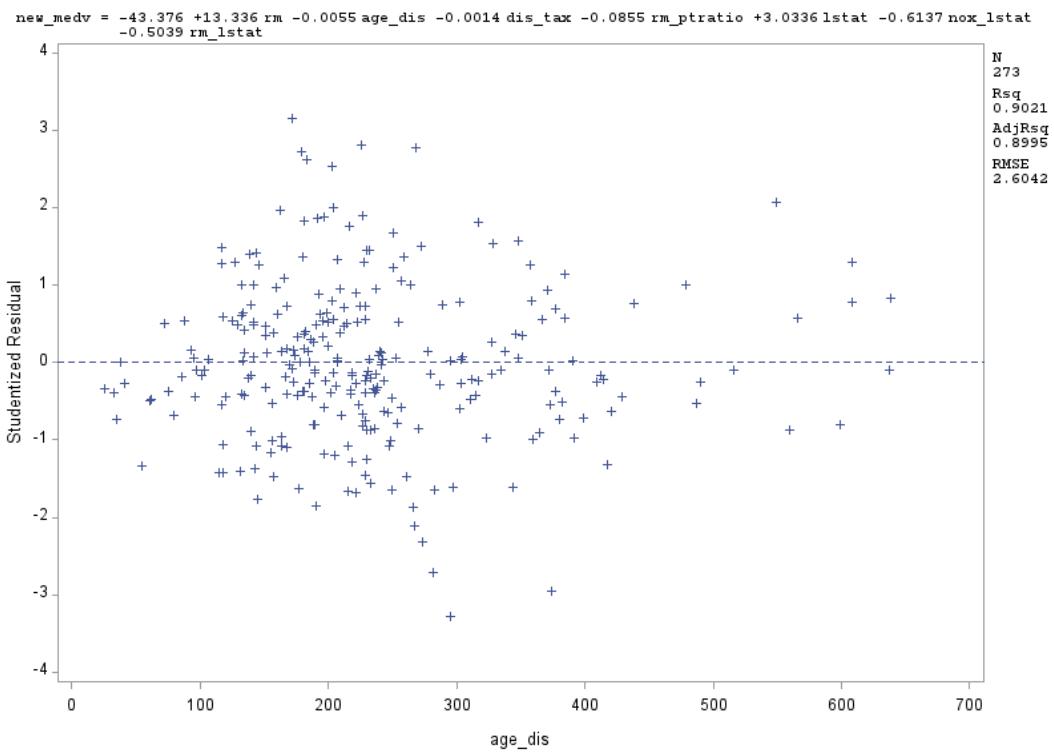


Fig. 29

Model 2 after removing outliers

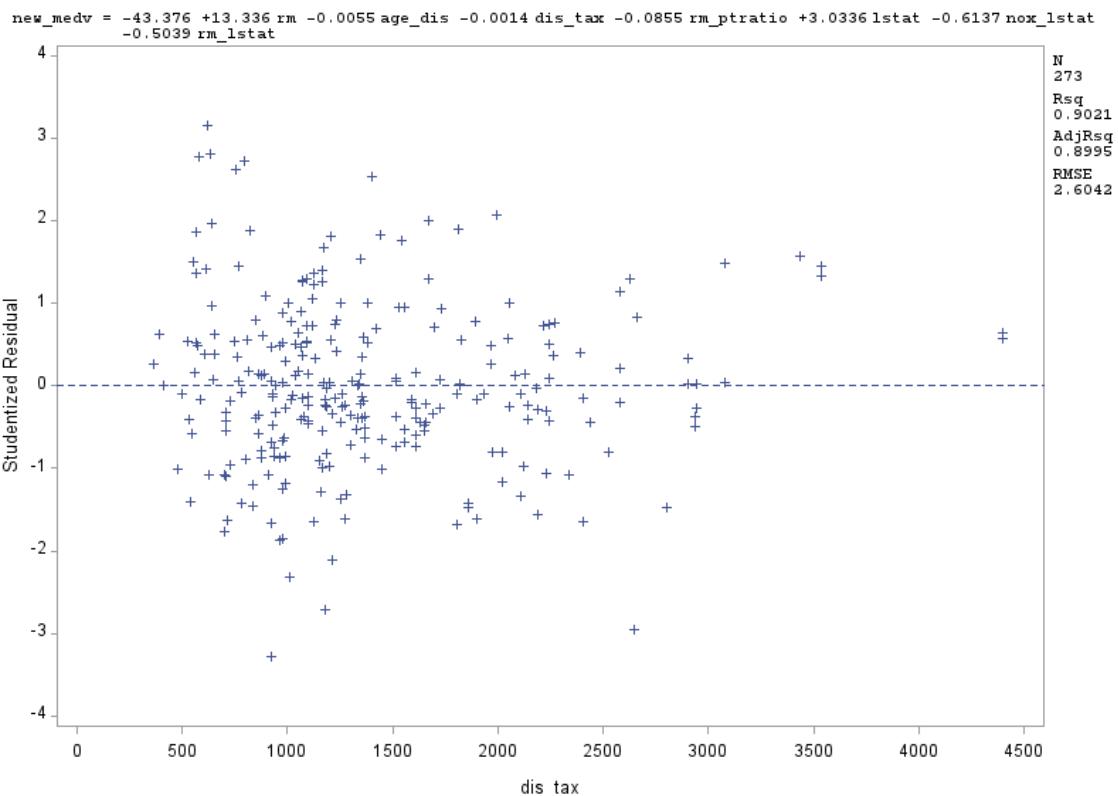


Fig. 30
Model 2 after removing outliers

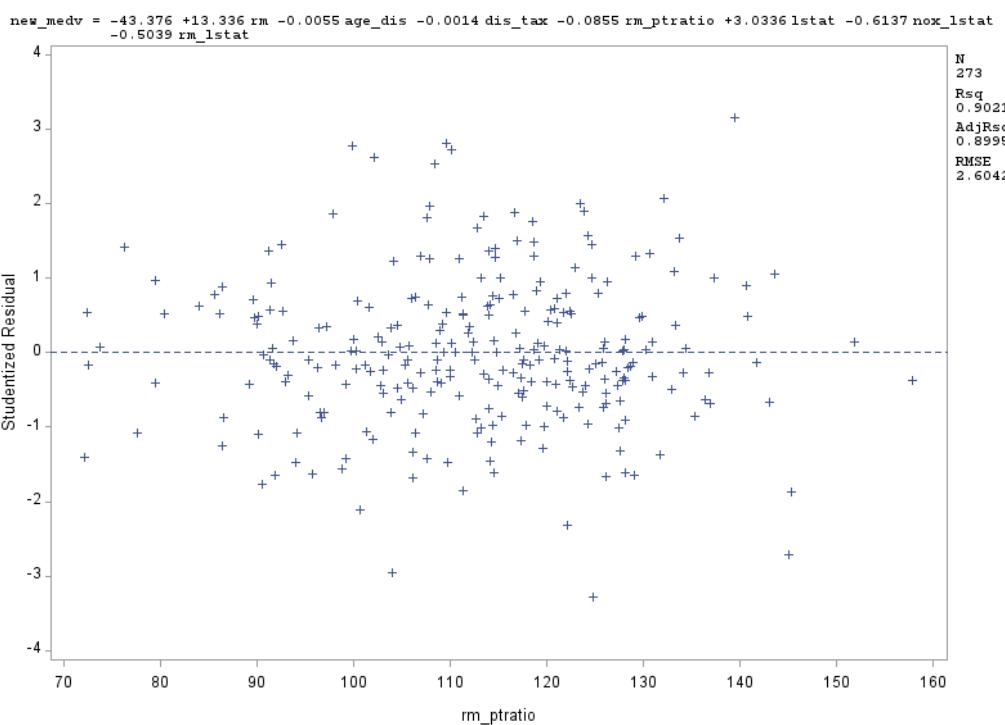


Fig. 31

Model 2 after removing outliers

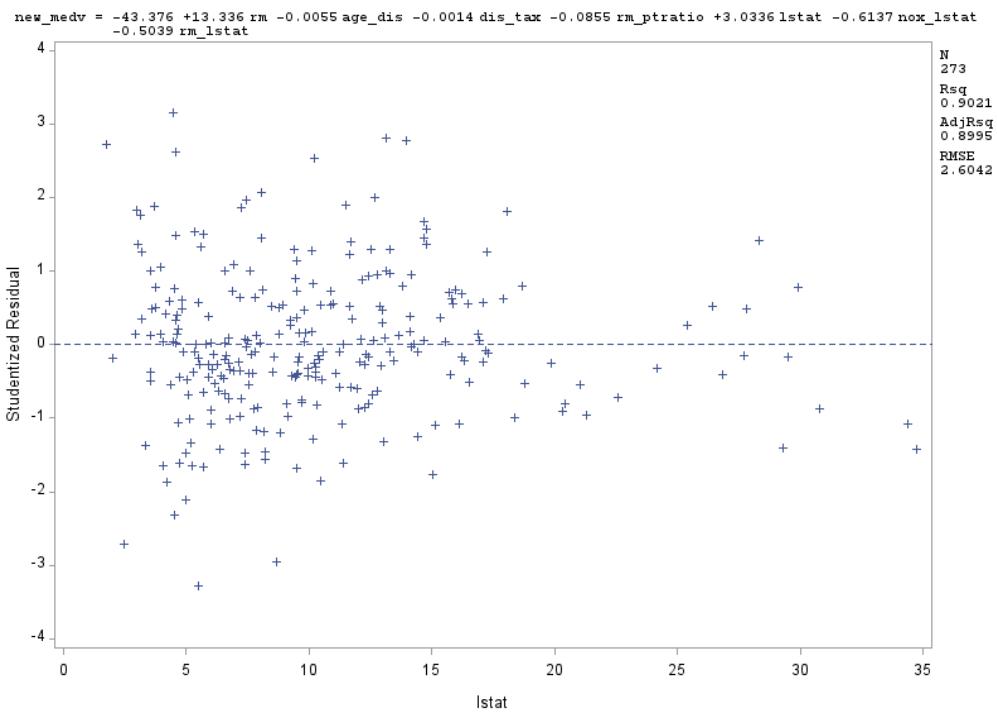


Fig. 32

Model 2 after removing outliers

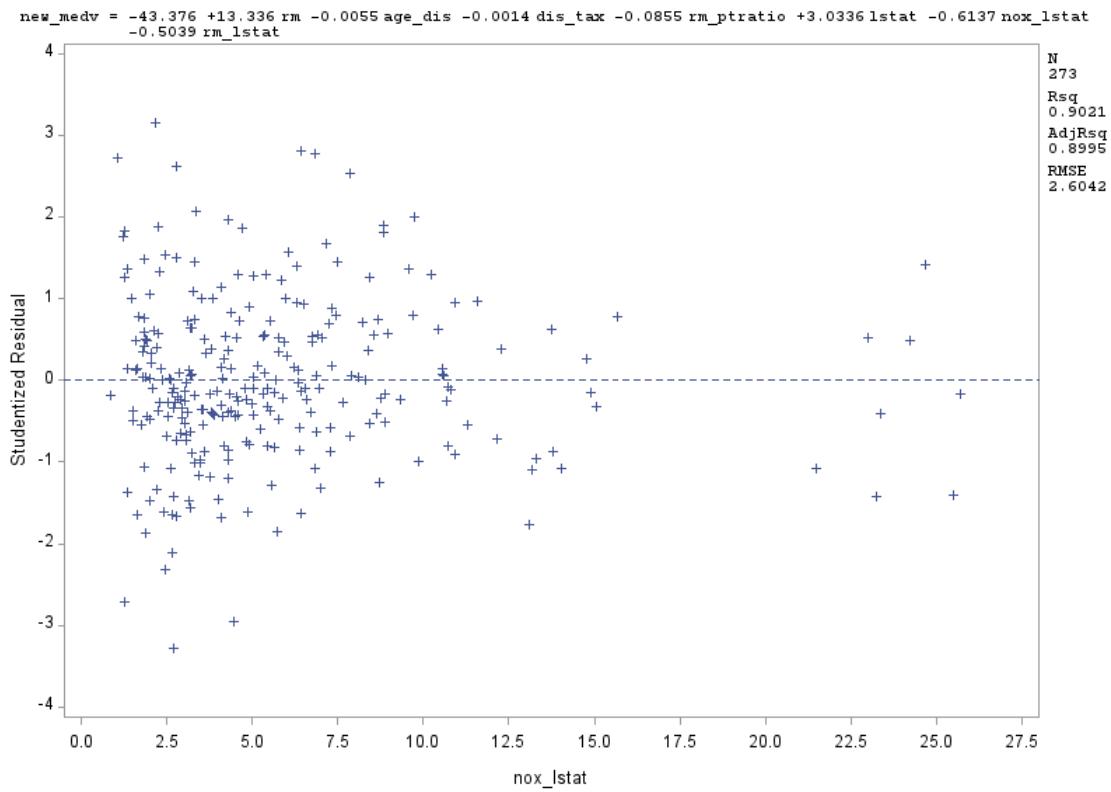


Fig. 33

Model 2 after removing outliers

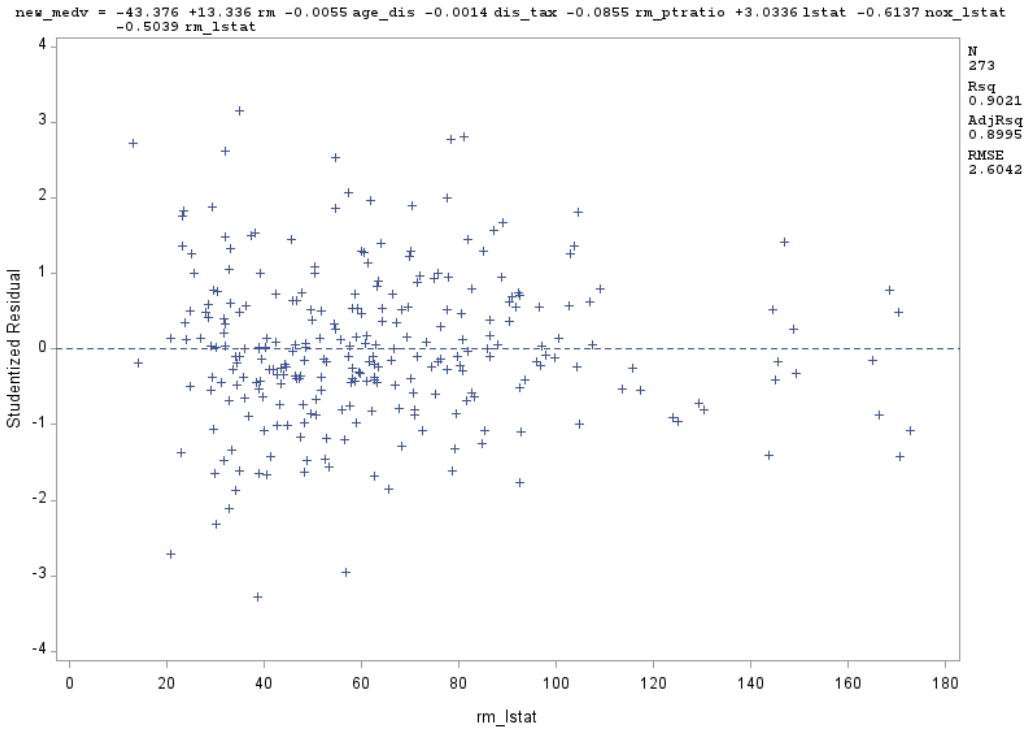


Fig. 34

Model 2 after removing outliers

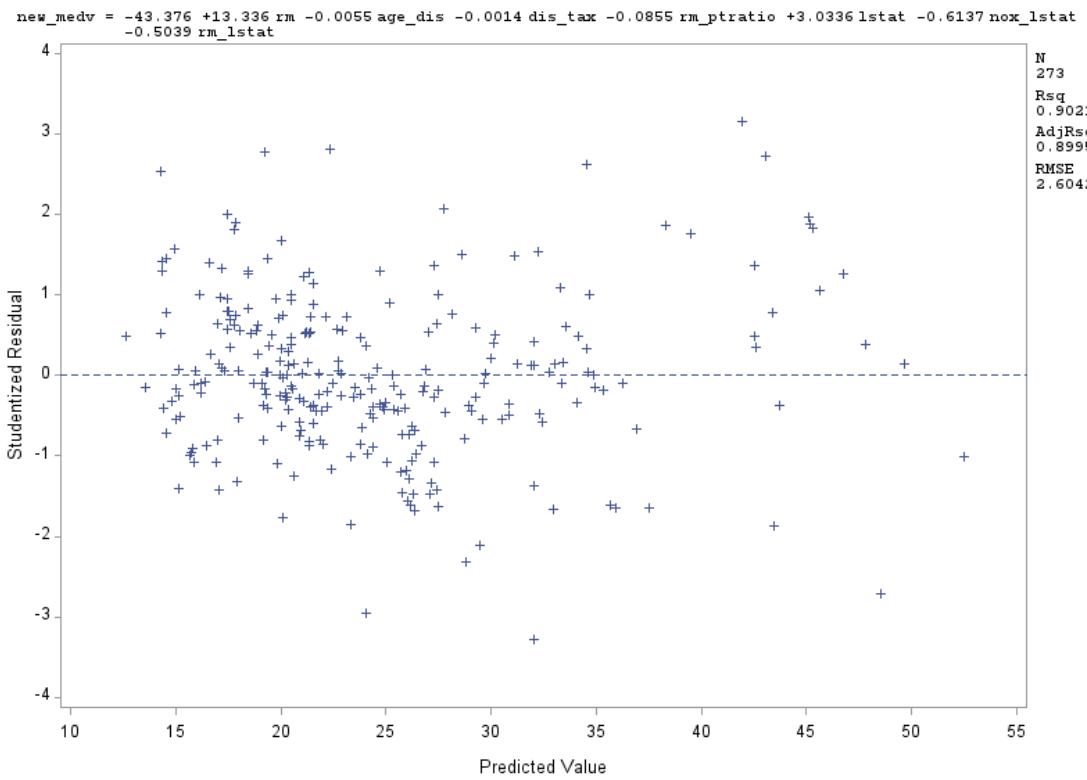
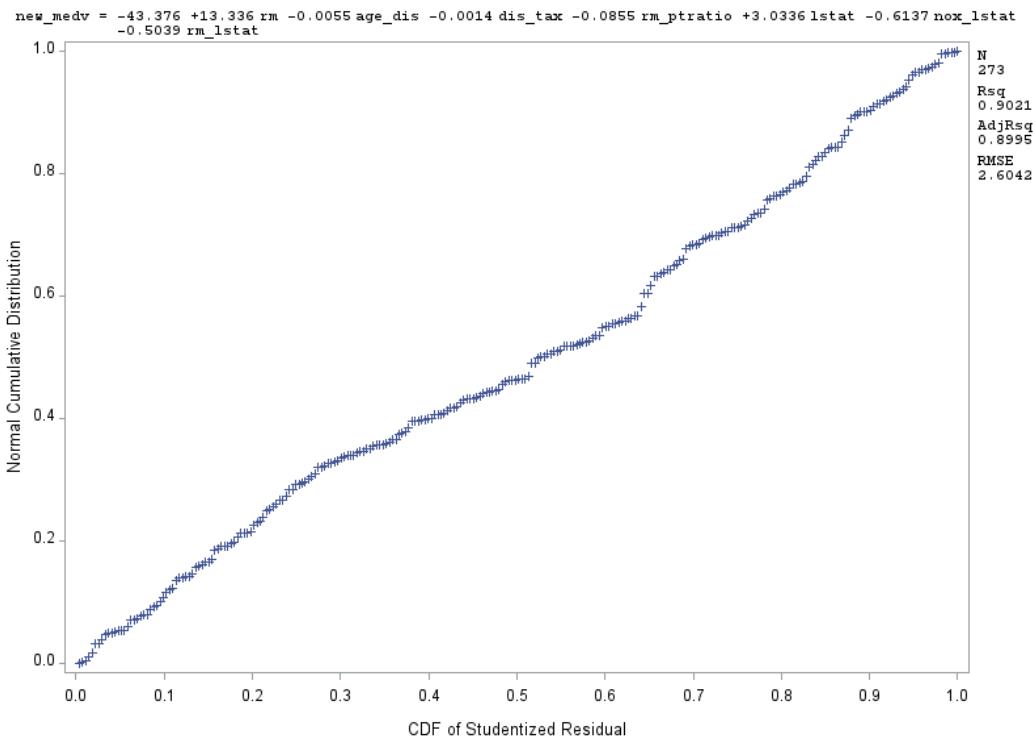


Fig. 35**Model 2 after removing outliers****Table 6**

Performance test parameters - Model 2

Obs	_TYPE_	_FREQ_	rmse	mae
1	0	91	3.36385	2.22168

Performance test parameters - Model 2

The CORR Procedure

2 Variables: medv yhat

Simple Statistics							
Variable	N	Mean	Std Dev	Sum	Minimum	Maximum	Label
medv	91	25.85055	9.13302	2352	14.30000	50.00000	
yhat	91	25.38279	8.81051	2310	8.51833	50.24374	Predicted Value of new_medv

Pearson Correlation Coefficients, N = 91		
Prob > r under H0: Rho=0		
	medv	yhat
medv	1.00000	0.93033 <.0001
yhat Predicted Value of new_medv	0.93033 <.0001	1.00000

APPENDIX E – FINAL MODEL & PREDICTIONS

Table 1

Final Model - The log transformation

The REG Procedure
 Model: MODEL1
 Dependent Variable: logmedv

Number of Observations Read	367
Number of Observations Used	367

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	12	30.71085	2.55924	210.28	<.0001
Error	354	4.30833	0.01217		
Corrected Total	366	35.01918			

Root MSE	0.11032	R-Square	0.8770
Dependent Mean	3.17054	Adj R-Sq	0.8728
Coeff Var	3.47952		

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	2.50502	0.15491	16.17	<.0001	0	0
crime	1	0.04650	0.00828	5.62	<.0001	0.19259	3.38398

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
zn	1	0.00070974	0.00032830	2.16	0.0313	0.05996	2.21345
chas	1	0.03864	0.02144	1.80	0.0723	0.03529	1.10302
nox	1	-0.52470	0.10277	-5.11	<.0001	-0.18340	3.71298
rm	1	0.24922	0.01319	18.90	<.0001	0.55401	2.47260
age	1	-0.00094060	0.00035229	-2.67	0.0079	-0.08674	3.03690
dis	1	-0.03668	0.00489	-7.50	<.0001	-0.25148	3.23198
rad	1	0.01153	0.00270	4.27	<.0001	0.14234	3.20443
tax	1	-0.00052892	0.00009845	-5.37	<.0001	-0.15733	2.46745
ptratio	1	-0.02459	0.00323	-7.62	<.0001	-0.17623	1.53953
minor	1	0.00060319	0.00016167	3.73	0.0002	0.08158	1.37591
lstat	1	-0.01494	0.00170	-8.81	<.0001	-0.28368	2.98341

Table 2

Final Model - Removing chas

The REG Procedure
Model: MODEL1

Dependent Variable: logmedv

Number of Observations Read	367
Number of Observations Used	367

Analysis of Variance					
Source	DF	Sum of Squares	Mean Square	F Value	Pr > F
Model	11	30.67131	2.78830	227.66	<.0001
Error	355	4.34787	0.01225		
Corrected Total	366	35.01918			

Root MSE	0.11067	R-Square	0.8758
Dependent Mean	3.17054	Adj R-Sq	0.8720

Coeff Var	3.49053		
------------------	---------	--	--

Parameter Estimates							
Variable	DF	Parameter Estimate	Standard Error	t Value	Pr > t	Standardized Estimate	Variance Inflation
Intercept	1	2.50664	0.15540	16.13	<.0001	0	0
crime	1	0.04657	0.00831	5.61	<.0001	0.19288	3.38391
zn	1	0.00071733	0.00032931	2.18	0.0300	0.06060	2.21308
nox	1	-0.52204	0.10309	-5.06	<.0001	-0.18247	3.71222
rm	1	0.25034	0.01321	18.95	<.0001	0.55651	2.46710
age	1	-0.00092748	0.00035333	-2.63	0.0090	-0.08553	3.03560
dis	1	-0.03722	0.00489	-7.60	<.0001	-0.25517	3.21991
rad	1	0.01231	0.00268	4.60	<.0001	0.15203	3.12132
tax	1	-0.00054065	0.00009854	-5.49	<.0001	-0.16082	2.45667
ptratio	1	-0.02506	0.00323	-7.77	<.0001	-0.17964	1.52922
minor	1	0.00060768	0.00016217	3.75	0.0002	0.08219	1.37558
lstat	1	-0.01482	0.00170	-8.72	<.0001	-0.28150	2.97921

Fig. 1

Predictions

The REG Procedure

Model: MODEL1

Dependent Variable: logmedv

Output Statistics								
Obs	Dependent Variable	Predicted Value	Std Error Mean Predict	95% CL Mean		95% CL Predict		Residual
1	.	2.1043	0.0418	2.0220	2.1865	1.8716	2.3370	.
2	.	3.1416	0.1361	2.8739	3.4094	2.7966	3.4867	.
3	3.18	3.3085	0.0153	3.2784	3.3386	3.0888	3.5282	-0.1305
4	3.07	3.1659	0.0137	3.1390	3.1928	2.9466	3.3852	-0.0932
5	3.55	3.4469	0.0148	3.4178	3.4761	3.2273	3.6665	0.0998
6	3.51	3.3973	0.0164	3.3650	3.4296	3.1773	3.6174	0.1112
7	3.59	3.3945	0.0164	3.3623	3.4267	3.1745	3.6145	0.1945
8	3.36	3.2091	0.0157	3.1782	3.2400	2.9893	3.4289	0.1478
9	3.13	3.0512	0.0144	3.0228	3.0795	2.8317	3.2706	0.0800