Build a Housing Simple Linear Regression Model Using SAS
Studio

Didem Bulut Aykurt

MIS500-1 – Foundations of Data Analytics

Colorado State University-Global Campus


Dr. Steve Chung

August 21, 2022

**Introduction**

In this case, I will work on The Boston Housing Dataset to apply a simple linear regression model by a dependent variable "MEDV" with countries variable "CRIM," using SAS Studio, which describes concerns housing values in the suburbs of Boston. The data sources originate from the StatLib library at Carnegie Mellon University on July 7, 1993, as follows housing data set available at the UCI Machine Learning repository. Five hundred six intakes represent aggregated data and 13 attributes(features) with the dependent variable(price).
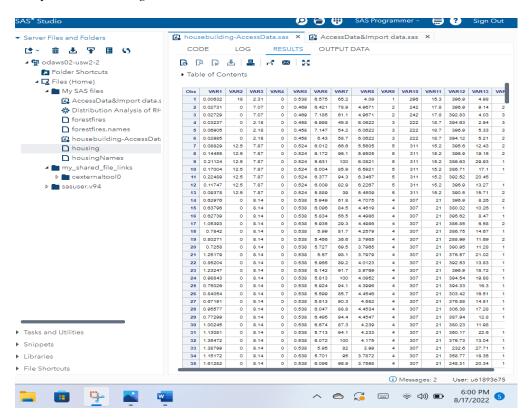
- CRIM - per capita crime rate by town

- ZN – the proportion of residential land zoned for over 25,000 sq. ft.

- INDUS – the proportion of non-retail business acres per town.

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

- NOX - nitric oxides concentration (parts per 10 million)

- RM - average number of rooms per dwelling

- AGE – the proportion of owner-occupied units built before 1940

- DIS - weighted distances to five Boston employment centers

- RAD - index of accessibility to radial highways

- TAX - full-value property-tax rate per $10,000

- PTRATIO - pupil-teacher ratio by town

- B - 1000(Bk - 0.63)^2 where Bk is the proportion of blacks by town

- LSTAT - % lower status of the population

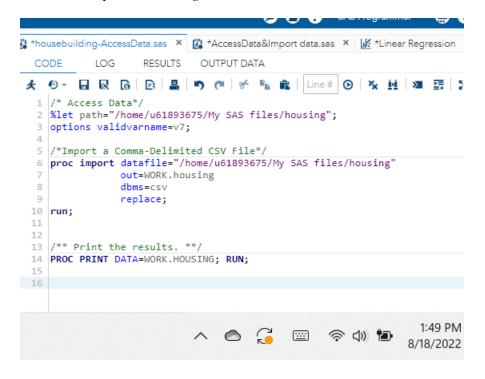- MEDV - Median value of owner-occupied homes in $1000's

**Uploading and importing the housing. Data CSV(Comma-delimited) file to SAS Studio**

The first step to uploading CSV data from my computer into SAS Studio is right click on the existing folder "My SAS files" under the main list "Files(Home)" and upload a file from my laptop. Now is the "housing. data" dataset online on SAS Studio; change the file name right, click then rename "housing" Screenshots 1 and 2 have details.

**Screenshot 1**
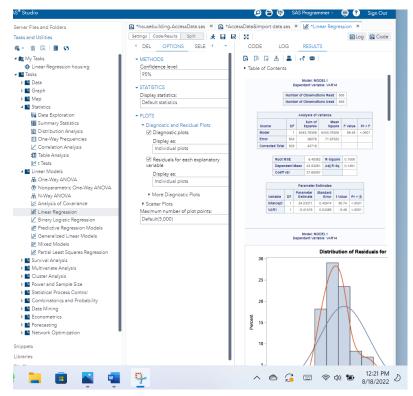*To upload the housing data set in SAS Studio.*

**Screenshot 2**
*Access and import the housing data set in SAS Studio.*



**To apply a simple linear regression model to the "housing"**

**dataset.**

Now housing dataset is ready to analyze, and I will do the first

step to click "Tasks and Utilities" on the main menu and then

double click "Linear Regression" under the "Linear Model." It will

pop up a new Linear regression page. The middle of the page has a

"DATA" organizer to choose the dependent variable

VAR14(MEDV) and the continuous variable VAR1(CRIM). The

next step is to edit the "Model" effect to a variable of VAR1. The

final part is "OPTION" to create diagnostic plots and residuals for

each variable with individual stories. Screenshot 3 shows the

detail.

**Screenshot 3**

*The simple linear regression model of housing dataset in SAS Studio.*



**The result for linear regression model of housing data set in SAS Studio.**

The first result table has statistical detail as we can say one Way ANOVA table. That table has great points like R-square, P-value, F-value, Sum of squares, and RMSE.

Look at each result to see how our data fit a linear and normal distribution. **SS**'s best result is zero, the case result of 6440. A higher number to tell not to do the housing dataset to the linear regression model.

**F value,** one-way ANOVA measures how a group of variables is jointly significant. The result from the critical of the F value table of 2.47.

**R-square** also tells us how strong the relation between two variables is. The result of 0.15 is 15%, pretty low to apply these two variables. That other way to say VAR1 is not a prediction for this model result.

**The P-value** shows how strongly this model supported the housing data set. The result of 0.001 is low. The simple linear regression model strongly supports the housing dataset as the linear regression model is statistically significant—all the table detail in screenshot 4.

The following table has the parameter estimate of VAR1 -0.42, which tells us VAR14 and VAR 1 have a negative correlation as VAR14 increases every unit, and VAR1 decreases by -0.42.

**Screenshot 4**

Mod 5-Critical Thinking; Option 1

**Model: MODEL1**
**Dependent Variable: VAR14**

| Number of Observations Read | 506 |
|---|---|
| Number of Observations Used | 506 |

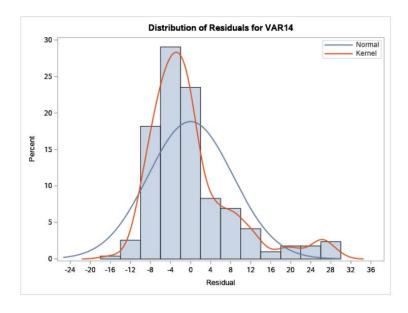| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 1 | 6440.78306 | 6440.78306 | 89.49 | <.0001 |
| Error | 504 | 36276 | 71.97522 | | |
| Corrected Total | 505 | 42716 | | | |

| Root MSE | 8.48382 | R-Square | 0.1508 |
|---|---|---|---|
| Dependent Mean | 22.53281 | Adj R-Sq | 0.1491 |
| Coeff Var | 37.65097 | | |

| Parameter Estimates | | | | | |
|---|---|---|---|---|---|
| Variable | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| |
| Intercept | 1 | 24.03311 | 0.40914 | 58.74 | <.0001 |
| VAR1 | 1 | -0.41519 | 0.04389 | -9.46 | <.0001 |

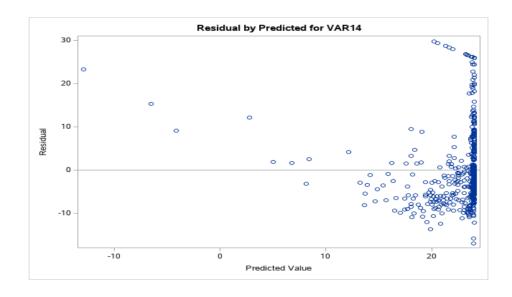# Screenshot 5

*The outcome of distribution of residual for VAR14*

# Screenshot 6

Mod 5-Critical Thinking; Option 1

Residual by Predicted for VAR14

**Screenshot 7**

RStudent by Predicted for VAR14

Mod 5-Critical Thinking; Option 1

## Screenshot 8

Observed by Predicted for VAR14

## Screenshot 9

Cook's D for VAR14

Mod 5-Critical Thinking; Option 1

## Screenshot 10

**Outlier and Leverage Diagnostics for VAR14**



## Screenshot 11

**Q-Q Plot of Residuals for VAR14**

Mod 5-Critical Thinking; Option 1

**Screenshot 12**

**Screenshot 13**

**Screenshot 14**

As a result of the definition, **screenshot 5** tells us the residual of the distribution plot for VAR14 has a heavy-tailed residual, which means there are many positive and negative residuals.

**Screenshot 11** shows a quantile plot of residuals for VAR14. The relationship between the residual and quantile is not linear.

**Screenshot 8** tells how VAR14 and VAR1 don't have a linear relationship that does not have a linear connection.

**Screenshots 6, 7, and 13** tell the linear regression model is not a good choice because all the variables are not distributed

homogeneously around the residual line or are not a random scatter against the predicted value as the residual should equal zero but not in this case. That means our model's not picking up all the signals of the dataset. It also implies that variability increases as the expected value increases.

**Screenshot 14, the** final chart, is the most interesting. The blue line represents the confidence interval for the average and is mainly created using y-intercept and that slope with a blue area. Those confidence intervals allow us to be confident in what we can say for the individual to predict value-dependent value. That also shows outlines out the outside lines.

**Conclusion**

Thus, I used the sample linear regression model for VAR1 and VAR14. An essential point for the linear model that is dependent and independent should be linear, independency, normal distribution, and equal error. In this case, two variables are not a linear and normal distribution. Other, not an independency and equal variance. That result might be outliers affecting.

## Reference

Housing Dataset

Index of /ml/machine-learning-databases/housing (uci.edu)

Normal Probability Plot of Residuals

4.6 - Normal Probability Plot of Residuals | STAT 462 (psu.edu)