# Regression Analysis of Air Pollution and Mortality

## STA-6013 Regression Analysis

Yahui Peng  @ifg583

Date: 11/29/2021

# 1. Problem Description

The objective of this project is to seek to find the association between air pollution and mortality as well as to obtain a prediction equation for mortality as a function of the predictor variables. The dataset used is named 'Pollution Data 12rev n57 Version 2018' with all numeric data collected from 57 US cities. Below is the variable description and the header of the dataset:

$y$ = Total age-adjusted mortality from all causes in deaths per 100,000 population
$x_1$ = Mean annual precipitation (in inches)
$x_2$ = Median number of school years completed for persons of age 25 years or older
$x_3$ = Percentage of the population that is nonwhite
$x_4$ = Relative pollution potential of oxides of nitrogen
$x_5$ = Relative pollution potential of sulfur dioxide
**Note:** Relative pollution potential is the product of the tons emitted per day per square
kilometer and a factor correcting for the dimensions of and exposure to the given area.

| Obs | City | y | x1 | x2 | x3 | x4 | x5 |
|---|---|---|---|---|---|---|---|
| 1 | San Jose, CA | 790.73 | 13 | 12.2 | 3 | 32 | 3 |
| 2 | Wichita, KS | 823.76 | 28 | 12.1 | 7.5 | 2 | 1 |
| 3 | San Diego, CA | 839.71 | 10 | 12.1 | 5.9 | 66 | 20 |
| 4 | Lancaster, PA | 844.05 | 43 | 9.5 | 2.9 | 7 | 32 |
| 5 | Minneapolis, MN | 857.62 | 25 | 12.1 | 3 | 11 | 26 |

We follow the $PISEAS$ (Planning, Investigating, Specification, Estimation, Assess, and Selection) rule, a general system for performing a regression analysis except for '$P$' since the data is designated here. MLR (multiple linear regression) is started with the most common OLS (ordinary least square) model. After investigation of the data, an increased power transformation on $y$ is conducted based on Box-Cox test and estimate parameters are then obtained for specification of the fitted model. However, MSE (mean squared error) and $R^2_{Adj}$ of the model are not desirable. Also, adequacy checking using diagnostic measures indicate the assumptions on error term do not hold well. To improve the fit, possible outliers and influential points are identified by influential analyses. After removing one outlier, the model is refitted, making the assumptions on error term hold while MSE and $R^2_{Adj}$ do not improve much. Another transformation on $y$ did not help much, likewise. Hence, we conclude that OLS model does not perform well on the data.

Since a model is supposed to be as informative as possible, deleting an outlier is not preferred or recommended. We try to accommodate it, such as downweighting it to a less impactful point, or almost to zero, by using WLS (weighted least square) model. Thus, the WLS model is fitted instead of the OLS one. To our delight, MSE, $R^2_{Adj}$ and even CV (coefficient variance) improved significantly with all observations kept. Finally, variable selection measures were performed to obtain the optimal model, and the prediction equation is

$$\hat{y} = 803.64449 + 2.33451x_1 + 2.28908x_3 - 0.12425x_4 + 0.42937x_5$$

We conclude $Mortality$ is positively associated with mean annual precipitation, percentage of the population that is nonwhite, relative pollution potential of sulfur dioxide, and is negatively related to relative pollution potential of oxides of nitrogen.

Fortunately, multicollinearity was not observed throughout the whole regression analysis.
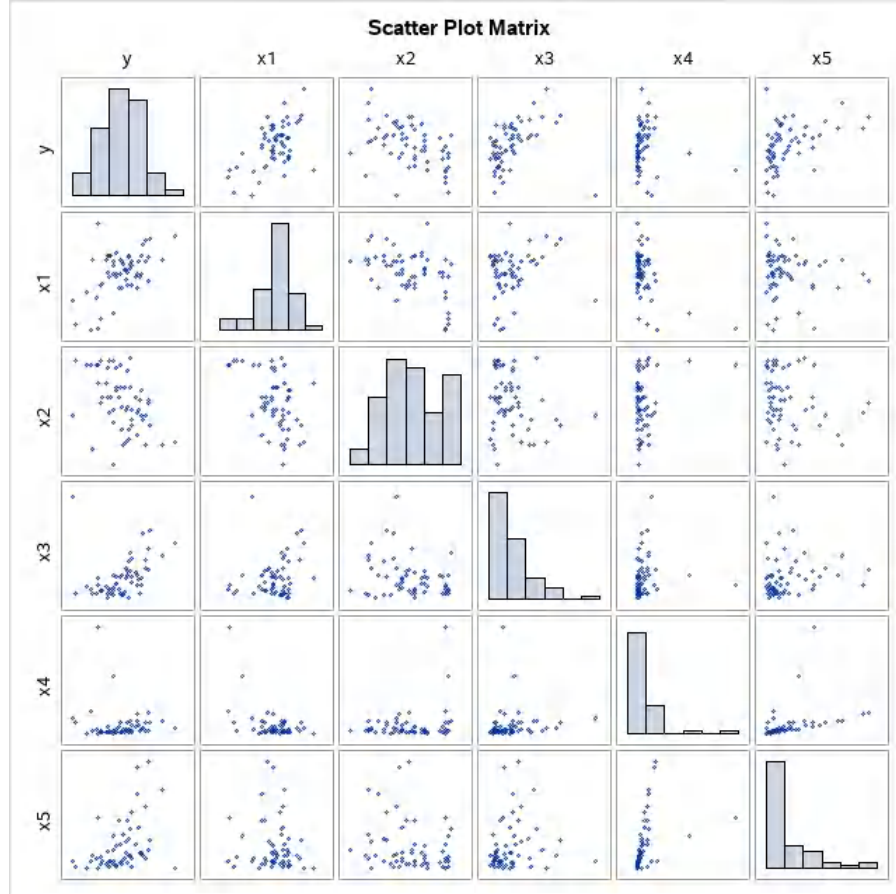
1

# 2. Investigation of the Data

The summary statistics for response and all predictors are shown in Table 1. $x_3$, $x_4$ and $x_5$ show large standard deviance compared to the mean.

### Table 1 Simple Summary Statistics

| | | | | Simple Statistics | | | |
|---|---|---|---|---|---|---|---|
| Variable | N | Mean | Std Dev | Sum | Minimum | Maximum | Label |
| y | 57 | 938.89702 | 65.79181 | 53517 | 782.37000 | 1113 | y |
| x1 | 57 | 36.96491 | 10.21023 | 2107 | 10.00000 | 60.00000 | x1 |
| x2 | 57 | 10.93333 | 0.85802 | 623.20000 | 9.00000 | 12.30000 | x2 |
| x3 | 57 | 12.18246 | 10.72048 | 694.40000 | 0.80000 | 57.60000 | x3 |
| x4 | 57 | 24.59649 | 47.40210 | 1402 | 1.00000 | 319.00000 | x4 |
| x5 | 57 | 56.45614 | 63.89726 | 3218 | 1.00000 | 278.00000 | x5 |

The pairwise scatter plots are shown in Fig.1, suggesting $y$ is positively associated with $x_1$,$x_3$ and $x_5$, negatively associated with $x_2$, and not obviously related to $x_4$. However, considering the two leverage points in the plot of $y$ vs $x_4$, there might be association between the two variables. Moreover, there appear to no or little linear pairwise association among the five predictors.
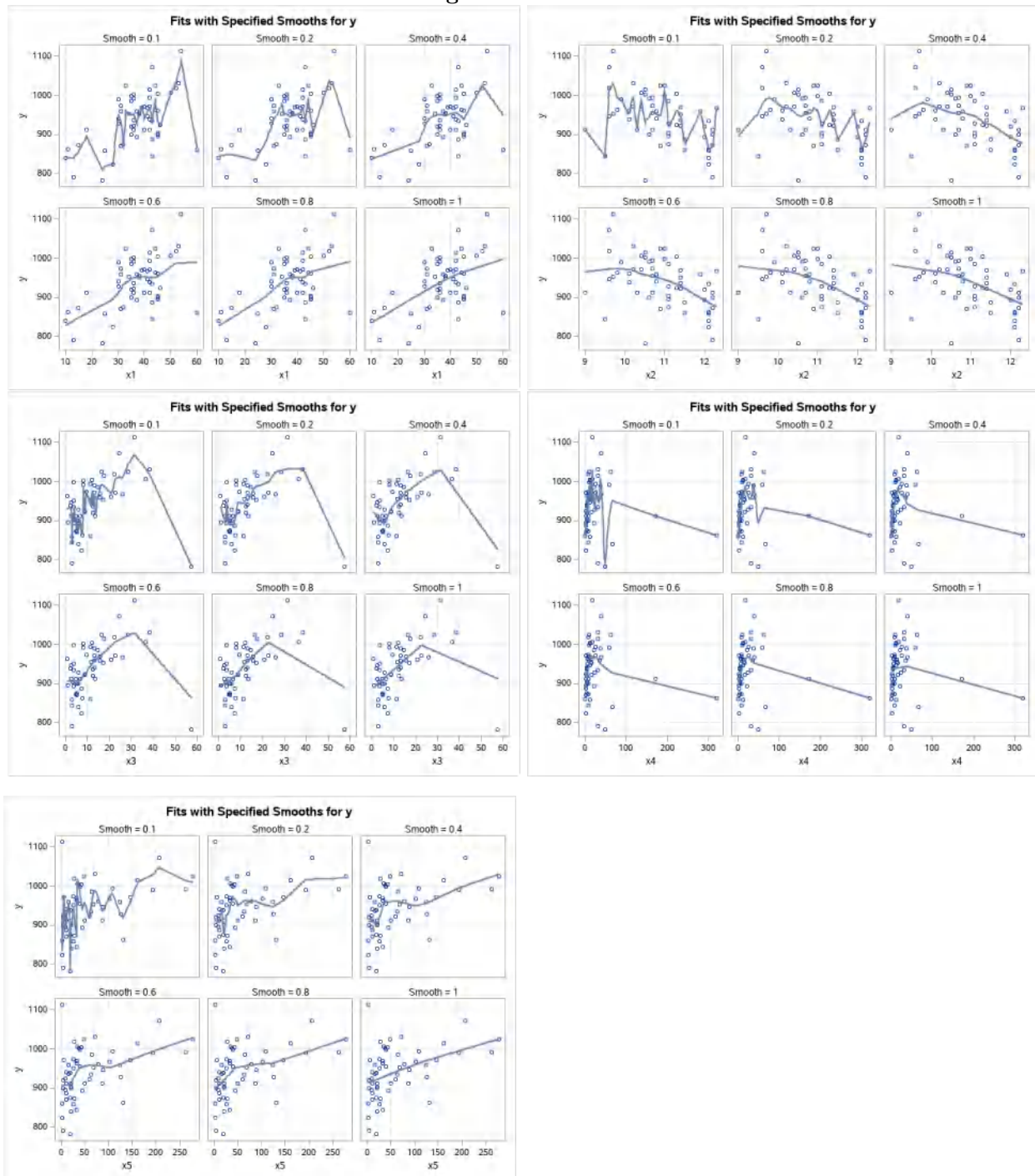
### Fig. 1 Paired Scatter Plots



Loess plots are obtained for all predictors using various smoothing parameters of 0.1, 0.2, 0.4, 0.6, 0.8 and 1. As seen in Fig.2, all predictors achieve the best Loess fit with *smooth* equals 0.6 or 0.8. However, the strong nonlinear trend with $x_3$ and $x_4$ suggests transformation of the original data, commonly on the power of response. To obtain the appropriate $\lambda$ for $y$ in the transformed form $y^\lambda$, Box-Cox analysis is recommended.

Observation 7, 17, 57 should be carefully scrutinized as possible outliers since they have a large $x_4$ or $x_3$ value and are off the trend of the majority data.
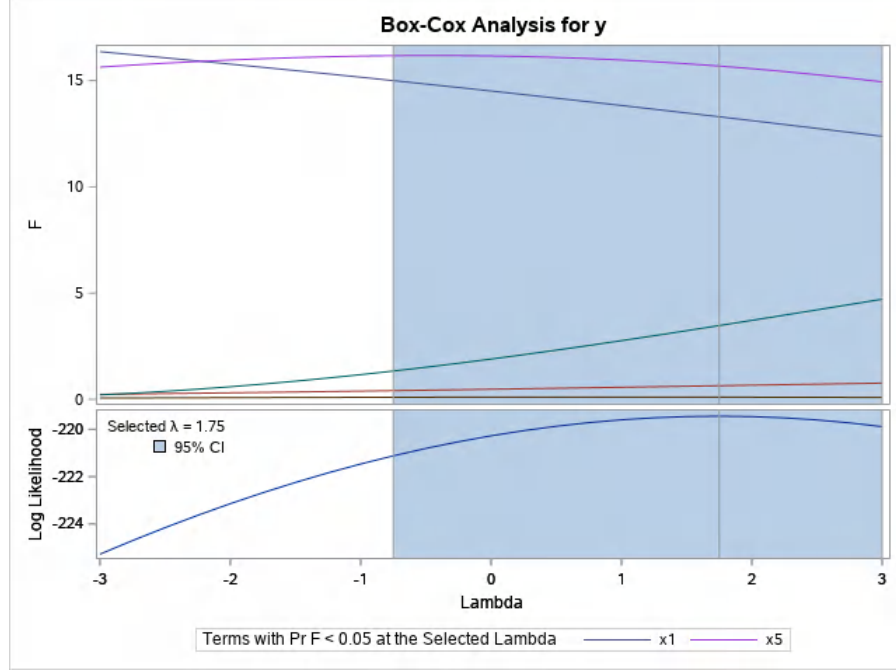
**Fig. 2 Loess Plots**

# 3. Specification of the Model

Base on preliminary judgments from the Investigation stage in Section 2 above, Box-Cox analysis suggests $y^{1.75}$ is a good transformation of the response. To simplify the prediction equation, we take value of 2, the nearest integer around 1.75, to be the proper order transformation on the response variable.

**Fig. 3 Box-Cox Analysis for $y$**



Therefore, the initial proposed MLR model is

$$y_i^2 = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3 x_{i3} + \beta_4 x_{i4} + \beta_5 x_{i5} + \epsilon_i, \ \ i = 1, ..., n$$

For $i$th observation, $y_i^2$ is the response variable, all $x_{ij}s$ are predictors for $j = 1, 2, 3, 4, 5$, and $\epsilon_i$ is the error term. $\beta_{ij}s$ are parameters where $\beta_{0j}$ is the response value when all predictors are set to zero, and other $\beta_{ij}s$ represent the change in $y_i$ per unit change in $x_{ij}$.

The model assumptions are:
1. Linear relationship between response and predictors in coefficient $\beta_{ij}$, or $\beta_{ij} \neq 0$ for at least one $j$
2. Error term $\epsilon_i$ follows $N(0, \sigma^2)$. In other words, the probability distribution of $\epsilon_i$ shows normality and constant variance.
3. No or light multicollinearity among all predictors $x_{ij}$.

# 4. Estimation of the Appropriate Model

According to Table 2, the prediction equation is obtained with the parameter estimates as

$$\hat{y}_2 = 139629 + 866.83842x_1 - 2122.80478x_2 + 343.27633x_3 - 16.63512x_4 + 133.20095x_5$$

The ANOVA analysis result given in Table 2 shows the $p$-value for $F$-test is smaller than 0.0001, indicating rejection of the $H_0$ that all $\beta_{ij}s = 0$. The assumption (1) in Section 3 holds. A good CV of 8.71% shows the ratio of standard deviation to mean is low, suggesting a low variability in the data. However, a $R^2_{Adj}$ of 49.04% means less than 50% of the variability in the data can be explained by the model. Also, the MSE is $13922^2$, which is extremely large.

### Table 2 ANOVA Analysis for the Initial Model

| | | Analysis of Variance | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 11413879112 | 2282775822 | 11.78 | <.0001 |
| Error | 51 | 9884331014 | 193810412 | | |
| Corrected Total | 56 | 21298210126 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 13922 | R-Square | 0.5359 |
| Dependent Mean | 159755 | Adj R-Sq | 0.4904 |
| Coeff Var | 8.71432 | | |

Read from the parameter estimates Table 3, the $p$-values obtained for $x_2$ and $x_4$ in $t$-test are greater than 0.05, indicating insignificant of these two variables at 0.05 level. The absolute value of $\hat{\beta}_2$ is extremely large while that of $\hat{\beta}_4$ is too small in magnitude compared with other predictor coefficients, suggesting the model does not fit well.

All VIFs (variance inflation factors) are smaller than 2, meaning only moderate correlations are observed among the predictors. Therefore, the assumption (3) in Section 3 holds. Multicollinearity is not a concern for the data.

### Table 3 Parameter Estimates for the Initial Model

| | | | Parameter Estimates | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Variable | Label | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Standardized Estimate | Variance Inflation | 95% Confidence Limits | |
| Intercept | Intercept | 1 | 139629 | 33991 | 4.11 | 0.0001 | 0 | 0 | 71389 | 207870 |
| x1 | x1 | 1 | 866.83942 | 237.71280 | 3.65 | 0.0006 | 0.45383 | 1.70211 | 389.61006 | 1344.06678 |
| x2 | x2 | 1 | -2122.80478 | 2642.57367 | -0.80 | 0.4255 | -0.09340 | 1.48544 | -7427.99279 | 3182.38322 |
| x3 | x3 | 1 | 343.27633 | 184.15581 | 1.86 | 0.0681 | 0.18870 | 1.12619 | -26.43189 | 712.98455 |
| x4 | x4 | 1 | -16.63512 | 50.81421 | -0.33 | 0.7447 | -0.04043 | 1.67640 | -118.64690 | 85.37665 |
| x5 | x5 | 1 | 133.20095 | 33.63467 | 3.96 | 0.0002 | 0.43643 | 1.33460 | 65.67653 | 200.72537 |

Also, the collinearity diagnostics table with intercept adjusted in Table 4 further confirms it as

$$\kappa = \frac{\lambda_{max}}{\lambda_{min}} = \frac{1.90269}{0.37954} = 1.32 < 10$$

### Table 4 Collinearity Diagnostics (intercept adjusted)

| | | | Collinearity Diagnostics (intercept adjusted) | | | | |
|---|---|---|---|---|---|---|---|
| | | | | Proportion of Variation | | | |
| Number | Eigenvalue | Condition Index | x1 | x2 | x3 | x4 | x5 |
| 1 | 1.90269 | 1.00000 | 0.12208 | 0.09084 | 0.02132 | 0.08798 | 0.00850 |
| 2 | 1.44511 | 1.14746 | 0.00071289 | 0.06619 | 0.14254 | 0.07239 | 0.23271 |
| 3 | 0.83232 | 1.51195 | 0.00146 | 0.08551 | 0.71547 | 0.00625 | 0.19161 |
| 4 | 0.44034 | 2.07869 | 0.33870 | 0.75730 | 0.00188 | 0.04948 | 0.36311 |
| 5 | 0.37954 | 2.23899 | 0.53705 | 0.00015162 | 0.11880 | 0.78356 | 0.20406 |

Overall, the assumption (2) in Section 3 does not hold with the initial model. The poor $R^2_{Adj}$ and MSE indicate poor adequacy of the model.

# 5. Assessment of the Chosen Prediction Equation

Shapiro-Wilk test result and QQ-plot of R-student residual vs predicted value in Table 5 suggests rejection of $H_0$ in normality test and we conclude the normality assumption on the error term does not hold.

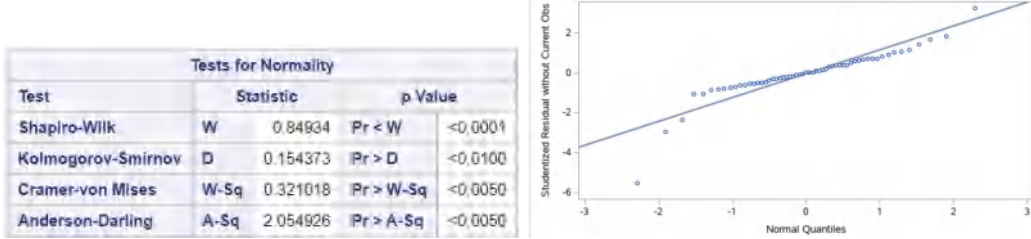**Table 5 Shapiro-Wilk Test on Normality**

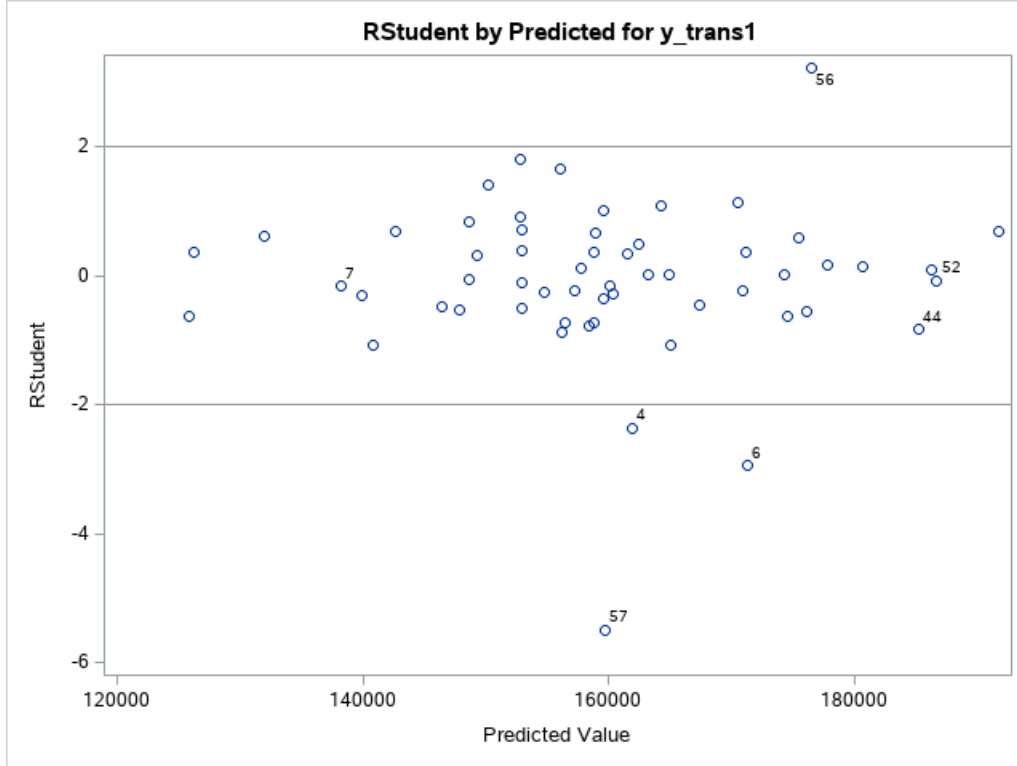| Tests for Normality | | | | |
|---|---|---|---|---|
| Test | Statistic | | p Value | |
| Shapiro-Wilk | W | 0.84934 | Pr < W | <0.0001 |
| Kolmogorov-Smirnov | D | 0.154373 | Pr > D | <0.0100 |
| Cramer-von Mises | W-Sq | 0.321018 | Pr > W-Sq | <0.0050 |
| Anderson-Darling | A-Sq | 2.054926 | Pr > A-Sq | <0.0050 |



Fig. 4 shows the of R-student residual vs predicted value. The variance on error term has an outward-funnel pattern, likely caused by observation 57, a possible outlier with R-student residual value of less than -4. Therefore, we conclude the constant variance assumption on the error term does not hold, and the assumption (2) in Section 3 does not hold.

**Fig. 4 R-student Residual vs Predicted Value**



Influential analysis is summarized in Table 6. According to *RStudent* values in Fig. 4, observation 57 is an outlier. Next, $H_{ii}$, $Cook's D$, $DFFITS$, and $DFBETAS$ are evaluated to find out possible influential points.
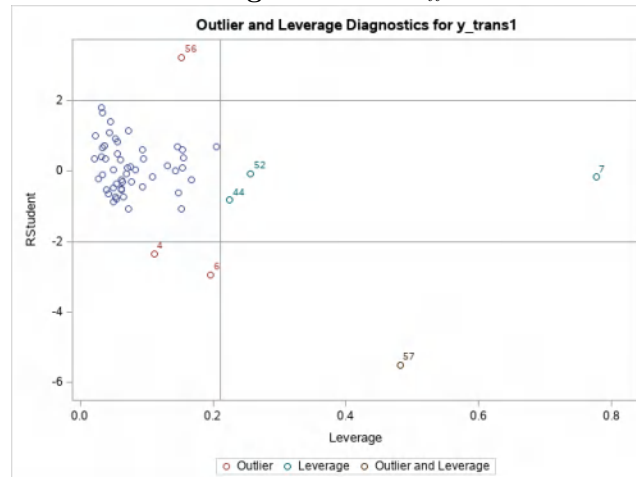
**Table 6 Influential Analysis Summary**

Output Statistics



Leverage points are determined by $h_{ii} > 2p/n = (2)(6)/57 = 0.2105$, suggesting point 7, 44, 52 and 57 to be leverage points. Combined with previous finding the point 57 is an outlier, then it is likely to be influential.

**Fig. 5 Plot of $h_{ii}$**



Outlier and Leverage Diagnostics for y_trans1

$Cook's D$ value greatly exceeds $Di > 1$ for observation 57, suggesting it is influential.

**Fig. 6 Plot of $CookD$**



The cutoff for $DFFITS_i$ is $|DFFITS_i| > 2\sqrt{p/n} = 2\sqrt{5/57} = 0.5923$. Observations 4, 6, 56, and 57 have values of $|DFFITS_i|$ that exceed this value, therefore are most likely influential.

**Fig. 7 Plot of $DFFITS$**



The cutoff for $DFBETAS_i$ is $|DFBETAS_i| > 2/\sqrt{n} = 2/\sqrt{57} = 0.2649$, we immediately noticed that observation 57 has effect on all five parameters and intercept, and its effect on $\hat{\beta}_1$ and $\hat{\beta}_3$ is large. Point 56 has effect on all five parameters other than intercept. Point 4 has effect on intercept, $\hat{\beta}_2$, $\hat{\beta}_3$ and $\hat{\beta}_5$, especially for $\hat{\beta}_2$. Point 6 has effect on intercept, $\hat{\beta}_1$, $\hat{\beta}_2$ and $\hat{\beta}_4$, especially for intercept and $\hat{\beta}_1$. Point 18 has small effect on intercept and $\hat{\beta}_2$. Point 7 has small effect on $\hat{\beta}_4$. Point 44 has small effect on $\hat{\beta}_5$.

The cutoff for $COVRATIO_i$ is $1 \pm 3p/n = 1 \pm (3)(6)/57$, or 1.316 and 0.684. Note that the values of observation 6, 8, 15, 52, 56, and 57 clearly exceed these limits, indicating these points are influential. However, point 4 barely exceed the cutoff, so the influence of it, from a practically point of view, is fairly small.

**Fig. 8 Plot of $DFBETA$**

8

**Influence Diagnostics for y_trans1**

Adopting a diagnostic view, point 57 is clearly influential since it has effect on $H_{ii}$, $RStudent$, $Cook'sD$, $\hat{\beta}$, $\hat{Y}_i$ and $COVRATIO_i$, point 56 have effect on $\hat{\beta}$, $\hat{Y}_i$ and $COVRATIO_i$, and point 4, 6 have moderate effect on $\hat{\beta}$ and $\hat{Y}_i$. Hence, we conclude that point 57 are influential.

After removing the influential point 57, the model specified in Section 3 is refitted and analyzed as follows.

The ANOVA analysis result given in Table 7 shows the $p$-value for $F$-test is smaller than 0.0001, indicating rejection of the $H_0$ that all $\beta_{ij}s = 0$. The assumption (1) in Section 3 holds. A slightly decreased CV of 6.91% shows the ratio of standard deviation to mean is low, suggesting a low variability in the data. However, a slightly higher $R^2_{Adj}$ of 64.94% means only less than 70% of the variability in the data can be explained by the model. Also, the MSE of $11100^2$ remains extremely large.

**Table 7 ANOVA Analysis after Deleting the Outlier**

| | | | Analysis of Variance | | | |
|---|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 13165586677 | 2633117335 | 21.37 | <.0001 |
| Error | 50 | 6160393749 | 123207875 | | |
| Corrected Total | 55 | 19325980426 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 11100 | R-Square | 0.6812 |
| Dependent Mean | 160541 | Adj R-Sq | 0.6494 |
| Coeff Var | 6.91405 | | |

Read from the parameter estimates Table 8, the $p$-values obtained for $x_1$, $x_2$ and $x_4$ in $t$-test are greater than 0.05, indicating insignificant of these three variables at 0.05 level, which will cause the model to be less informative due to great loss of predictors. The absolute value of $\hat{\beta}_2$ is extremely large while that of $\hat{\beta}_4$ and $\hat{\beta}_5$ are too small in magnitude compared with other predictor coefficients, suggesting the model does not fit well.

All VIFs are smaller than 3, meaning only moderate correlations are observed among the predictors. Therefore, the assumption (3) in Section 3 holds. Multicollinearity is not a concern for the data.

9

## Table 8 Parameter Estimates after Deleting the Outlier

| | | | Parameter Estimates | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Variable | Label | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Standardized Estimate | Variance Inflation | 95% Confidence Limits | |
| Intercept | Intercept | 1 | 173851 | 27808 | 6.25 | <.0001 | 0 | 0 | 117998 | 229705 |
| x1 | x1 | 1 | 378.71431 | 209.29784 | 1.81 | 0.0764 | 0.20507 | 2.01480 | -41.67278 | 799.10140 |
| x2 | x2 | 1 | -4013.61098 | 2134.85271 | -1.88 | 0.0659 | -0.18495 | 1.51795 | -8301.58885 | 274.36688 |
| x3 | x3 | 1 | 1036.70608 | 193.56682 | 5.36 | <.0001 | 0.49108 | 1.31872 | 647.91569 | 1425.49647 |
| x4 | x4 | 1 | -38.62336 | 40.71193 | -0.95 | 0.3473 | -0.09833 | 1.68523 | -120.39568 | 43.14896 |
| x5 | x5 | 1 | 99.10320 | 27.52533 | 3.60 | 0.0007 | 0.33981 | 1.39720 | 43.81695 | 154.38945 |

Shapiro-Wilk test result and QQ-plot of R-student residual vs predicted value in Table 9 indicates we fail to reject $H_0$ in normality test and we conclude the normality assumption on the error term hold.

## Table 9 Shapiro-Wilk Test on Normality

| Tests for Normality | | | | |
|---|---|---|---|---|
| Test | | Statistic | | p Value |
| Shapiro-Wilk | W | 0.963682 | Pr < W | 0.0898 |
| Kolmogorov-Smirnov | D | 0.099041 | Pr > D | >0.1500 |
| Cramer-von Mises | W-Sq | 0.129012 | Pr > W-Sq | 0.0453 |
| Anderson-Darling | A-Sq | 0.760501 | Pr > A-Sq | 0.0460 |



Fig. 9 shows the of R-student residual vs predicted value. The variance on error term has an decent even pattern along size the zero line, suggesting the constant variance assumption on the error term holds. Therefore, the assumption (2) in Section 3 holds.

**Fig. 9 R-student Residual vs Predicted Value**



Overall, the assumptions (1), (2) and (3) in Section 3 all hold with the refitted model after deleting the outlier. And there is no sign of strong multicollinearity among the predictors. However, the decent $R^2_{Adj}$ of 64.94% and the poor MSE of $11100^2$ indicate lack of adequacy of the refitted model.

Box-Cox analysis is performed again, and the result in Fig. 10 suggests $y^{1.5}$ is a proper transformation of the response after deleting the outlier. However, the $R^2_{Adj}$ and MSE do not improve much as shown in Table 10.

**Fig. 10 Box-Cox Analysis after Deleting the Outlier**

**Table 10 ANOVA Analysis after Second Transformation on Response**

| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 314016453 | 62803291 | 21.27 | <.0001 |
| Error | 50 | 147645256 | 2952905 | | |
| Corrected Total | 55 | 461661709 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 1718.40191 | R-Square | 0.6802 |
| Dependent Mean | 28945 | Adj R-Sq | 0.6482 |
| Coeff Var | 5.93675 | | |

Hence, we conclude MLR with an OLS model does not provide a desirable fit for the data, especially for the complete data. To make the model as informative as possible, we want to keep all the observations and consider WLS as a method to accommodate the outlier via downweighting it to a less impactful point, or almost to zero.

To our delight, read from Tabel 11 and Table 12, the ANOVA analysis shows the $p$-value for $F$-test is smaller than 0.0001, indicating rejection of the $H_0$ that all $\beta_{ij}s = 0$. A smaller CV of 0.17%, a higher $R^2_{Adj}$ of 71.62%, and a profoundly decreased MSE of $1.5448^2$ are obtained. All VIFs are smaller than 3, meaning only moderate correlations are observed among the predictors. Therefore, multicollinearity is not a concern for the data. We conclude the WLS model provides a better fit to the data than the OLS model.

**Table 11 ANOVA Analysis of the WLS Model**

Weight: wt2

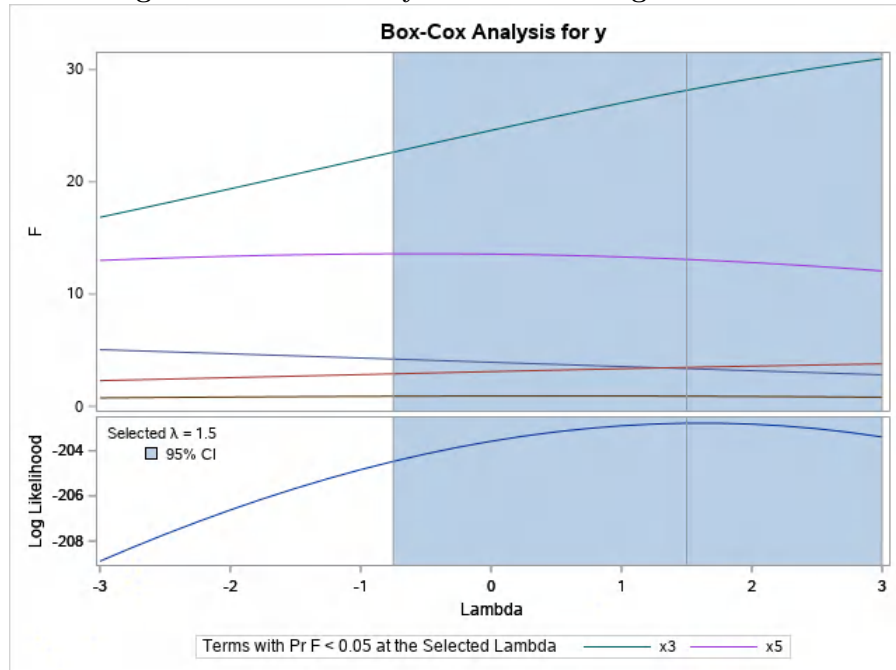| Analysis of Variance | | | | | |
|---|---|---|---|---|---|
| Source | DF | Sum of Squares | Mean Square | F Value | Pr > F |
| Model | 5 | 349.11200 | 69.82240 | 29.26 | <.0001 |
| Error | 51 | 121.70630 | 2.38640 | | |
| Corrected Total | 56 | 470.81830 | | | |

| | | | |
|---|---|---|---|
| Root MSE | 1.54480 | R-Square | 0.7415 |
| Dependent Mean | 934.69189 | Adj R-Sq | 0.7162 |
| Coeff Var | 0.16527 | | |

**Table 12 Parameter Estimates of the WLS Model**

| Parameter Estimates | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Variable | Label | DF | Parameter Estimate | Standard Error | t Value | Pr > \|t\| | Standardized Estimate | Variance Inflation | 95% Confidence Limits | |
| Intercept | Intercept | 1 | 823.10440 | 50.03136 | 16.45 | <.0001 | 0 | 0 | 722.66226 | 923.54654 |
| x1 | x1 | 1 | 2.41085 | 0.48239 | 5.00 | <.0001 | 0.61271 | 2.96529 | 1.44242 | 3.37929 |
| x2 | x2 | 1 | -1.85437 | 4.28649 | -0.43 | 0.6671 | -0.04280 | 1.93109 | -10.45986 | 6.75112 |
| x3 | x3 | 1 | 2.15583 | 0.77624 | 2.78 | 0.0076 | 0.28730 | 2.11127 | 0.59747 | 3.71419 |
| x4 | x4 | 1 | -0.11054 | 0.06866 | -1.61 | 0.1137 | -0.19447 | 2.88061 | -0.24842 | 0.02735 |
| x5 | x5 | 1 | 0.42484 | 0.06398 | 6.64 | <.0001 | 0.85790 | 1.93689 | 0.29639 | 0.55329 |

The weights obtained for the WLS model are listed in the below Table 13.

Table 13 Weights of the WLS Model

Output Statistics

| Obs | Weight | Dependent Variable | Predicted value | Std Error Mean Predict | 95% CL Mean | | 95% CL Predict | | Residual | Std Error Residual | Student Residual | Cook's D |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 9.90E-04 | 791 | 836.0270 | 15.2658 | 805.3796 | 866.6744 | 732.8233 | 939.2307 | 45.2970 | 46.864 | 0.971 | 0.017 |
| 2 | 9.84E-04 | 824 | 804.5425 | 10.6464 | 963.1700 | 905.9157 | 703.7504 | 985.3367 | 60.7029 | 47.895 | -1.289 | 0.012 |
| 3 | 9.84E-04 | 840 | 838.9850 | 14.9068 | 806.7851 | 868.8084 | 725.9140 | 952.0775 | 1.0142 | 52.395 | 0.019 | 0.000 |
| 4 | 1.58E-03 | 844 | 926.2275 | 11.0597 | 906.6042 | 950.4308 | 842.3407 | 1009 | 84.1775 | 37.131 | 2.287 | 0.026 |
| 5 | 2.54E-02 | 855 | 877.2250 | 10.1030 | 866.3619 | 897.5185 | 812.4088 | 942.0016 | -19.8152 | 26.939 | -0.677 | 0.029 |
| 6 | 2.50E-03 | 861 | 957.6067 | 11.3160 | 942.8168 | 994.2646 | 900.1605 | 1043 | 110.0967 | 31.765 | 3.467 | 0.254 |
| 7 | 1.97E-02 | 862 | 853.9691 | 11.8079 | 840.2741 | 887.6641 | 810.2280 | 897.7101 | 8.1391 | 1.864 | 1.089 | 7.104 |
| 8 | 9.79E-04 | 871 | 952.5097 | 9.1221 | 876.2159 | 906.8455 | 791.9117 | 950.1077 | 51.1997 | 48.790 | -0.435 | 0.001 |
| 9 | 1.22E-03 | 872 | 657.7874 | 19.9249 | 627.2022 | 688.3628 | 761.9337 | 951.5711 | 13.9826 | 41.513 | 0.007 | 0.003 |
| 10 | 1.55E-03 | 874 | 897.6523 | 7.2506 | 880.0600 | 912.2245 | 809.8050 | 985.4995 | 23.3721 | 42.557 | 0.549 | 0.001 |
| 11 | 1.80E-02 | 887 | 904.8846 | 3.3635 | 314.1174 | 935.6517 | 857.4655 | 998.3036 | -17.4140 | 35.776 | 1.048 | 0.004 |
| 12 | 1.14E-02 | 889 | 914.1861 | 3.5788 | 907.0013 | 921.3700 | 884.2469 | 944.1262 | 20.1961 | 14.026 | 1.440 | 0.022 |
| 13 | 5.05E-03 | 890 | 916.2022 | 5.7699 | 904.8468 | 927.6159 | 871.0923 | 961.3721 | 20.5322 | 26.662 | 0.960 | 0.012 |
| 14 | 2.86E-03 | 899 | 904.8722 | 6.6698 | 891.0426 | 918.1017 | 845.3052 | 964.4391 | 5.6122 | 28.189 | 0.199 | 0.000 |
| 15 | 1.44E-03 | 900 | 934.6191 | 7.3946 | 909.7746 | 939.4640 | 841.4360 | 1006 | 26.0896 | 40.085 | 0.626 | 0.012 |
| 16 | 4.01E-03 | 904 | 926.2907 | 4.9500 | 916.4516 | 936.1250 | 876.4175 | 976.3094 | 23.9237 | 23.099 | 0.501 | 0.009 |
| 17 | 1.44E-03 | 912 | 891.0454 | 20.0383 | 671.0857 | 910.9974 | 806.9578 | 975.1236 | 20.0544 | 39.467 | 0.520 | 0.001 |
| 18 | 1.76E-03 | 912 | 937.6616 | 13.2788 | 913.3046 | 962.5684 | 844.4587 | 1031 | 26.1396 | 33.212 | 0.605 | 0.005 |
| 19 | 1.55E-03 | 912 | 910.2903 | 6.2058 | 890.7904 | 926.7620 | 810.0146 | 990.5450 | 1.9196 | 30.236 | 0.050 | 0.000 |
| 20 | 1.06E-03 | 912 | 911.2204 | 6.8096 | 896.5496 | 925.6910 | 815.9637 | 1000 | 0.1296 | 46.069 | 0.000 | 0.000 |
| 21 | 7.56E-04 | 920 | 912.6502 | 10.5306 | 849.6114 | 934.7905 | 790.5000 | 1020 | 8.0775 | 55.185 | 0.110 | 0.000 |
| 22 | 2.35E-03 | 922 | 911.7340 | 4.9179 | 921.2609 | 941.0070 | 806.0760 | 996.1090 | 9.2640 | 31.050 | 0.290 | 0.000 |
| 23 | 1.46E-03 | 920 | 935.0676 | 6.6615 | 921.0747 | 948.3008 | 863.7668 | 1017 | 11.857E | 39.995 | 0.297 | 0.000 |
| 24 | 7.18E-03 | 925 | 937.9126 | 3.3302 | 935.2003 | 950.6250 | 913.3490 | 962.4768 | -6.7929 | 0.337 | -1.061 | 0.106 |
| 25 | 2.67E+04 | 925 | 924.6813 | 0.2399 | 904.3010 | 935.2917 | 913.9325 | 935.5305 | 0.0107 | 0.06782 | 2.447 | 1554.925 |
| 26 | 7.29E-04 | 908 | 931.8420 | 6.9544 | 900.8867 | 939.6190 | 805.6160 | 1088 | 14.3872 | 58.492 | 0.255 | 0.000 |
| 27 | 1.94E-03 | 908 | 917.5790 | 6.9269 | 899.6573 | 935.5004 | 841.1605 | 991.9973 | 20.9218 | 34.833 | 0.600 | 0.014 |
| 28 | 4.38E-04 | 941 | 907.6752 | 10.7012 | 886.3917 | 929.3680 | 759.1163 | 1056 | 33.3048 | 23.046 | 0.456 | 0.002 |
| 29 | 6.16E-03 | 948 | 948.5620 | 10.1118 | 926.2624 | 968.8802 | 905.1697 | 990.9759 | 0.3020 | 18.863 | 0.020 | 0.000 |
| 30 | 3.10E-02 | 951 | 911.8074 | 4.8304 | 902.5416 | 921.1302 | 856.3800 | 968.2960 | 38.8020 | 27.365 | 1.420 | 0.040 |
| 31 | 4.39E-04 | 954 | 951.3972 | 10.3911 | 900.5963 | 972.2682 | 801.6668 | 1101 | 2.1828 | 73.009 | 0.000 | 0.000 |
| 32 | 8.58E-03 | 954 | 935.2000 | 7.2810 | 920.8726 | 949.6471 | 826.2161 | 1042 | 18.2100 | 52.300 | 0.307 | 0.003 |
| 33 | 3.59E-04 | 969 | 923.8601 | 9.7017 | 904.3613 | 943.3260 | 816.2778 | 1091 | 24.9819 | 51.801 | 0.675 | 0.003 |
| 34 | 1.43E-03 | 969 | 930.6869 | 7.4510 | 945.7295 | 975.6450 | 877.9449 | 1044 | -1.4689 | 40.905 | 0.036 | 0.000 |
| 35 | 5.52E-04 | 961 | 959.7254 | 12.3012 | 985.0303 | 1014 | 856.9687 | 1121 | 28.7159 | 64.005 | 0.449 | 0.001 |
| 36 | 3.59E-03 | 962 | 936.0905 | 10.8554 | 904.8925 | 947.2870 | 865.7924 | 996.3956 | 26.2590 | 26.085 | 1.391 | 0.061 |
| 37 | 9.77E-04 | 960 | 994.1097 | 16.2379 | 960.5059 | 1020 | 869.0537 | 1099 | 28.3507 | 48.670 | 0.586 | 0.007 |
| 38 | 7.29E-04 | 989 | 942.8647 | 9.8660 | 923.0579 | 960.6715 | 826.2997 | 1099 | 25.7952 | 56.361 | 0.458 | 0.001 |
| 39 | 4.46E-03 | 970 | 987.8021 | 6.6267 | 974.7004 | 1001 | 939.4829 | 1085 | 17.3021 | 23.226 | 0.760 | 0.009 |
| 40 | 2.65E-04 | 971 | 956.8047 | 15.7740 | 924.1371 | 967.4724 | 762.6952 | 1149 | 16.3153 | 93.624 | 0.164 | 0.000 |
| 41 | 6.59E-04 | 972 | 905.3087 | 6.3903 | 891.5024 | 925.1509 | 760.9436 | 1028 | 54.1332 | 58.473 | 1.101 | 0.004 |
| 42 | 9.12E-04 | 966 | 943.4590 | 6.0847 | 927.2291 | 959.6905 | 709.4320 | 1047 | 42.9902 | 50.510 | 0.841 | 0.003 |
| 43 | 6.46E-03 | 969 | 983.2169 | 7.8331 | 967.4913 | 998.9425 | 939.4738 | 1029 | 6.0521 | 19.536 | 0.313 | 0.003 |
| 44 | 2.34E-03 | 981 | 1013 | 10.8763 | 991.0766 | 1035 | 945.2460 | 1081 | 21.6216 | 29.691 | 0.721 | 0.011 |
| 45 | 4.29E-03 | 995 | 971.8079 | 5.2155 | 961.4173 | 982.3684 | 923.3899 | 1020 | 22.7021 | 22.014 | 0.909 | 0.000 |
| 46 | 2.75E-03 | 990 | 910.0849 | 5.8949 | 598.8819 | 921.5079 | 845.6749 | 973.5140 | 87.7951 | 26.999 | 3.027 | 0.052 |
| 47 | 1.05E-03 | 1002 | 922.2590 | 6.7257 | 909.7767 | 935.7513 | 826.6152 | 1016 | 79.6215 | 47.201 | 1.697 | 0.010 |
| 48 | 1.57E-03 | 1004 | 953.3512 | 6.9483 | 909.4019 | 967.3005 | 873.6392 | 1083 | 50.1490 | 30.866 | 1.007 | 0.009 |
| 49 | 1.76E-04 | 1006 | 1010 | 26.3907 | 984.3018 | 1067 | 776.1094 | 1254 | 9.4456 | 113.9 | -0.084 | 0.000 |
| 50 | 3.76E-03 | 1015 | 1037 | 0.1049 | 991.0463 | 1023 | 954.4412 | 1061 | 7.5419 | 23.747 | 0.318 | 0.001 |
| 51 | 3.35E-04 | 1010 | 969.1126 | 6.6715 | 955.6431 | 1023 | 816.4468 | 1163 | 28.4974 | 82.712 | 0.345 | 0.000 |
| 52 | 1.15E-02 | 1025 | 1021 | 10.4390 | 1003 | 1060 | 983.3235 | 1084 | -5.8410 | 9.650 | -0.388 | 0.031 |
| 53 | 3.08E-04 | 1026 | 989.3800 | 16.6010 | 952.4945 | 1027 | 809.2747 | 1170 | 35.8828 | 86.003 | 0.415 | 0.001 |
| 54 | 2.11E-04 | 1000 | 1042 | 25.9990 | 988.0205 | 1094 | 721.9826 | 1263 | 11.6057 | 103.2 | 0.113 | 0.000 |
| 55 | 1.84E-02 | 1071 | 1045 | 13.6696 | 1019 | 1070 | 964.2927 | 1125 | 26.4024 | 35.931 | 0.735 | 0.011 |
| 56 | 1.76E-04 | 1113 | 1002 | 23.7037 | 953.9546 | 1049 | 764.5821 | 1239 | 111.5181 | 113.2 | 0.985 | 0.007 |
| 57 | 5.70E-05 | 782 | 866.4081 | 41.1731 | 905.7776 | 1071 | 589.3064 | 1406 | 208.9087 | 200.5 | 1.028 | 0.007 |

# 6. Selection of Variable Subset

Stepwise selection and all possible regressions are used for variable selection of the WLS model at 0.05 level, the outputs are shown in Table 14 and Table 15, respectively. All of the three statistics, $R_p^2$ , $C_p$ (close to $p + 1 = 6$), and $MS_{Res}$ optimize when $x_1$, $x_3$, $x_4$ and $x_5$ are included in the model, indicating we should remove $x_2$ in the final model.

**Table 14 Backward Selection**

Weight: wt2

| | | | | Summary of Stepwise Selection | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Step | Variable Entered | Variable Removed | Label | Number Vars In | Partial R-Square | Model R-Square | C(p) | F Value | Pr > F |
| 1 | x5 | | x5 | 1 | 0.2545 | 0.2545 | 94.0797 | 18.78 | <.0001 |
| 2 | x1 | | x1 | 2 | 0.3988 | 0.6533 | 17.4088 | 62.10 | <.0001 |
| 3 | x3 | | x3 | 3 | 0.0662 | 0.7195 | 6.3446 | 12.51 | 0.0009 |
| 4 | x4 | | x4 | 4 | 0.0211 | 0.7406 | 4.1872 | 4.22 | 0.0449 |

**Table 15 All Possible Regressions**

C(p) Selection Method

| Number of Observations Read | 57 |
|---|---|
| Number of Observations Used | 57 |

Weight: wt2

| Number in Model | C(p) | R-Square | Adjusted R-Square | MSE | Variables in Model |
|---|---|---|---|---|---|
| 4 | 4.1872 | 0.7406 | 0.7206 | 2.34909 | x1 x3 x4 x5 |
| 5 | 6.0000 | 0.7415 | 0.7162 | 2.38640 | x1 x2 x3 x4 x5 |
| 3 | 6.3446 | 0.7195 | 0.7036 | 2.49197 | x1 x3 x5 |
| 4 | 6.5901 | 0.7284 | 0.7075 | 2.45937 | x1 x2 x3 x5 |
| 4 | 11.7133 | 0.7024 | 0.6795 | 2.69449 | x1 x2 x4 x5 |
| 3 | 11.9304 | 0.6912 | 0.6737 | 2.74347 | x1 x2 x5 |
| 3 | 12.5085 | 0.6882 | 0.6706 | 2.76950 | x1 x4 x5 |
| 2 | 17.4088 | 0.6533 | 0.6404 | 3.02316 | x1 x5 |
| 4 | 28.9774 | 0.6149 | 0.5853 | 3.48677 | x2 x3 x4 x5 |
| 3 | 29.2262 | 0.6035 | 0.5811 | 3.52224 | x3 x4 x5 |

# 7. Conclusion

After using the OLS and WLS models, WLS is chosen as the best estimation method. And the final prediction equation is

$$\hat{y} = 803.64449 + 2.33451x_1 + 2.28908x_3 - 0.12425x_4 + 0.42937x_5$$

We conclude *Mortality* is positively associated with mean annual precipitation ($x_1$), percentage of the population that is nonwhite ($x_3$), relative pollution potential of sulfur dioxide ($x_5$), and is negatively related to relative pollution potential of oxides of nitrogen ($x_4$). However, *Mortality* is not associated with Median number of school years completed for persons of age 25 years or older ($x_2$). An advantage of the WLS model is relatively even magnitudes of the predictor coefficients. A concern is the negative association between $y$ and $x_4$, but a reasonable explanation maybe the existed large composition of nitrogen oxides in the air we breathe. Further investigation should be focused on whether the outlier data is reliable, because it is influential to modeling.

# 8. Appendix

```
/* 1. import database */
filename reffile '/home/u49400069/Regression Analysis/Final_Project/Pollution_Mortality.xls';


proc import datafile=reffile
        dbms=xls
        out=work.pollution;
        getnames=yes;
run;
proc print data=work.pollution;run;
proc contents data=work.pollution; run;



/* 2. investigation of the data */
/* Boxplots */
data long;
   set pollution;
   array tm(*) x1 - x5;

   do i=1 to dim(tm);
      hour=compress(vname(tm(i)), 'kd');
      Value=tm(i);
      output;
   end;
   keep hour value;
run;


proc sgplot data=long;
```

```
    vbox value / group=hour;

run;


proc sgplot data=have;

    vbox x1 x2 x3 x4 x5;

run;


* correlation table & scatter plots ;

proc corr data=pollution plots=matrix(histogram nvar=all);

run;


* loess plots ;

proc sgplot data=pollution;

        reg x=x1 y=y / clm cli;

run;


proc loess data=pollution;

        model y=x1/smooth=0.1 0.2 0.4 0.6 0.8 1.0;

run;


proc sgplot data=pollution;

        reg x=x2 y=y / clm cli;

run;


proc loess data=pollution;

        model y=x2/smooth=0.1 0.2 0.4 0.6 0.8 1.0;

run;


proc sgplot data=pollution;
```

```
        reg x=x3 y=y / clm cli;
run;


proc loess data=pollution;
        model y=x3/smooth=0.1 0.2 0.4 0.6 0.8 1.0;
run;


proc sgplot data=pollution;
        reg x=x4 y=y / clm cli;
run;


proc loess data=pollution;
        model y=x4/smooth=0.1 0.2 0.4 0.6 0.8 1.0;
run;


proc sgplot data=pollution;
        reg x=x5 y=y / clm cli;
run;


proc loess data=pollution;
        model y=x5/smooth=0.1 0.2 0.4 0.6 0.8 1.0;
run;



/* 3. specification of the model */
* box-cox analysis and transformation on y ;
proc transreg data=pollution;
        model boxcox(y)=identity(x1 x2 x3 x4 x5);
run;
```

```
data trans1;

        set pollution;

        y_trans1=y**2;

run;




/* 4. estimation of the appropriate model */

proc reg data=trans1 plots(label)=(cooksd RSTUDENTBYPREDICTED dfbetas dffits diagnostics

                observedbypredicted);

        model y_trans1=x1 x2 x3 x4 x5/

        alpha=.05 r p clb cli clm stb vif partial influence collinoint collin;

        output out=one r=resid student=sresid p=pred rstudent=rs r=y_res;

        run;

proc univariate data=one plot normal;

        var rs; * qqplot r-student residual vs predicted ;

run;




/* 5. assessment of the chosen prediction equation */

* delete an outlier ;

data pollution_new;

        set trans1 end=last;

        if not last then

                output;

run;


* refit the model ;

proc reg data=pollution_new plots(label)=(cooksd RSTUDENTBYPREDICTED dfbetas dffits diagnostics
```

```
                    observedbypredicted);

          model y_trans1=x1 x2 x3 x4 x5/alpha=.05 r p clb cli clm stb vif partial

                    influence collinoint collin;

          output out=one r=resid student=sresid p=pred rstudent=rs r=y_res;

          run;


proc univariate normal plot data=one;

          var rs; * qqplot r-student vs predicted ;

run;


* box-cox analysis and transformation on y ;

proc transreg data=pollution_new;

          model boxcox(y)=identity(x1 x2 x3 x4 x5);

run;


data trans2;

          set pollution_new;

          y_trans2=y**1.5;

run;


* refit the model ;

proc reg data=trans2 plots(label)=(cooksd dfbetas dffits RSTUDENTBYPREDICTED diagnostics
observedbypredicted);

          model y_trans2=x1 x2 x3 x4 x5/alpha=.05 r p clb cli clm stb vif partial

                    influence collinoint collin;

          output out=one r=resid student=sresid p=pred rstudent=rs r=y_res;

          run;


proc univariate normal plot data=one;
```

```
        var rs; * qqplot r-student vs predicted ;
run;


/* weighted ls est */
* step0: initial step ;
proc reg data=pollution;
        model y=x1 x2 x3 x4 x5 / clb;
        output out=result1 p=yhat r=resid;
        run;


* step1: estimate standard dev. function ;
data result1;
        set result1;
        absres=abs(resid);
run;


proc reg data=result1;
        model absres=x1 x2 x3 x4 x5;
        output out=step1 p=preds1 r=ress;
        run;


data step1;
        set step1;
        wt1=1/(preds1)**2;
run;


proc reg data=step1;
        model y=x1 x2 x3 x4 x5 /p clb;
        weight wt1;
```

```
        output out=result2 p=wyhat r=wres;

        run;


* step2: estimate the standard dev. function ;
data result2;

        set result2;

        abswres=abs(wres);
run;


proc reg data=result2;

        model abswres=x1 x2 x3 x4 x5;

        output out=step2 p=preds2;

        run;


data step2;

        set step2;

        wt2=1/(preds2)**2;
run;


proc reg data=step2 plots(label)=(RSTUDENTBYPREDICTED);

        model y=x1 x2 x3 x4 x5 /p r clb cli clm stb vif partial collinoint collin;

        weight wt2;

        output out=one r=resid student=sresid p=pred rstudent=rs r=y_res;

        run;


/* 6. variable selection */
* forward selection ;
proc reg data=step2;
```

```
        model y=x1 x2 x3 x4 x5 / selection=forward slentry=0.25;

        weight wt2;

        run;


* backward selection ;

proc reg data=step2;

        model y=x1 x2 x3 x4 x5 / selection=backward slstay=0.1;

        weight wt2;

        run;


* stepwise selection ;

proc reg data=step2;

        model y=x1 x2 x3 x4 x5 / selection=stepwise slentry=0.15 slstay=0.15;

        weight wt2;

        run;


* all possible selection ;

proc reg data=step2;

        model y=x1 x2 x3 x4 x5 / selection=cp rsquare mse adjrsq p clm cli best=10;

        weight wt2;

        run;
```