
STAT 5302 Applied Regression Analysis

Course Project

By Chunni Zhao

December 5, 2017

1 DATA PREPARATION

First transform variable *Use* into factor. From the summary statistics in Table 4.1, we can see that all of the regressors have large range. Besides, in the regressor *Use*, no category has extremely small number of observations. From the scatter plot matrix Figure 4.1 we can see that there might be a point with high regressor values but relatively low response value. By checking the point with $FTE > 2000$, $ImprV > 1.5 * 10^7$, $LandV > 2 * 10^6$, $Size > 6 * 10^5$, we can find out that it is a single point, the case 77. The case 77 has $Wst = 109.2$, which is relatively small compared with the regressors.

We may use power transformation to make regressors linearly related. From the output in Table 4.2 we can see that all of the confidence intervals include 0. And the likelihood ratio tests in Table 4.3 shows that at the significance level 0.05, we may reject the null hypothesis of making no power transformation, but can not reject making all log transformations. As for a comparison, I have also tried using the power transformation for the data after dropping the case 77. The results are in Table 4.4 and Table 4.5. We can see that the result after dropping case 77 is in accordance with the result using the full data. From the scatter plot matrix Figure 4.2 for the response and log-transformed regressors, we can see that the regressors are generally linearly related. However, the relations between response and regressors are not linear. The scatter plot matrix after dropping case 77 (Figure 4.3) has similar pattern.

2 MODELING

We may first consider fitting the regression of the *Wst* on the transformed regressors. Since for the linearly related regressors and assuming that $E(Y|X = x) = g(\beta'x)$, then the OLS estimation $\hat{\beta}$ is a consistent estimate of $c\beta$ with a constant $c[1]$.

In model1, I first considered all two way interactions of the transformed regressors and categorical variable *Use*. Then I use backward selection with AIC for variables screening. The

selected model using wilkinson-rogers notation is:

$$Wst \sim \log(FTE) + \log(ImprV) + \log(LandV) + \log(Size) + \\ Use + \log(FTE) : \log(ImprV) + \log(FTE) : \log(LandV) + \log(ImprV) : \log(LandV) + \\ \log(ImprV) : \log(Size) + \log(ImprV) : Use$$

Figure 4.4 shows the scatter plot of the response on the fitted value. We can see that the data points are concentrated on the lower left but spreads out on the right.

For the next step, I considered transforming the fitted value to get the function $g(\beta'x)$. From the result from Figure 4.5 and Table 4.6, I choose 4 as the power and fit the polynomial model:

$$E(Wst|fitted) = \beta_0 + \beta_1 fitted + \beta_2 fitted^2 + \beta_3 fitted^3 + \beta_4 fitted^4 \quad (2.1)$$

The residual plot for the fitted model is shown in Figure 4.6. We can see that there is still a cluster on the left. For the non-constant variance test, it rejects the null hypothesis of constant variance with $\chi^2 = 74.46437$, $df = 1$ and $p = 6.2e - 18$.

To solve this, I considered the weighted regression. I have tried to fit the squared residual from Equation 2.1 on the fitted value to get a weight:

$$resid^2 = \theta_0 + \theta_1 fitted$$

The summary table Table 4.7 shows that we can take

$$\frac{1}{72.6325 + 25.4773 fitted}$$

as the weight. Refit the model and from the Pearson residual plot Figure 4.6 for the weighted regression, we can see that the small trend in the cluster on the left gets smaller compared with Figure 4.7. Besides, non-constant variance test gives $\chi^2 = 0.4139518$, $df = 1$ and $p = 0.51$. We can not reject the null hypothesis of constant variance at the significance level 0.05 this time.

The outlier test in Table 4.8 shows that there are two possible outliers: case 81 and case 61. However, the Cook's distance in Figure 4.8 only shows case 41 is the most influential. The confidence region scatter plot matrix of the model using original data(using solid line) and the model removing case 41(using dashed line) is shown in Figure 4.9. We can see that although some points are outside the sold ellipse, the coefficients before and after dropping case 41 all have the same sign and close to each other to some extent. So that we may claim that the finally fitted model is kind of robust.

3 FINAL MODEL

The final model fitted is:

$$\begin{aligned}
fitted(FTE, ImprV, LandV, Size, Use) = & 1564.6 + 3.8\log(FTE) - 163\log(ImprV) - 311.2\log(LandV) \\
& + 174.2\log(Size) + 942.9Use3 + 674.2Use4 + 483.8Use5 \\
& + 547.8Use6 + 11.8\log(FTE) : \log(ImprV) \\
& - 13\log(FTE) : \log(LandV) + 28.8\log(ImprV) : \log(LandV) \\
& - 14.3\log(ImprV) : \log(Size) - 78.6\log(ImprV) : Use3 \\
& - 55.1\log(ImprV) : Use4 - 40.4\log(ImprV) : Use5 \\
& - 46.2\log(ImprV) : Use6
\end{aligned}$$

$$\hat{E}(Wst|fitted) = 25.6 + 567.3fitted + 307.4fitted^2 + 164.6fitted^3 + 89.1fitted^4$$

$$Weight(fitted) = 72.6 + 25.5fitted$$

This model might be hard for interpretation, but it provides an estimation of conditional mean for the waste given the regressors(without transforming the response). So it might be competent for our purpose of setting the tax which is a monotonically increasing function of the waste.

4 APPENDIX

4.1 TABLES

Table 4.1: Summary statistics for regressors

| | FTE | ImprV | LandV | Size | Use | Wst |
|---|----------------|-----------------|-----------------|----------------|------|----------------|
| 1 | Min. : 1.00 | Min. : 11700 | Min. : 6300 | Min. : 913 | 2:24 | Min. : 0.09 |
| 2 | 1st Qu.: 5.00 | 1st Qu.: 84700 | 1st Qu.: 29650 | 1st Qu.: 3664 | 3:27 | 1st Qu.: 2.60 |
| 3 | Median : 10.00 | Median : 200600 | Median : 54100 | Median : 6680 | 4:14 | Median : 9.70 |
| 4 | Mean : 49.86 | Mean : 542341 | Mean : 127215 | Mean : 19407 | 5:39 | Mean : 25.73 |
| 5 | 3rd Qu.: 26.00 | 3rd Qu.: 393700 | 3rd Qu.: 123850 | 3rd Qu.: 14732 | 6:43 | 3rd Qu.: 20.80 |
| 6 | Max. :3000.00 | Max. :22996600 | Max. :2530900 | Max. :788876 | | Max. :520.00 |

Table 4.2: Power transformation for regressors

| | Est.Power | Std.Err. | Wald Lower Bound | Wald Upper Bound |
|-------|-----------|----------|------------------|------------------|
| FTE | -0.03 | 0.04 | -0.11 | 0.05 |
| ImprV | 0.03 | 0.03 | -0.03 | 0.10 |
| LandV | -0.10 | 0.06 | -0.22 | 0.02 |
| Size | -0.04 | 0.04 | -0.12 | 0.05 |

Table 4.3: Likelihood ratio tests for the power for regressors

| | LRT | df | pval |
|-----------------------------|---------|----|------|
| LR test, lambda = (0 0 0 0) | 7.98 | 4 | 0.09 |
| LR test, lambda = (1 1 1 1) | 1712.24 | 4 | 0.00 |

Table 4.4: Power transformation for regressors after dropping case 77

| | Est.Power | Std.Err. | Wald Lower Bound | Wald Upper Bound |
|-------|-----------|----------|------------------|------------------|
| FTE | -0.02 | 0.05 | -0.11 | 0.08 |
| ImprV | 0.06 | 0.04 | -0.01 | 0.13 |
| LandV | -0.11 | 0.07 | -0.24 | 0.03 |
| Size | -0.03 | 0.05 | -0.12 | 0.07 |

Table 4.5: Likelihood ratio tests for the power for regressors after dropping case 77

| | LRT | df | pval |
|-----------------------------|---------|----|------|
| LR test, lambda = (0 0 0 0) | 7.47 | 4 | 0.11 |
| LR test, lambda = (1 1 1 1) | 1298.50 | 4 | 0.00 |

Table 4.6: Power transformation on the fitted value

| | lambda | RSS |
|---|-----------|----------|
| 1 | 4.050555 | 124577.3 |
| 2 | -1.000000 | 561501.6 |
| 3 | 0.000000 | 467963.1 |
| 4 | 1.000000 | 231108.9 |

Table 4.7: Summary table of fitting the squared residual on fitted

| | Estimate | Std. Error | t value | Pr(> t) |
|-------------|----------|------------|---------|----------|
| (Intercept) | 72.6325 | 208.8363 | 0.35 | 0.7285 |
| fitted(s1) | 25.4773 | 3.9088 | 6.52 | 0.0000 |

Table 4.8: result for outlier test

| | rstudent | unadjusted p-value | Bonferonni p |
|----|----------|--------------------|--------------|
| 81 | 5.965307 | 1.8742e-08 | 2.7551e-06 |
| 61 | 5.958074 | 1.9417e-08 | 2.8542e-06 |

4.2 GRAPHS

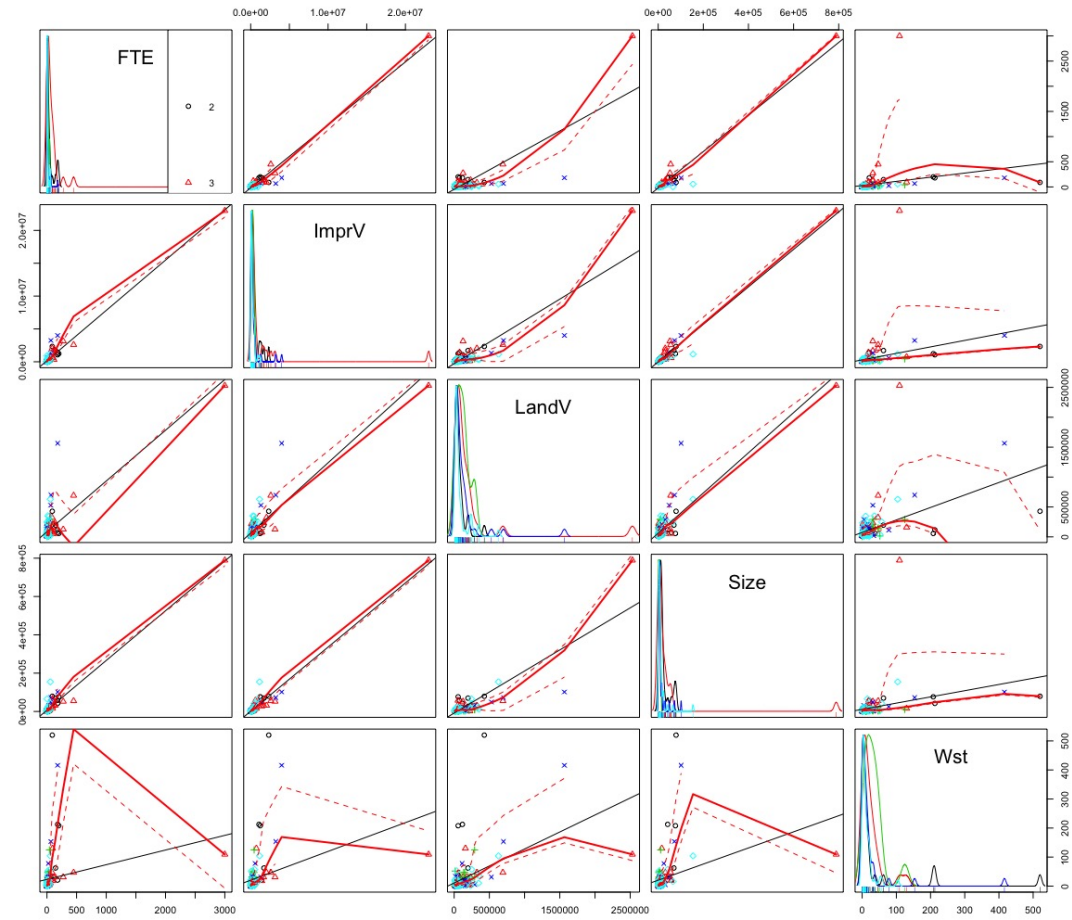


Figure 4.1: scatter plot matrix of original response and regressors

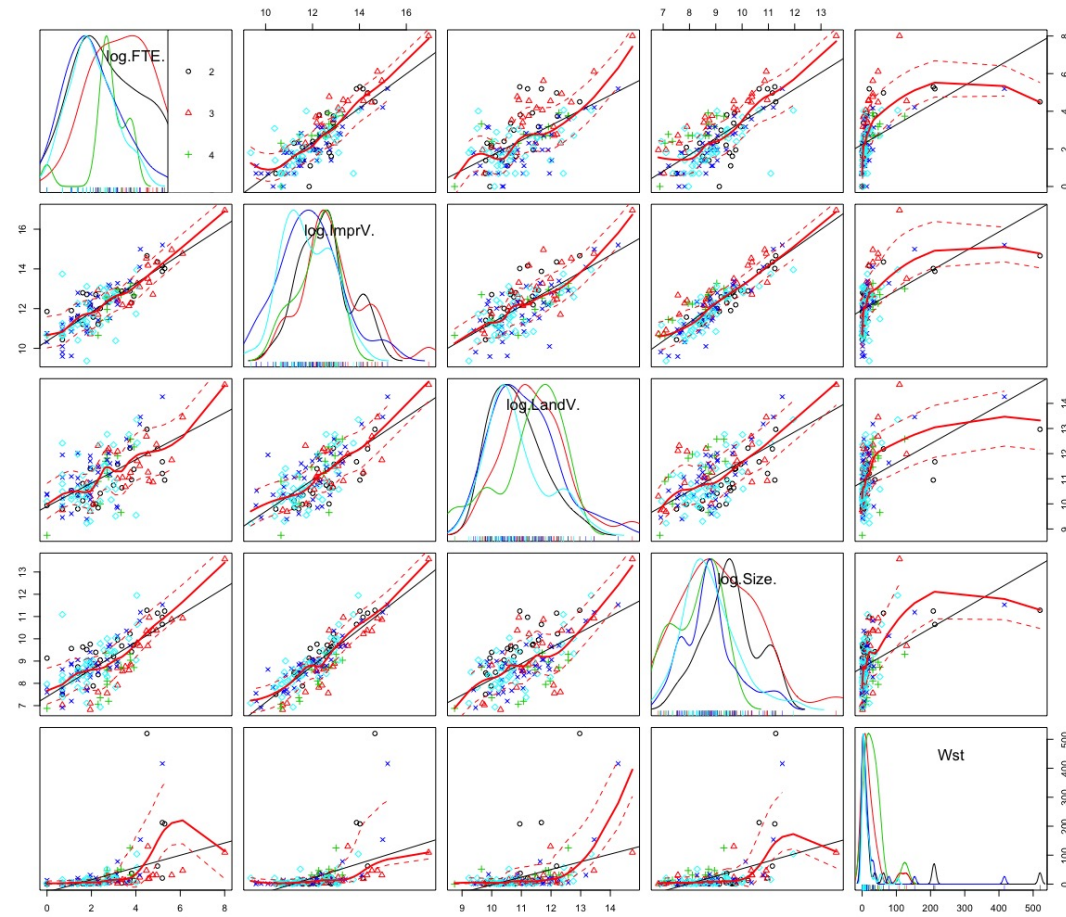


Figure 4.2: scatter plot matrix of response and regressors after using log-transformation on the regressors

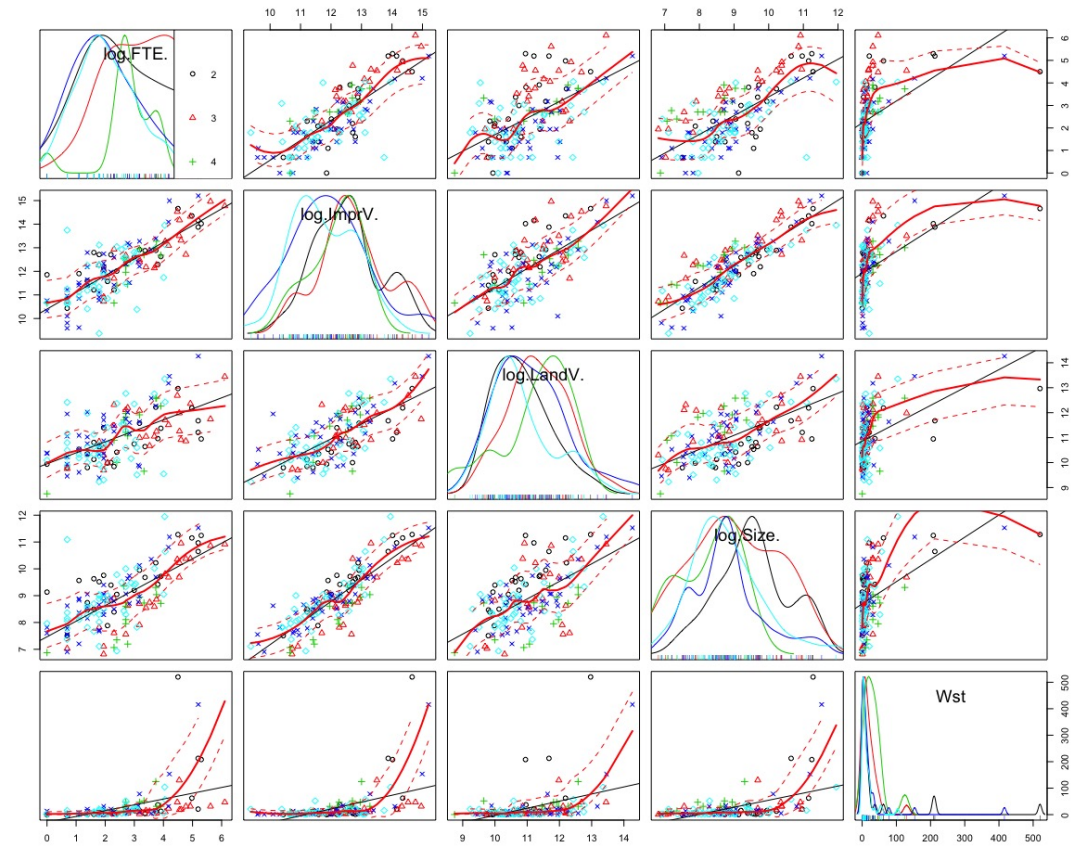


Figure 4.3: scatter plot matrix of response and regressors after using log-transformation on the regressors after dropping case 77

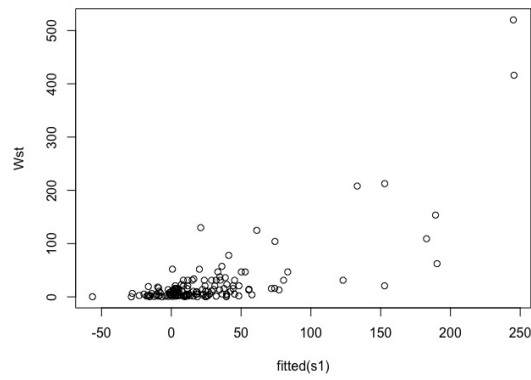


Figure 4.4: scatter plot between the response Wst and the fitted value from model selected by stepwise regression

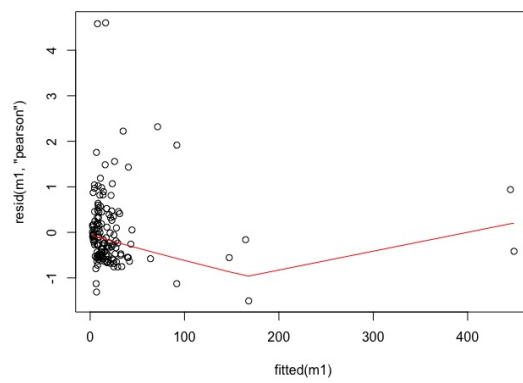


Figure 4.5: Power transforming the fitted value

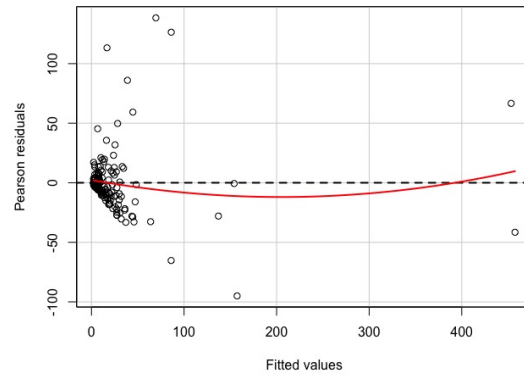


Figure 4.6: the residual plot of the model using 4th order polynomial of fitted value

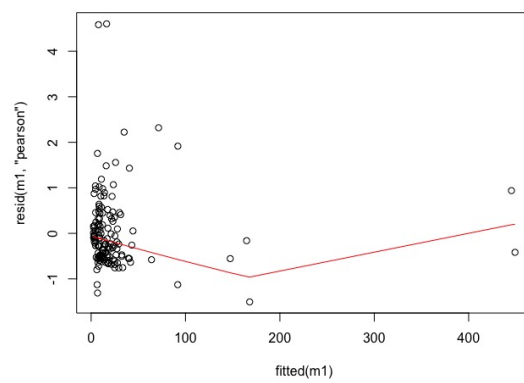


Figure 4.7: the residual plot of the weighted regression model using 4th order polynomial of fitted value

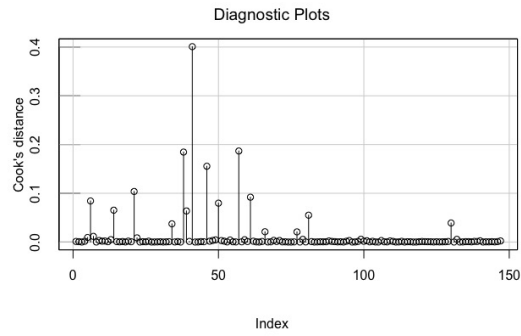


Figure 4.8: Cook's distance of the weighted regression model using 4th order polynomial of fitted value

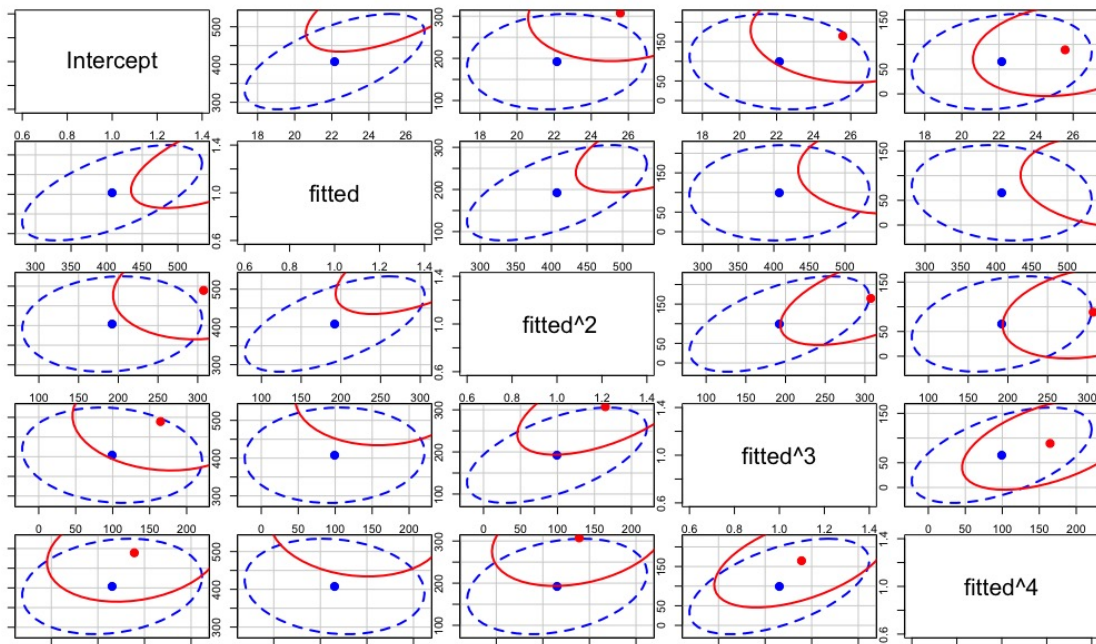


Figure 4.9: The confidence region scatter plot matrix of the model using original data(solid line) and the model removing case 41(dashed line).

4.3 CODES

```
1 #read in data
2 assDat <- read.table("/Users/mikezhang/Documents/17fall/STAT5302_TA/lab12/Waste.
   txt", header = FALSE)
3
4 #name the variables
5 nam <- c("FTE", "ImprV", "LandV", "Size", "Use", "Wst")
6 names(assDat) <- nam
7
8 #change the variable Use into factor
9 assDat$Use <- as.factor(assDat$Use)
10
11 #get the summary of the data
12 summary(assDat)
13
14 #draw the scatterplot matrix of the data
15 library(car)
16 scatterplotMatrix(~ FTE + ImprV + LandV + Size
17                   + Wst | Use, data = assDat)
18
19 #there is a point that has huge FTE, ImprV, LandV
20 #and Size but small Wst
21 which(assDat$FTE > 2000)
22 which(assDat$ImprV > 1.5*10^7)
23 which(assDat$LandV > 2*10^6)
24 which(assDat$Size > 6*10^5)
25 assDat$Wst[77]
26
27 #####
28 #modeling
29 #####
30
31
32 #try the power transformation with that point (77)
33 pt1 <- powerTransform(cbind(FTE, ImprV, LandV, Size) ~ 1,
34                       data = assDat)
35 #power transformation suggests using all log transformation
36 summary(pt1)
37
38
39 #try the power transformation without observation 77
40 pt2 <- powerTransform(cbind(FTE, ImprV, LandV, Size) ~ 1,
41                       data = assDat[-77, ])
42 #It still suggests using all log transformation
43 summary(pt2)
44
45 #scatterplot matrix of variables after log transformation
46 #with observation 77
47 scatterplotMatrix(~ log(FTE) + log(ImprV) + log(LandV)
48                   + log(Size) + Wst | Use, data = assDat)
49
50 #scatterplot matrix of variables after log transformation
51 #without observation 77
```

```

52 scatterplotMatrix(~ log(FTE) + log(ImprV) + log(LandV)
53                   + log(Size) + Wst | Use, data = assDat[-77, ])
54
55 #comment: we can see that whether dropping the observation
56 #77 makes little difference in the making the regressors
57 #linearly related
58
59 #variable screening using backward selection
60 s1 <- step(lm(Wst ~ (log(FTE) + log(ImprV) + log(LandV)
61                + log(Size) + Use)^2, data = assDat), direction = "backward")
62
63 #plot the response versus the fitted value
64 plot(Wst ~ fitted(s1), data = assDat)
65
66 #Try to do a power transformation on the fitted value
67 #the result shows that we may choose 4 as the power
68 with(assDat, invTranPlot(abs(fitted(s1)), Wst))
69
70
71
72 #choose 4 as a power and transform the fitted value
73 s2 <- lm(Wst ~ poly(I(fitted(s1)), 4), data = assDat)
74
75 #plot the residual plot
76 #there is a huge cluster on the left
77 residualPlot(s2)
78
79 #test nonconstant variance
80 #the result shows that it is significant that
81 #the model has a non constant variance
82 ncvTest(s2)
83
84 #we can do the regression of square of residuals on
85 #the fitted value to model the variance
86 we <- lm(resid(s2)^2 ~ fitted(s1), data = assDat)
87 #the result turns out that it is significant
88 summary(we)
89
90 #fit the model with the reciprocal of fitted we as the weight
91 ml <- lm(Wst ~ poly(I(fitted(s1)), 4), data = assDat,
92          weights = 1/abs(fitted(we)))
93 #plot the residual plot, we can see that
94 #there is still a cluster of points on the left.
95 plot(resid(ml, "pearson") ~ fitted(ml))
96 lines(lowess(fitted(ml), f = 1,
97             resid(ml, "pearson")), col="red")
98
99 #the summary shows that almost all of the terms are very
100 #significant
101 summary(ml)
102 #ncvTest shows that we can not reject the null hypothesis
103 #that there is no nonconstant variance
104 ncvTest(ml)
105

```

```

106 #the outlier test shows that there are 2 possible outliers
107 outlierTest(ml)
108
109 #plot the values
110 influenceIndexPlot(ml, vars = c("Cook"))
111 #find out the observation with largest Cook's distance
112 which(cooks.distance(ml) == max(cooks.distance(ml)))
113
114
115 #####
116 #plot the confidence region matrix of retaining that
117 #point and removing that point
118 #####
119
120 #create function cE to plot confidence ellipses:
121 #betas is a vector with length 2 to store the indices
122 #for betas.
123 #c is the confidence level
124 #dp is the indice for the observation to drop
125 ce <- function(betas, c, dp){
126
127   #confidence region after dropping case dp
128   confidenceEllipse(lm(Wst ~ poly(I(fitted(s1)[-dp]), 4),
129                        data = assDat[-dp, ],
130                        weights = 1/abs(fitted(we)[-dp])),
131                     which.coef = betas,
132                     levels = c,
133                     col = "blue", lty = 2)
134   #confidence region before dropping case dp for betaas
135   confidenceEllipse(ml, which.coef = betas,
136                     levels = c ,add=TRUE)
137
138 }
139
140
141 #there are 6 coefficients
142 d <- 5
143 #set the arrangement of the plots
144 par(mfrow = c(d,d))
145 #change the figure margins setting
146 par(mar=c(1,1,1,1))
147
148 #create a vector to store variable names
149 vname <- c("Intercept", "fitted", "fitted^2", "fitted^3",
150            "fitted^4")
151 for (i in c(1:d)){
152   for (j in c(1:d)){
153     if (i == j){
154       plot(1, type = "n")
155       #text the variable name
156       text(1,1, vname[i], cex = 2)
157     } else {
158       ce(c(i,j), 0.95, 41)
159     }
160   }
161 }

```

```
160 }  
161 }  
162  
163 #set the figure margins to the default  
164 par(mar=c(5.1, 4.1, 4.1, 2.1))  
165 #set the figure arrangement to default  
166 par(mfrow = c(1,1))
```

REFERENCES

- [1] Weisberg, Sanford. Applied linear regression. Vol. 528. John Wiley & Sons, 2005, p194.