

CS 498 AML HW6

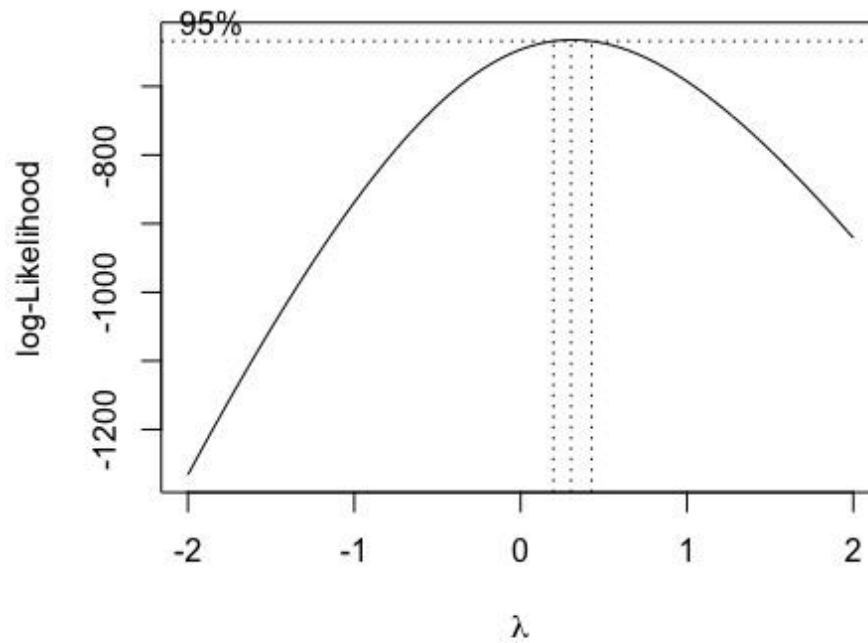
Yidi Yang(yyang160)

Huiyun Wu(hwu63)

1. Removed Row Numbers and BoxCox Curve.

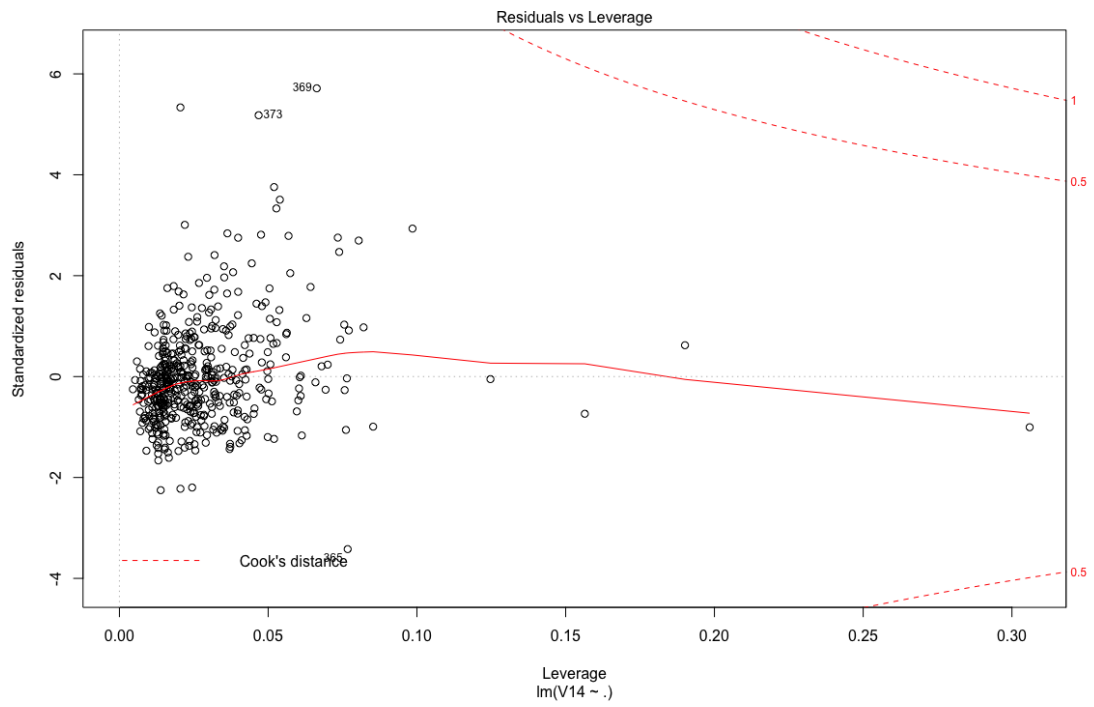
a. row numbers removed : 365, 366, 369, 372, 373, 368, 370, 371, 413

b. Box-Cox transformation curve

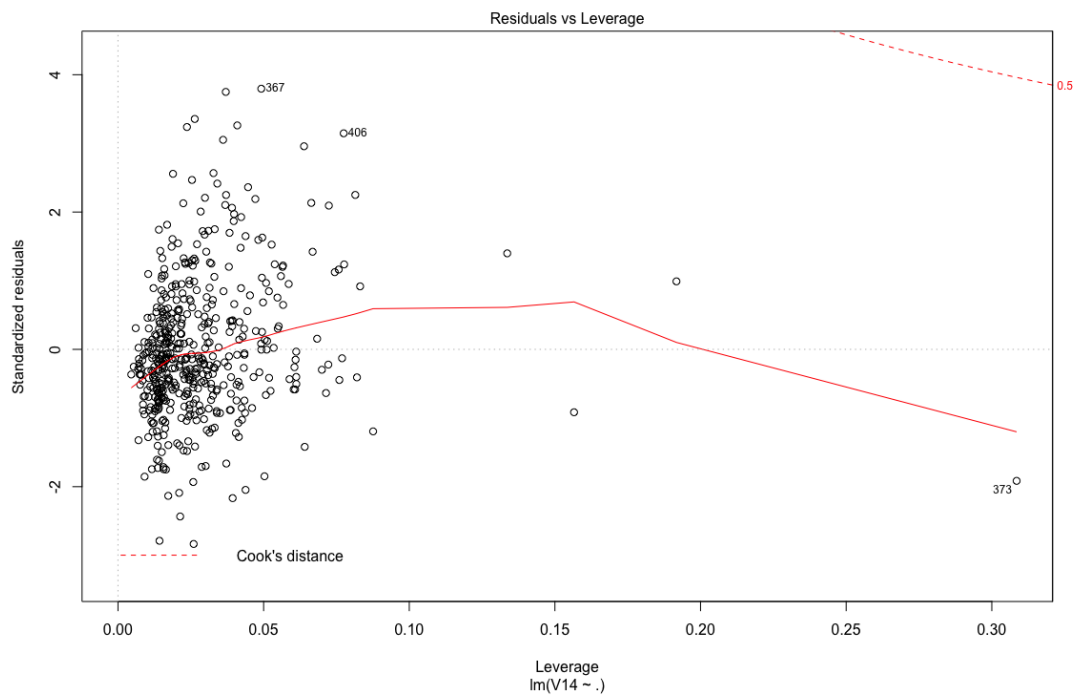


Best value of the parameter: 0.303030303

2. Plot for Identification of Outliers (3 metrics in 1 plot).

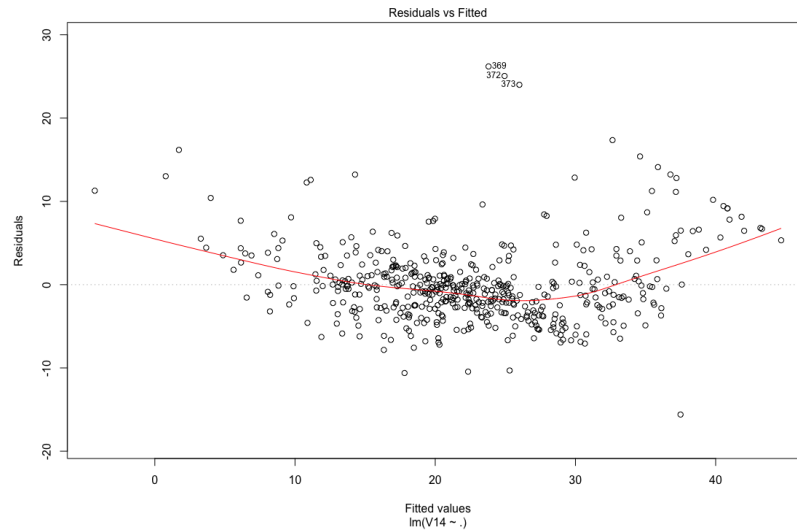


Final diagnostic plot after removing all outliers (but before transforming y).

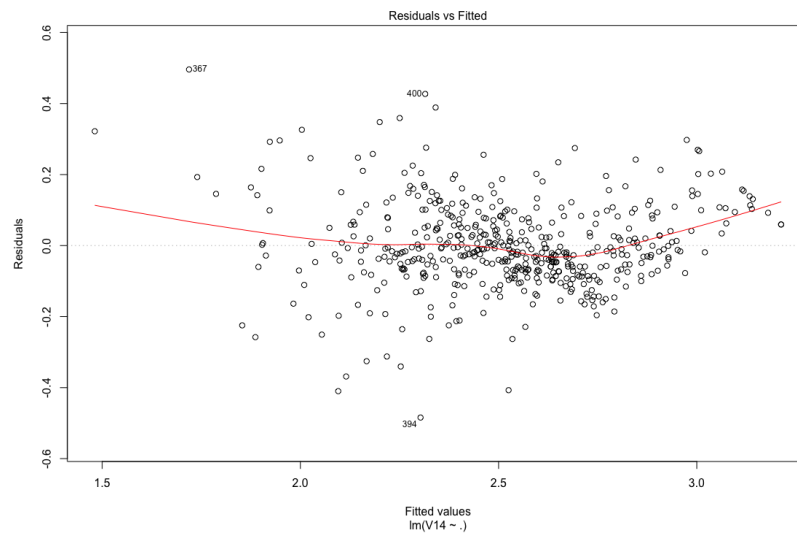


3. Plots of Standardized Residuals VS Fitted Values.

- a. Plot of Standardized residuals vs Fitted values for the linear regression model obtained without any transforms (like removing outliers or transforming dependent variables).



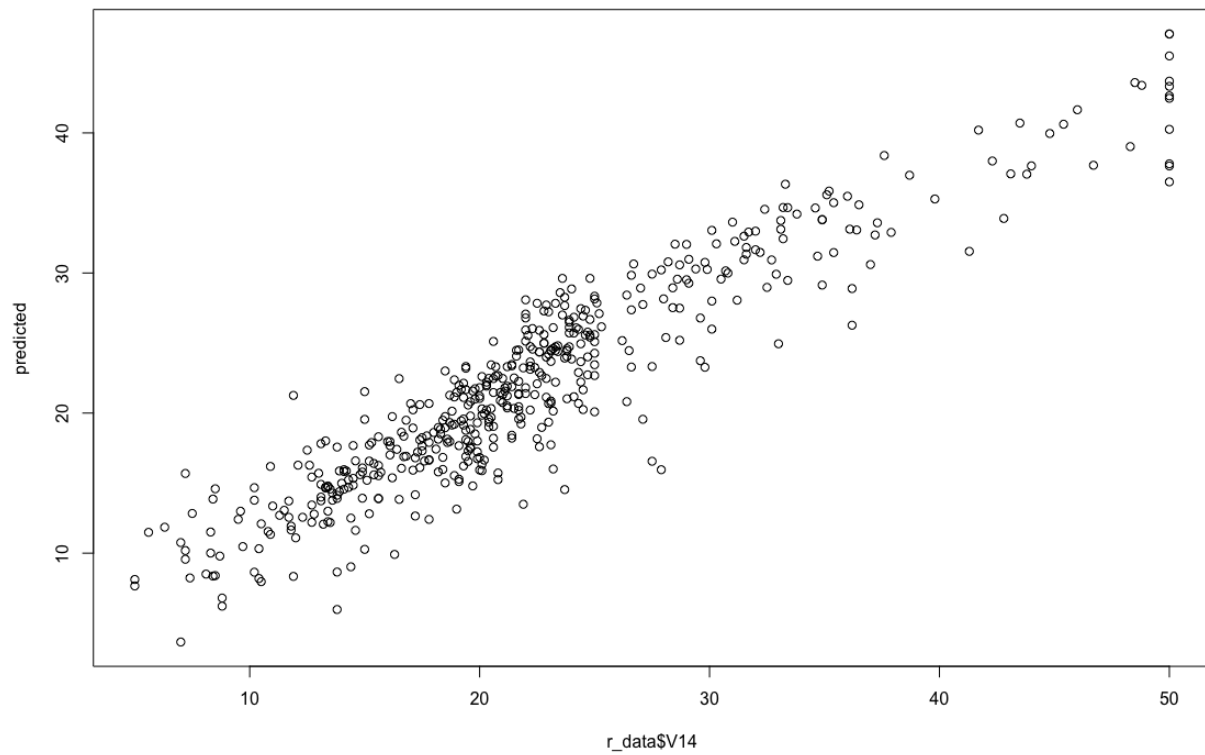
- b. Plot of Standardized residuals vs Fitted values for the final linear regression model obtained after removing all outliers and transforming the dependent variable.



- c. Compare the two plots. What do you observe ?

We observed that the second plot looks more noisy. The residual seems to be uncorrelated to the predicted value, and the mean of the residual seems to be zero, which are both good signs as mentioned on the textbook.

4. Final plot of Fitted house price vs True house price.



The true price values and the predicted ones seem to be linearly correlated, indicating that removing the outliers helps predict the house price more accurately.

5. Code screenshot.

```
1 library(car)
2 library(MASS)
3 library(pracma)
4 Data <- read.table("~/Documents/AML_hw6/Boston_Housing_data", quote="\"", comment.char="")
5 linearMod <- lm(V14 ~ ., data=Data)
6 summary(linearMod)
7 plot(linearMod)
8
9 r_data <- Data[-c(365,366,369,372,373, 368, 370, 371, 413),]
10 row.names(r_data) <- 1:nrow(r_data)
11 fit <- lm(V14 ~ ., data= r_data)
12 summary(fit)
13 plot(fit)
14 cd <- cooks.distance(fit)
15 cd <- data.frame(1:nrow(r_data), cd)
16 colnames(cd) <- c("idx", "value")
17 plot(cd$idx ~ cd$value)
18 text(cd$idx ~ cd$value, labels=idx, data=cd, cex=0.9, font=2)
19
20 bc <- boxcox(V14 ~ ., data=r_data)
21 lambda <- bc$x[which.max(bc$y)]
22 print(lambda)
23
24 new_data <- r_data
25 new_data$V14 <- new_data$V14 ** lambda
26 new_fit <- lm(V14 ~ ., data=new_data)
27 summary(new_fit)
28 plot(new_fit)
29 plot(rstandard(new_fit))
30 predicted <- fitted(new_fit) ** (1/lambda)
31
32 plot(predicted ~ r_data$V14)
33 |
34
```