

- The index of the outliers removed: 365, 369, 370, 372, 381, 406, 411, 415, 419.
- Following is a box-cox plot of the lambda parameter and its log-likelihood. The best lambda is around **0.182**.

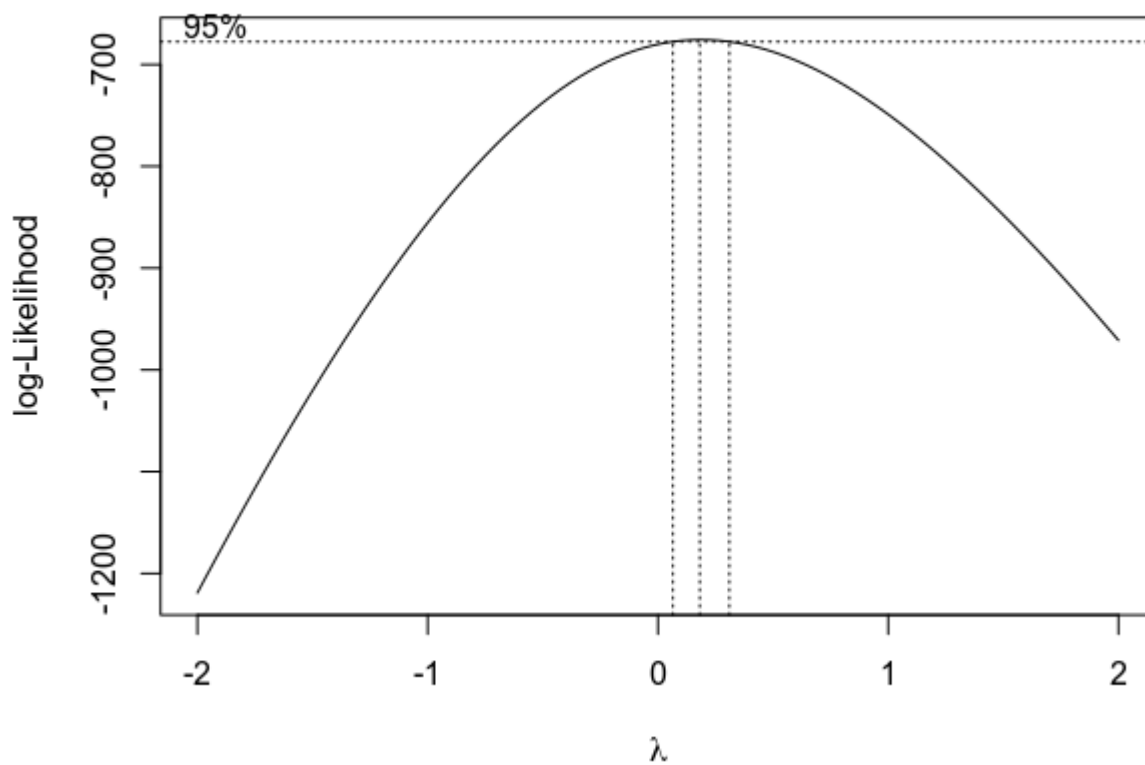


Figure 1: The log-likelihood vs. parameter of box-cox transformation.

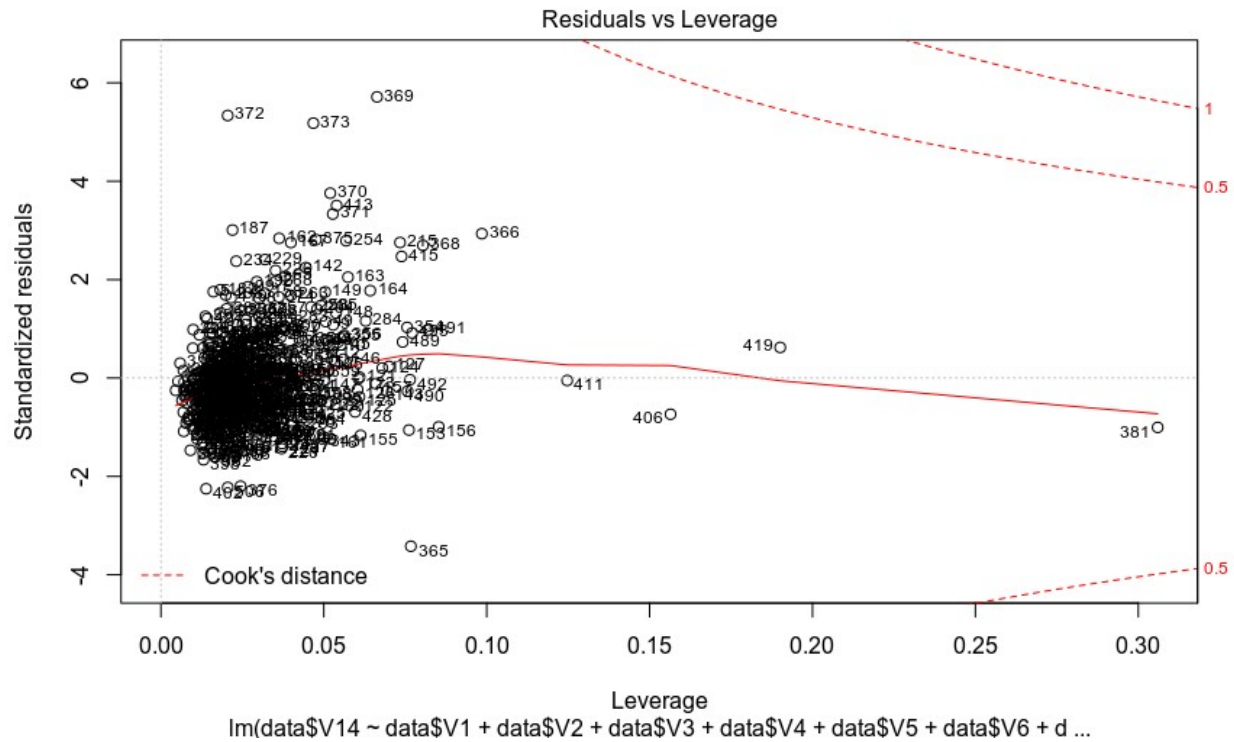


Figure 2: Standardized Residual vs. Leverage vs. Cook's Distance plot with outliers.

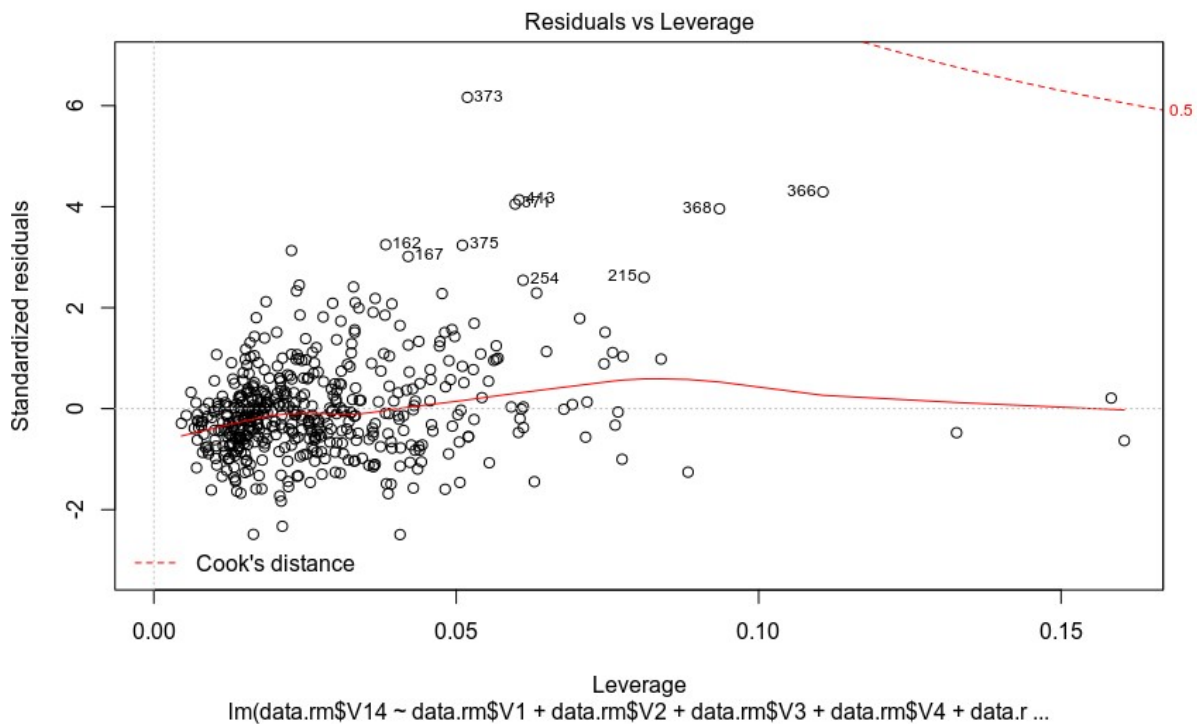


Figure 3: Standardized Residual vs. Leverage vs. Cook's Distance plot without outliers.

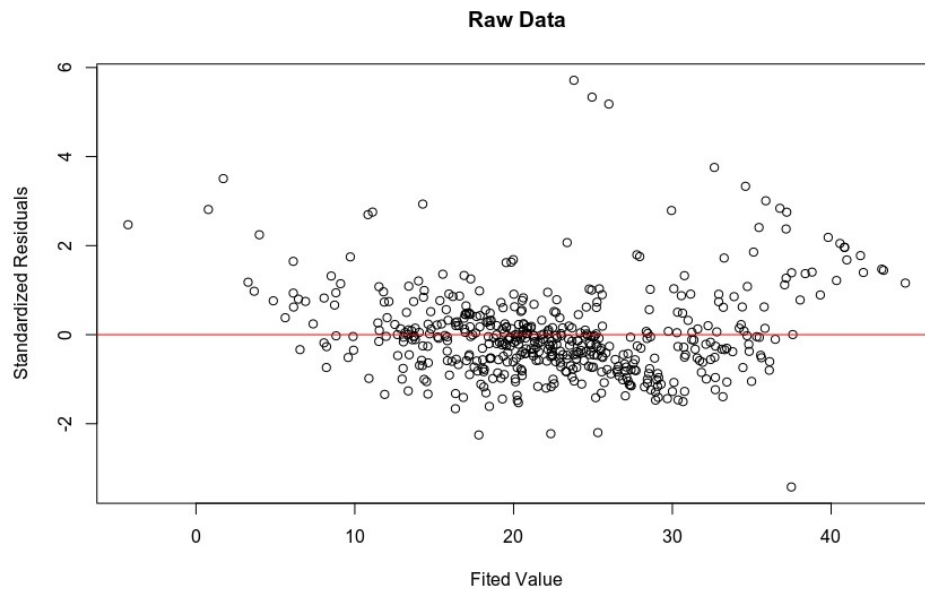


Figure 4: Standardized Residual vs. Fitted Values **without** transformation

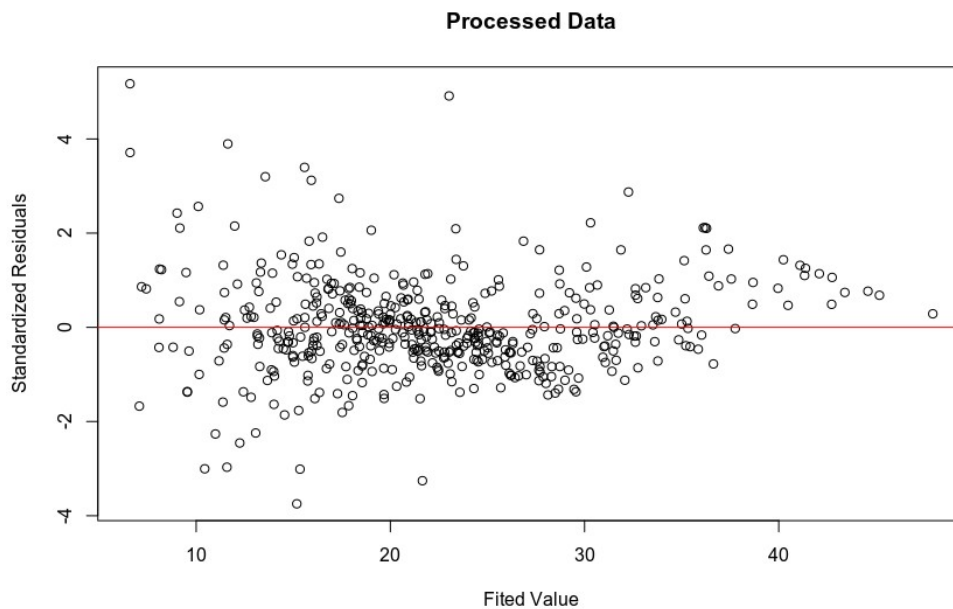


Figure 5: Std Residual vs. Fitted Values **with** 10 outliers removed & transforming the dependent variable

Comparing the above two graphs, we can observe that with transformation, the linear regression can yield results with standardized residuals that are closely center to 0 compared to the result generated by using original data.

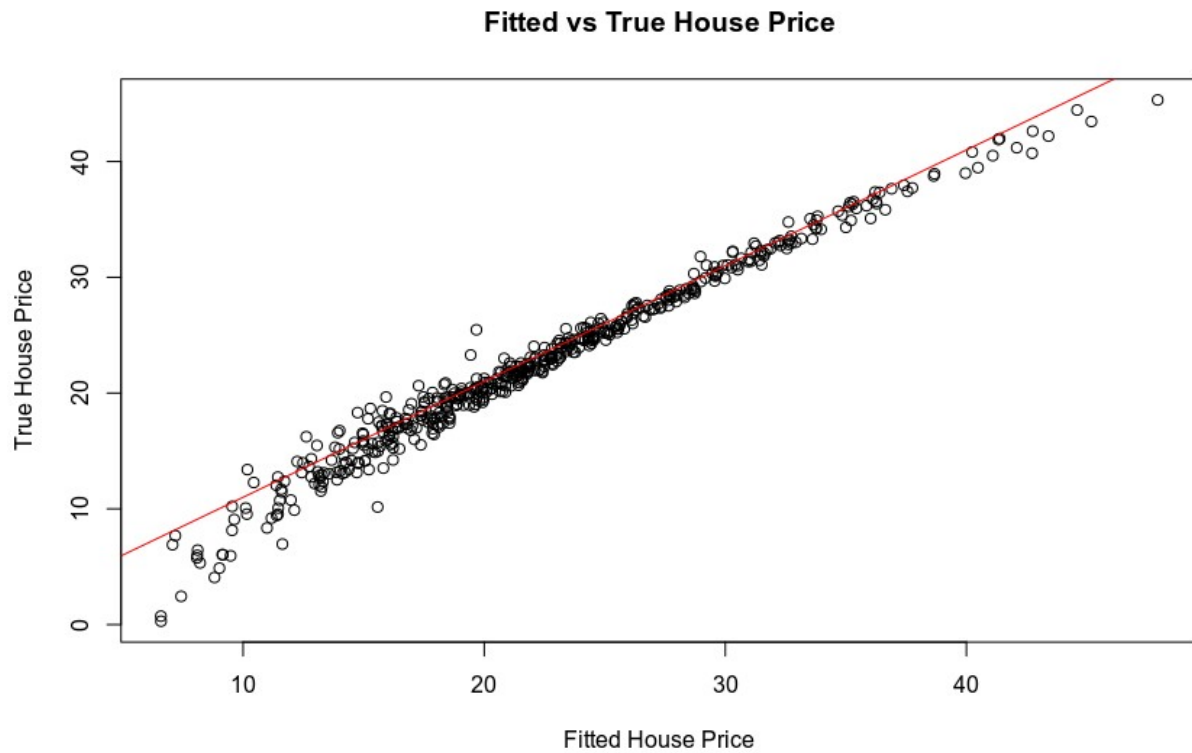


Figure 6: The plot of fitted house price vs true house price

From the above graph, we can observe that the plot of fitted house price vs true house price follows closely as the line of $y = x$. This indicates that our linear regression has yielded desirable result.

```

1 setwd('~/Documents/Applied-Machine-Learning/HW-7/')
2 library(MASS)
3 data <- read.table('housing_data.txt', header = FALSE)
4
5 # Problem A, starting fitting using the data and plot the graphs
6 data.lm <-
7   lm(data$V14~data$V1+data$V2+data$V3+data$V4+data$V5+data$V6+data$V7+data$V8+data
8     $V9+data$V10+data$V11+data$V12+data$V13, data = data)
9   # par(mfrow=c(2,2), oma = c(0, 0, 2, 0)) # Used for formatting, not necessary
10  # here
11  par(c(0, 0, 2, 0))
12  plot(data.lm, id.n = 506) # Use this line to show all the indices of all the
13  # data
14  sum = summary(data.lm)
15
16 # Problem B, performing data removal
17 # Simply re-run many times of the reduction to finalized the following 9 points
18 data.rm <- data[-c(365, 369, 370, 372, 372, 381, 406, 411, 415, 419),]
19 data.rm.lm <-
20   lm(data.rm$V14~data.rm$V1+data.rm$V2+data.rm$V3+data.rm$V4+data.rm$V5+data.rm$V6
21     +data.rm$V7+data.rm$V8+data.rm$V9+data.rm$V10+data.rm$V11+data.rm$V12+data.rm$V1
22     3, data = data.rm)
23   # par(mfrow=c(2,2), oma = c(0, 0, 2, 0)) # Used for formatting, not necessary
24   # here
25   par(c(0, 0, 2, 0))
26   plot(data.rm.lm, id.n = 10)
27
28 # Problem C
29 box_cox_result <- boxcox(data.rm.lm)
30 best_lam <- box_cox_result$x[which(box_cox_result$y == max(box_cox_result$y))]
31
32 # Problem D
33 data.box <- (data.rm$V14^best_lam-1)/best_lam
34 box.lm <-
35   lm(data.box~data.rm$V1+data.rm$V2+data.rm$V3+data.rm$V4+data.rm$V5+data.rm$V6+da
36     ta.rm$V7+data.rm$V8+data.rm$V9+data.rm$V10+data.rm$V11+data.rm$V12+data.rm$V13,
37     data = data.rm)
38   # Plot the original graph
39   data.stdres = rstandard(data.lm)
40   plot(data.lm$fitted.values, data.stdres, ylab="Standardized Residuals",
41     xlab="Fited Value", main="Raw Data")
42   abline(0, 0, col='red') # the horizon
43
44 # Plot the processed data
45 box.stdres = rstandard(box.lm)
46 box.lm.fitted_retrans = (box.lm$fitted.values*best_lam+1)^(1./best_lam)
47 plot(box.lm.fitted_retrans, box.stdres, ylab="Standardized Residuals",
48   xlab="Fited Value", main="Processed Data")
49 abline(0, 0, col='red') # the horizon
50
51 # Problem E
52 plot(box.lm.fitted_retrans, data.rm.lm$fitted.values, ylab="True House Price",
53   xlab="Fitted House Price", main="Fitted vs True House Price")
54 abline(1, 1, col='red')

```