

Class Activity#4

Write Your Name

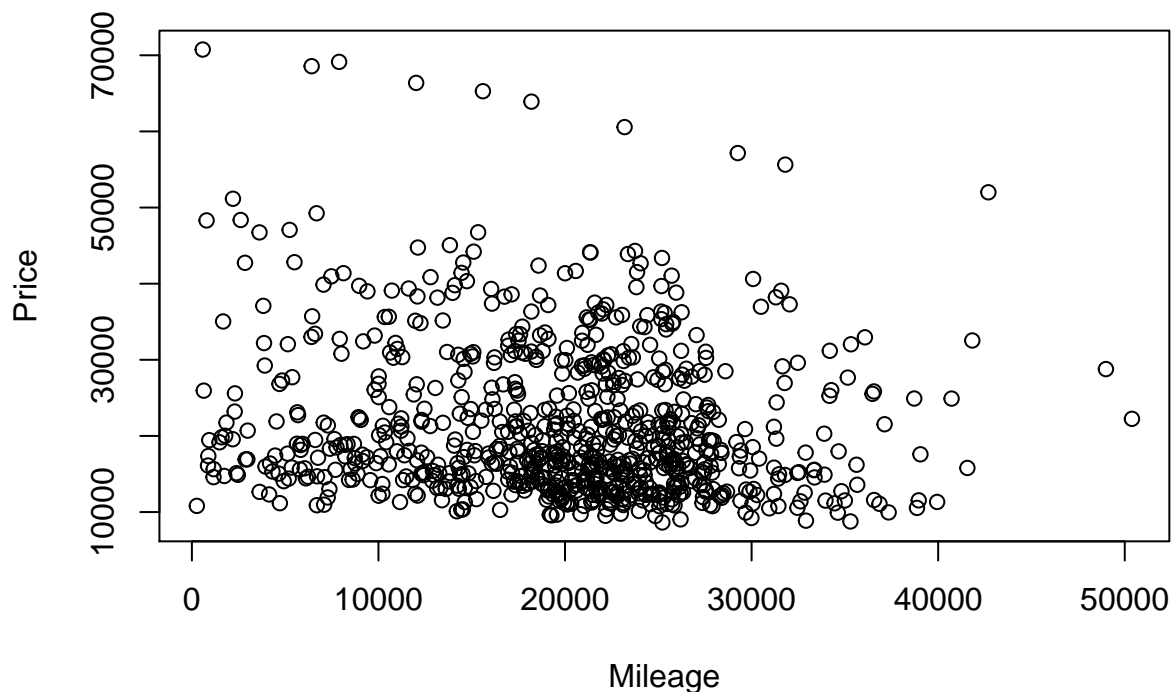
2021-02-08

Multiple Linear Regression:

Use the cars data in order to understand multiple linear regression:

- 1) produce a scatterplot from the cars data set to display the relationship between mileage(Mileage) and suggested retail price (Price). Does the scatterplot show a strong relationship between Mileage and Price?

```
cars <- read.csv("C3 Cars.csv")  
  
plot(Price~Mileage, data = cars)
```



From the scatterplot above, there definitely is not a strong relationship between Mileage and Price.

- 2) Calculate the least squares regression line, $Price = b_0 + b_1(mileage)$. Report the regression model, the R^2 value, the correlations coefficient, the t-statistics, and p-values for the estimated model coefficients (the intercept and slope). Based on these statistics, can you conclude that Mileage is a strong indicator

of price? Explain your reasoning in a few sentences.

```
M1 <- lm(Price~Mileage, data = cars)

summary(M1)

##
## Call:
## lm(formula = Price ~ Mileage, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13905   -7254   -3520    5188   46091
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.476e+04  9.044e+02  27.383  < 2e-16 ***
## Mileage      -1.725e-01  4.215e-02  -4.093  4.68e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9789 on 802 degrees of freedom
## Multiple R-squared:  0.02046,    Adjusted R-squared:  0.01924
## F-statistic: 16.75 on 1 and 802 DF,  p-value: 4.685e-05

cor(cars$Price, cars$Mileage)

## [1] -0.1430505
```

The correlations coefficient can be found using the `cor()` function, and when we pass Price and Mileage into `cor()`, we get -0.1430505, which is quite close to zero, meaning Price and Mileage do not have a very high linear correlation. For our t-statistic, we get 27.383 for our β_0 and -4.093 for our β_1 . From the summary above, the intercept and Mileage β 's have quite low p values thus, we conclude β_0 is a significant predictor of Price as well as Mileage. However, looking at the R^2 value, we find that our model captures a very low amount of the variation in the data given. Thus, while we can conclude Mileage is significant in predicting Price, our model does not perform well.

- 3) The first car in this data set is a Buick Century with 8221 miles. Calculate the residual value for this car (observed retail price minus the expected price calculated from the regression line).

```
M1$residuals[1]
```

```
##           1
## -6032.165
```

By using our linear model, we can find the residual of the first car by accessing the first index of the residuals.

- 4) Use the Cars data to conduct a stepwise regression analysis.
- a) Calculate the seven regression models, each with one of the following explanatory variables: Cyl, Liter, Doors, Cruise, Sound, Leather, and Mileage. Identify the explanatory variable that corresponds to the model with the largest R^2 value. Call this variable X_1 .

```
A1 <- lm(Price~Cyl, data = cars)
summary(A1)

##
## Call:
## lm(formula = Price ~ Cyl, data = cars)
```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11216  -5230  -2749   2773  38339
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -17.06    1126.94  -0.015   0.988
## Cyl           4054.20     206.85  19.600 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8133 on 802 degrees of freedom
## Multiple R-squared:  0.3239, Adjusted R-squared:  0.323
## F-statistic: 384.1 on 1 and 802 DF,  p-value: < 2.2e-16
```

```
A2 <- lm(Price~Liter, data = cars)
summary(A2)
```

```
##
## Call:
## lm(formula = Price ~ Liter, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10186  -5128  -3172   3032  41614
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   6185.8      846.7    7.306 6.66e-13 ***
## Liter         4990.4      262.0   19.050 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8207 on 802 degrees of freedom
## Multiple R-squared:  0.3115, Adjusted R-squared:  0.3107
## F-statistic: 362.9 on 1 and 802 DF,  p-value: < 2.2e-16
```

```
A3 <- lm(Price~Doors, data = cars)
summary(A3)
```

```
##
## Call:
## lm(formula = Price ~ Doors, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13018  -7052  -2800   5420  46948
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  27033.6     1475.2   18.325 < 2e-16 ***
## Doors        -1613.2      406.6   -3.968 7.91e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 9795 on 802 degrees of freedom
## Multiple R-squared:  0.01925,    Adjusted R-squared:  0.01803
## F-statistic: 15.74 on 1 and 802 DF,  p-value: 7.906e-05

A4 <- lm(Price~Cruise, data = cars)
summary(A4)

##
## Call:
## lm(formula = Price ~ Cruise, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14913  -6020  -1454   3634  46971
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  13921.9      632.7   22.00  <2e-16 ***
## Cruise       9862.3      729.4   13.52  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8926 on 802 degrees of freedom
## Multiple R-squared:  0.1856, Adjusted R-squared:  0.1846
## F-statistic: 182.8 on 1 and 802 DF,  p-value: < 2.2e-16

A5 <- lm(Price~Sound, data = cars)
summary(A5)

##
## Call:
## lm(formula = Price ~ Sound, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -14491  -6874  -3184   5014  50257
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  23130.1      611.0  37.856  < 2e-16 ***
## Sound       -2631.4      741.4   -3.549 0.000409 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9814 on 802 degrees of freedom
## Multiple R-squared:  0.01546, Adjusted R-squared:  0.01423
## F-statistic: 12.6 on 1 and 802 DF,  p-value: 0.0004092

A6 <- lm(Price~Leather, data = cars)
summary(A6)

##
## Call:
## lm(formula = Price ~ Leather, data = cars)
##
```

```
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13260  -7435  -2691   5422  48453
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  18828.8      655.6   28.720 < 2e-16 ***
## Leather      3473.5       770.5    4.508 7.53e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9768 on 802 degrees of freedom
## Multiple R-squared:  0.02471,    Adjusted R-squared:  0.02349
## F-statistic: 20.32 on 1 and 802 DF,  p-value: 7.526e-06
```

```
A7 <- lm(Price~Mileage, data = cars)
summary(A7)
```

```
##
## Call:
## lm(formula = Price ~ Mileage, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13905  -7254  -3520   5188  46091
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  2.476e+04  9.044e+02  27.383 < 2e-16 ***
## Mileage      -1.725e-01  4.215e-02  -4.093 4.68e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 9789 on 802 degrees of freedom
## Multiple R-squared:  0.02046,    Adjusted R-squared:  0.01924
## F-statistic: 16.75 on 1 and 802 DF,  p-value: 4.685e-05
```

```
#how to create a linear regression without intercept
A1 <- lm(Price~Cyl-1, data = cars)
summary(A1)
```

```
##
## Call:
## lm(formula = Price ~ Cyl - 1, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11209  -5232  -2745   2780  38346
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## Cyl  4051.17      52.62      77 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Residual standard error: 8128 on 803 degrees of freedom
## Multiple R-squared:  0.8807, Adjusted R-squared:  0.8806
## F-statistic: 5928 on 1 and 803 DF,  p-value: < 2.2e-16
```

We find that Cyl has the highest R^2 value among the rest, meaning for a single variable, Cyl captures the greatest amount of variation.

b) Calculate six regression models. Each model should have two explanatory variables, X_1 and one of the six explanatory variables. Find the two-variable model that has highest R^2 value. How much did R^2 improve when this variable was included?

```
B1 <- lm(Price~Cyl+Liter, data = cars)
summary(B1)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Liter, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10479   -5182   -2944    3034   39076
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1372.4     1434.5   0.957   0.339
## Cyl           2976.4       719.8   4.135 3.92e-05 ***
## Liter         1412.2       903.4   1.563   0.118
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8126 on 801 degrees of freedom
## Multiple R-squared:  0.3259, Adjusted R-squared:  0.3242
## F-statistic: 193.6 on 2 and 801 DF,  p-value: < 2.2e-16
```

```
B2 <- lm(Price~Cyl+Doors, data = cars)
summary(B2)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Doors, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -12093   -5565   -2888    3085   35847
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5713.3     1614.9   3.538 0.000426 ***
## Cyl           4056.4       204.0  19.888 < 2e-16 ***
## Doors        -1627.8       332.9  -4.890 1.22e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8019 on 801 degrees of freedom
## Multiple R-squared:  0.3435, Adjusted R-squared:  0.3418
## F-statistic: 209.5 on 2 and 801 DF,  p-value: < 2.2e-16
```

```
B3 <- lm(Price~Cyl+Cruise, data = cars)
summary(B3)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Cruise, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11724  -5695  -1961   3555  38661
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -1046.4     1082.7  -0.967   0.334
## Cyl           3392.6       211.3  16.058 <2e-16 ***
## Cruise       6000.4       678.8   8.839 <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7768 on 801 degrees of freedom
## Multiple R-squared:  0.3839, Adjusted R-squared:  0.3824
## F-statistic: 249.6 on 2 and 801 DF,  p-value: < 2.2e-16
```

```
B4 <- lm(Price~Cyl+Sound, data = cars)
summary(B4)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Sound, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -11946  -5429  -2607   2792  38970
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1293.7     1235.7   1.047  0.2955
## Cyl           4007.0       207.0  19.359 <2e-16 ***
## Sound        -1563.7       614.8  -2.543  0.0112 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8106 on 801 degrees of freedom
## Multiple R-squared:  0.3293, Adjusted R-squared:  0.3276
## F-statistic: 196.6 on 2 and 801 DF,  p-value: < 2.2e-16
```

```
B5 <- lm(Price~Cyl+Leather, data = cars)
summary(B5)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Leather, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -11748 -5318 -2838 3078 37807
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1528.9      1179.4  -1.296   0.195
## Cyl          3992.4       205.5  19.423 < 2e-16 ***
## Leather      2538.3       637.5   3.981 7.47e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8059 on 801 degrees of freedom
## Multiple R-squared:  0.337, Adjusted R-squared:  0.3353
## F-statistic: 203.6 on 2 and 801 DF, p-value: < 2.2e-16
B6 <- lm(Price~Cyl+Mileage, data = cars)
summary(B6)
```

```
##
## Call:
## lm(formula = Price ~ Cyl + Mileage, data = cars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -10264  -5121  -2838   3102   35477
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)
## (Intercept) 3145.75032 1325.93436   2.372   0.0179 *
## Cyl          4027.67463  204.61180  19.684 < 2e-16 ***
## Mileage      -0.15243    0.03464  -4.401 1.22e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 8042 on 801 degrees of freedom
## Multiple R-squared:  0.3398, Adjusted R-squared:  0.3382
## F-statistic: 206.2 on 2 and 801 DF, p-value: < 2.2e-16
```

From the summary above, we find that Cyl + Cruise has the highest R^2 value, meaning Cyl and Cruise captures the greatest amount of variation with two predictors. The R^2 value increased by 0.0609.

- c) Instead of continuing this process to identify more variables, use the software instructions provided to conduct a stepwise regression analysis. List each of the explanatory variables in the model suggested by the stepwise regression procedure.

```
library(leaps)
attach(cars)
X <- cbind(Cyl, Liter, Sound, Mileage, Cruise, Doors, Leather)
Best.model <- leaps(X, Price, method = "Cp")
cbind(Best.model$Cp, Best.model$which)
```

```
##           1 2 3 4 5 6 7
## 1 171.958749 1 0 0 0 0 0
## 1 189.686532 0 1 0 0 0 0
## 1 370.659915 0 0 0 0 1 0
## 1 601.986973 0 0 0 0 0 1
## 1 608.092654 0 0 0 1 0 0
```



```
## 1 609.834887 0 0 0 0 0 1 0
## 1 615.281723 0 0 1 0 0 0 0
## 2 87.578694 1 0 0 0 1 0 0
## 2 110.439808 0 1 0 0 1 0 0
## 2 145.781417 1 0 0 0 0 1 0
## 2 151.013703 1 0 0 1 0 0 0
## 2 155.097238 1 0 0 0 0 0 1
## 2 166.172861 1 0 1 0 0 0 0
## 2 166.384702 0 1 0 1 0 0 0
## 2 171.002543 1 1 0 0 0 0 0
## 2 174.648560 0 1 0 0 0 0 1
## 2 178.764912 0 1 0 0 0 1 0
## 3 61.038972 1 0 0 0 1 0 1
## 3 63.091938 1 0 0 1 1 0 0
## 3 66.219495 1 0 0 0 1 1 0
## 3 83.864809 0 1 0 1 1 0 0
## 3 84.753451 1 0 1 0 1 0 0
## 3 85.688900 0 1 0 0 1 0 1
## 3 89.416157 1 1 0 0 1 0 0
## 3 100.730101 0 1 0 0 1 1 0
## 3 105.210991 0 1 1 0 1 0 0
## 3 123.962780 1 0 0 1 0 1 0
## 4 36.150805 1 0 0 1 1 0 1
## 4 40.977639 1 0 0 1 1 1 0
## 4 42.974873 1 0 0 0 1 1 1
## 4 53.188852 1 0 1 0 1 0 1
## 4 58.747764 0 1 0 1 1 0 1
## 4 59.643905 1 0 1 1 1 0 0
## 4 61.859720 1 0 1 0 1 1 0
## 4 63.037481 1 1 0 0 1 0 1
## 4 64.790527 1 1 0 1 1 0 0
## 4 67.247794 1 1 0 0 1 1 0
## 5 17.403533 1 0 0 1 1 1 1
## 5 27.353986 1 0 1 1 1 0 1
## 5 33.463503 1 0 1 0 1 1 1
## 5 35.865078 1 0 1 1 1 1 0
## 5 38.117245 1 1 0 1 1 0 1
## 5 42.246805 1 1 0 1 1 1 0
## 5 43.423588 1 1 0 0 1 1 1
## 5 46.774957 0 1 1 1 1 0 1
## 5 50.367042 0 1 0 1 1 1 1
## 5 55.119836 1 1 1 0 1 0 1
## 6 6.824315 1 0 1 1 1 1 1
## 6 18.159331 1 1 0 1 1 1 1
## 6 29.176771 1 1 1 1 1 0 1
## 6 34.361320 1 1 1 0 1 1 1
## 6 36.814371 0 1 1 1 1 1 1
## 6 37.403986 1 1 1 1 1 1 0
## 6 97.352706 1 1 1 1 0 1 1
## 7 8.000000 1 1 1 1 1 1 1
```

```
Best.lm <- lm(Price~Cyl+Sound+Mileage+Cruise+Doors+Leather)
summary(Best.lm)
```

```
##
```

```
## Call:
## lm(formula = Price ~ Cyl + Sound + Mileage + Cruise + Doors +
##      Leather)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13104  -5566  -1544   3877   33349
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.323e+03  1.771e+03   4.135 3.92e-05 ***
## Cyl          3.200e+03  2.030e+02  15.765 < 2e-16 ***
## Sound       -2.024e+03  5.707e+02  -3.547 0.000412 ***
## Mileage     -1.705e-01  3.186e-02  -5.352 1.14e-07 ***
## Cruise       6.206e+03  6.515e+02   9.525 < 2e-16 ***
## Doors       -1.463e+03  3.083e+02  -4.747 2.45e-06 ***
## Leather      3.327e+03  5.971e+02   5.572 3.45e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7387 on 797 degrees of freedom
## Multiple R-squared:  0.4457, Adjusted R-squared:  0.4415
## F-statistic: 106.8 on 6 and 797 DF,  p-value: < 2.2e-16
Model_adj <- lm(Price~Cyl+Sound+Mileage+Cruise+Doors+Leather)
summary(Model_adj)

##
## Call:
## lm(formula = Price ~ Cyl + Sound + Mileage + Cruise + Doors +
##      Leather)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13104  -5566  -1544   3877   33349
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  7.323e+03  1.771e+03   4.135 3.92e-05 ***
## Cyl          3.200e+03  2.030e+02  15.765 < 2e-16 ***
## Sound       -2.024e+03  5.707e+02  -3.547 0.000412 ***
## Mileage     -1.705e-01  3.186e-02  -5.352 1.14e-07 ***
## Cruise       6.206e+03  6.515e+02   9.525 < 2e-16 ***
## Doors       -1.463e+03  3.083e+02  -4.747 2.45e-06 ***
## Leather      3.327e+03  5.971e+02   5.572 3.45e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 7387 on 797 degrees of freedom
## Multiple R-squared:  0.4457, Adjusted R-squared:  0.4415
## F-statistic: 106.8 on 6 and 797 DF,  p-value: < 2.2e-16
```

From the above, we find that our leaps function finds that out of our 7 variables, only six are needed. We came to this conclusion by taking the lowest value in the first column of the cbind function which happens to be where there are ones for every significant explanatory variable. The significant explanatory variables were:

Cyl, Sound, Mileage, Cruise, Doors, Leather.

```
##### Cross validation using model selection and subsets
Data_selection <- data.frame(Cyl, Liter, Sound, Mileage, Cruise, Doors, Leather, Price)
Reg_subset <- regsubsets(Price~., data = Data_selection)
summary(Reg_subset)
```

```
## Subset selection object
## Call: regsubsets.formula(Price ~ ., data = Data_selection)
## 7 Variables (and intercept)
##      Forced in Forced out
## Cyl      FALSE      FALSE
## Liter     FALSE     FALSE
## Sound     FALSE     FALSE
## Mileage   FALSE     FALSE
## Cruise    FALSE     FALSE
## Doors     FALSE     FALSE
## Leather   FALSE     FALSE
## 1 subsets of each size up to 7
## Selection Algorithm: exhaustive
##      Cyl Liter Sound Mileage Cruise Doors Leather
## 1 ( 1 ) "*" " " " " " " " " " "
## 2 ( 1 ) "*" " " " " " " "*" " " "
## 3 ( 1 ) "*" " " " " " " "*" " "*"
## 4 ( 1 ) "*" " " " " "*" "*" " " "*"
## 5 ( 1 ) "*" " " " " "*" "*" "*" "*"
## 6 ( 1 ) "*" " " "*" "*" "*" "*" "*"
## 7 ( 1 ) "*" "*" "*" "*" "*" "*" "*"

```

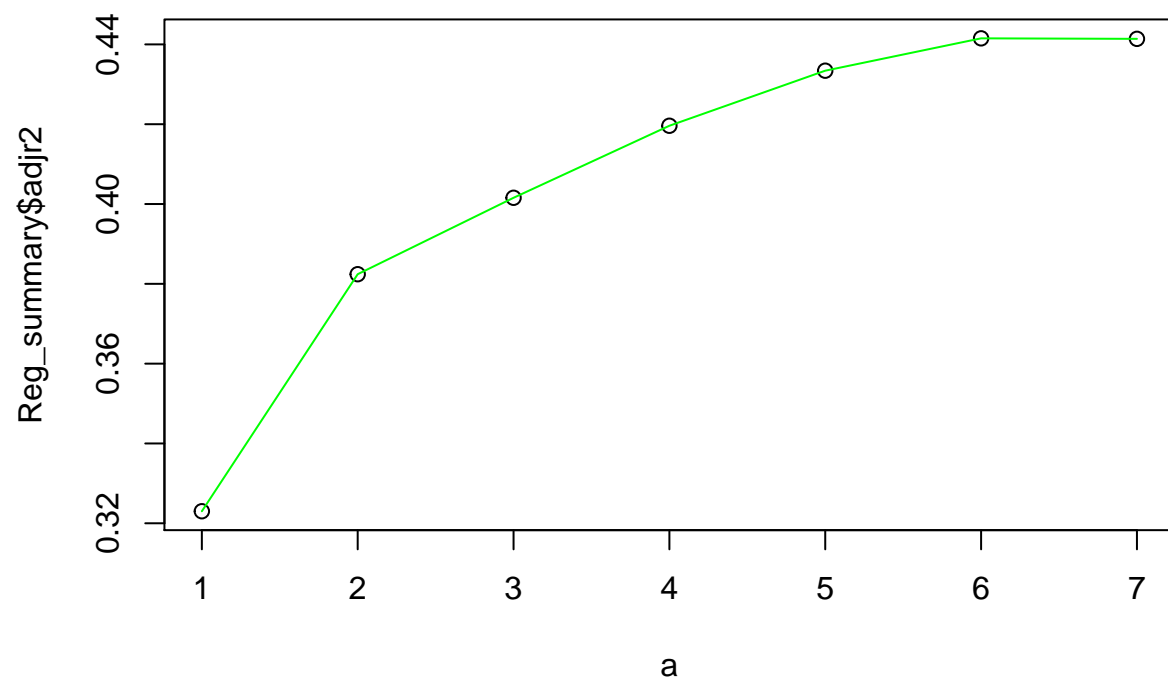
```
Reg_summary <- summary(Reg_subset)
Reg_summary$adjr2
```

```
## [1] 0.3230160 0.3824109 0.4015670 0.4196163 0.4334123 0.4415180 0.4413948
```

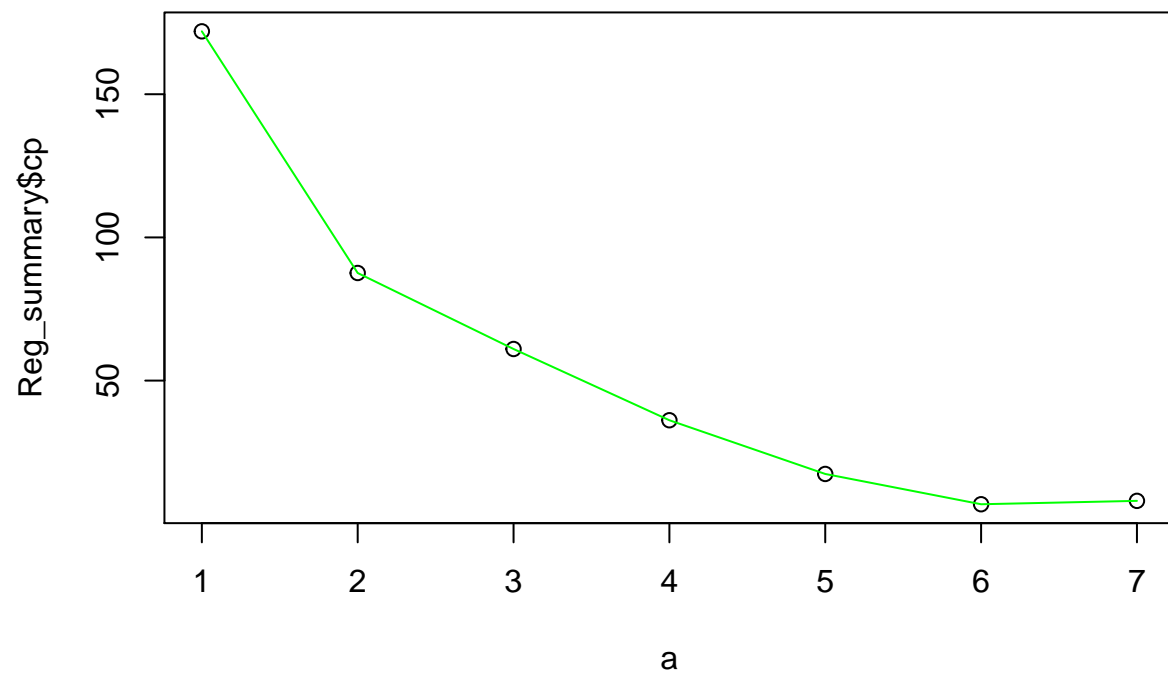
```
Reg_summary$cp
```

```
## [1] 171.958749 87.578694 61.038972 36.150805 17.403533 6.824315 8.000000
```

```
a<-c(1, 2, 3, 4, 5, 6, 7)
plot(Reg_summary$adjr2~a)
points(Reg_summary$adjr2, type = "l", col = "green")
```



```
plot(Reg_summary$cp~a)
points(Reg_summary$cp, type = "l", col = "green")
```



```
Reg_summary$rss
```

```
## [1] 53050954921 48336202986 46778462504 45310866304 44178449357 43491856868  
## [7] 43446864552
```