

Multivariate statistics – Regression and regression diagnostics

Suckwon Hong

Ulsan National Institute of Science and Technology, 50 UNIST-gil Ulsan, Republic of Korea

amoeba94@unist.ac.kr

Question 1. Dealing with missing values and outlier

1.1. Missing values

A variable *Horsepower* contains six missing values. We assigned mean value of the variable to the missing ones.

1.2. Outliers

We employed two approaches to handle outliers existing in continuous variables (i.e., *mpg*, *cylinders*, *displacement*, *horsepower*, *weight*, *acceleration*, *model.year*). First, for each continuous variable, we replace data points above ($75\% \text{ quantile} + 1.5 * \text{interquartile}$) or ($25\% \text{ quantile} - 1.5 * \text{interquartile}$) with the value of ($70\% \text{ quantile}$) and ($30\% \text{ quantile}$), respectively. Second, we made simple linear regression model, and calculated Cook's distance. According to the rule of thumb, we checked whether there exist data points with the value above 0.03 . As a result, we eliminated four data points from the original.

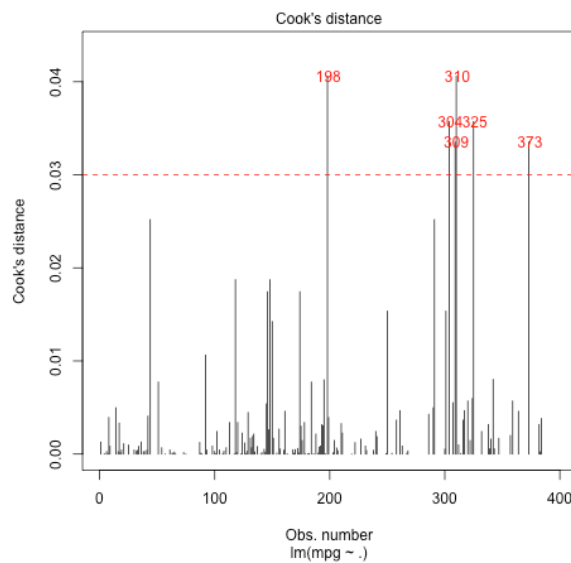


Figure 1. Cook's distance

Question 2. Testing assumptions for a linear regression

2.1. Normality assumption

Three approaches are employed to test the normality assumption: histogram, qqplot, Shapiro-Wilk test. Histogram and qqplot for residuals show that residuals generally follow a normal distribution. After, we conducted Shapiro-Wilk normality test. However, the result showed that residuals do not follow a normal distribution ($p\text{-value} < 2.2e-16$).

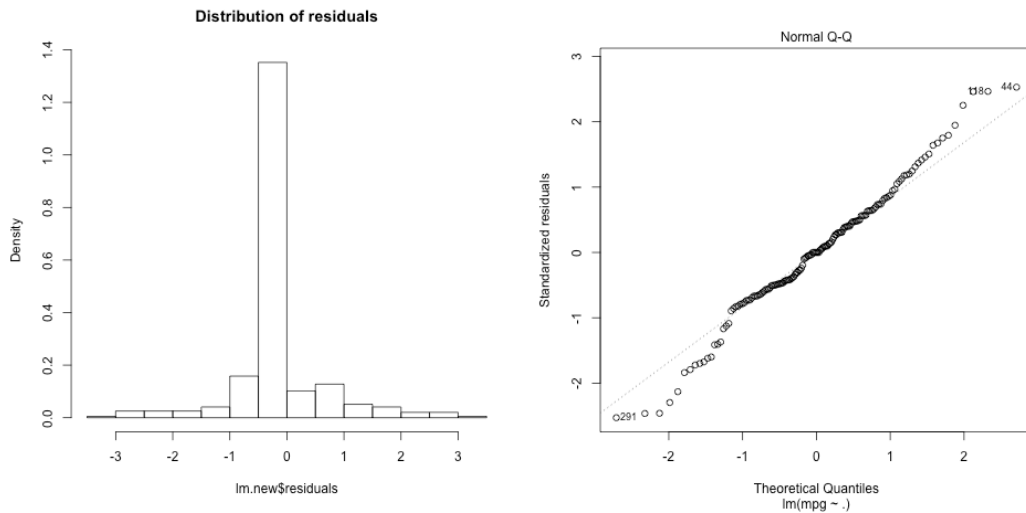
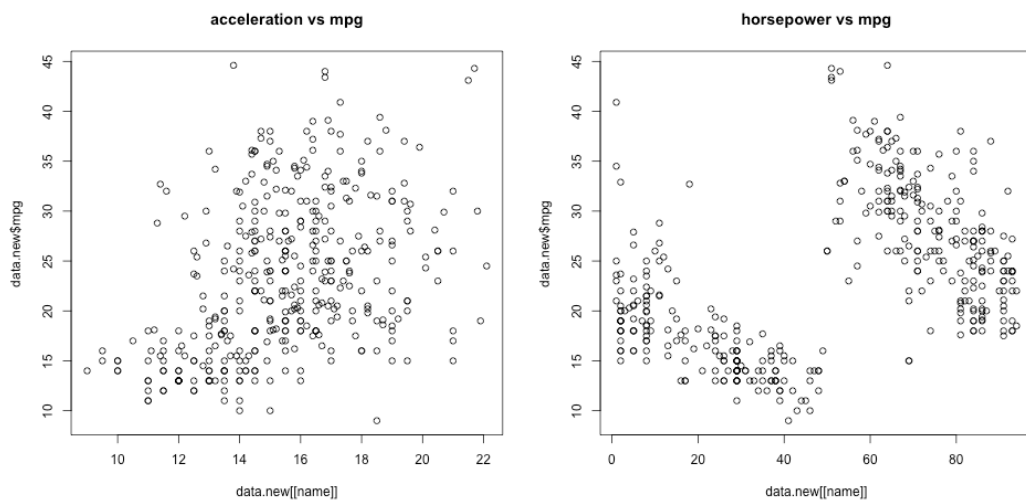


Figure 2. Histogram and qqplot for continuous variables

2.2. Linearity assumption

To test a linearity assumption, we drew scatterplots between dependent variable and different independent variables. It seems that the variable *horsepower* does not meet an assumption, so we omit the variable in running a regression.



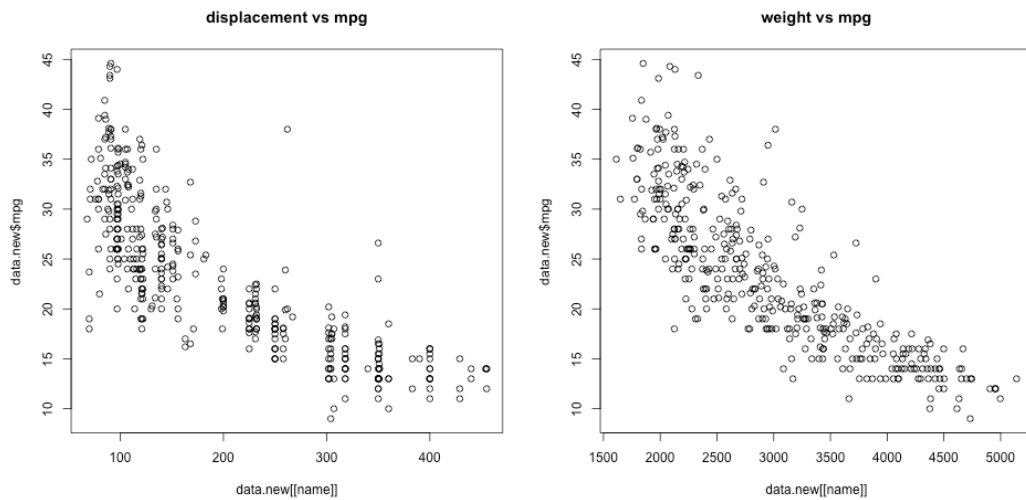


Figure 3. Scatterplots vs dependent variable

2.3. Homoscedasticity assumption

We drew residual plot for the last assumption (i.e. homoscedasticity). Since we observed no pattern from the plot, the data seem to meet the assumption.

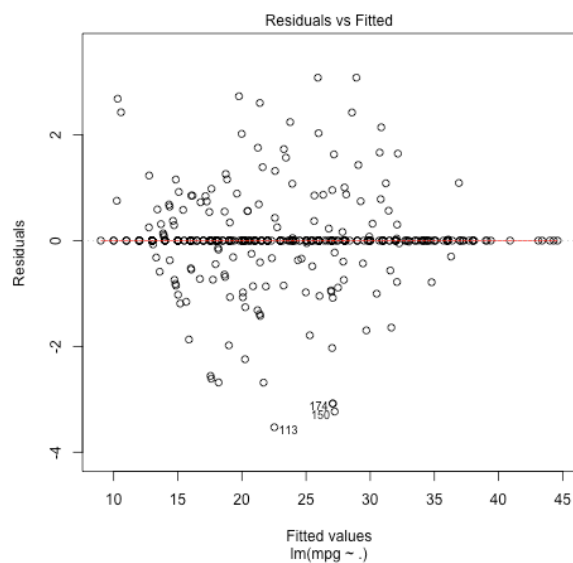


Figure 4. Residuals vs fitted plot

Question 3. Evaluate multicollinearity

We computed correlation between continuous variables and variance inflation factor (VIF) to check an existence of multicollinearity. Since we found that variables *weight* and *displacement* are highly correlated from the correlation matrix, we computed VIF. The VIF of the model was 93.3245, and it is found that multiple values of the variable *car.name* are highly correlated. Thus, we omit the variable.

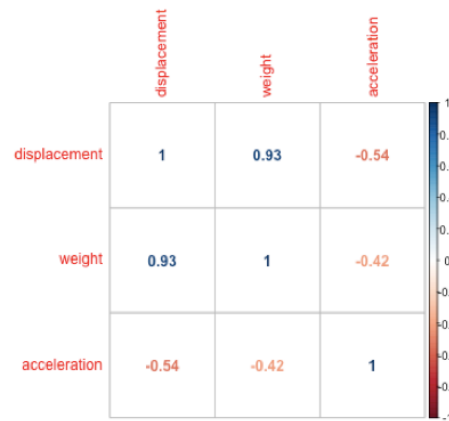


Figure 5. Correlation matrix

Question 4. Regression results

As a result, we run regression on *mpg* using variables of *cylinders*, *displacement*, *weight*, *acceleration*, *model.year*, *origin*. We tested the assumptions again, and it seems that the data meet the assumptions. Additionally, the value of VIF was 7.5330 implying that there is no multicollinearity problem.

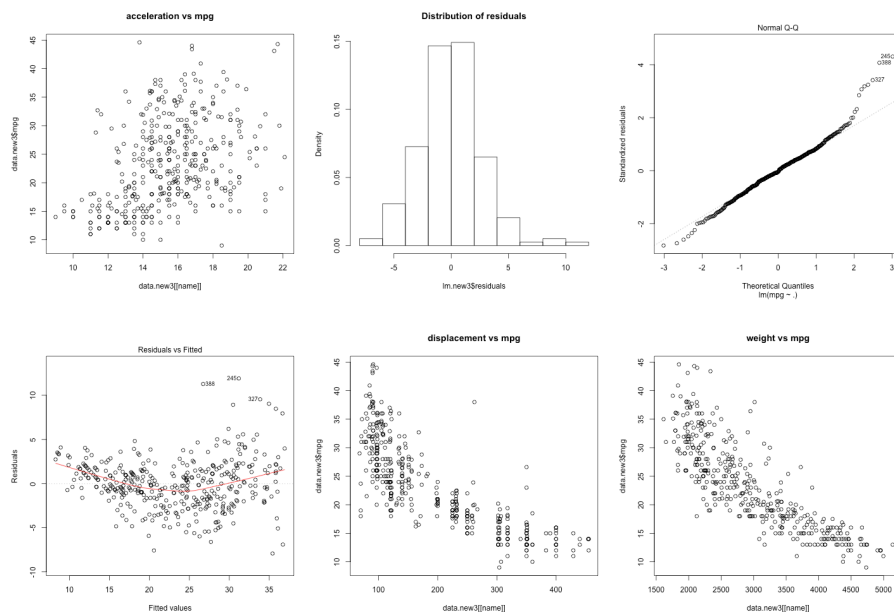


Figure 6. Regression diagnostics

The value of R-squared of our final regression model was 0.8673, which means that the independent variables explained about 87% of total variance of dependent variable *mpg*. Furthermore, we concluded that the model is significantly significant because the p-value of the model was small enough. Most values of categorical variables (i.e., *cylinders*, *model.year*, *origin*) were shown to be statistically significant while only one continuous variable (i.e., *weight*) was significant. From the coefficients, we can say that for every additional unit of *weight*, we can expect *mpg* to decrease by average of 0.006.

```
Call:
lm(formula = mpg ~ ., data = data.new3)

Residuals:
    Min       1Q   Median       3Q      Max
-7.9205 -1.6678 -0.1363  1.5740 11.9114

Coefficients:
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  28.4268767   2.1001197   13.536 < 2e-16 ***
cylinders4    6.8947847   1.5480193    4.454 1.12e-05 ***
cylinders5    7.4197881   2.3561718    3.149 0.001771 **
cylinders6    4.8748593   1.7074968    2.855 0.004547 **
cylinders8    6.7294035   1.9744644    3.408 0.000726 ***
displacement  0.0059664   0.0065765    0.907 0.364873
weight       -0.0061118   0.0005379  -11.363 < 2e-16 ***
acceleration  0.0913918   0.0775128    1.179 0.239134
model.year71  1.6854487   0.7810048    2.158 0.031567 *
model.year72  0.1737705   0.7921370    0.219 0.826483
model.year73 -0.0163790   0.7112849   -0.023 0.981641
model.year74  2.1902900   0.8146049    2.689 0.007496 **
model.year75  1.8900314   0.7925717    2.385 0.017597 *
model.year76  2.4481421   0.7708549    3.176 0.001619 **
model.year77  3.8036186   0.7969196    4.773 2.62e-06 ***
model.year78  3.6821787   0.7618101    4.833 1.97e-06 ***
model.year79  5.7138255   0.8000996    7.141 4.93e-12 ***
model.year80  8.9162413   0.8367824   10.655 < 2e-16 ***
model.year81  7.5835506   0.8162609    9.291 < 2e-16 ***
model.year82  8.7312281   0.8243937   10.591 < 2e-16 ***
origin2       1.8148616   0.5158109    3.518 0.000488 ***
origin3       1.9194970   0.4894114    3.922 0.000105 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.865 on 370 degrees of freedom
Multiple R-squared:  0.8673,    Adjusted R-squared:  0.8597
F-statistic: 115.1 on 21 and 370 DF,  p-value: < 2.2e-16
```

Figure 7. Regression results

R Code.

```
1  ### Multivariate statistics - assignment 2
2  rm(list = ls())
3  ## install packages
4  install.packages("car")
5  install.packages("dplyr")
6  install.packages("corrplot")
7  install.packages("fmsb")
8  library(car)
9  library(dplyr)
10 library(corrplot)
11 library(fmsb)
12 ## load data
13 data <- read.csv("mpg.csv", header = TRUE)
14 ## data manipulation
15 data$origin <- as.factor(data$origin)
16 data[["model.year"]] <- as.factor(data[["model.year"]])
17 data[["cylinders"]] <- as.factor(data[["cylinders"]])
18 data[["horsepower"]] <- as.numeric(data[["horsepower"]])
19 str(data)
20 ## handling missing data
21 data[data=="?"] <- NA
22 na_count <- sapply(data, function(x) sum(length(which(is.na(x)))))
23 data[is.na(data)] <- mean(data$horsepower, na.rm = TRUE) # mean implementation
24
25 ## evaluate outliers
26 # using IQR
27 for(name in names(data)){
28   if((name == "car.name") || (name == "origin") || (name == "cylinders") || (name == "model.year")){
29     }
30   else{
31     x <- data[[name]]
32     qnt <- quantile(x, probs=c(.25, .75), na.rm = T)
33     caps <- quantile(x, probs=c(.3, .7), na.rm = T)
34     H <- 1.5 * IQR(x, na.rm = T)
35     x[x < (qnt[1] - H)] <- caps[1]
36     x[x > (qnt[2] + H)] <- caps[2]
37     data[[name]] <- x
38   }
39 }
40
41 # using cook's distance
42 lm.pre <- lm(mpg~., data = data)
43 cooks.d <- cooks.distance(lm.pre)
44 png(filename = "cookd.png")
45 plot(lm.pre, which=4, labels.id = "")
46 abline(h=0.03, lty=2, col="red")
47 text(x=1:length(cooks.d)+1, y=cooks.d, labels=ifelse(cooks.d>0.03, names(cooks.d), ""), col="red")
48 dev.off()
49
50 data$cooks.d <- cooks.d
51 data[data=="NaN"] <- 0
52 data$cooky <- ifelse(data$cooks.d<0.03, "keep", "no")
53 data.new <- subset(data, cooky=="keep")
54 data.new <- data.new[-c(10,11)]
55 lm.new <- lm(mpg~., data = data.new)
56
57 ## test regression assumptions
58 #normality using histogram
59 png(filename = "histogram.png")
60 hist(lm.new$residuals, main = "Distribution of residuals", freq = FALSE)
61 dev.off()
62 #normality using qqplot
63 png(filename = "qqplot.png")
64 plot(lm.new, which = 2)
65 dev.off()
66 #normality using Shapiro-Wilk W test
67 print(shapiro.test(lm.new$residuals))
68 #linearity using scatterplot
69 for(name in names(data.new)){
70   if((name == "car.name") || (name == "origin") || (name == "cylinders") || (name == 'mpg') || (name == "model.year")){
71     }
72   else{
73     png(filename = paste("scatterplot", name, ".png"))
74     plot(data.new[[name]], data.new$mpg, main=paste(name, "vs mpg"))
75     dev.off()
76   }
77 }
78 data.new2 <- subset(data.new, select = -c(horsepower))
79 lm.new2 <- lm(mpg~., data = data.new2)
```

```

79 ##Homoscedasticity using residual plot
80 png(filename = "residual vs fitted.png")
81 plot(lm.new2,which=1)
82 dev.off()
83
84 ## evaluate multicollinearity
85 ##correlation matrix
86 cor<-cor(select(data.new2, displacement, weight, acceleration)) #only continuous
87 cor <- as.matrix(cor)
88 png(filename = "correlation matrix.png")
89 corrplot(as.matrix(cor), method = "number")
90 dev.off()
91
92 fmsb::VIF(lm.new2)
93 ld.vars <- attributes(alias(lm.new2)$Complete)
94 data.new3 <- subset(data.new2,select = -c(car.name))
95
96 ## run final regression
97 lm.new3 <- lm(mpg~., data = data.new3)
98 # test assumptions again
99 png(filename = "histogram.png")
100 hist(lm.new3$residuals, main = "Distribution of residuals", freq = FALSE)
101 dev.off()
102 png(filename = "qqplot.png")
103 plot(lm.new3, which = 2)
104 dev.off()
105 print(shapiro.test(lm.new3$residuals))
106
107 ~ for(name in names(data.new3)){
108 ~   if((name == "car.name") || (name=="origin")||(name=="cylinders")||(name=="mpg")||(name=="model.year")){
109 ~     }
110 ~   else{
111     png(filename = paste("scatterplot",name,".png"))
112     plot(data.new3[[name]],data.new3$mpg, main=paste(name,"vs mpg"))
113     dev.off()
114   }
115 }
116 png(filename = "residual vs fitted.png")
117 plot(lm.new3,which=1)
118 dev.off()
119 fmsb::VIF(lm.new3)
120 summary(lm.new3)

```