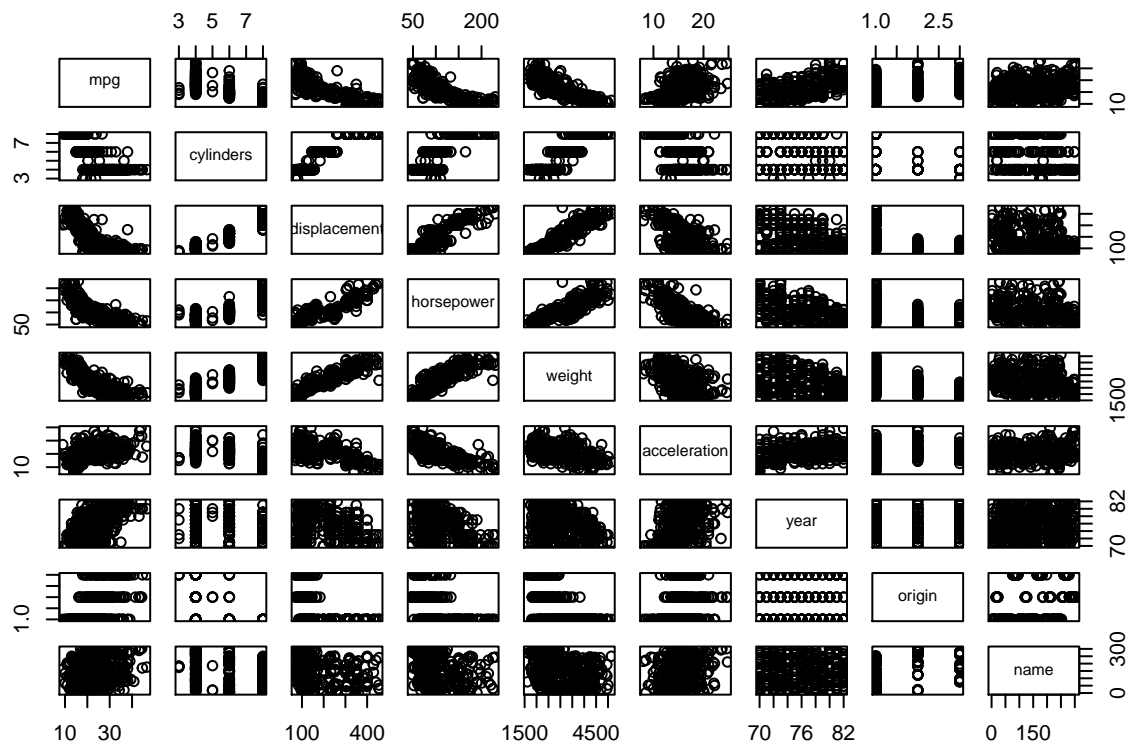# MA685 HW1#9

*Jiayuan Shi*

*01/25/2016*

```r
library(ISLR)
```

## 9a.

```r
pairs(Auto)
```



## 9b.

```r
cor(subset(Auto, select=-name))
```

```
##                     mpg  cylinders displacement horsepower     weight
## mpg           1.0000000 -0.7776175   -0.8051269 -0.7784268 -0.8322442
## cylinders    -0.7776175  1.0000000    0.9508233  0.8429834  0.8975273
## displacement -0.8051269  0.9508233    1.0000000  0.8972570  0.9329944
## horsepower   -0.7784268  0.8429834    0.8972570  1.0000000  0.8645377
## weight       -0.8322442  0.8975273    0.9329944  0.8645377  1.0000000
## acceleration  0.4233285 -0.5046834   -0.5438005 -0.6891955 -0.4168392
```

```
## year           0.5805410 -0.3456474   -0.3698552 -0.4163615 -0.3091199
## origin         0.5652088 -0.5689316   -0.6145351 -0.4551715 -0.5850054
##              acceleration      year     origin
## mpg             0.4233285  0.5805410  0.5652088
## cylinders      -0.5046834 -0.3456474 -0.5689316
## displacement   -0.5438005 -0.3698552 -0.6145351
## horsepower     -0.6891955 -0.4163615 -0.4551715
## weight         -0.4168392 -0.3091199 -0.5850054
## acceleration    1.0000000  0.2903161  0.2127458
## year            0.2903161  1.0000000  0.1815277
## origin          0.2127458  0.1815277  1.0000000
```
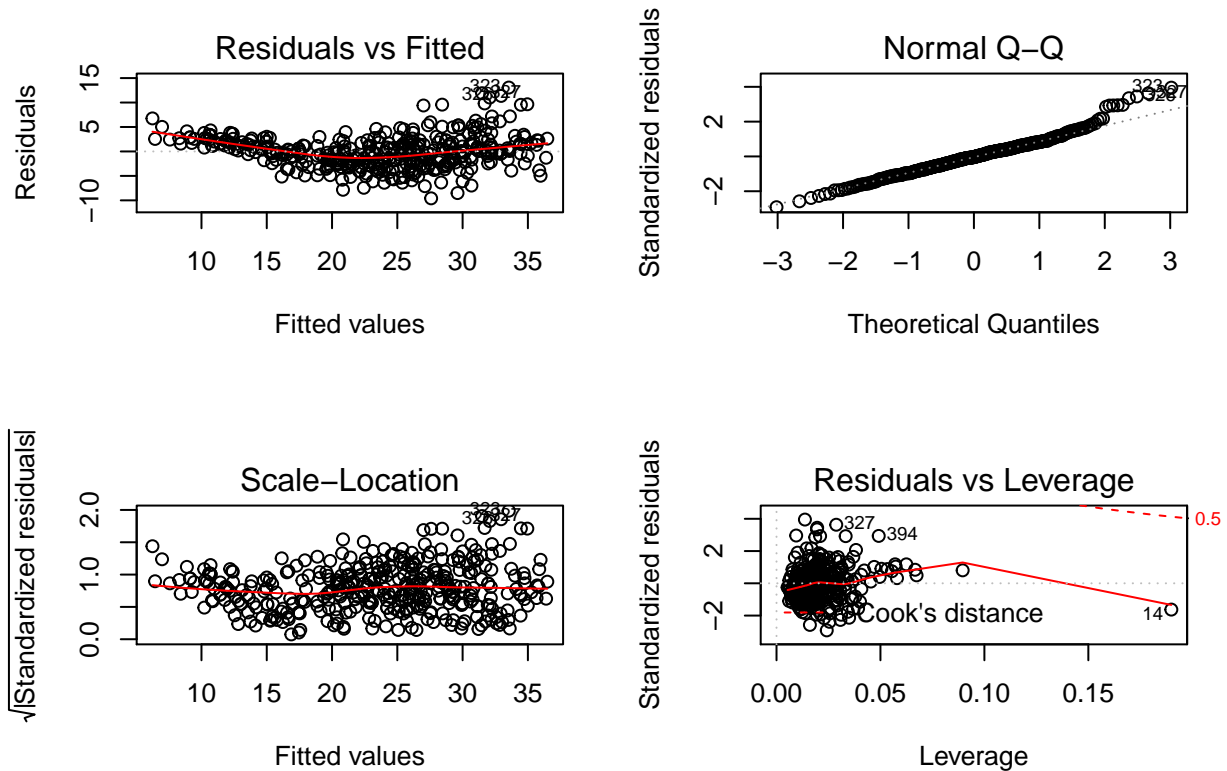
## 9c.

```
lm1 = lm(mpg~.-name, data=Auto)
summary(lm1)
```

```
##
## Call:
## lm(formula = mpg ~ . - name, data = Auto)
##
## Residuals:
##     Min      1Q  Median      3Q     Max
## -9.5903 -2.1565 -0.1169  1.8690 13.0604
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -17.218435   4.644294  -3.707  0.00024 ***
## cylinders     -0.493376   0.323282  -1.526  0.12780
## displacement   0.019896   0.007515   2.647  0.00844 **
## horsepower    -0.016951   0.013787  -1.230  0.21963
## weight        -0.006474   0.000652  -9.929  < 2e-16 ***
## acceleration   0.080576   0.098845   0.815  0.41548
## year           0.750773   0.050973  14.729  < 2e-16 ***
## origin         1.426141   0.278136   5.127 4.67e-07 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 3.328 on 384 degrees of freedom
## Multiple R-squared:  0.8215, Adjusted R-squared:  0.8182
## F-statistic: 252.4 on 7 and 384 DF,  p-value: < 2.2e-16
```

i. Yes, there is a relatioship between the predictors and the response by testing the null hypothesis of whether all the regression coefficients are zero. The F -statistic is has a small p-value, less than 0.01, indicating that we can reject the null hypothesis. Therefore, there is a significant relatioship between the predictors and the response.

ii. From the p-values of each predictor, displacement, weight, year, and origin have a statistically significant relationship to the response(since they have small p-values $< 0.01$), while cylinders, horsepower, and acceleration do not.

iii. The coefficient for year, 0.7508, suggests that for every one year increase, mpg also increases by 0.7508.
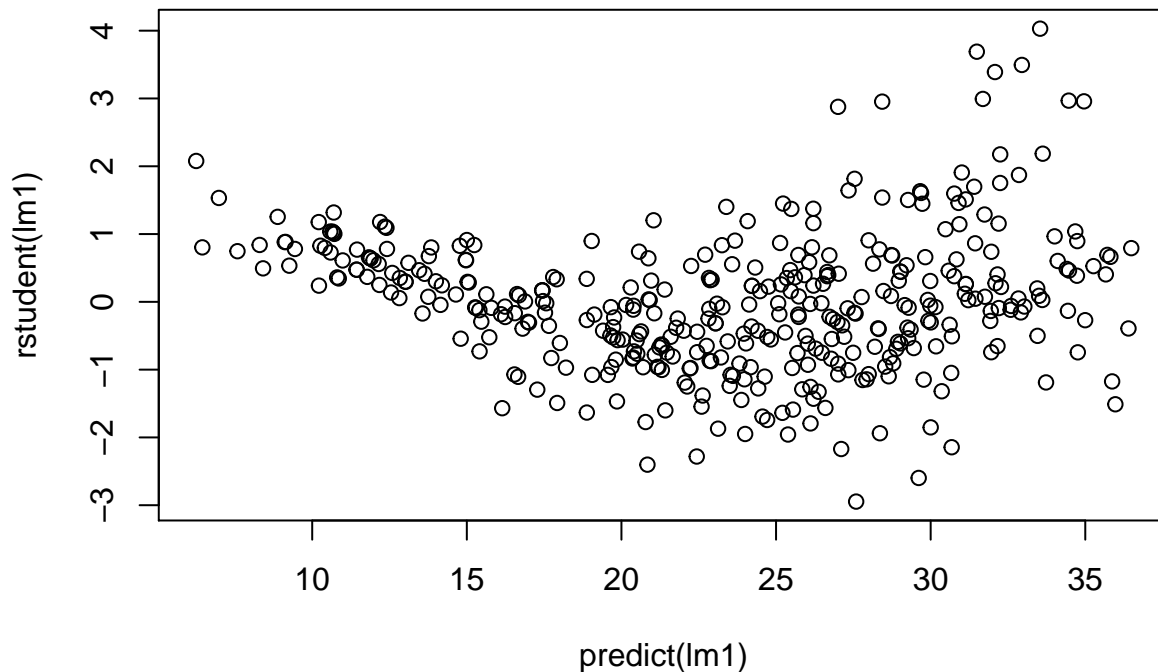
## 9d.

```
par(mfrow=c(2,2))
plot(lm1)
```



There are three problems with the fit. One is that the curve pattern in the residuals plot indicates non-linearity in the data. The other is from the leverage plot, point 14 appears to have high leverage, although not a high magnitude residual.

```
plot(predict(lm1), rstudent(lm1))
```

## 9d.

```
par(mfrow=c(2,2))
plot(lm1)
```



There are three problems with the fit. One is that the curve pattern in the residuals plot indicates non-linearity in the data. The other is from the leverage plot, point 14 appears to have high leverage, although not a high magnitude residual.

```
plot(predict(lm1), rstudent(lm1))
```

Also there may be some outliers in the plot of studentized residuals because there exist some values greater than 3.

## 9e.

```
lm2 = lm(mpg~cylinders*displacement+displacement*weight, data=Auto)
summary(lm2)
```

```
##
## Call:
## lm(formula = mpg ~ cylinders * displacement + displacement *
##     weight, data = Auto)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -13.2934  -2.5184  -0.3476   1.8399  17.7723
##
## Coefficients:
##                         Estimate Std. Error t value Pr(>|t|)
## (Intercept)            5.262e+01  2.237e+00  23.519  < 2e-16 ***
## cylinders              7.606e-01  7.669e-01   0.992    0.322
## displacement          -7.351e-02  1.669e-02  -4.403 1.38e-05 ***
## weight                -9.888e-03  1.329e-03  -7.438 6.69e-13 ***
## cylinders:displacement -2.986e-03  3.426e-03  -0.872    0.384
## displacement:weight    2.128e-05  5.002e-06   4.254 2.64e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.103 on 386 degrees of freedom
```

4

```
## Multiple R-squared:  0.7272, Adjusted R-squared:  0.7237
## F-statistic: 205.8 on 5 and 386 DF,  p-value: < 2.2e-16
```
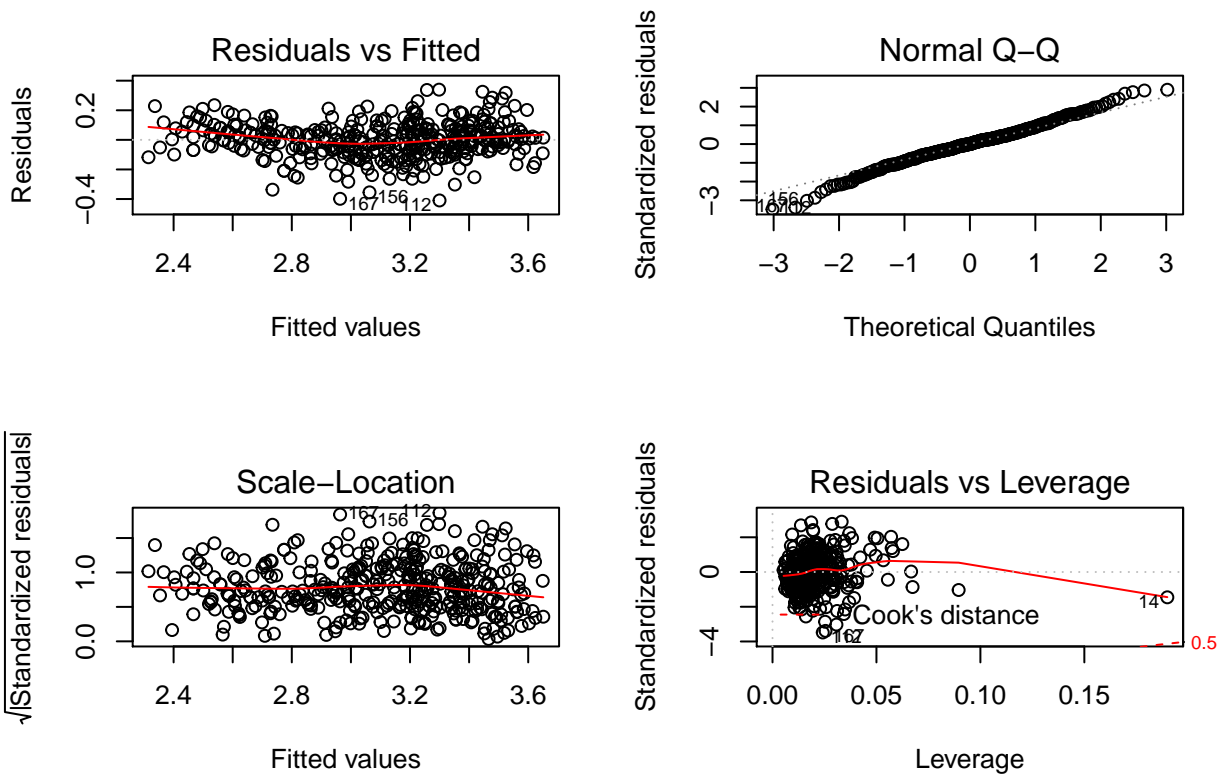
From the correlation matrix, I got the two highest correlated pairs($>0.9$) and used them in picking my interaction effects. From the p-values, we can see that the interaction between displacement and weight is statistically signifcant with p-value $< 0.01$, while the interactiion between cylinders and displacement is not.
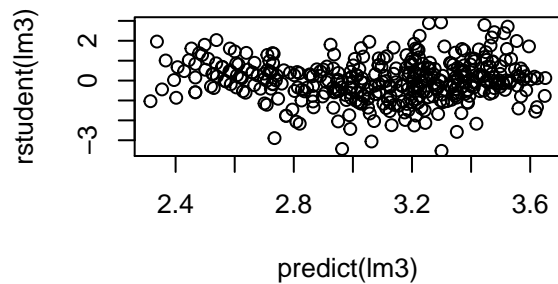
## 9f.

```
lm3 = lm(log(mpg)~.-name, data=Auto)
summary(lm3)
```

```
##
## Call:
## lm(formula = log(mpg) ~ . - name, data = Auto)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.40955 -0.06533  0.00079  0.06785  0.33925
##
## Coefficients:
##                Estimate Std. Error t value Pr(>|t|)
## (Intercept)   1.751e+00  1.662e-01  10.533  < 2e-16 ***
## cylinders    -2.795e-02  1.157e-02  -2.415  0.01619 *
## displacement  6.362e-04  2.690e-04   2.365  0.01852 *
## horsepower   -1.475e-03  4.935e-04  -2.989  0.00298 **
## weight       -2.551e-04  2.334e-05 -10.931  < 2e-16 ***
## acceleration -1.348e-03  3.538e-03  -0.381  0.70339
## year          2.958e-02  1.824e-03  16.211  < 2e-16 ***
## origin        4.071e-02  9.955e-03   4.089 5.28e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1191 on 384 degrees of freedom
## Multiple R-squared:  0.8795, Adjusted R-squared:  0.8773
## F-statistic: 400.4 on 7 and 384 DF,  p-value: < 2.2e-16
```

```
par(mfrow=c(2,2))
plot(lm3)
```

```
plot(predict(lm3), rstudent(lm3))
```



From the correlation matrix in 9a., displacement, horsepower and weight show a similar nonlinear pattern against the response mpg. This nonlinear pattern is very close to a log form. So we use log(mpg) as our response variable.

The outputs show that log transform of mpg yield better model fitting. For example, R^2 increases from 0.8215 to 0.8795. Residuals tend to be more normal. The curve pattern in the residuals plot indicates more linearity in the data. Also, there are less outliers in the plot of studentized residuals because all values are less than 3.