

# Peer graded assignment regression models

*yinshu zhang*

*October 24, 2018*

## Introduction

This study is conducted by Motor Trend magazine to explore the relationship between a set of variables and the fuel consumption. The focus of this paper is to address two questions.

“Is an automatic or manual transmission better for MPG”

“Quantify the MPG difference between automatic and manual transmissions”

## Data Exploration

We will use “mtcars” data to draw conclusions, mtcars is a 32 by 11 data frame, with columns defined as:

column name	definition
mpg	miles per gallon(US)
cyl	number of cylinders
displacement	engine displacement(cu. in.)
hp	Gross horsepower
drat	Rear axle ratio
wt	Weight (1000 lbs)
qsec	zero to 1/4 mile time in seconds
vs	Engine shape (0 = V-shaped, 1 = straight)
am	Transmission (0 = automatic, 1 = manual)
gear	Number of forward gears
carb	Number of carburetors

A quick look at each variables

## Relationship between MPG and AM

### Correlation test

From summary figure, we can see the MPG and AM variables has 0.5998 correlation value, it is an indication MPG and AM has relationship, but not strong enough. Therefore, to answer first question, we need better evidences.

We will test the correlation between MPG and AM with alternative hypothesis “true correlation is greater than zero”

```
cor.test(mpg, am, alternative = "greater")
```

```
##  
## Pearson's product-moment correlation  
##
```

```
## data: mpg and am
## t = 4.1061, df = 30, p-value = 0.0001425
## alternative hypothesis: true correlation is greater than 0
## 95 percent confidence interval:
## 0.3691544 1.0000000
## sample estimates:
## cor
## 0.5998324
```

From the output, p-value is significantly away from 95% confidence interval, this is a good evidence *supporting* MPG and AM are correlated.

## One way ANOVA test

To further prove the opposite is false, we can use one-way ANOVA. The null hypothesis of ANOVA is designed as the mean MPG are the same between auto and manual transmission cars, ie. either those two samples are from either from same population (or different populations with same mean), or when two sample show difference, it is due to chance.

```
anova1 <- aov(mpg ~ am, data = mtcars)
summary(anova1)

##           Df Sum Sq Mean Sq F value    Pr(>F)
## am           1  405.2   405.2    16.86 0.000285 ***
## Residuals    30  720.9    24.0
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

f_crit <- round(qf(.95, 1, 30), 2)
```

F value of 16.86 is significantly larger than the 95% confidence interval critical value 4.17, the result shows strong evident *rejecting* null hypothesis.

## Regression

### simple regression

For the sake of comparison, let's do a single variable (factor actually) linear model.

```
mod.sim <- lm(mpg ~ am, data = mtcars)
summary(mod.sim)

##
## Call:
## lm(formula = mpg ~ am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.3923 -3.0923 -0.2974  3.2439  9.5077
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   17.147      1.125   15.247 1.13e-15 ***
## am              7.245      1.764    4.106 0.000285 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.902 on 30 degrees of freedom
## Multiple R-squared:  0.3598, Adjusted R-squared:  0.3385
## F-statistic: 16.86 on 1 and 30 DF,  p-value: 0.000285
```

To interpret the output, the coefficient of AM is 7.245, which means changing from auto to manual, will result on average 7.245 mile/gallon fuel consumption increase. However, as expected, the model quality is low, we can't use it to predict the MPG, when only around 36% of residuals variance are explained by this model, AM is a dummy variable after all.

## selection of independent variables

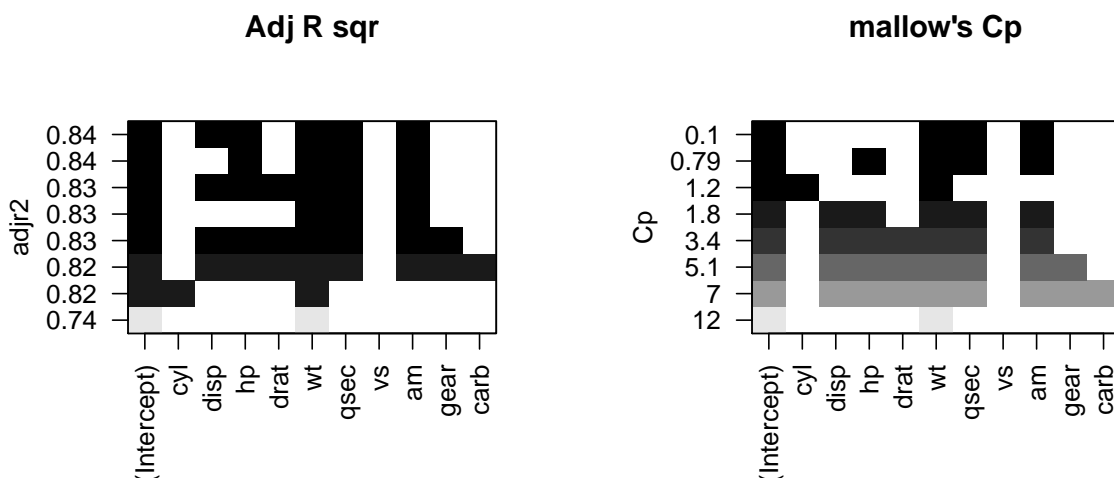
To make a multi-variable linear regression model, there are ten variables to choose from. we need to avoid those highly correlated variables, take a quick look at first summary plot, we can see quite some variables are highly related.

This is not a big surprise, with basic car knowledge we can understand why some of independent variables are correlated. For example, more cylinders of an engine normally means bigger displacement, higher the horsepower normally results lower quarter mile time.

We can categorize those ten variables into following group

- 1, drive train specs: cyl, disp, drat, vs, carb, drat, am, gear
- 2, other vehicle specs: wt
- 3, performance measurements: qsec, hp

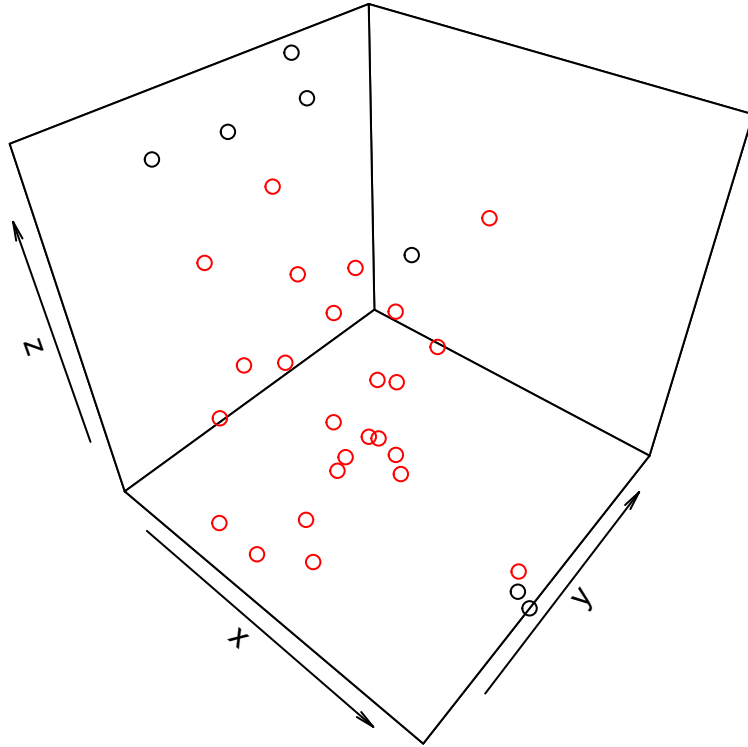
we can choose ourselves and compare the results but better option is to get help from package "leaps".



Here we did an exhausted subset comparison, with one best result from each number of variable combination. from the right plot, we can see the adjusted R square does not show big difference after 0.7, the Mallows' Cp however has one very low value compared to other combinations. Therefore, we will create a regression model with independent variables of weight, quarter mile time, and auto/manual transmission.

```
##
```

```
## Call:
## lm(formula = mpg ~ wt + qsec + am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4811 -1.5555 -0.7257  1.4110  4.6610
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.6178     6.9596   1.382 0.177915
## wt          -3.9165     0.7112  -5.507 6.95e-06 ***
## qsec         1.2259     0.2887   4.247 0.000216 ***
## am           2.9358     1.4109   2.081 0.046716 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.459 on 28 degrees of freedom
## Multiple R-squared:  0.8497, Adjusted R-squared:  0.8336
## F-statistic: 52.75 on 3 and 28 DF,  p-value: 1.21e-11
## Warning in `[<-.factor`(`*tmp*`, is.na(Col), value = "white"): invalid
## factor level, NA generated
```



## Executive Summary