

# RegMods Assignment

## Summary

Manual transmission better for MPG (miles per gallon) than automatic transmission. We did a multivariate regression to improve estimate of transmission types on MPG.

## Data Processing

```
data(mtcars)

mtcars$cyl <- as.factor(mtcars$cyl)
mtcars$vs <- as.factor(mtcars$vs)
mtcars$gear <- factor(mtcars$gear)
mtcars$carb <- factor(mtcars$carb)
mtcars$am <- factor(mtcars$am, labels=c('Automatic', 'Manual'))
str(mtcars)

## 'data.frame': 32 obs. of 11 variables:
## $ mpg : num 21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
## $ cyl : Factor w/ 3 levels "4","6","8": 2 2 1 2 3 2 3 1 1 2 ...
## $ disp: num 160 160 108 258 360 ...
## $ hp : num 110 110 93 110 175 105 245 62 95 123 ...
## $ drat: num 3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 3.92 3.92 ...
## $ wt : num 2.62 2.88 2.32 3.21 3.44 ...
## $ qsec: num 16.5 17 18.6 19.4 17 ...
## $ vs : Factor w/ 2 levels "0","1": 1 1 2 2 1 2 1 2 2 2 ...
## $ am : Factor w/ 2 levels "Automatic","Manual": 2 2 2 1 1 1 1 1 1 1 ...
## $ gear: Factor w/ 3 levels "3","4","5": 2 2 2 1 1 1 1 2 2 2 ...
## $ carb: Factor w/ 6 levels "1","2","3","4",...: 4 4 1 1 2 1 4 2 2 4 ...
```

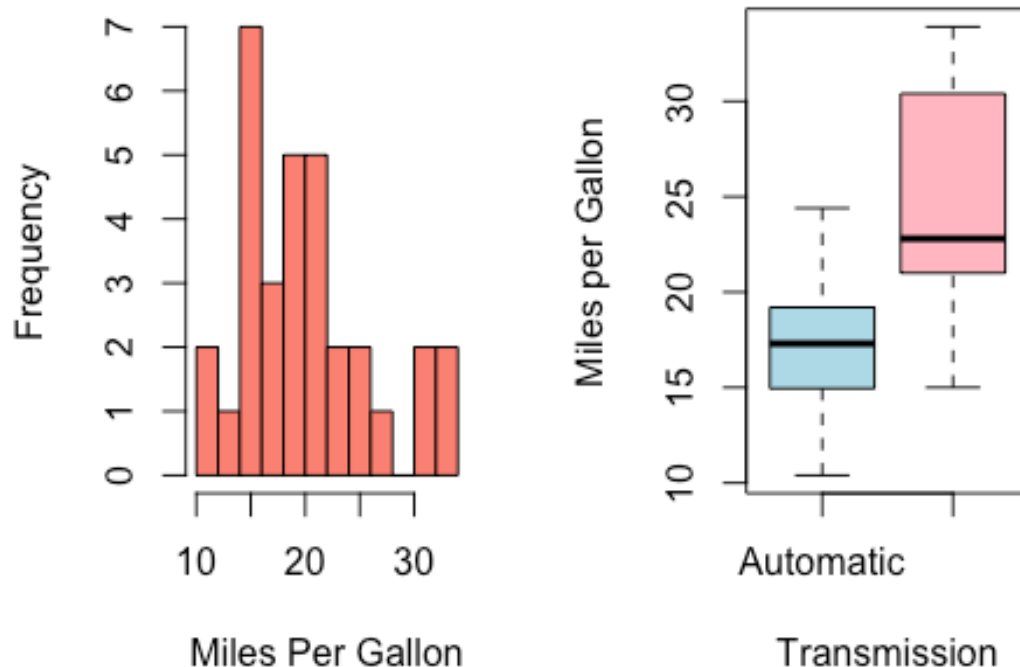
## EDA

Do a boxplot to examine car transmission types on mpg. We can say there is increase in mpg for manual transmission vs automatic transmission. Also, plot histogram to check for normal curve.

```
par(mfrow = c(1, 2))
# Histogram with Normal Curve
x <- mtcars$mpg
h <- hist(x, breaks=10, col="salmon", xlab="Miles Per Gallon",
          main="Histogram of Miles per Gallon")

boxplot(mpg~am, data = mtcars,
        col = c("light blue", "light pink"),
        xlab = "Transmission",
        ylab = "Miles per Gallon",
        main = "MPG by Transmission Type")
```

## Histogram of Miles per Ga MPG by Transmission Ty



## Hypotheses Testing

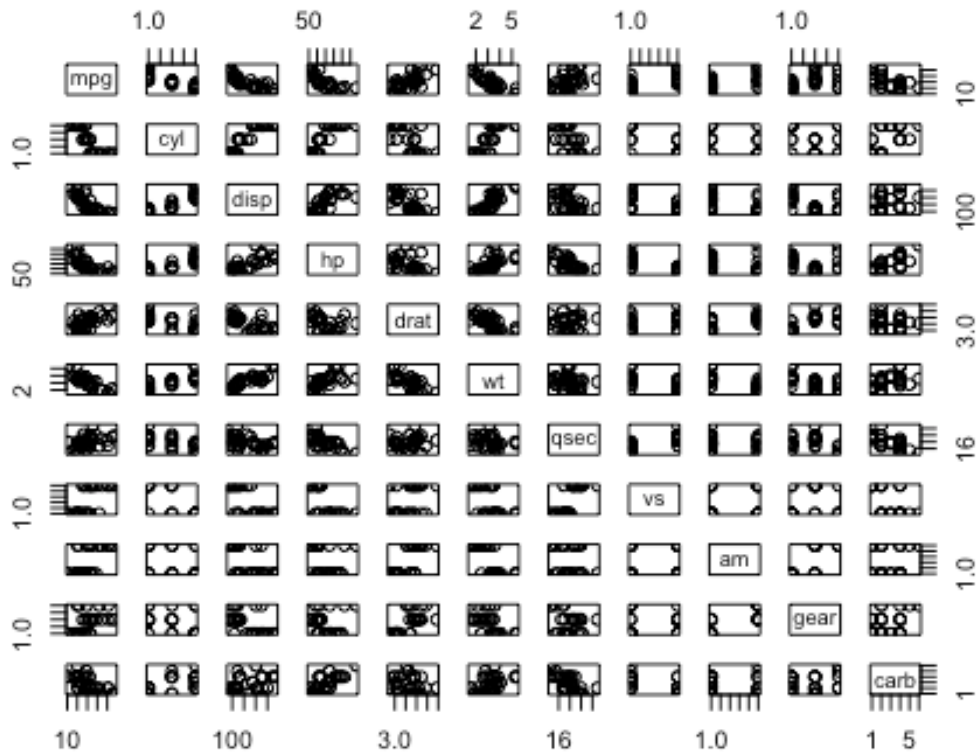
Seems that mean MPG of manual transmission cars is 7.24 MPGs higher than that of automatic transmission cars. We need to check whether this is a significant difference. Set alpha-value at 0.5 and run a t-test to find out.

```
autoData <- mtcars[mtcars$am == "Automatic",]  
manualData <- mtcars[mtcars$am == "Manual",]  
t.test(autoData$mpg, manualData$mpg)
```

```
##  
## Welch Two Sample t-test  
##  
## data: autoData$mpg and manualData$mpg  
## t = -3.7671, df = 18.332, p-value = 0.001374  
## alternative hypothesis: true difference in means is not equal to 0  
## 95 percent confidence interval:  
## -11.280194 -3.209684  
## sample estimates:  
## mean of x mean of y  
## 17.14737 24.39231
```

## Regression Analysis

```
pairs(mtcars)
```



```
data(mtcars)
sort(cor(mtcars)[1,])

##          wt          cyl          disp          hp          carb          qsec
## -0.8676594 -0.8521620 -0.8475514 -0.7761684 -0.5509251  0.4186840
##          gear          am          vs          drat          mpg
##  0.4802848  0.5998324  0.6640389  0.6811719  1.0000000
```

Based on pairwise correlation of variables with mpg, we see that there is little linear correlation between mpg and the variables qsec, gear, and carb.

## Model building and selection

Model1 explains variation less (adjusted R Square: 0.3385) compared to Model2 (adjusted R Square: 0.8066). However, Model1's am variable give lower p-value (less than 0.05) whereas no variable in the model2 gives lower p-value than 0.05 (due to overfitting). Therefore, we use the step method to iterate over the variables and obtain the best model.

```
model1 <- lm(mpg~am, data = mtcars)
summary(model1)
```

```
##
## 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

model2 <- lm(mpg~., data = mtcars)
summary(model2)

##
## Call:
## lm(formula = mpg ~ ., data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4506 -1.6044 -0.1196  1.2193  4.6271
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  12.30337   18.71788    0.657  0.5181
## cyl          -0.11144    1.04502   -0.107  0.9161
## disp           0.01334    0.01786    0.747  0.4635
## hp            -0.02148    0.02177   -0.987  0.3350
## drat           0.78711    1.63537    0.481  0.6353
## wt            -3.71530    1.89441   -1.961  0.0633 .
## qsec           0.82104    0.73084    1.123  0.2739
## vs             0.31776    2.10451    0.151  0.8814
## am             2.52023    2.05665    1.225  0.2340
## gear           0.65541    1.49326    0.439  0.6652
## carb          -0.19942    0.82875   -0.241  0.8122
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.65 on 21 degrees of freedom
## Multiple R-squared:  0.869, Adjusted R-squared:  0.8066
## F-statistic: 13.93 on 10 and 21 DF,  p-value: 3.793e-07
```

```
best <- step(model2, direction = "both", trace = FALSE)
summary(best)

##
## 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
```

## Residual diagnostics

Residual are normally distributed and homoskedastic. So, we can report the estimates from our report.

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

