

# Motor Trend Car Road Tests - Effects of transmission on MPG

## Executive Summary

This detailed analysis has been performed to fulfill the requirements of the course project for the course Regression Models (<https://www.coursera.org/course/regmods>) offered by the Johns Hopkins University (<https://www.coursera.org/jhu>) on Coursera (<https://www.coursera.org/>). In this project, we will analyze the `mtcars` data set and explore the relationship between a set of variables and miles per gallon (MPG) which will be our outcome.

The main objectives of this research are as follows

- Is an automatic or manual transmission better for MPG?
- Quantifying how different is the MPG between automatic and manual transmissions?

The key takeaway from our analysis was

- Manual transmission is better for MPG by a factor of 1.8 compared to automatic transmission.
- Means and medians for automatic and manual transmission cars are significantly different.

## Data processing and transformation

We load in the data set, perform the necessary data transformations by factoring the necessary variables and look at the data, in the following section.

```
data(mtcars)
mtcars$cyl <- factor(mtcars$cyl)
mtcars$vs <- factor(mtcars$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 ...
```

# Exploratory Data Analysis

In this section, we dive deeper into our data and explore various relationships between variables of interest. Initially, we plot the relationships between all the variables of the dataset (see Figure 2 in the appendix). From the plot, we notice that variables like `cyl`, `disp`, `hp`, `drat`, `wt`, `vs` and `am` seem to have some strong correlation with `mpg`. But we will use linear models to quantify that in the regression analysis section.

Since we are interested in the effects of car transmission type on `mpg`, we plot boxplots of the variable `mpg` when `am` is `Automatic` or `Manual` (see Figure 3 in the appendix). This plot clearly depicts an increase in the `mpg` when the transmission is `Manual`.

# Regression Analysis

In this section, we start building linear regression models based on the different variables and try to find out the best model fit and compare it with the base model which we have using `anova`. After model selection, we also perform analysis of residuals.

## Model building and selection

Like we mentioned earlier, based on the pairs plot where several variables seem to have high correlation with `mpg`, We build an initial model with all the variables as predictors, and perform stepwise model selection to select significant predictors for the final model which is the best model. This is taken care by the `step` method which runs `lm` multiple times to build multiple regression models and select the best variables from them using both **forward selection** and **backward elimination** methods by the `AIC` algorithm. The code is depicted in the section below, you can run it to see the detailed computations if required.

```
init_model <- lm(mpg ~ ., data = mtcars)
best_model <- step(init_model, direction = "both")
```

The best model obtained from the above computations consists of the variables, `cyl`, `wt` and `hp` as confounders and `am` as the independent variable. Details of the model are depicted below.

```
summary(best_model)
```

```
##
## Call:
## lm(formula = mpg ~ cyl + hp + wt + am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.939 -1.256 -0.401  1.125  5.051
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  33.7083     2.6049   12.94  7.7e-13 ***
## cyl6         -3.0313     1.4073   -2.15  0.0407 *
## cyl8         -2.1637     2.2843   -0.95  0.3523
## hp           -0.0321     0.0137   -2.35  0.0269 *
## wt           -2.4968     0.8856   -2.82  0.0091 **
## amManual      1.8092     1.3963    1.30  0.2065
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.41 on 26 degrees of freedom
## Multiple R-squared:  0.866, Adjusted R-squared:  0.84
## F-statistic: 33.6 on 5 and 26 DF, p-value: 1.51e-10
```

From the above model details, we observe that the adjusted  $(R^2)$  value is **0.84** which is the maximum obtained considering all combinations of variables. Thus, we can conclude that more than 84% of the variability is explained by the above model.

In the following section, we compare the base model with only `am` as the predictor variable and the best model which we obtained earlier containing confounder variables also.

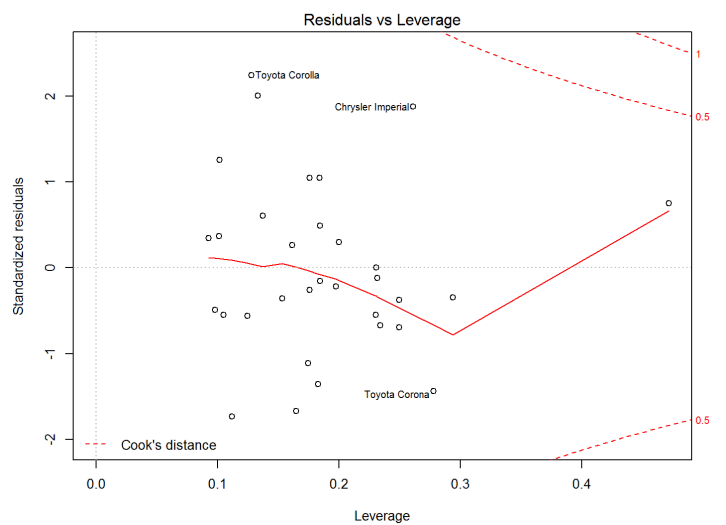
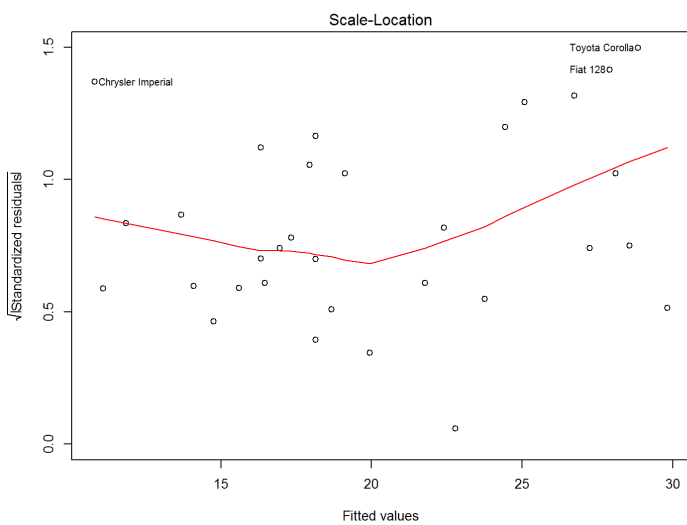
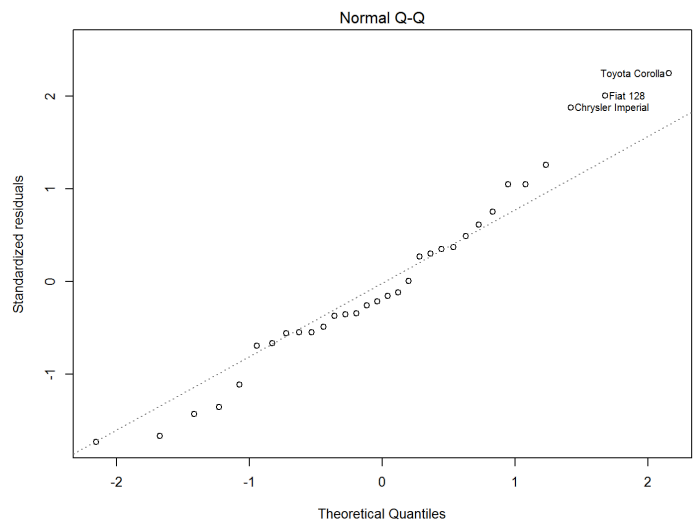
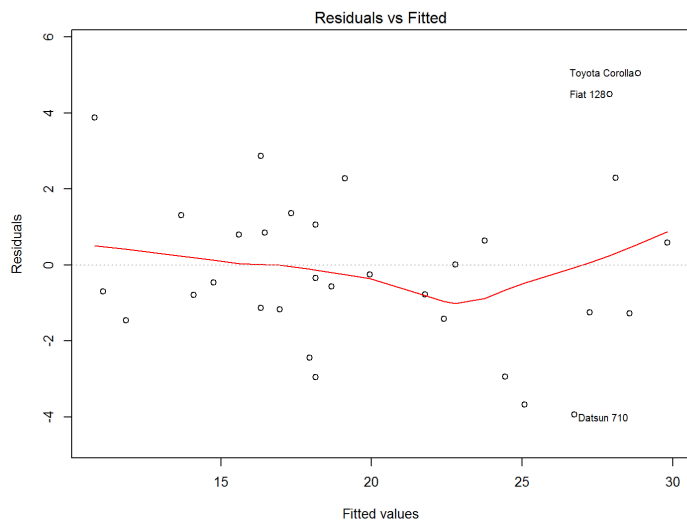
```
base_model <- lm(mpg ~ am, data = mtcars)
anova(base_model, best_model)
```

```
## Analysis of Variance Table
##
## Model 1: mpg ~ am
## Model 2: mpg ~ cyl + hp + wt + am
##   Res.Df  RSS Df Sum of Sq    F  Pr(>F)
## 1      30 721
## 2      26 151  4      570 24.5 1.7e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Looking at the above results, the ***p-value*** obtained is highly significant and we reject the null hypothesis that the confounder variables `cyl`, `hp` and `wt` don't contribute to the accuracy of the model.

## Residuals and Diagnostics

In this section, we shall study the residual plots of our regression model and also compute some of the regression diagnostics for our model to find out some interesting leverage points (often called as outliers) in the data set.



From the above plots, we can make the following observations,

- The points in the **Residuals vs. Fitted** plot seem to be randomly scattered on the plot and verify the independence condition.
- The **Normal Q-Q** plot consists of the points which mostly fall on the line indicating that the residuals are normally distributed.
- The **Scale-Location** plot consists of points scattered in a constant band pattern, indicating constant variance.
- There are some distinct points of interest (outliers or leverage points) in the top right of the plots.

We now compute some regression diagnostics of our model to find out these interesting leverage points as shown in the following section. We compute top three points in each case of influence measures.

```
leverage <- hatvalues(best_model)
tail(sort(leverage), 3)
```

##	Toyota Corona	Lincoln Continental	Maserati Bora
##	0.2778	0.2937	0.4714

```
influential <- dfbetas(best_model)
tail(sort(influential[,6]),3)
```

## Chrysler Imperial	Fiat 128	Toyota Corona
## 0.3507	0.4292	0.7305

Looking at the above cars, we notice that our analysis was correct, as the same cars are mentioned in the residual plots.

## Inference

We also perform a t-test assuming that the transmission data has a normal distribution and we clearly see that the manual and automatic transmissions are significantly different.

```
t.test(mpg ~ am, data = mtcars)
```

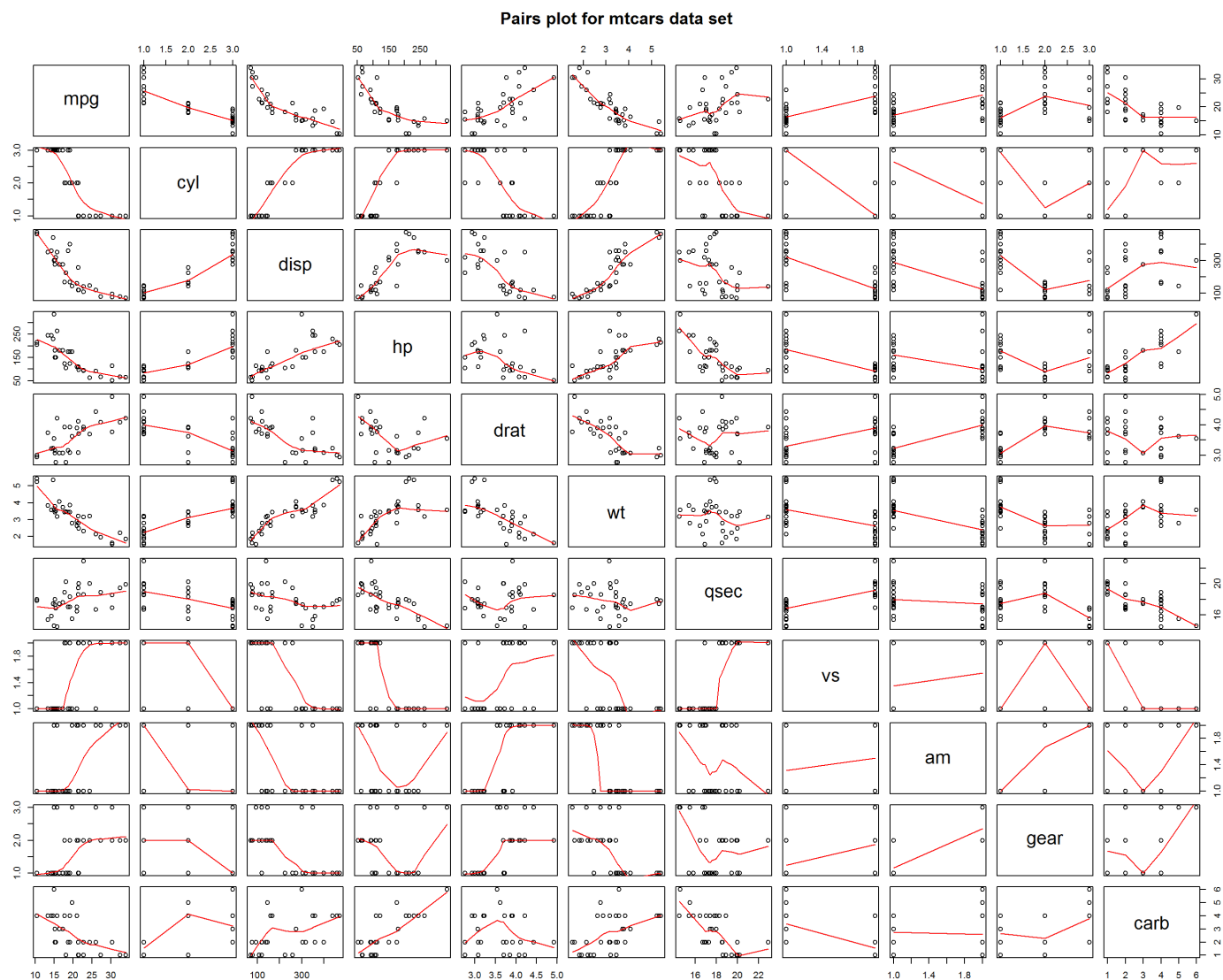
```
##
## Welch Two Sample t-test
##
## data: mpg by am
## t = -3.767, df = 18.33, p-value = 0.001374
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -11.28 -3.21
## sample estimates:
## mean in group Automatic mean in group Manual
## 17.15 24.39
```

## Conclusion

Based on the observations from our best fit model, we can conclude the following,

- Cars with `Manual` transmission get more miles per gallon `mpg` compared to cars with `Automatic` transmission. (1.8 adjusted by hp, cyl, and wt).
- `mpg` will decrease by 2.5 (adjusted by hp, cyl, and am) for every 1000 lb increase in `wt`.
- `mpg` decreases negligibly with increase of `hp`.
- If number of cylinders, `cyl` increases from 4 to 6 and 8, `mpg` will decrease by a factor of 3 and 2.2 respectively (adjusted by hp, wt, and am).

# Appendix





Car MPG by transmission type

