

Motor Trends Car Study

Jason Battles

October 15, 2016

Executive Summary

Motor Trend, a magazine about the automobile industry, requested this sample study to investigate the impact of several variables on miles per gallon (MPG). This study has a few primary points of inquiry.

1. Is an automatic or manual transmission better for MPG?
2. What is the quantifiable effect on MPG between automatic and manual transmissions?
3. Is there a more effective model for MPG prediction than simply transmission type?

This study reveals that vehicles with manual transmissions do indeed have a slightly better fuel efficiency than those with automatic transmissions. We also do indeed find a better model of MPG prediction than just using transmission type.

Environment Preparation

We configure the environment, load the data set, perform some necessary data transformations, and verify that the data is loaded correctly. This working data set has **32 observations of 11 variables**.

```
library(ggplot2)
library(corrplot)
data(mtcars)
mtcars.orig <- mtcars # Maintain original copy for Correlation Matrix
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

The focus of our study is `mpg` so let us first determine the nature of its distribution. We find that the median and mean are close to each other, so we can assume a normal distribution of data.

```
summary(mtcars$mpg)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##    10.40   15.42   19.20   20.09   22.80   33.90
```

In this analysis, we are interested in the effects of car transmission type (Automatic or Manual) on MPG so we start our exploration by creating a box plot (*Appendix 1 - Figure 1*). **This plot clearly depicts that manual transmissions tend to have higher MPG than automatic.**

Independent Variable Selection

We now more closely explore the available variables to identify those which are most significant to any improved regression models that we may build. We create a correlation matrix and learn that the variables `cyl`, `disp`, `hp`, `drat`, `wt`, `vs` and `am` (*transmission type*) have a strong correlation with dependent variable `mpg` (*sig.level = 0.0005*). Variables `qsec`, `gear`, and `carb` are NOT significant and will be excluded from further analysis. (**Appendix 1 - Figure 2**).

Regression Analysis

This data is further analyzed using regression analysis and fitted linear models. We then compare these fitted linear models using `anova` tables.

Simple Linear Regression

We begin our regression analysis with our original assumption that `mpg` is mostly dependent on `am` (transmission type).

```
first.mdl <- lm(mpg ~ factor(am), data = mtcars)
summary(first.mdl)
```

```
##
## Call:
## lm(formula = mpg ~ factor(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 ***
## factor(am)Manual    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
```

We see that the adjusted R squared value of this first model is only 0.338 which means that only **33.8%** of the regression variance can be explained with just transmission type. However, as outlined earlier by the Correlation Matrix, we must not forget there are several other predictor variables that we should investigate to see if any may lead to a better model for MPG.

Multivariate Regression Model

Although we identified the most highly correlated variables in the Correlation Matrix, we would like to sharpen our model further by conducting an **AIC** analysis using the **step** function.

Using this method, we build an initial model with all variables as possible predictors and then perform an iterative model selection process using the **step** method. This method runs **lm** multiple times to build multiple regression models and selects the best variables using both forward selection and backward elimination methods using the AIC algorithm. The full output is included in *Appendix 1 - Figure 3*.

```
init.mdl <- lm(mpg ~ ., data = mtcars)
best.mdl <- step(init.mdl, direction = "both")
```

The best model resulting in the lowest **AIC** value (AIC=61.65) is **lm(mpg ~ cyl + hp + wt + am, data = mtcars)**. This model depicts the variables **cyl**, **hp**, and **wt** as confounders and **am** as the independent variable. Details of the best model are depicted below.

```
summary(best.mdl)
```

```
##
## Call:
## lm(formula = mpg ~ cyl + hp + wt + am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.9387 -1.2560 -0.4013  1.1253  5.0513
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  33.70832    2.60489   12.940 7.73e-13 ***
## cyl6         -3.03134    1.40728   -2.154  0.04068 *
## cyl8         -2.16368    2.28425   -0.947  0.35225
## hp           -0.03211    0.01369   -2.345  0.02693 *
## wt           -2.49683    0.88559   -2.819  0.00908 **
## amManual      1.80921    1.39630    1.296  0.20646
## ---
## 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:  0.8659, Adjusted R-squared:  0.8401
## F-statistic: 33.57 on 5 and 26 DF,  p-value: 1.506e-10
```

We observe that the R squared value of this updated model is 0.8401 which means that **84%** of the variability is now explained.

We now compare the first model (only transmission type as predictor) with the best model.

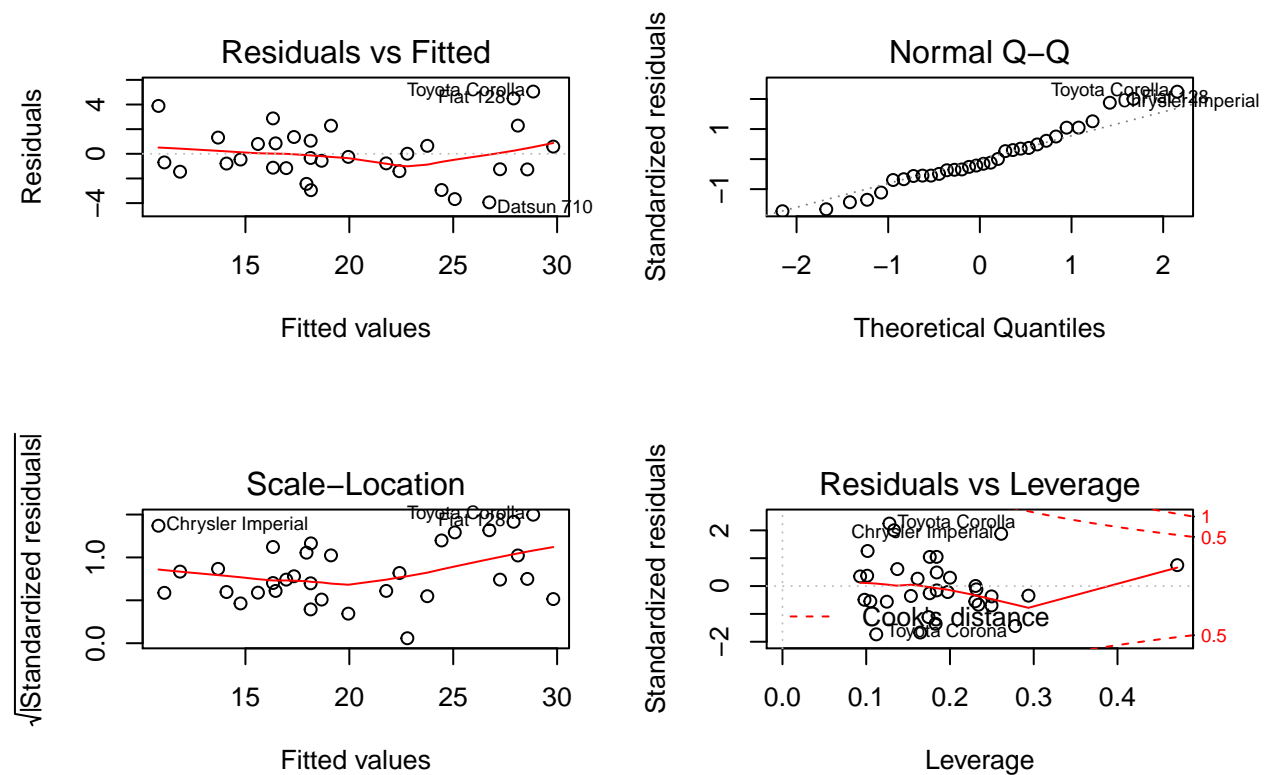
```
anova(first.mdl, best.mdl)
```

```
## Analysis of Variance Table
##
## Model 1: mpg ~ factor(am)
## Model 2: mpg ~ cyl + hp + wt + am
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      30 720.90
## 2      26 151.03  4    569.87 24.527 1.688e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

The P-Value is highly significant so we may reject the null hypothesis that the confounder variables `cyl`, `hp`, and `wt` do not contribute to the accuracy of the model. Therefore, the best model is indeed a better model than the first model.

Model Residuals and Diagnostics

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



1. The points in the Residuals vs. Fitted plot are randomly scattered on the plot which verifies the independence condition.

2. The Normal Q-Q plot consists of the points which mostly fall on the line indicating that the residuals are normally distributed.
3. The Scale-Location plot consists of points scattered in a constant band pattern, indicating constant variance.
4. There are some distinct points of interest (outliers or leverage points) in the top right of the plots that may indicate values of increased leverage of outliers.

Conclusions

1. Manual transmissions result in higher MPG than automatic transmissions.
2. Manual transmissions will achieve an additional 1.8mpg compared to automatic transmissions.
3. The optimal regression model includes `cyl`, `hp`, `wt`, `am` (transmission type)
 - Baseline for the model is a 4 cylinders engine. 6 cylinders will decrease `mpg` by 3.03.
 - 8 cylinders will decrease `mpg` by an additional 2.16.
 - `mpg` will decrease by 2.49 for every 1000lb increase in `wt`
 - `mpg` will slightly decrease by 0.32 with every increase in 10hp

Appendix 1

Figure 1 - Box Plot (MPG vs. Transmission Type)

```
p <- ggplot(mtcars, aes(factor(am), mpg))
p + geom_boxplot() + geom_jitter() +
  labs(title="Box Plot of MPG vs. Transmission Type", x = "Transmission")
```

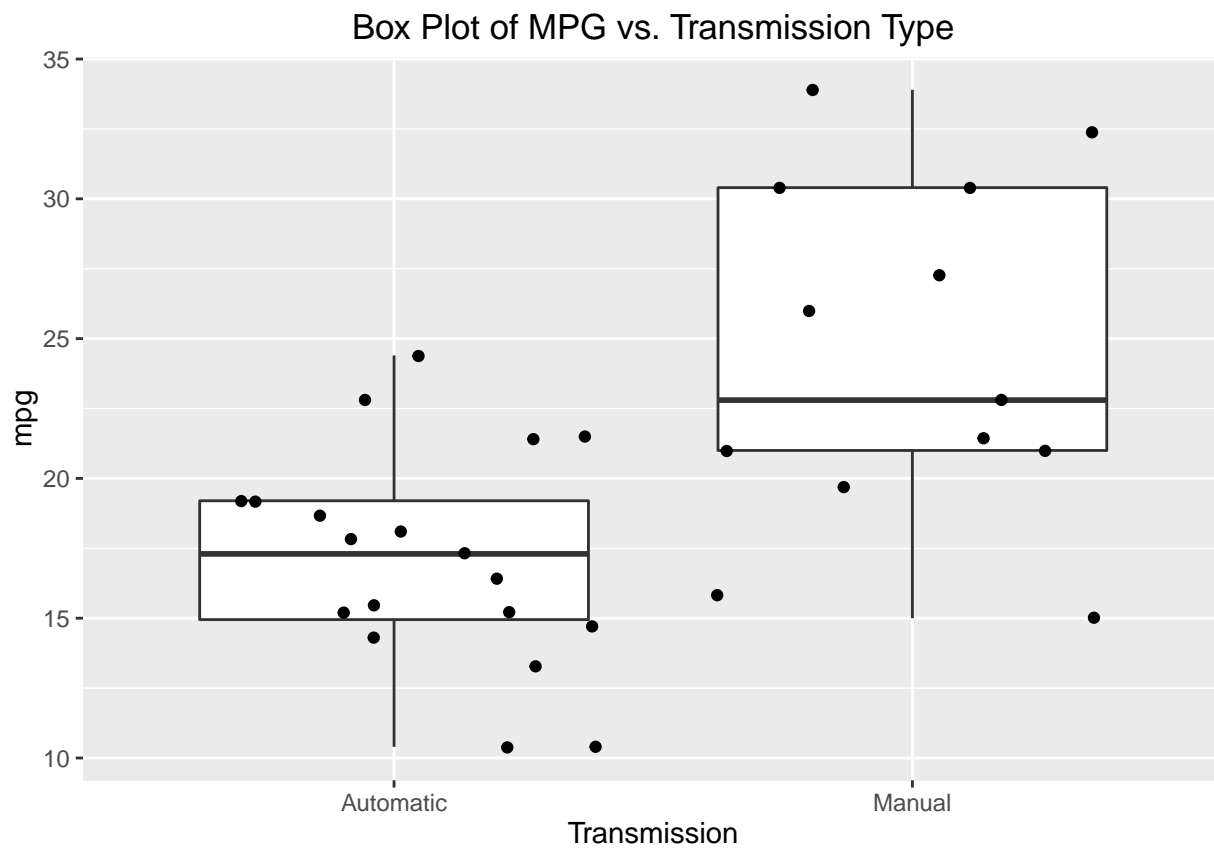


Figure 2 - Multivariate Correlation Matrix

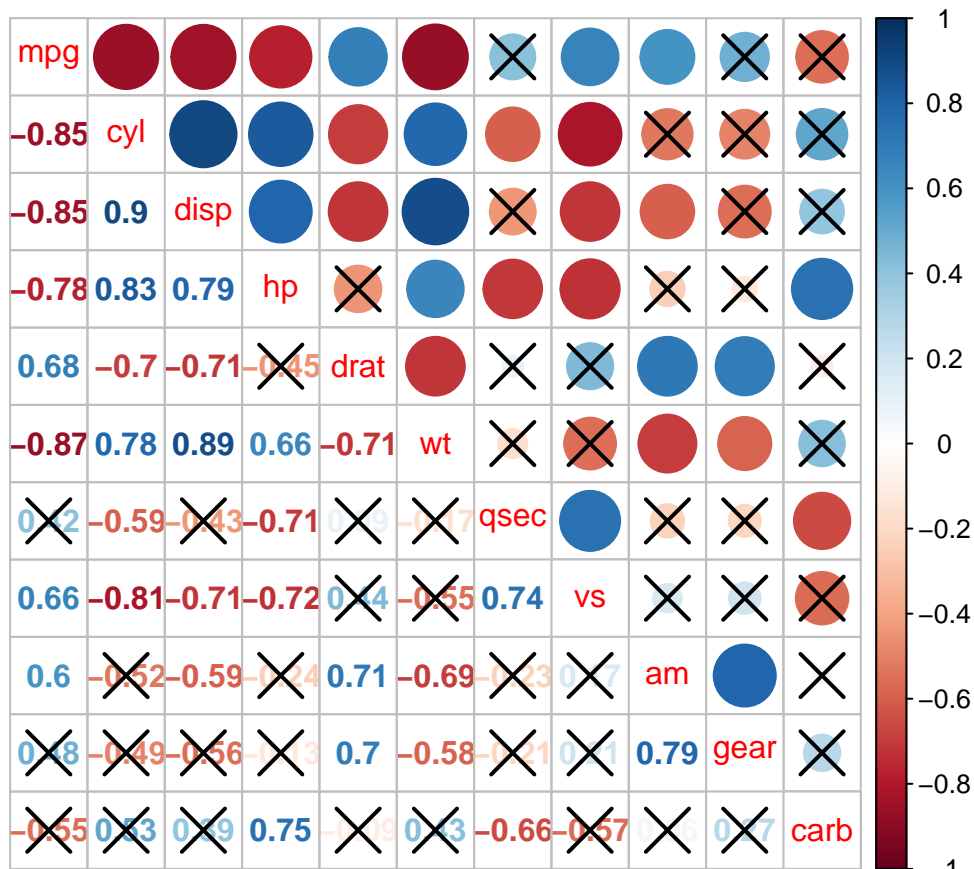


Figure 3 - Output from AIC Step Analysis

```
init.mdl <- lm(mpg ~ ., data = mtcars)
best.mdl <- step(init.mdl, direction = "both")

## Start:  AIC=76.4
## mpg ~ cyl + disp + hp + drat + wt + qsec + vs + am + gear + carb
##
##      Df Sum of Sq  RSS   AIC
## - carb  5   13.5989 134.00 69.828
## - gear  2    3.9729 124.38 73.442
## - am    1    1.1420 121.55 74.705
## - qsec  1    1.2413 121.64 74.732
## - drat  1    1.8208 122.22 74.884
## - cyl   2   10.9314 131.33 75.184
## - vs    1    3.6299 124.03 75.354
## <none>             120.40 76.403
## - disp  1    9.9672 130.37 76.948
## - wt    1   25.5541 145.96 80.562
## - hp    1   25.6715 146.07 80.588
##
```

```

## Step: AIC=69.83
## mpg ~ cyl + disp + hp + drat + wt + qsec + vs + am + gear
##
##      Df Sum of Sq  RSS   AIC
## - gear  2    5.0215 139.02 67.005
## - disp  1    0.9934 135.00 68.064
## - drat  1    1.1854 135.19 68.110
## - vs    1    3.6763 137.68 68.694
## - cyl   2   12.5642 146.57 68.696
## - qsec  1    5.2634 139.26 69.061
## <none>          134.00 69.828
## - am    1   11.9255 145.93 70.556
## - wt    1   19.7963 153.80 72.237
## - hp    1   22.7935 156.79 72.855
## + carb  5   13.5989 120.40 76.403
##
## Step: AIC=67
## mpg ~ cyl + disp + hp + drat + wt + qsec + vs + am
##
##      Df Sum of Sq  RSS   AIC
## - drat  1    0.9672 139.99 65.227
## - cyl   2   10.4247 149.45 65.319
## - disp  1    1.5483 140.57 65.359
## - vs    1    2.1829 141.21 65.503
## - qsec  1    3.6324 142.66 65.830
## <none>          139.02 67.005
## - am    1   16.5665 155.59 68.608
## - hp    1   18.1768 157.20 68.937
## + gear  2    5.0215 134.00 69.828
## - wt    1   31.1896 170.21 71.482
## + carb  5   14.6475 124.38 73.442
##
## Step: AIC=65.23
## mpg ~ cyl + disp + hp + wt + qsec + vs + am
##
##      Df Sum of Sq  RSS   AIC
## - disp  1    1.2474 141.24 63.511
## - vs    1    2.3403 142.33 63.757
## - cyl   2   12.3267 152.32 63.927
## - qsec  1    3.1000 143.09 63.928
## <none>          139.99 65.227
## + drat  1    0.9672 139.02 67.005
## - hp    1   17.7382 157.73 67.044
## - am    1   19.4660 159.46 67.393
## + gear  2    4.8033 135.19 68.110
## - wt    1   30.7151 170.71 69.574
## + carb  5   13.0509 126.94 72.095
##
## Step: AIC=63.51
## mpg ~ cyl + hp + wt + qsec + vs + am
##
##      Df Sum of Sq  RSS   AIC
## - qsec  1    2.442 143.68 62.059
## - vs    1    2.744 143.98 62.126

```



```

## - cyl    2    18.580 159.82 63.466
## <none>                141.24 63.511
## + disp   1      1.247 139.99 65.227
## + drat   1      0.666 140.57 65.359
## - hp     1     18.184 159.42 65.386
## - am     1     18.885 160.12 65.527
## + gear   2      4.684 136.55 66.431
## - wt     1     39.645 180.88 69.428
## + carb   5      2.331 138.91 72.978
##
## Step:  AIC=62.06
## mpg ~ cyl + hp + wt + vs + am
##
##      Df Sum of Sq    RSS    AIC
## - vs   1      7.346 151.03 61.655
## <none>                143.68 62.059
## - cyl   2     25.284 168.96 63.246
## + qsec   1      2.442 141.24 63.511
## - am     1     16.443 160.12 63.527
## + disp   1      0.589 143.09 63.928
## + drat   1      0.330 143.35 63.986
## + gear   2      3.437 140.24 65.284
## - hp     1     36.344 180.02 67.275
## - wt     1     41.088 184.77 68.108
## + carb   5      3.480 140.20 71.275
##
## Step:  AIC=61.65
## mpg ~ cyl + hp + wt + am
##
##      Df Sum of Sq    RSS    AIC
## <none>                151.03 61.655
## - am     1      9.752 160.78 61.657
## + vs     1      7.346 143.68 62.059
## + qsec   1      7.044 143.98 62.126
## - cyl    2     29.265 180.29 63.323
## + disp   1      0.617 150.41 63.524
## + drat   1      0.220 150.81 63.608
## + gear   2      1.361 149.66 65.365
## - hp     1     31.943 182.97 65.794
## - wt     1     46.173 197.20 68.191
## + carb   5      5.633 145.39 70.438

```