

# Regression Models : Automatic|Manual Transmission better for MPG

*Aysegul Sonmez*

*March 4, 2018*

## R Markdown

### Motor Trends Regression Model Project Week 4 Executive Summary

In this report, we will analyze mtcars data set and explore the relationship between a set of variables and miles per gallon (MPG). The data was extracted from the 1974 Motor Trend US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973–74 models). We use regression models and exploratory data analyses to mainly explore how automatic ( $am = 0$ ) and manual ( $am = 1$ ) transmissions features affect the MPG feature. The t-test shows that the performance difference between cars with automatic and manual transmission. And it is about 7 MPG more for cars with manual transmission than those with automatic transmission. Then, we fit several linear regression models and select the one with highest Adjusted R-squared value. So, given that weight and 1/4 mile time are held constant, manual transmitted cars are  $14.079 + (-4.141) \cdot \text{weight}$  more MPG (miles per gallon) on average better than automatic transmitted cars. Thus, cars that are lighter in weight with a manual transmission and cars that are heavier in weight with an automatic transmission will have higher MPG values.

## Exploratory Data Analysis

First, we load the data set mtcars and change some variables from numeric class to factor class.

```
library(ggplot2)
data(mtcars)
mtcars[1:3, ] # Sample Data
```

##		mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
##	Mazda RX4	21.0	6	160	110	3.90	2.620	16.46	0	1	4	4
##	Mazda RX4 Wag	21.0	6	160	110	3.90	2.875	17.02	0	1	4	4
##	Datsun 710	22.8	4	108	93	3.85	2.320	18.61	1	1	4	1

```
dim(mtcars)

## [1] 32 11

mtcars$cyl <- as.factor(mtcars$cyl)
mtcars$vs <- as.factor(mtcars$vs)
mtcars$am <- factor(mtcars$am)
mtcars$gear <- factor(mtcars$gear)
mtcars$carb <- factor(mtcars$carb)
attach(mtcars)

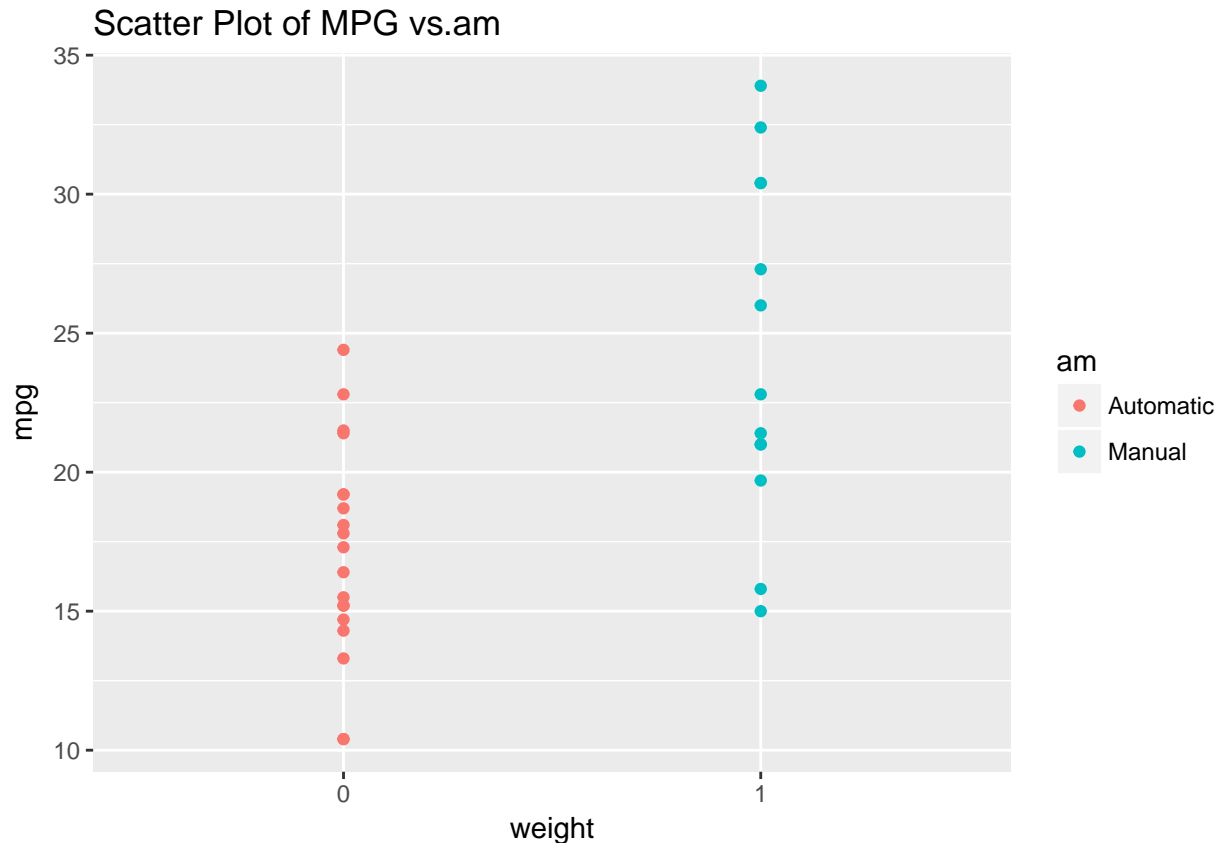
## The following object is masked from package:ggplot2:
##
## mpg
```

```

par(mfrow = c(1,2))
data_matrix <- data.matrix(mtcars, rownames.force = NA)
barplot(data_matrix, main="Mtcars", ylab=data_matrix[,0], beside=TRUE, col=rainbow(5), pch=22, lwd=1)
legend("bottom", rownames(mtcars), cex=0.6, bty="n", fill=rainbow(5), lwd=2, box.lwd=2);

ggplot(mtcars, aes(x=am, y=mpg, group="rownames(mtcars)", am, color=am, height=3, width=3)) + geom_point
scale_colour_discrete(labels=c("Automatic", "Manual")) +
xlab("weight") + ggtitle("Scatter Plot of MPG vs.am ")

```



Then, we do some basic exploratory data analyses. Please refer to the Appendix: Figures section for the plots. According to the box plot, we see that manual transmission yields higher values of MPG in general. And as for the pair graph, we can see some higher correlations between variables like “wt”, “disp”, “cyl” and “hp”.

### Inference

At this step, we make the null hypothesis as the MPG of the automatic and manual transmissions are from the same population (assuming the MPG has a normal distribution). We use the two sample T-test to show it.

```
## [1] 0.001373638
```

```
## mean in group 0 mean in group 1
##      17.14737      24.39231
```

Since the p-value is 0.00137, we reject our null hypothesis. So, the automatic and manual transmissions are from different populations. And the mean for MPG of manual transmitted cars is about 7 more than that of automatic transmitted cars.

## Regression Analysis

First, we fit the full model as the following.

```
##
## Call:
## lm(formula = mpg ~ ., data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5087 -1.3584 -0.0948  0.7745  4.6251
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  23.87913    20.06582   1.190   0.2525
## cyl6         -2.64870     3.04089  -0.871   0.3975
## cyl8         -0.33616     7.15954  -0.047   0.9632
## disp         0.03555     0.03190   1.114   0.2827
## hp          -0.07051     0.03943  -1.788   0.0939 .
## drat         1.18283     2.48348   0.476   0.6407
## wt          -4.52978     2.53875  -1.784   0.0946 .
## qsec         0.36784     0.93540   0.393   0.6997
## vs1          1.93085     2.87126   0.672   0.5115
## am1          1.21212     3.21355   0.377   0.7113
## gear4        1.11435     3.79952   0.293   0.7733
## gear5        2.52840     3.73636   0.677   0.5089
## carb2       -0.97935     2.31797  -0.423   0.6787
## carb3        2.99964     4.29355   0.699   0.4955
## carb4        1.09142     4.44962   0.245   0.8096
## carb6        4.47757     6.38406   0.701   0.4938
## carb8        7.25041     8.36057   0.867   0.3995
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.833 on 15 degrees of freedom
## Multiple R-squared:  0.8931, Adjusted R-squared:  0.779
## F-statistic:  7.83 on 16 and 15 DF,  p-value: 0.000124
```

This model has the Residual standard error as 2.833 on 15 degrees of freedom. And the Adjusted R-squared value is 0.779, which means that the model can explain about 78% of the variance of the MPG variable. However, none of the coefficients are significant at 0.05 significant level.

Then, we use backward selection to select some statistically significant variables.

```
stepModel <- step(fullModel, k=log(nrow(mtcars)))
summary(stepModel) # results hidden
```

This model is “mpg ~ wt + qsec + am”. It has the Residual standard error as 2.459 on 28 degrees of freedom. And the Adjusted R-squared value is 0.8336, which means that the model can explain about 83% of the variance of the MPG variable. All of the coefficients are significant at 0.05 significant level.

Please refer to the Appendix: Figures section for the plots again. According to the scatter plot, it indicates that there appear to be an interaction term between “wt” variable and “am” variable, since automatic cars tend to weigh heavier than manual cars. Thus, we have the following model including the interaction term:

```
amIntWtModel<-lm(mpg ~ wt + qsec + am + wt:am, data=mtcars)
summary(amIntWtModel) # results hidden
```

This model has the Residual standard error as 2.084 on 27 degrees of freedom. And the Adjusted R-squared value is 0.8804, which means that the model can explain about 88% of the variance of the MPG variable. All of the coefficients are significant at 0.05 significant level. This is a pretty good one.

Next, we fit the simple model with MPG as the outcome variable and Transmission as the predictor variable.

```
amModel<-lm(mpg ~ am, data=mtcars)
summary(amModel) # results hidden
```

It shows that on average, a car has 17.147 mpg with automatic transmission, and if it is manual transmission, 7.245 mpg is increased. This model has the Residual standard error as 4.902 on 30 degrees of freedom. And the Adjusted R-squared value is 0.3385, which means that the model can explain about 34% of the variance of the MPG variable. The low Adjusted R-squared value also indicates that we need to add other variables to the model.

Finally, we select the final model.

```
anova(amModel, stepModel, fullModel, amIntWtModel)
confint(amIntWtModel) # results hidden
```

We end up selecting the model with the highest Adjusted R-squared value, “mpg ~ wt + qsec + am + wt:am”.

```
{r , echo=TRUE) summary(amIntWtModel)$coef
```

Thus, the result shows that when “wt” (weight lb/1000) and “qsec” (1/4 mile time) remain constant, cars with manual transmission add  $14.079 + (-4.141) \cdot \text{wt}$  more MPG (miles per gallon) on average than cars with automatic transmission. That is, a manual transmitted car that weighs 2000 lbs have 5.797 more MPG than an automatic transmitted car that has both the same weight and 1/4 mile time.

## Residual Analysis and Diagnostics

Please refer to the Appendix: Figures section for the plots. According to the residual plots, we can verify the following underlying assumptions:

- 1.The Residuals vs. Fitted plot shows no consistent pattern, supporting the accuracy of the independence.
- 2.The Normal Q-Q plot indicates that the residuals are normally distributed because the points lie close to the line.
- 3.The Scale-Location plot confirms the constant variance assumption, as the points are randomly distributed around the line.
- 4.The Residuals vs. Leverage argues that no outliers are present, as all values fall well within the 0.5 Cook's distance.

As for the Dfbetas, the measure of how much an observation has effected the estimate of a regression coefficient, we get the following result:

```
{r , echo=TRUE) sum((abs(dfbetas(amIntWtModel)))>1)
```

Therefore, the above analyses meet all basic assumptions of linear regression and well answer the questions. Appendix: Figures

- 1.Boxplot of MPG vs. Transmission

```
{r , echo=TRUE) boxplot(mpg ~ am, xlab="Transmission (0 = Automatic, 1 = Manual)", ylab="MPG",
main="Boxplot of MPG vs. Transmission")
```

## 2. Pair Graph of Motor Trend Car Road Tests

```
“{r , echo=TRUE)
```

```
pairs(mtcars, panel=panel.smooth, main="Pair Graph of Motor Trend Car Road Tests") “
```

## 3. Scatter Plot of MPG vs. Weight by Transmission

```
{r , echo=TRUE) ggplot(mtcars, aes(x=wt, y=mpg, group=am, color=am, height=3, width=3)) +
geom_point() + scale_colour_discrete(labels=c("Automatic", "Manual")) + xlab("weight")
+ ggtitle("Scatter Plot of MPG vs. Weight by Transmission")
```

## 4. Residual Plots

