

Manual or automatic: how does affect fuel performance?

Javier Santibáñez

Saturday, January 24, 2015

Summary

In this report we present the results of a comparison between manual and automatic vehicles, measured in miles per gallon of fuel (mpg). The results shows that automatic cars has a better performance than manual cars. The expected difference was estimated in XXX mpg.

Introduction

We used a data set that consists of the measurements of 11 variables from observations of cars from 32 cars. The variables are the following:

- `mpg`, miles(US)/gallon
- `cyl`, number of cylinders
- `dis`, displacement (cu.in.)
- `hp`, gross horsepower
- `drat`, rear axie ratio
- `wt`, weigth (lb/1000)
- `qsec` 1/4 mile time
- `vs` V/S
- `am` transmission (automatic, manual)
- `gear`, number of forward gears
- `carb`, number of carburetors

Additionally, we created a new variable, called `am2`, which is the factor version of `am` with levels *automatic* and *manual*.

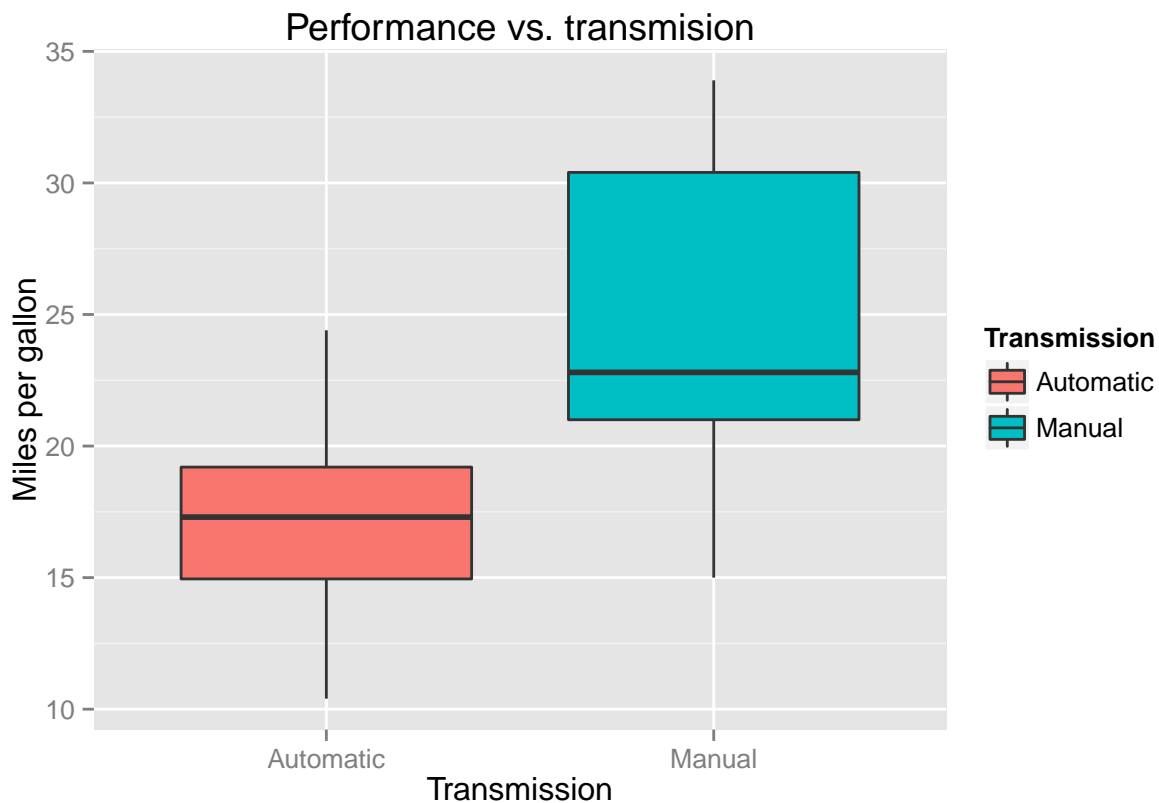
In this report we have to address the questions:

1. **Is an automatic or manual transmission better for performance?**
2. **What is the difference in performance between automatic and manual transmissions?**

To deal with this questions we used a statistical approach, in specific, we used regression models.

The first qustion can be easily answered. Figure *Performance vs. transmission* suggests that performance, measured in miles per gallon, depends of the car transmissions.

```
suppressWarnings(suppressMessages(library(ggplot2)))
suppressWarnings(suppressMessages(library(knitr)))
suppressWarnings(suppressMessages(library(corrplot)))
mtcars$am2<-factor(mtcars$am,labels=c("Automatic","Manual"))
g<-ggplot(data=mtcars,aes(factor(am2),mpg))
g+geom_boxplot(aes(fill=am2))+
  xlab("Transmission")+ylab("Miles per gallon")+labs(fill="Transmission")+
  ggtitle("Performance vs. transmsion")
```



Despite the previous graphic, we must verify that manual transmission is better than automatic transmission. Also, it is important to consider other variables, because some of them can explain part of the variability. Finally, we used the software statistical environment R and the packages `ggplot2` and `corrplot` for graphics.

Methodology

The following is a brief description of the analysis performed:

1. Fit a regression model with `mpg` (performance) as output and `am` (transmission) as input.
2. Select other covariables from the data set, considering their correlation with `mpg`.
3. Fit a regression model with `mpg` as output and the set of variables from the preceding step as input.
4. Evaluate and compare models fitted in steps 1 and 4.

Results

Step 1. Basic model.

The basic model consists in only one input (`mpg`). We used the next code to fit the model.

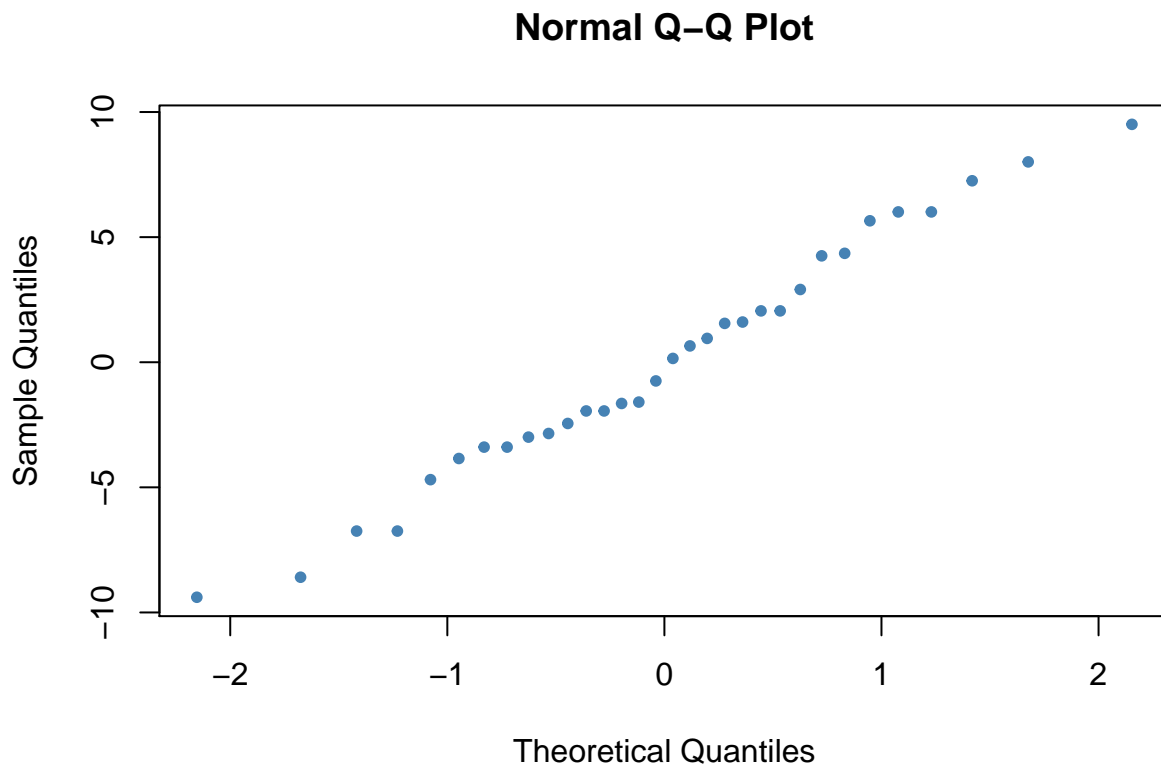
```
fitMod1<-lm(mpg~am2,data=mtcars)
summary(fitMod1)
```

```
##
## Call:
## lm(formula = mpg ~ am2, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -9.392 -3.092 -0.297  3.244  9.508
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    17.15      1.12    15.25 1.1e-15 ***
## am2Manual        7.24      1.76     4.11 0.00029 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 4.9 on 30 degrees of freedom
## Multiple R-squared:  0.36,    Adjusted R-squared:  0.338
## F-statistic: 16.9 on 1 and 30 DF,  p-value: 0.000285
```

The previous results indicates that there is a highly significant difference in performance between automatic and manual transmission, estimated as 7.2449 miles per gallon, on average. With this information we can answer our two questions, but we improve them with the following steps.

We verify the accuracy of the normal assumption with a *qq-plot*, it shows that the residuals are approximately normal distributed.

```
qqnorm(fitMod1$residual,pch=20,col="steelblue")
```



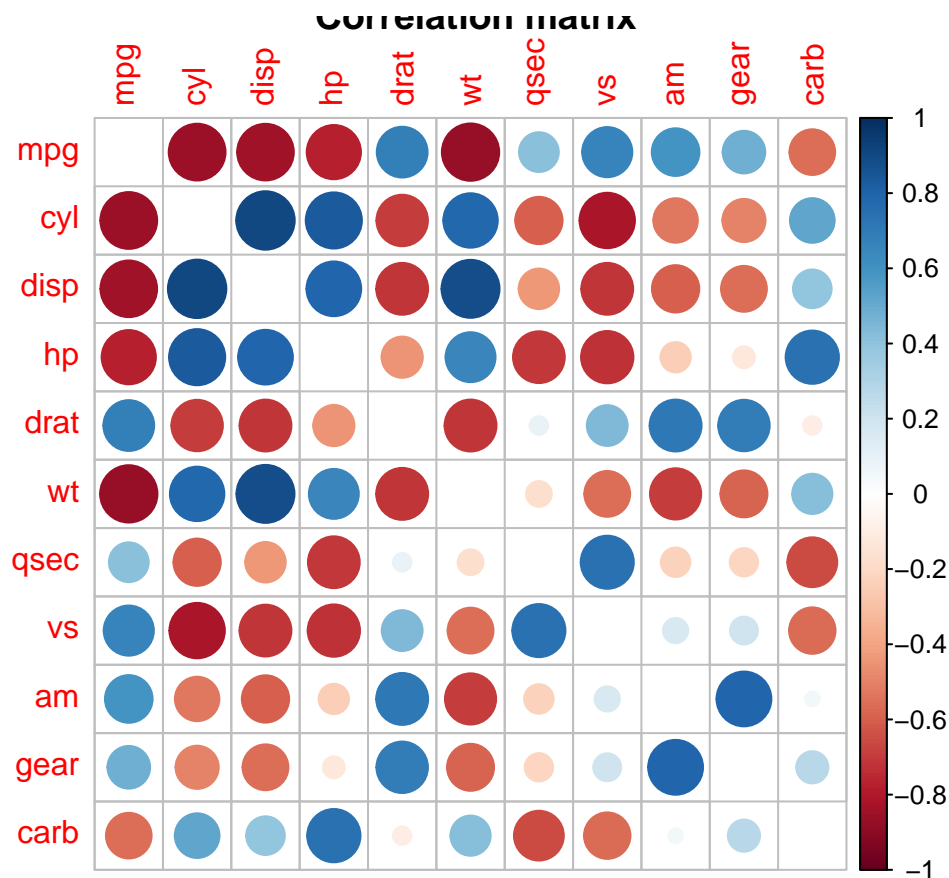
Step 2. Select new variables.

To select new variables to include in the model, we used the correlation matrix via the `corrplot` function. From Figure *Correlation matrix* we can identify a set of variables highly correlated:

- 6, 6, 4, 6, 8, 6, 8, 4, 4, 6, 6, 8, 8, 8, 8, 8, 8, 4, 4, 4, 4, 8, 8, 8, 8, 4, 4, 4, 8, 6, 8, 4
- 160, 160, 108, 258, 360, 225, 360, 146.7, 140.8, 167.6, 167.6, 275.8, 275.8, 275.8, 472, 460, 440, 78.7, 75.7, 71.1, 120.1, 318, 304, 350, 400, 79, 120.3, 95.1, 351, 145, 301, 121
- 110, 110, 93, 110, 175, 105, 245, 62, 95, 123, 123, 180, 180, 180, 205, 215, 230, 66, 52, 65, 97, 150, 150, 245, 175, 66, 91, 113, 264, 175, 335, 109
- 3.9, 3.9, 3.85, 3.08, 3.15, 2.76, 3.21, 3.69, 3.92, 3.92, 3.92, 3.07, 3.07, 3.07, 2.93, 3, 3.23, 4.08, 4.93, 4.22, 3.7, 2.76, 3.15, 3.73, 3.08, 4.08, 4.43, 3.77, 4.22, 3.62, 3.54, 4.11

Despite this result, we cannot include the whole set of variables because there are a high correlation between them. Then, we have to choose one variable from that set. We select the variable `hp` because we can use it as a confounder in an extended model, to reduce variability and to improve the estimation of difference in `mpg` due to transmission.

```
cor_mtcars<-cor(mtcars[, -12])  
corrplot(cor_mtcars, diag=F, title="Correlation matrix")
```



Step 3. Extended model.

Now we extend our basic model adding the variable `hp`. First we fit a model with an interaction term with the following code.

```
fitMod2<-lm(mpg~am2+hp+am2*hp,data=mtcars)
summary(fitMod2)
```

```
##
## Call:
## lm(formula = mpg ~ am2 + hp + am2 * hp, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.382 -2.270  0.134  1.706  5.875
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  26.624848   2.182943   12.20   1e-12 ***
## am2Manual     5.217653   2.665093    1.96    0.06 .
## hp           -0.059137   0.012945   -4.57   9e-05 ***
## am2Manual:hp  0.000403   0.016460    0.02    0.98
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.96 on 28 degrees of freedom
## Multiple R-squared:  0.782, Adjusted R-squared:  0.759
## F-statistic: 33.5 on 3 and 28 DF, p-value: 2.11e-09
```

The previous results suggest that there is evidence to support the hipotesis that the interaction term is not necessary.

We fit a new model wihtout the interaction term between `am` and `hp`, these are the results:

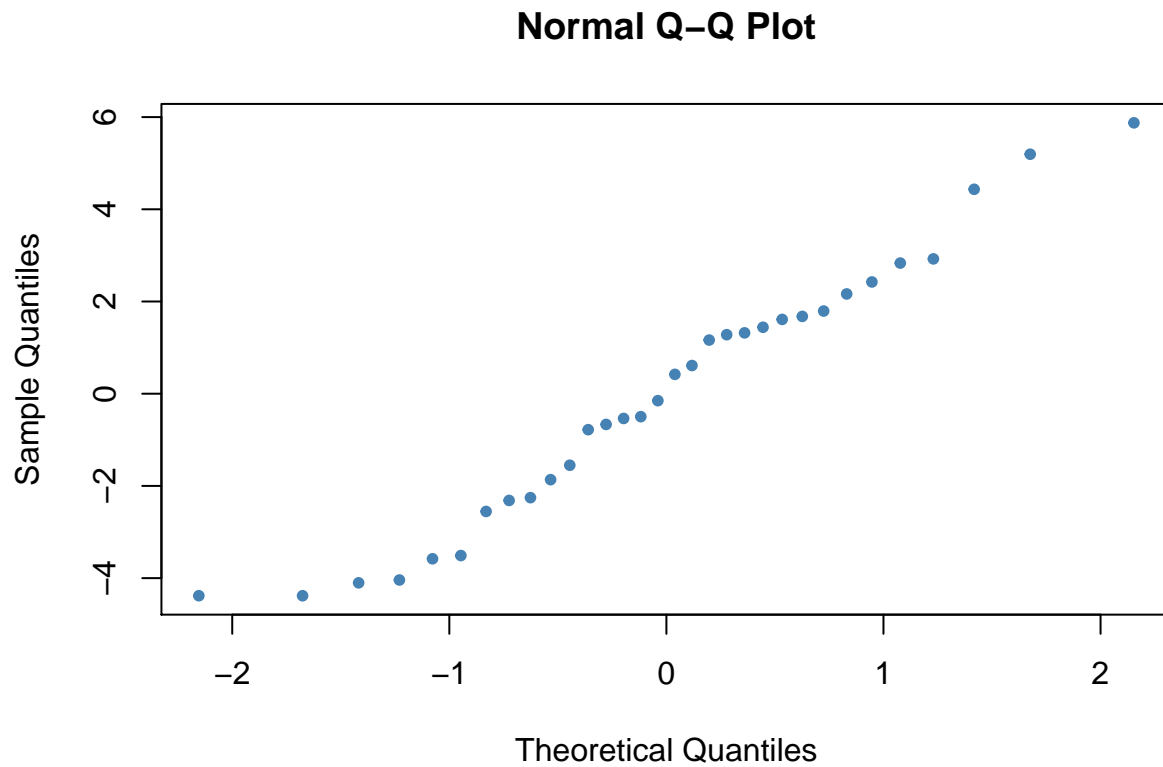
```
fitMod3<-lm(mpg~am2+hp,data=mtcars)
summary(fitMod3)
```

```
##
## Call:
## lm(formula = mpg ~ am2 + hp, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.384 -2.264  0.137  1.697  5.866
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  26.58491   1.42509   18.65 < 2e-16 ***
## am2Manual     5.27709   1.07954    4.89 3.5e-05 ***
## hp           -0.05889   0.00786   -7.50 2.9e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.91 on 29 degrees of freedom
## Multiple R-squared:  0.782, Adjusted R-squared:  0.767
## F-statistic: 52 on 2 and 29 DF, p-value: 2.55e-10
```

These results suggest support again that there is a difference in performance between automatic and manual transmission, estimated in 5.2771 miles per gallon, but in the case, we take into account a confounder that explains some of the residual variability.

We evaluated the assumption of normality with a *qq-plot* and it seems that there is no problems with normality.

```
qqnorm(fitMod2$residual,pch=20,col="steelblue")
```



Step 4. Compare models.

We estimated two models. An analysis of variance (ANOVA) shows that

```
kable(anova(fitMod1,fitMod3))
```

```
##
##
## | Res.Df |   RSS | Df | Sum of Sq |    F | Pr(>F) |
## |-----:|-----:|---:|-----: |-----:|-----: |
## |     30 | 720.9 | NA |    NA |    NA |    NA |
## |     29 | 245.4 |  1 |   475.5 | 56.18 |    0 |
```

Conclusions