

# Motor Trend magazine - Data analysis of influence on MPG for Automatic vs. Manual Transmission.

*by jmvilaverde*

*Thursday, June 18, 2015*

## Context

You work for Motor Trend, a magazine about the automobile industry. Looking at a data set of a collection of cars, they are interested in exploring the relationship between a set of variables and miles per gallon (MPG) (outcome). They are particularly interested in the following two questions:

“Is an automatic or manual transmission better for MPG” “Quantify the MPG difference between automatic and manual transmissions”

## Question

Take the mtcars data set and write up an analysis to answer their question using regression models and exploratory data analyses.

Your report must be:

Written as a PDF printout of a compiled (using knitr) R markdown document. Brief. Roughly the equivalent of 2 pages or less for the main text. Supporting figures in an appendix can be included up to 5 total pages including the 2 for the main report. The appendix can only include figures. Include a first paragraph executive summary. Upload your PDF by clicking the Upload button below the text box.

## Peer Grading

The criteria that your classmates will use to evaluate and grade your work are shown below. Each criteria is binary: (1 point = criteria met acceptably; 0 points = criteria not met acceptably) Your Course Project score will be the sum of the points and will count as 40% of your final grade in the course.

## Evaluation/feedback on the above work

- Did the student interpret the coefficients correctly?
- Did the student do some exploratory data analyses?
- Did the student fit multiple models and detail their strategy for model selection?
- Did the student answer the questions of interest or detail why the question(s) is (are) not answerable?
- Did the student do a residual plot and some diagnostics?
- Did the student quantify the uncertainty in their conclusions and/or perform an inference correctly?
- Was the report brief (about 2 pages long) for the main body of the report and no longer than 5 with supporting appendix of figures?
- Did the report include an executive summary?
- Was the report done in Rmd (knitr)?

## Report brief (about 2 pages)

### Executive summary (First paragraph)

Automatic/Manual transmission is better...

### Data adquisition

#### 1. Data adquisition and Initial structure analysis

```
data(mtcars)

#?mtcars

#Show colnames of mtcars
names(mtcars)
```

##### 1.1. Variables

```
## [1] "mpg" "cyl" "disp" "hp" "drat" "wt" "qsec" "vs" "am" "gear"
## [11] "carb"
```

**mtcars**

Format: A data frame with 32 observations on 11 variables.

Variable	Units
<b>mpg</b>	<b>Miles/(US) gallon</b>
cyl	Number of cylinders
disp	Displacement (cu.in.)
hp	Gross horsepower
drat	Rear axle ratio
wt	Weight (lb/1000)
qsec	1/4 mile time
vs	V/S
<b>am</b>	<b>Transmission (0 = automatic, 1 = manual)</b>
gear	Number of forward gears
carb	Number of carburetors

```
#Show structure of mtcars
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 : num 6 6 4 6 8 6 8 4 4 6 ...
## $ 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 : num 0 0 1 1 0 1 0 1 1 1 ...
## $ am : num 1 1 1 0 0 0 0 0 0 0 ...
## $ gear: num 4 4 4 3 3 3 3 4 4 4 ...
## $ carb: num 4 4 1 1 2 1 4 2 2 4 ...
```

**2.Initial analysis** Selected variables:

- Predictor:  $X = \text{am}$ , Transmission with values 0 for automatic, 1 for manual.
- Outcome:  $Y = \text{mpg}$ , Miles/(US) gallon

Linear regression model formula:

- $Y = \beta_0 + \beta_1 X \rightarrow \text{mpg} = \beta_0 + \beta_1 \text{am}$

Create initial model:

Calculate  $\beta_0$  (Intercept):  $\hat{\beta}_0 = \bar{Y} + \hat{\beta}_1 \bar{X}$

Calculate  $\beta_1$  (Slope):  $\hat{\beta}_1 = \text{Cor}(Y, X) \frac{Sd(Y)}{Sd(X)}$

Fitting the best line:

- $\sum_{i=1}^n (Y_i - (\beta_0 + \beta_1 X_i))^2$  to minimize the distance.

```
#Theoretical formula
y <- mtcars$mpg
x <- mtcars$am
beta1 <- cor(y,x)*(sd(y)/sd(x))
beta0 <- mean(y) - beta1 * mean(x)
#R function
model.Initial <- lm(y ~ x)
coeffs <- coef(model.Initial)

#Comparison
rbind(c(beta0, beta1), coeffs)
```

```
##      (Intercept)      x
##      17.14737 7.244939
## coeffs 17.14737 7.244939
```

For linear model  $Y = \beta_0 + \beta_1 X \rightarrow \text{mpg} = \beta_0 + \beta_1 \text{am}$

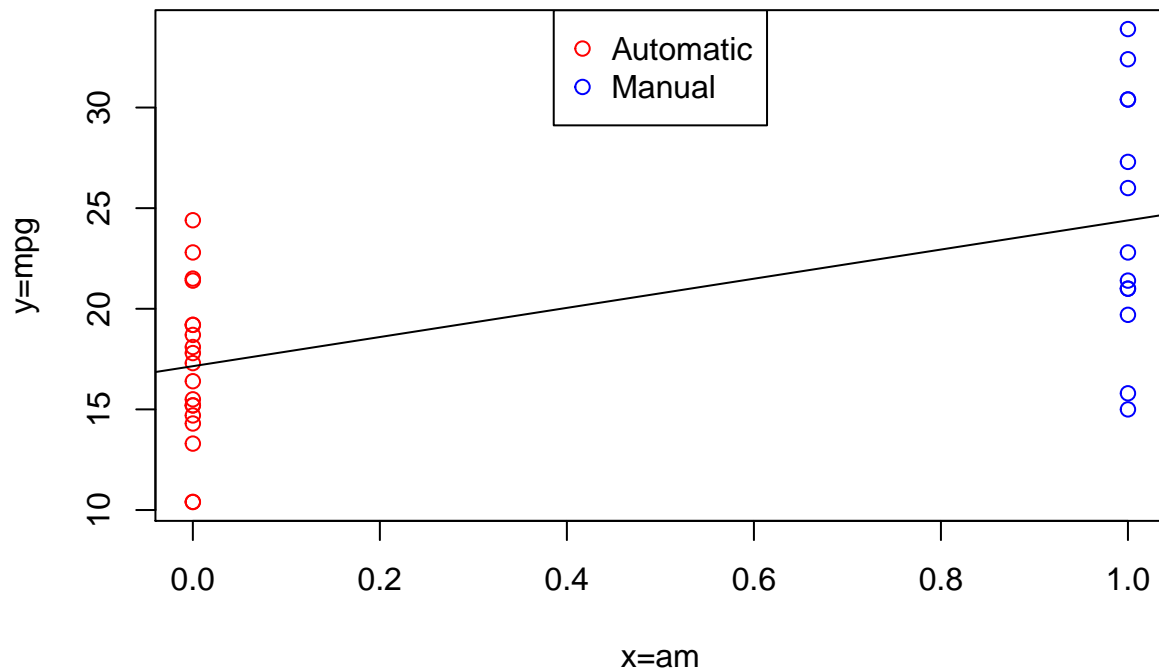
- Intercept  $\beta_0 = 17.1473684$
- Slope  $\beta_1 = 7.2449393$

For this linear regression model, the expected difference in MPG is 7.2449393 when the car have manual transmission in comparison to the same car with automatic transmission. The Intercept 17.1473684 is the expected MPG of a automatic transmission car.

```

color <- ifelse(x==0,"red","blue")
plot(x,y,col=color,xlab="x=am",ylab="y=mpg")
legend("top",c("Automatic","Manual"),col=c("red","blue"),pch=1)
abline(model.Initial)

```



View Figure 1 for plot with regression line.

### Normalizing data????

**3.Basic regression model with additive Gaussian errors.** Into point 2 is obtained a linear regression model that fits the data, but doesn't take in consideration the impact of the others variables. We are going to analyze gaussian errors for Initial Model.

Selected variables:

- **Predictor:** X = am, Transmission with values 0 for automatic, 1 for manual.
- **Outcome:** Y = mpg, Miles/(US) gallon

Probabilistic model for linear regression:

- $Y_i = \beta_0 + \beta_1 X_i + \epsilon_i \rightarrow mpg_i = \beta_0 + \beta_1 am_i + \epsilon_i$
- $\epsilon_i$  are assumed iid  $N(\mu_i, \sigma^2)$ .
- Note,  $E[Y_i|X_i = x_i] = \mu_i = \beta_0 + \beta_1 x_i$  and  $Var(Y_i|X_i = x_i) = \sigma^2$ .
- $\hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X}$  and  $\hat{\beta}_1 = Cor(Y, X) \frac{Sd(Y)}{Sd(X)}$ .

Residuals analysis:

- Observed outcome  $i$  is  $Y_i$  at a predictor value  $X_i$ .
- Predicted outcome  $i$  is  $\hat{Y}_i$  at a predictor value  $X_i$  is  $\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$ .
- Residual is the difference between observed and predicted:  $e_i = Y_i - \hat{Y}_i$ , the vertical distance between the observed data point and the regression line.
- Least squares minimizes  $\sum_{i=1}^n e_i^2$ .
- $e_i$  can be thought of as estimates of the  $\epsilon_i$ .

```
#Calculate residuals
```

```
e.ModelInitial <- y-(beta0+beta1*x)
e.ModelInitial
```

```
## [1] -3.3923077 -3.3923077 -1.5923077  4.2526316  1.5526316  0.9526316
## [7] -2.8473684  7.2526316  5.6526316  2.0526316  0.6526316 -0.7473684
## [13]  0.1526316 -1.9473684 -6.7473684 -6.7473684 -2.4473684  8.0076923
## [19]  6.0076923  9.5076923  4.3526316 -1.6473684 -1.9473684 -3.8473684
## [25]  2.0526316  2.9076923  1.6076923  6.0076923 -8.5923077 -4.6923077
## [31] -9.3923077 -2.9923077
```

```
#R function for residuals
```

```
resid(model.Initial)
```

```
##          1          2          3          4          5          6
## -3.3923077 -3.3923077 -1.5923077  4.2526316  1.5526316  0.9526316
##          7          8          9         10         11         12
## -2.8473684  7.2526316  5.6526316  2.0526316  0.6526316 -0.7473684
##         13         14         15         16         17         18
##  0.1526316 -1.9473684 -6.7473684 -6.7473684 -2.4473684  8.0076923
##         19         20         21         22         23         24
##  6.0076923  9.5076923  4.3526316 -1.6473684 -1.9473684 -3.8473684
##         25         26         27         28         29         30
##  2.0526316  2.9076923  1.6076923  6.0076923 -8.5923077 -4.6923077
##         31         32
## -9.3923077 -2.9923077
```

```
#Calculate predicted y (mpg)
```

```
yhat <- predict(model.Initial)
```

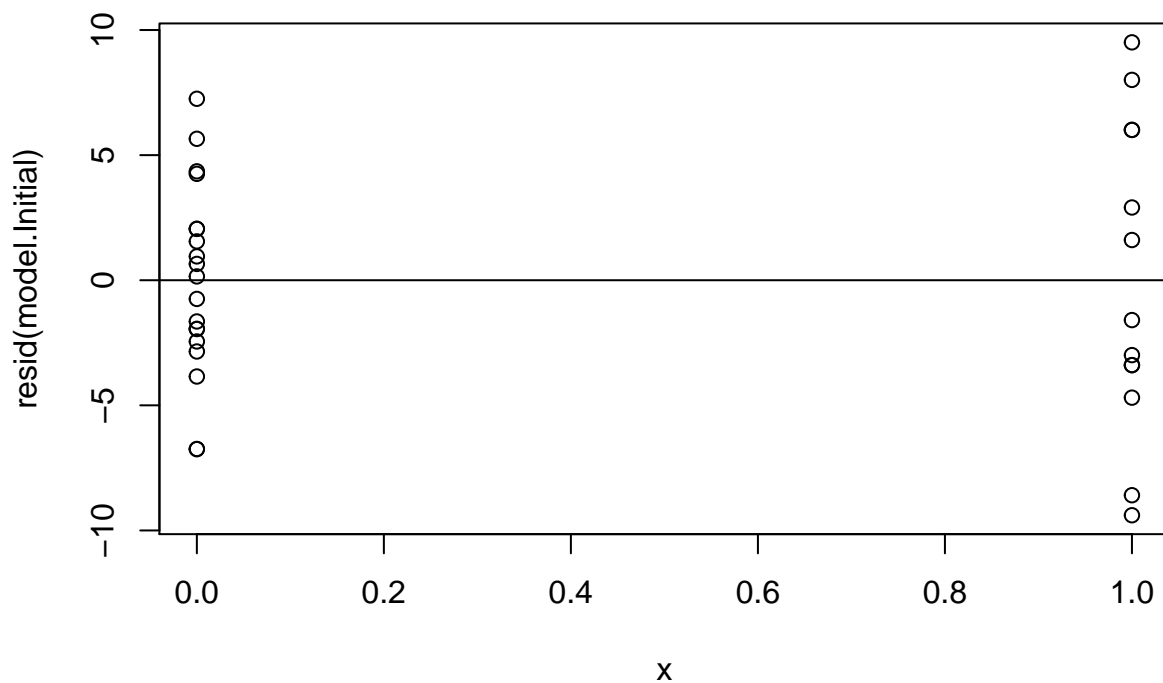
```
#Calculate max difference between residual and observed Y - predicted Y (y hat)
```

```
max(abs(e.ModelInitial -(y - yhat)))
```

```
## [1] 3.552714e-15
```

Plot of residuals:

```
plot(x, resid(model.Initial))
abline(h=0)
```



Formula to estimate residual variation:

- ML estimate of  $\sigma^2$  is  $\frac{1}{n} \sum_{i=1}^n e_i^2$
- For  $E[\hat{\sigma}^2] = \sigma^2$  most people use  $\frac{1}{n-2} \sum_{i=1}^n e_i^2$

```
#Calculate variation
n <- length(y)
var.e.Model.Initial <- sqrt((1/(n-2))*sum((e.ModelInitial^2)))
var.e.Model.Initial
```

```
## [1] 4.902029
```

```
#R function to calculate residual variation
summary(model.Initial)$sigma
```

```
## [1] 4.902029
```

Model Initial has a residual variation of 4.9020288.

**Total variation.** Formula Total variation = Residual variation + Regression Variation:

- $\sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 + \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$

Define the percent of total variation described by the model as:

- $R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = 1 - \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2}$

Relation between  $R^2$  and  $r$  (the correlation):

Recall that:  $(\hat{Y}_i - \bar{Y}) = \hat{\beta}_1(X_i - \bar{X})$  so that

- $R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \hat{\beta}_1^2 \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = Cor(Y, X)^2$

```
#Calculate R2
```

```
R2.Model.Initial <- sum((yhat - mean(y))^2)/sum((y-mean(y))^2)
R2.Model.Initial
```

```
## [1] 0.3597989
```

```
#R function for R2
```

```
cor(y,x)^2
```

```
## [1] 0.3597989
```

Inference in regression

Create confidence intervals and perform hypothesis tests.

```
#Calculation of coefficients
```

```
sigma <- var.e.Model.Initial
ssx <- sum((x - mean(x))^2)
seBeta0 <- (1/n + mean(x)^2/ssx) ^ 0.5 * sigma
seBeta1 <- sigma / sqrt(ssx)
tBeta0 <- beta0 / seBeta0
tBeta1 <- beta1 / seBeta1
pBeta0 <- 2 * pt(abs(tBeta0), df=n-2, lower.tail=FALSE)
pBeta1 <- 2 * pt(abs(tBeta1), df=n-2, lower.tail=FALSE)
coefTable <- rbind(c(beta0, seBeta0, tBeta0, pBeta0), c(beta1, seBeta1, tBeta1, pBeta1))
colnames(coefTable) <- c("Estimate", "Std.Error", "t value", "P(>|t|)")
rownames(coefTable) <- c("(Intercept)", "x")
coefTable
```

```
##              Estimate Std.Error   t value      P(>|t|)
## (Intercept) 17.147368  1.124603 15.247492 1.133983e-15
## x           7.244939  1.764422  4.106127 2.850207e-04
```

```
#R function for calculate coefficients
```

```
summary(model.Initial)$coef
```

```
##              Estimate Std. Error   t value      Pr(>|t|)
## (Intercept) 17.147368  1.124603 15.247492 1.133983e-15
## x           7.244939  1.764422  4.106127 2.850207e-04
```

Getting a confidence interval

```
sumCoef <- summary(model.Initial)$coef
confInterval.Model.Initial <- sumCoef[1,1] + c(-1,1) * qt(0.975, df=model.Initial$df) * sumCoef[1,2]
confInterval.Model.Initial
```

```
## [1] 14.85062 19.44411
```

With 95% confidence, we estimate that a conversion from automatic transmission to manual transmission results in a 14.8506236 to 19.4441132 increase in MPG for the Initial Model.

## Final Analysis for Model Initial

*This analysis doesn't take in consideration other variables that MUST be evaluated due the influence on relation between MPG and the type of transmission AM. Is required a Multivariable regression model.*

### Is an automatic or manual transmission better for MPG?

The manual transmission is better for MPG.

### Quantify the MPG difference between automatic and manual transmissions.

With 95% confidence, we estimate that a manual transmission results in a 14.8506236 to 19.4441132 increase in MPG comparing to use of automatic transmission for the Model Initial.

- Predictor:  $X = am$ , Transmission with values 0 for automatic, 1 for manual.
- Outcome:  $Y = mpg$ , Miles/(US) gallon

Linear regression model formula:  $Y = \beta_0 + \beta_1 X \rightarrow mpg = \beta_0 + \beta_1 am$

### Use of other models

This model lacks of ?precision?. is required to use additional variables from the original dataset.

For this, is needed the use of multivariable regression analyses.

**“Is an automatic or manual transmission better for MPG”**

**Main Body + Appendix only figures (not more than 5)**