

# Regression Model Project

sourabh jain

14/08/2020

## 1. Exploratory Analysis

```
data(mtcars)
head(mtcars)
```

```
##           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
## Hornet 4 Drive  21.4   6  258  110 3.08 3.215 19.44 1   0    3    1
## Hornet Sportabout 18.7   8  360  175 3.15 3.440 17.02 0   0    3    2
## Valiant        18.1   6  225  105 2.76 3.460 20.22 1   0    3    1
```

```
str(mtcars)
```

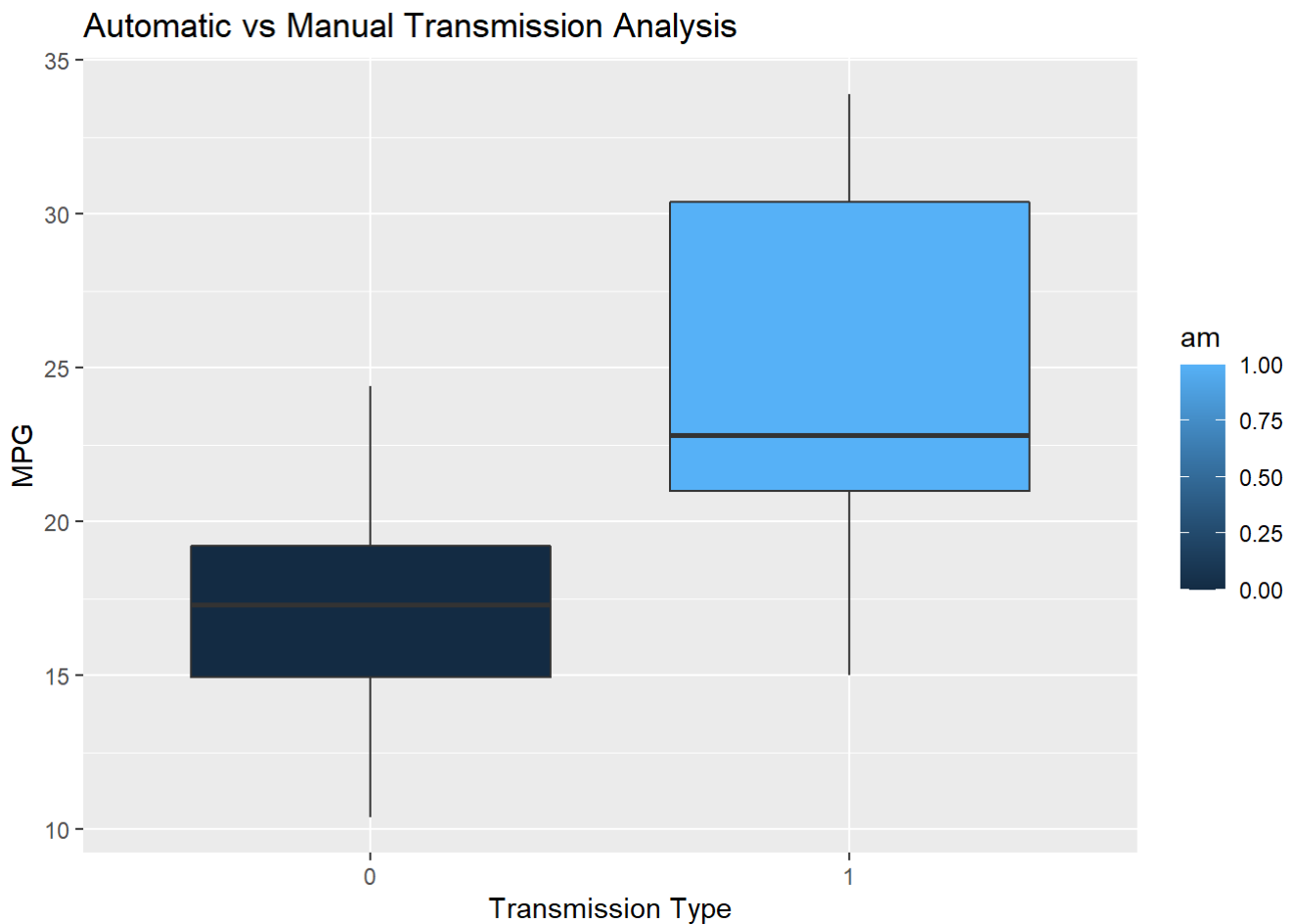
```
## 'data.frame':   32 obs. of  11 variables:
## $ mpg : num  21 21 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 ...
## $ cyl : num  6 6 4 6 8 6 8 4 4 6 ...
## $ disp: num  160 160 108 258 360 ...
## $ hp : num  110 110 93 110 175 105 245 62 95 123 ...
## $ drat: num  3.9 3.9 3.85 3.08 3.15 2.76 3.21 3.69 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 ...
```

```
summary(mtcars)
```

##	mpg	cyl	disp	hp
##	Min. :10.40	Min. :4.000	Min. : 71.1	Min. : 52.0
##	1st Qu.:15.43	1st Qu.:4.000	1st Qu.:120.8	1st Qu.: 96.5
##	Median :19.20	Median :6.000	Median :196.3	Median :123.0
##	Mean :20.09	Mean :6.188	Mean :230.7	Mean :146.7
##	3rd Qu.:22.80	3rd Qu.:8.000	3rd Qu.:326.0	3rd Qu.:180.0
##	Max. :33.90	Max. :8.000	Max. :472.0	Max. :335.0
##	drat	wt	qsec	vs
##	Min. :2.760	Min. :1.513	Min. :14.50	Min. :0.0000
##	1st Qu.:3.080	1st Qu.:2.581	1st Qu.:16.89	1st Qu.:0.0000
##	Median :3.695	Median :3.325	Median :17.71	Median :0.0000
##	Mean :3.597	Mean :3.217	Mean :17.85	Mean :0.4375
##	3rd Qu.:3.920	3rd Qu.:3.610	3rd Qu.:18.90	3rd Qu.:1.0000
##	Max. :4.930	Max. :5.424	Max. :22.90	Max. :1.0000
##	am	gear	carb	
##	Min. :0.0000	Min. :3.000	Min. :1.000	
##	1st Qu.:0.0000	1st Qu.:3.000	1st Qu.:2.000	
##	Median :0.0000	Median :4.000	Median :2.000	
##	Mean :0.4062	Mean :3.688	Mean :2.812	
##	3rd Qu.:1.0000	3rd Qu.:4.000	3rd Qu.:4.000	
##	Max. :1.0000	Max. :5.000	Max. :8.000	

## 2. Boxplot

```
library(ggplot2)
g<-ggplot(data=mtcars, aes(x=as.factor(am),y=mpg))+geom_boxplot(aes(fill=am))
g<-g+labs(title="Automatic vs Manual Transmission Analysis")+xlab("Transmission Type")
+ylab("MPG")
g
```



In this graph: 0 = Auto and 1 = Manual

### 3. Simple Regression Model Analysis

```
fit1<- lm(mpg~am, data=mtcars)
summary(fit1)
```

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

```
coef(fit1)
```

```
## (Intercept)          am  
## 17.147368      7.244939
```

```
confint(fit1)
```

```
##              2.5 %    97.5 %  
## (Intercept) 14.85062 19.44411  
## am          3.64151 10.84837
```

Given the results on this analysis we can conclude that the transmission type is statistically significant. Also we can say that a manual transmission will have on average a 7.2 MPG greater consumption compared to an automatic one. Finally the 95% confidence interval for this coefficient is (3.64, 10.84).

## 4. Multiple Regression Models

```
fit2<-lm(mpg~as.factor(am)+hp, data=mtcars)  
summary(fit2)
```

```
##  
## Call:  
## lm(formula = mpg ~ as.factor(am) + hp, data = mtcars)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -4.3843 -2.2642  0.1366  1.6968  5.8657   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept)  26.584914   1.425094  18.655 < 2e-16 ***  
## as.factor(am)1  5.277085   1.079541   4.888 3.46e-05 ***  
## hp          -0.058888   0.007857  -7.495 2.92e-08 ***  
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 2.909 on 29 degrees of freedom  
## Multiple R-squared:  0.782, Adjusted R-squared:  0.767   
## F-statistic: 52.02 on 2 and 29 DF, p-value: 2.55e-10
```

```
fit3<-lm(mpg~as.factor(am)+hp+wt, data=mtcars)  
summary(fit3)
```

```
##
## Call:
## lm(formula = mpg ~ as.factor(am) + hp + wt, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4221 -1.7924 -0.3788  1.2249  5.5317
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   34.002875    2.642659   12.867 2.82e-13 ***
## as.factor(am)1  2.083710    1.376420    1.514 0.141268
## hp            -0.037479    0.009605   -3.902 0.000546 ***
## wt            -2.878575    0.904971   -3.181 0.003574 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.538 on 28 degrees of freedom
## Multiple R-squared:  0.8399, Adjusted R-squared:  0.8227
## F-statistic: 48.96 on 3 and 28 DF,  p-value: 2.908e-11
```

It is worth mentioning that using the variable weight in the model makes the transmission type non significant, which suggests there is confounding between these variables

## 5. Model Selection

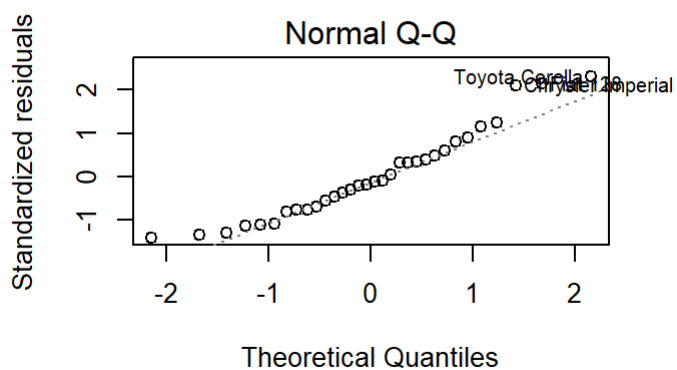
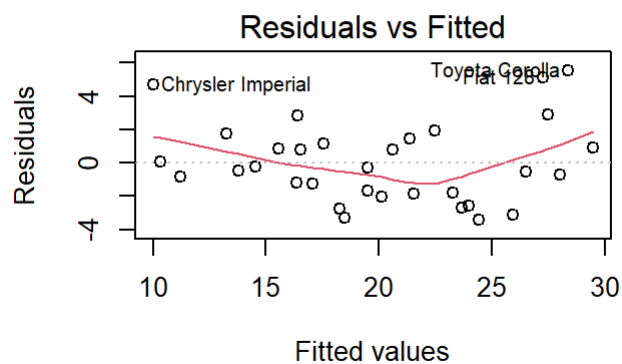
```
anova(fit1,fit2,fit3)
```

```
## Analysis of Variance Table
##
## Model 1: mpg ~ am
## Model 2: mpg ~ as.factor(am) + hp
## Model 3: mpg ~ as.factor(am) + hp + wt
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      30 720.90
## 2      29 245.44  1    475.46 73.841 2.445e-09 ***
## 3      28 180.29  1     65.15 10.118 0.003574 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

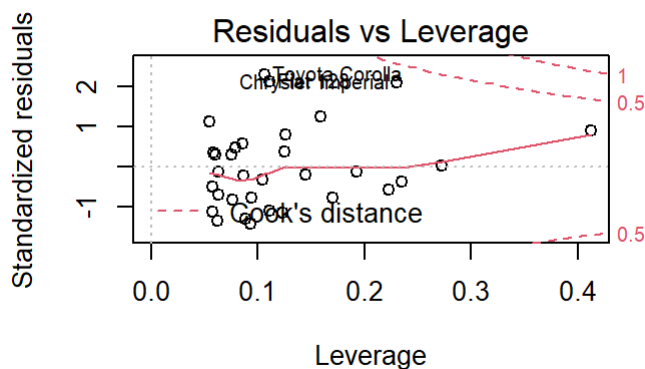
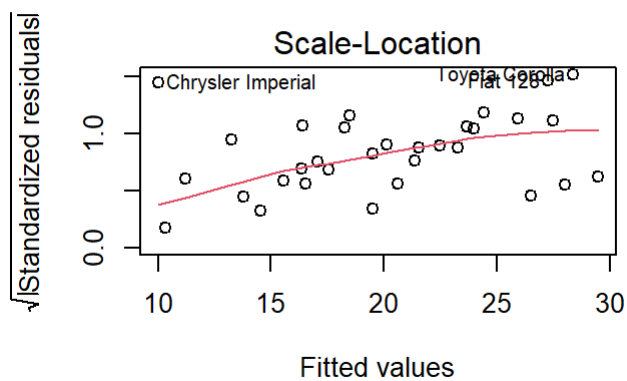
Using the RSS we can conclude that the model with the 3 variables is the best, but also we have to take into account that one of the variables is not significant. It is suggested that a new model is fit without the transmission type.

## 6. Residual Plots

```
par(mfrow=c(2,2))
plot(fit3)
```



This



analysis shows us that there is no heteroskedacity, but the residuals still show an underlying relationship in the data that is not taken into account in the model.