# Regression Model week-4 project

*Vasudha Singh*

*December 16, 2018*

## Overview/ Executive Summary

Here we have a dataset of Motor Trend, a magazine about the automobile industry. Looking at a data set of a collection of cars,we are interested in exploring the relationship between a set of variables and miles per gallon (MPG) (outcome). We are particularly interested in the following two questions:

1."Is an automatic or manual transmission better for MPG"
2."Quantify the MPG difference between automatic and manual transmissions"

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. A data frame consists 32 observations and 11 Variables

The analysis we have done, we can say that Manual Transmisson is better than Automatic transmission for MPG. the factors which included MPG can be multiple, by using regression method we can pick relatively right variables into our model .

## Loading And Processing the Data

```r
library(ggplot2)
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```r
library(datasets)
?mtcars
```

```
## starting httpd help server ...
```

```
##  done
```

```r
data(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 ...
```
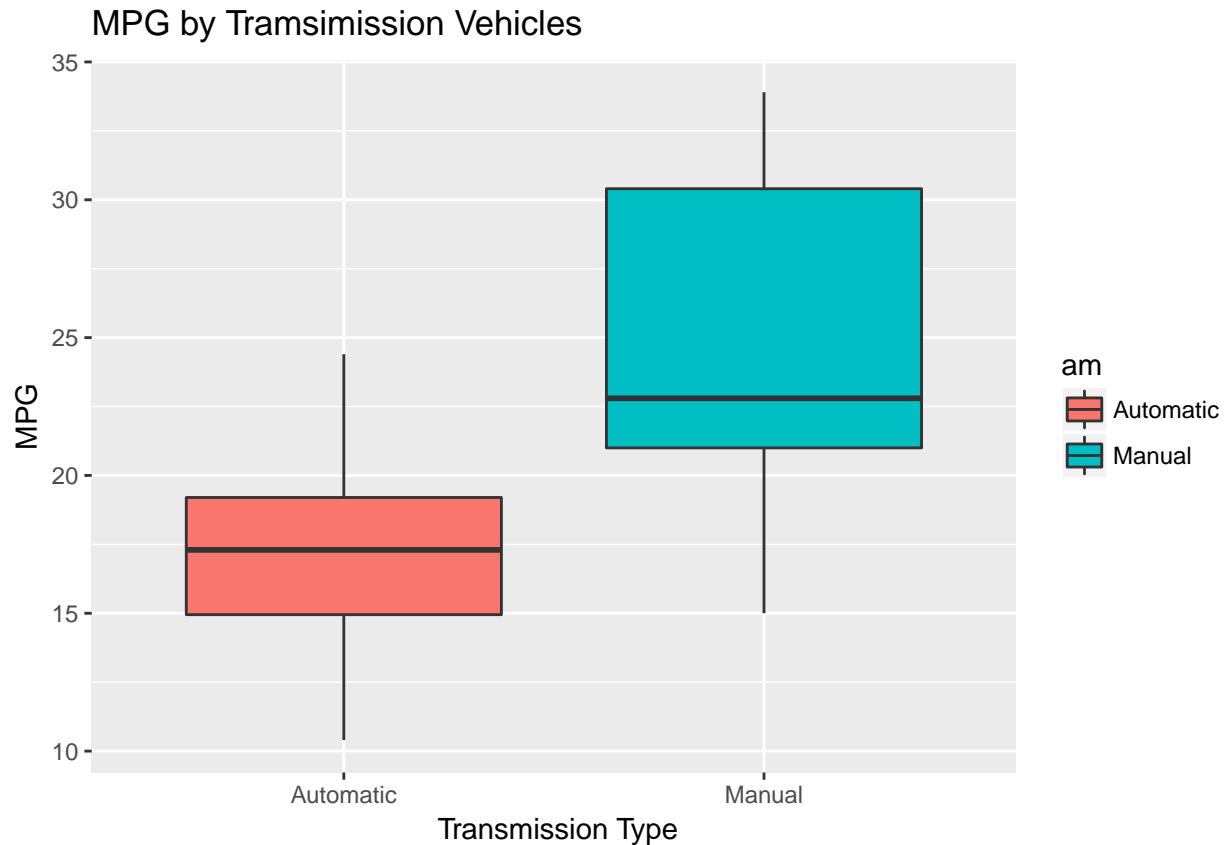
```r
summary(lm(mpg~. -1, mtcars))$coef
```

```
##          Estimate Std. Error     t value    Pr(>|t|)
## cyl   0.35082641 0.76292423   0.45984438 0.65014009
## disp  0.01354278 0.01762273   0.76848373 0.45037109
## hp   -0.02054767 0.02143989  -0.95838513 0.34828334
## drat  1.24158213 1.46276742   0.84878985 0.40513967
## wt   -3.82613150 1.86238084  -2.05443023 0.05200271
## qsec  1.19139689 0.45942323   2.59324480 0.01659185
## vs    0.18972068 2.06824861   0.09173011 0.92774262
## am    2.83222230 1.97512820   1.43394353 0.16564985
## gear  1.05426253 1.34668717   0.78285629 0.44205756
## carb -0.26321386 0.81235653  -0.32401273 0.74898869
```

```r
###Converting variables to Factor
mtcars$am<-factor(mtcars$am, labels=c("Automatic", "Manual"))
mtcars$cyl<-factor(mtcars$cyl)
mtcars$vs<- factor(mtcars$vs)
mtcars$gear<-factor(mtcars$gear)
mtcars$carb<-factor(mtcars$carb)
```

## Exploratory Analysis

```r
g<- ggplot(data=mtcars, aes(y=mpg, x=am, fill= am))
g<-g +geom_boxplot()+ xlab(" Transmission Type")+ ylab(" MPG ")+ labs(title="MPG by Tramsimission Vehicl
g
```

## MPG by Tramsimission Vehicles



The boxplot graph shows that there is a Significant increase in MPG(miles per gallon) with in vehicles via Automatic transmission Versus Manual transmission.

**Inference**

```
t_test<- t.test(mpg ~ am, mtcars)
t_test
```

```
##
##   Welch Two Sample t-test
##
## data:  mpg by am
## t = -3.7671, df = 18.332, p-value = 0.001374
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -11.280194  -3.209684
## sample estimates:
## mean in group Automatic    mean in group Manual
##               17.14737               24.39231
```

From the result we can say that, t_test **Reject** the **Null hypothesis** that there is a difference between transmission mean is equal to 0. i.e(p_value=0.001374)

And, the difference between the estimates of the two transmission is **7.245** MPG(miles per gallon), which tells the **Manual Transmission Vehicles** is **better** than Automatic Transmission Vehicles.
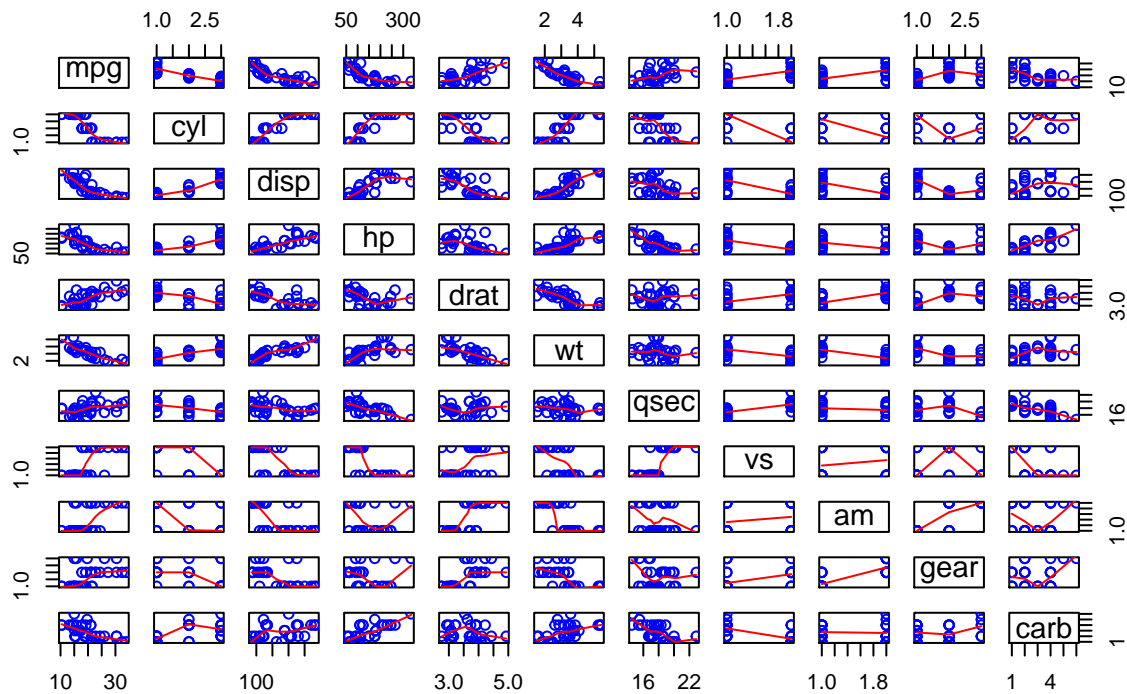
## Regression Analysis

**Fit the linear regression model for the Whole data.**

```
fit<- lm(mpg~. , mtcars)
summary(fit)$coeff
```

```
##                 Estimate  Std. Error      t value   Pr(>|t|)
## (Intercept) 23.87913244 20.06582026   1.19004018 0.25252548
## cyl6        -2.64869528  3.04089041  -0.87102622 0.39746642
## cyl8        -0.33616298  7.15953951  -0.04695316 0.96317000
## disp         0.03554632  0.03189920   1.11433290 0.28267339
## hp          -0.07050683  0.03942556  -1.78835344 0.09393155
## drat         1.18283018  2.48348458   0.47627845 0.64073922
## wt          -4.52977584  2.53874584  -1.78425732 0.09461859
## qsec         0.36784482  0.93539569   0.39325050 0.69966720
## vs1          1.93085054  2.87125777   0.67247551 0.51150791
## amManual     1.21211570  3.21354514   0.37718957 0.71131573
## gear4        1.11435494  3.79951726   0.29328856 0.77332027
## gear5        2.52839599  3.73635801   0.67670068 0.50889747
## carb2       -0.97935432  2.31797446  -0.42250436 0.67865093
## carb3        2.99963875  4.29354611   0.69863900 0.49546781
## carb4        1.09142288  4.44961992   0.24528452 0.80956031
## carb6        4.47756921  6.38406242   0.70136677 0.49381268
## carb8        7.25041126  8.36056638   0.86721532 0.39948495
```

```
pairs(mpg~. ,panel= panel.smooth, main="Mtcars Data", col=4, mtcars)
```

**Mtcars Data**

In this, we notice that none of the coefficient of the variables have p-value is less than 0.05, so we cannot say that there is which variable is Statistically Significant.To observe the intial relationship between the variables we plot scatterplot by plotting each variables against all other variables which tells the correlation between them.

**fitting the best model**

```
bestmodel<- step(fit, 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
```

```r
summary(bestmodel)$coeff
```

```
##               Estimate Std. Error   t value     Pr(>|t|)
## (Intercept) 33.70832390 2.60488618 12.940421 7.733392e-13
## cyl6        -3.03134449 1.40728351 -2.154040 4.068272e-02
## cyl8        -2.16367532 2.28425172 -0.947214 3.522509e-01
## hp          -0.03210943 0.01369257 -2.345025 2.693461e-02
## wt          -2.49682942 0.88558779 -2.819404 9.081408e-03
## amManual     1.80921138 1.39630450  1.295714 2.064597e-01
```
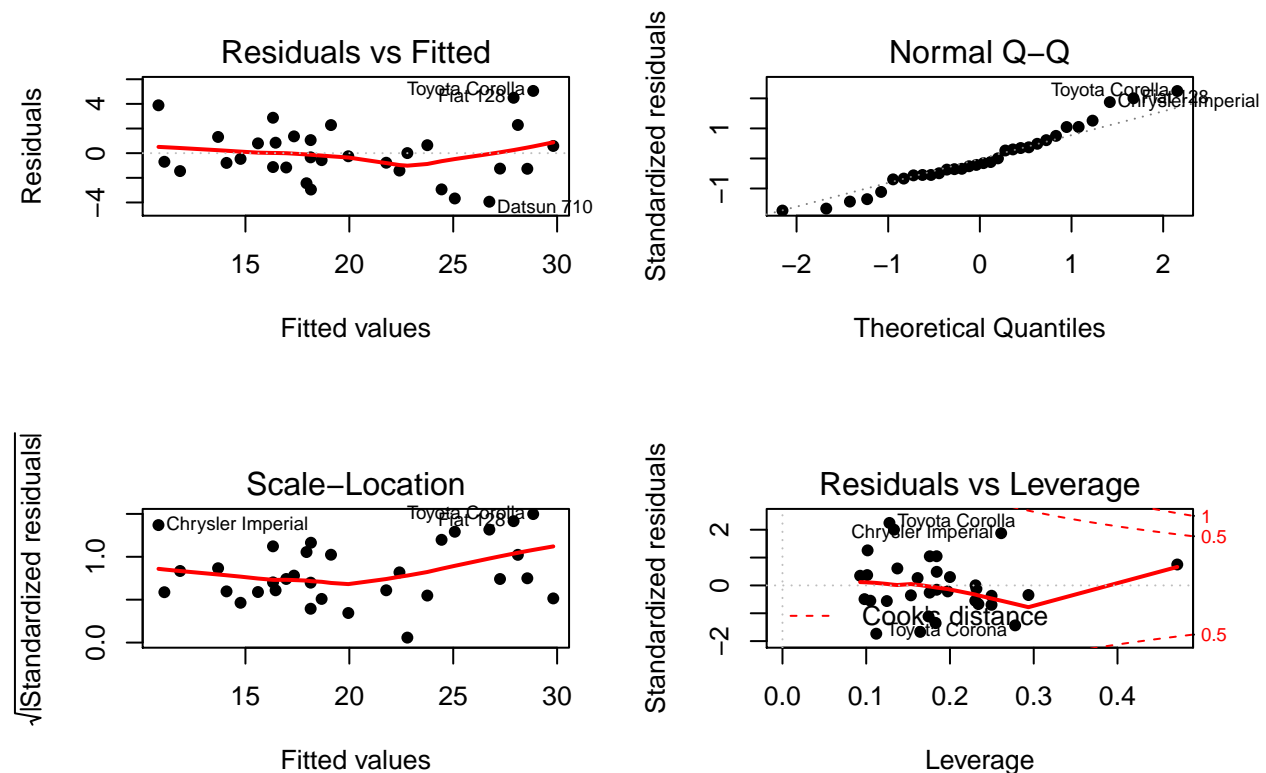
```r
confint(bestmodel,"amManual", level=0.95)
```

```
##                2.5 %   97.5 %
## amManual -1.060934 4.679356
```

In this model we can see, the bestmodel consists **4** variables(Cylinder, horsepower, weight, transmission) with R-squared value of **0.8659** which tells that best model consists the **87%** variance in MPG. in this also we have a P-Value which is statistically significant. The coefficient concluded that increasing a number of cylinder to 6 with decrease the MPG by **3.03** , increasing the number of cylinders from 6 to 8 define the decreases the MPG by **2.16** . Increases the horse power decreases the MPG by **3.21** per 100 horsepower. Weight decreases the MPG by **2.49** . Mannual transmission improves the MPG by **1.81** than Automatic transmission.

**Residual Plots and Diagnostics**

```r
par(mfrow=c(2,2))
plot(bestmodel, pch=16, lty=1, lwd=2)
```



1. Randomness of Residuals and fitted plot support the assumption of randomness.
2. The Normal Q-Q plot tells the distribution of residuals is normal.
3. Scale location plots tells the constant variance assumption.
4. The REsidual leverage plots tells that there is **no outliers** present.
   these models tells that there is no Hetroscadisticty.

## Conclusion

From the above analysis we can answer the above questions that,
1. Manual Transmisson Vehicles is better than Automatic Transmission Vehicles.

2. The result of multi variate model suggest that fuel efficeincyis 2.94 MPG higher for Manual over Automatic Transmission Vehicles.