

linear regression course

Mohamed Elashri

May 18, 2019

Task

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”

Data preparation

```
data(mtcars)
names(mtcars)
```

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

and this is the summary of dataset information

```
summary(cars)
```

```
##      speed      dist
##  Min.   : 4.0    Min.   : 2.00
##  1st Qu.:12.0    1st Qu.: 26.00
##  Median :15.0    Median : 36.00
##  Mean   :15.4    Mean   : 42.98
##  3rd Qu.:19.0    3rd Qu.: 56.00
##  Max.   :25.0    Max.   :120.00
```

There are eleven attributes in this dataset, we want to get what is the relation between MPG attribute and the other attributes.

Which is better ?

To get the relation we should first see the correlation between the MPG and the others, which is the first step. Thanks R that it is very easy to do this using built-in function `cor()`.

```
cor(mtcars$mpg,mtcars[, -1])
```

```
##           cyl      disp      hp      drat      wt      qsec
## [1,] -0.852162 -0.8475514 -0.7761684 0.6811719 -0.8676594 0.418684
##           vs      am      gear      carb
## [1,] 0.6640389 0.5998324 0.4802848 -0.5509251
```

we see that there are positive correlation which are (drat, qsec, vs, am, and gear) and negative correlation for the others (cyl, disp, hp, wt and carb).

Now to answer the transmission part we should first assign values to this categorical attribute (transmission type) as 0 -> automatic transmission type 1 -> manual transmission type

doing this with R as

```
mtcars$am <- as.factor(mtcars$am)
levels(mtcars$am) <-c("Automatic", "Manual")
```

we can see which is better from using boxplot (more about boxplot idea here <https://towardsdatascience.com/understanding-boxplots-5e2df7bcbd51>) and plot is in appendix 1

The plot tells us the manual have better MPG than the automatic .But we need to test this claim (we want to see if we can reject the null hypothesis).

perform t-test for example using R

```
t.test(mtcars$mpg~mtcars$am,conf.level=0.95)
```

```
##
##  Welch Two Sample t-test
##
## data:  mtcars$mpg by mtcars$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
```

We can see p-value of 0.001374 which gives us clue that we should reject the null hypothesis which is that there is not difference in MPG. This mean that the automatic cars has lower mpg than the manual cars.

The conculsion is based on the assumption that the automatic and manual cars other 10 characteristics are the same which is not accurate and to invistigate this we should apply multiple regression models to fit more attributes in our analysis (see Appendix 3 for details).

Now I want

MPG Difference Qualification

First, we can try to do a multivariate linear regression with all attributes (see appendix 3)

```
model = lm(data = mtcars, mpg ~ .)
summary(model)
```

```
##
## Call:
## lm(formula = mpg ~ ., data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4506 -1.6044 -0.1196  1.2193  4.6271
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  12.30337    18.71788   0.657   0.5181
## cyl         -0.11144     1.04502  -0.107   0.9161
## disp          0.01334     0.01786   0.747   0.4635
## hp          -0.02148     0.02177  -0.987   0.3350
## drat          0.78711     1.63537   0.481   0.6353
## wt          -3.71530     1.89441  -1.961   0.0633 .
## qsec          0.82104     0.73084   1.123   0.2739
```

```
## vs          0.31776    2.10451    0.151    0.8814
## amManual    2.52023    2.05665    1.225    0.2340
## gear        0.65541    1.49326    0.439    0.6652
## carb       -0.19942    0.82875   -0.241    0.8122
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.65 on 21 degrees of freedom
## Multiple R-squared:  0.869, Adjusted R-squared:  0.8066
## F-statistic: 13.93 on 10 and 21 DF,  p-value: 3.793e-07
```

If we look to the coefficients we observe that wt is the only attribute changing with our mpg so we don't use all variables to fit our model because it will probably result in overfitting problem. So we need to test different models with different exploratory variables.

We use R function step() which is this automatic model-choosing function that choose the best linear regression model

```
bestmodel = step(lm(data = mtcars, mpg ~ .), trace=0)
summary(bestmodel)
```

```
##
## Call:
## lm(formula = mpg ~ wt + qsec + am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.4811 -1.5555 -0.7257  1.4110  4.6610
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   9.6178     6.9596   1.382 0.177915
## wt           -3.9165     0.7112  -5.507 6.95e-06 ***
## qsec          1.2259     0.2887   4.247 0.000216 ***
## amManual      2.9358     1.4109   2.081 0.046716 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.459 on 28 degrees of freedom
## Multiple R-squared:  0.8497, Adjusted R-squared:  0.8336
## F-statistic: 52.75 on 3 and 28 DF,  p-value: 1.21e-11
```

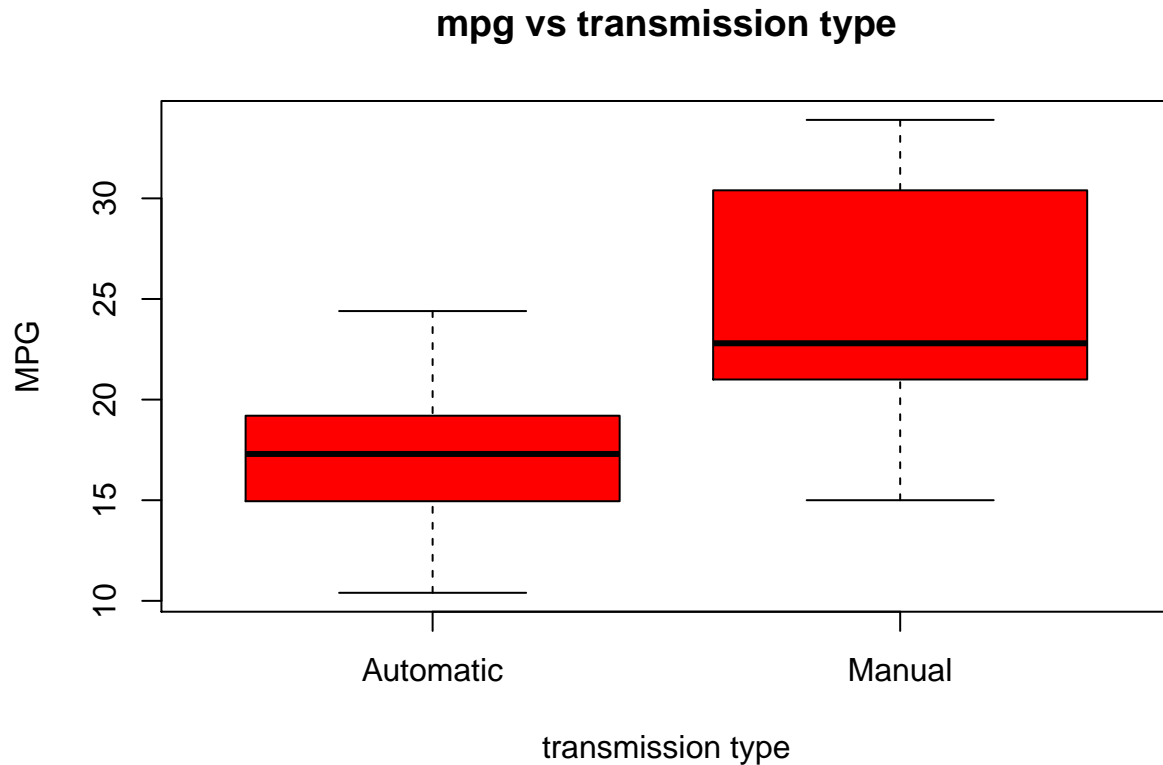
looks like the best model is the one that includes wt, qsec and am, So they also needs to be considered. wt negatively changes with mpg, and qsec and am positively changes. Every kg/1000 weight increase will cause a decrease of roughly 4 mpg, every increase of 1/4 mile time will cause an increase of 1.2 mpg, and on average, manual transmission is 2.9 mpg better than automatic transmission. The model is able to explain 85% of variance. The residual plots also seems to be randomly scattered (see appendix 4).

The colclusion

Based on the previous analysis, we can say that on average manual transmission is better than automatic transmission by 2.9 mpg but also transmission type is not the only factor accounting for MPG, weight, and acceleration (1/4 mile time) also needs to be considered.

Appendix 1 (Barplot and vilion plot)

```
boxplot(mtcars$mpg ~ mtcars$am, data = mtcars, outpch = 19, ylab="MPG", xlab="transmission type", main="mpg vs transmission type")
```

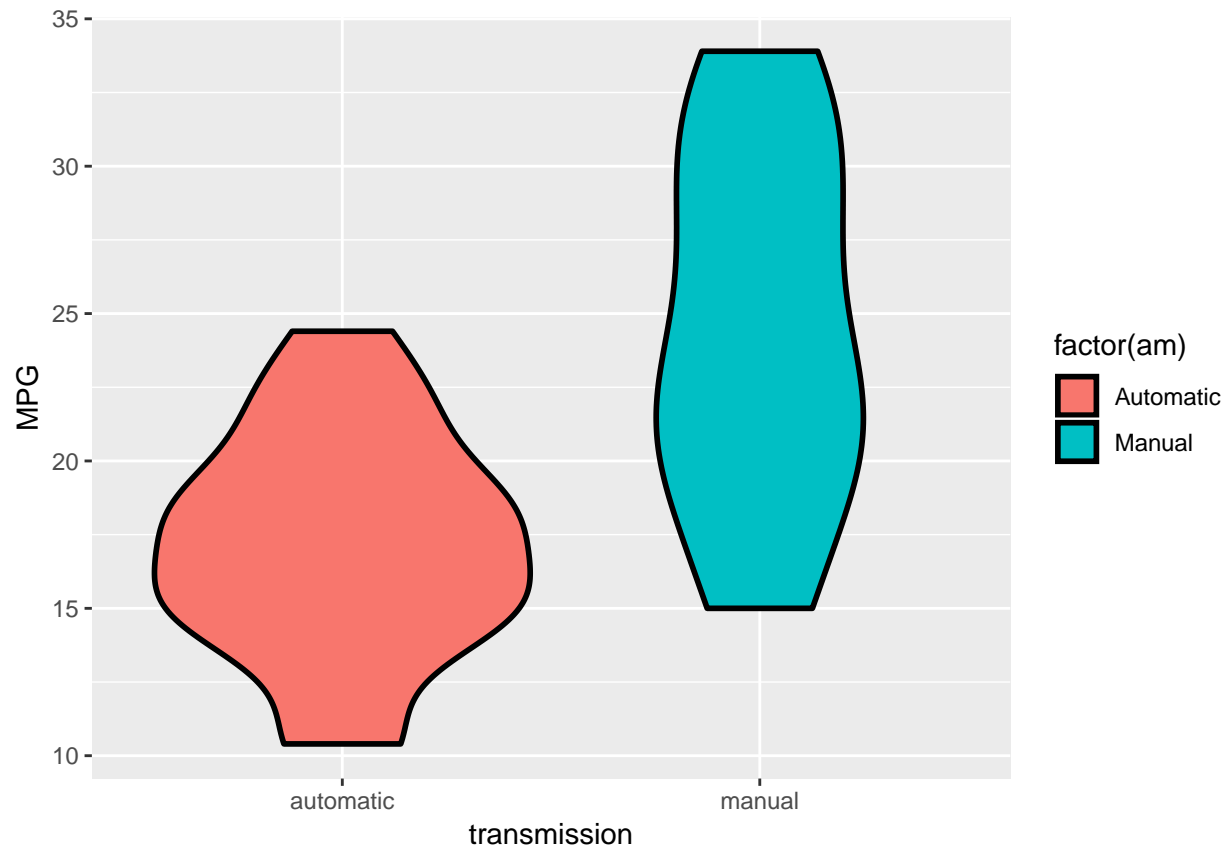


We can also get same information from vilion plot.

```
library(stats)
library(ggplot2)
```

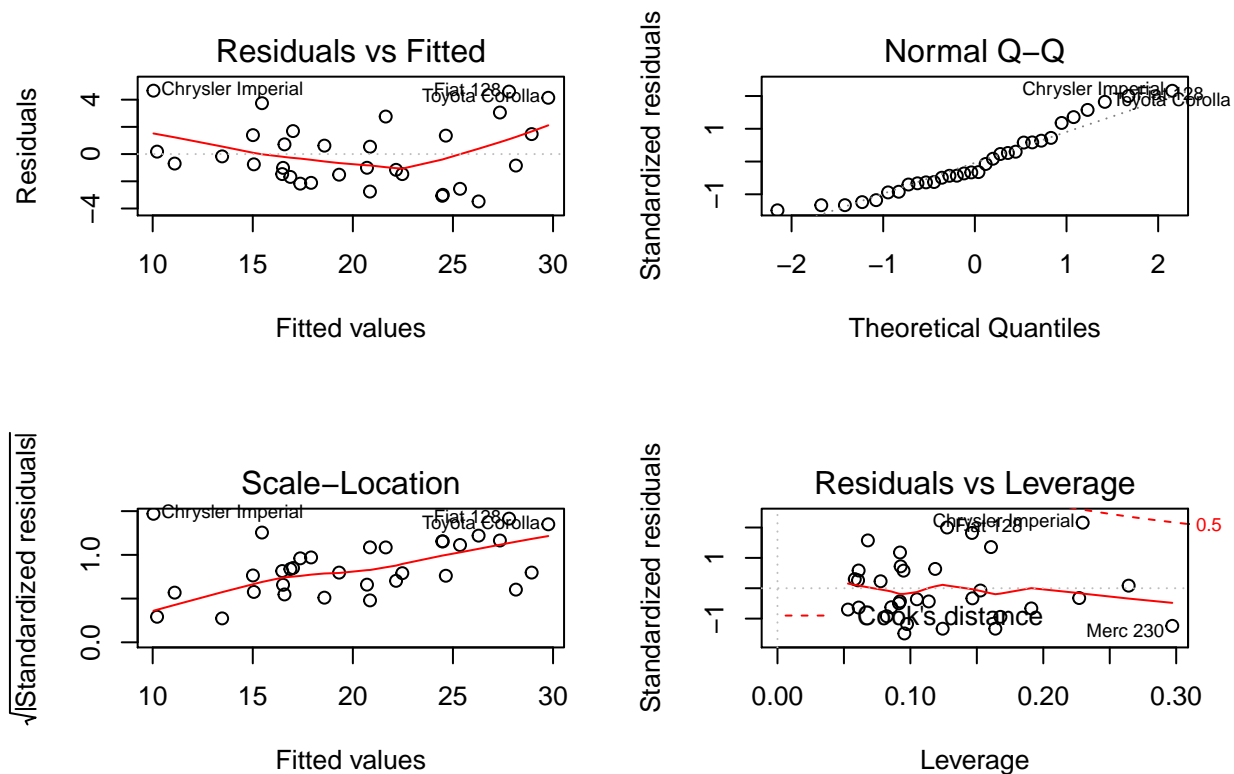
```
## Warning: package 'ggplot2' was built under R version 3.5.3
```

```
ggplot(mtcars, aes(y=mpg, x=factor(am, labels = c("automatic", "manual")), fill=factor(am)))+
  geom_violin(colour="black", size=1)+
  xlab("transmission") + ylab("MPG")
```



Appendix 2 (Residual Plots)

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



Appendix 3 (Regression Model Results)

comparison between best model to fit the dataset and simple regression model that we used.

```
full.model <- lm(mpg ~ ., data = mtcars)
best.model <- step(full.model, direction = "backward")

## Start: AIC=70.9
## mpg ~ cyl + disp + hp + drat + wt + qsec + vs + am + gear + carb
##
##      Df Sum of Sq  RSS   AIC
## - cyl    1    0.0799 147.57 68.915
## - vs     1    0.1601 147.66 68.932
## - carb    1    0.4067 147.90 68.986
## - gear    1    1.3531 148.85 69.190
## - drat    1    1.6270 149.12 69.249
## - disp    1    3.9167 151.41 69.736
## - hp      1    6.8399 154.33 70.348
## - qsec    1    8.8641 156.36 70.765
## <none>                 147.49 70.898
## - am      1   10.5467 158.04 71.108
## - wt      1   27.0144 174.51 74.280
##
## Step: AIC=68.92
## mpg ~ disp + hp + drat + wt + qsec + vs + am + gear + carb
##
##      Df Sum of Sq  RSS   AIC
```

```

## - vs      1      0.2685 147.84 66.973
## - carb    1      0.5201 148.09 67.028
## - gear     1      1.8211 149.40 67.308
## - drat     1      1.9826 149.56 67.342
## - disp     1      3.9009 151.47 67.750
## - hp       1      7.3632 154.94 68.473
## <none>                147.57 68.915
## - qsec     1     10.0933 157.67 69.032
## - am       1     11.8359 159.41 69.384
## - wt       1     27.0280 174.60 72.297
##
## Step:  AIC=66.97
## mpg ~ disp + hp + drat + wt + qsec + am + gear + carb
##
##      Df Sum of Sq    RSS    AIC
## - carb  1      0.6855 148.53 65.121
## - gear  1      2.1437 149.99 65.434
## - drat  1      2.2139 150.06 65.449
## - disp  1      3.6467 151.49 65.753
## - hp    1      7.1060 154.95 66.475
## <none>                147.84 66.973
## - am    1     11.5694 159.41 67.384
## - qsec  1     15.6830 163.53 68.200
## - wt    1     27.3799 175.22 70.410
##
## Step:  AIC=65.12
## mpg ~ disp + hp + drat + wt + qsec + am + gear
##
##      Df Sum of Sq    RSS    AIC
## - gear  1      1.565 150.09 63.457
## - drat  1      1.932 150.46 63.535
## <none>                148.53 65.121
## - disp  1     10.110 158.64 65.229
## - am    1     12.323 160.85 65.672
## - hp    1     14.826 163.35 66.166
## - qsec  1     26.408 174.94 68.358
## - wt    1     69.127 217.66 75.350
##
## Step:  AIC=63.46
## mpg ~ disp + hp + drat + wt + qsec + am
##
##      Df Sum of Sq    RSS    AIC
## - drat  1      3.345 153.44 62.162
## - disp  1      8.545 158.64 63.229
## <none>                150.09 63.457
## - hp    1     13.285 163.38 64.171
## - am    1     20.036 170.13 65.466
## - qsec  1     25.574 175.67 66.491
## - wt    1     67.572 217.66 73.351
##
## Step:  AIC=62.16
## mpg ~ disp + hp + wt + qsec + am
##
##      Df Sum of Sq    RSS    AIC

```

```

## - disp 1      6.629 160.07 61.515
## <none>      153.44 62.162
## - hp 1      12.572 166.01 62.682
## - qsec 1     26.470 179.91 65.255
## - am 1      32.198 185.63 66.258
## - wt 1      69.043 222.48 72.051
##
## Step: AIC=61.52
## mpg ~ hp + wt + qsec + am
##
##      Df Sum of Sq  RSS    AIC
## - hp  1      9.219 169.29 61.307
## <none>      160.07 61.515
## - qsec 1     20.225 180.29 63.323
## - am  1     25.993 186.06 64.331
## - wt  1     78.494 238.56 72.284
##
## Step: AIC=61.31
## mpg ~ wt + qsec + am
##
##      Df Sum of Sq  RSS    AIC
## <none>      169.29 61.307
## - am  1     26.178 195.46 63.908
## - qsec 1    109.034 278.32 75.217
## - wt  1    183.347 352.63 82.790

```