

Regression Models Course Project

aps2201

February 13, 2016

Overview

In this project we want to see two things:

- Is an automatic or manual transmission better for MPG?
- The MPG difference between automatic and manual transmissions

First, we need to look at the `mtcars` dataset

```
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          am
##  Min.   :2.760   Min.    :1.513   Min.    :14.50   V:18   automatic:19
## 1st Qu.:3.080   1st Qu.:2.581   1st Qu.:16.89   S:14   manual    :13
## Median :3.695   Median :3.325   Median :17.71
## Mean   :3.597   Mean    :3.217   Mean    :17.85
## 3rd Qu.:3.920   3rd Qu.:3.610   3rd Qu.:18.90
## Max.   :4.930   Max.    :5.424   Max.    :22.90
##      gear          carb
##  Min.   :3.000   Min.    :1.000
## 1st Qu.:3.000   1st Qu.:2.000
## Median :4.000   Median :2.000
## Mean   :3.688   Mean    :2.812
## 3rd Qu.:4.000   3rd Qu.:4.000
## Max.   :5.000   Max.    :8.000
```

The `vs` and `am` are not supposed to be numeric, since they are actually codes for V/S and automatic/manual respectively.

Lets fix that

```
mtcars$vs=sub("0","V",mtcars$vs)
mtcars$vs=sub("1","S",mtcars$vs)

mtcars$am=sub("0","automatic",mtcars$am)
mtcars$am=sub("1","manual",mtcars$am)

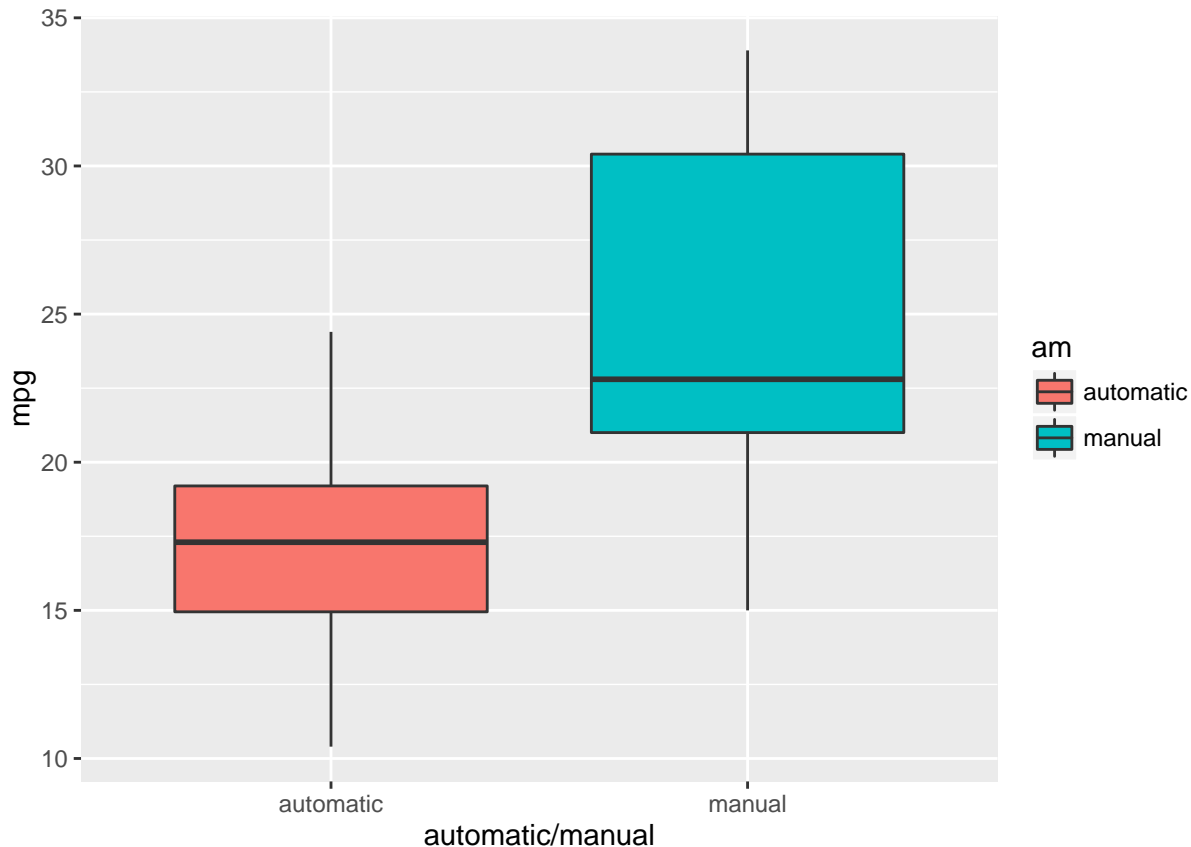
mtcars$vs=as.factor(mtcars$vs)
mtcars$am=as.factor(mtcars$am)
```

```
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          am
##  Min.    :2.760   Min.    :1.513   Min.    :14.50   S:14   automatic:19
## 1st Qu.:3.080   1st Qu.:2.581   1st Qu.:16.89   V:18   manual    :13
## Median :3.695   Median :3.325   Median :17.71
## Mean    :3.597   Mean    :3.217   Mean    :17.85
## 3rd Qu.:3.920   3rd Qu.:3.610   3rd Qu.:18.90
## Max.    :4.930   Max.    :5.424   Max.    :22.90
##      gear          carb
##  Min.    :3.000   Min.    :1.000
## 1st Qu.:3.000   1st Qu.:2.000
## Median :4.000   Median :2.000
## Mean    :3.688   Mean    :2.812
## 3rd Qu.:4.000   3rd Qu.:4.000
## Max.    :5.000   Max.    :8.000
```

ok, now we need to answer the question *Is an automatic or manual transmission better for MPG?*, to do this we can explore this by plotting the mpg with the transmission type.

```
amformpg=ggplot(aes(x=am,y=mpg),data=mtcars)+
  geom_boxplot(aes(fill=am))+
  xlab("automatic/manual")
amformpg
```



Now we can see that from miles per gallon (mpg) automatic seems to be the more gas guzzling compared to manual. Hence, manual is better compared to automatic regarding mpg.

Ok, this gives us a general idea, we need to look at the correlation table to see how it actually correlates, for this we need to convert the factors back to numeric for the `cor()` function to read.

```
mtcars$vs=as.numeric(mtcars$vs)
mtcars$am=as.numeric(mtcars$am)
cor(mtcars)
```

```
##           mpg           cyl           disp           hp           drat           wt
## mpg      1.0000000 -0.8521620 -0.8475514 -0.7761684  0.68117191 -0.8676594
## cyl     -0.8521620  1.0000000  0.9020329  0.8324475 -0.69993811  0.7824958
## disp    -0.8475514  0.9020329  1.0000000  0.7909486 -0.71021393  0.8879799
## hp      -0.7761684  0.8324475  0.7909486  1.0000000 -0.44875912  0.6587479
## drat     0.6811719 -0.6999381 -0.7102139 -0.4487591  1.00000000 -0.7124406
## wt      -0.8676594  0.7824958  0.8879799  0.6587479 -0.71244065  1.0000000
## qsec     0.4186840 -0.5912421 -0.4336979 -0.7082234  0.09120476 -0.1747159
## vs      -0.6640389  0.8108118  0.7104159  0.7230967 -0.44027846  0.5549157
## am       0.5998324 -0.5226070 -0.5912270 -0.2432043  0.71271113 -0.6924953
## gear     0.4802848 -0.4926866 -0.5555692 -0.1257043  0.69961013 -0.5832870
## carb    -0.5509251  0.5269883  0.3949769  0.7498125 -0.09078980  0.4276059
##           qsec           vs           am           gear           carb
## mpg      0.41868403 -0.6640389  0.59983243  0.4802848 -0.55092507
## cyl     -0.59124207  0.8108118 -0.52260705 -0.4926866  0.52698829
```

```
## disp -0.43369788  0.7104159 -0.59122704 -0.5555692  0.39497686
## hp   -0.70822339  0.7230967 -0.24320426 -0.1257043  0.74981247
## drat  0.09120476 -0.4402785  0.71271113  0.6996101 -0.09078980
## wt   -0.17471588  0.5549157 -0.69249526 -0.5832870  0.42760594
## qsec  1.00000000 -0.7445354 -0.22986086 -0.2126822 -0.65624923
## vs   -0.74453544  1.0000000 -0.16834512 -0.2060233  0.56960714
## am   -0.22986086 -0.1683451  1.00000000  0.7940588  0.05753435
## gear -0.21268223 -0.2060233  0.79405876  1.0000000  0.27407284
## carb -0.65624923  0.5696071  0.05753435  0.2740728  1.00000000
```

So, now that we have proof on our assumption, we need to fit a regression model to the correlation.

By the looks of of the correlation table there are some outstanding numbers, remember we are just looking for mpg difference for automatic and manual.

```
initialmodel <- lm(mpg ~ ., data = mtcars)
bestmodel <- step(initialmodel, direction = "both")
```

```
## 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
## + cyl      1     0.0799 147.49 70.898
## - 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
## + vs 1 0.2685 147.57 68.915
## + cyl 1 0.1883 147.66 68.932
## - 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
## + carb 1 0.685 147.84 66.973
## + vs 1 0.434 148.09 67.028
## + cyl 1 0.414 148.11 67.032
## - 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
## + gear 1 1.565 148.53 65.121
## + cyl 1 1.003 149.09 65.242
## + vs 1 0.645 149.45 65.319
## + carb 1 0.107 149.99 65.434
## - 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
## + drat 1 3.345 150.09 63.457
## + gear 1 2.977 150.46 63.535
## + cyl 1 2.447 150.99 63.648

```

```
## + vs      1      1.121 152.32 63.927
## + carb    1      0.011 153.43 64.160
## - 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
## + disp  1      6.629 153.44 62.162
## + carb  1      3.227 156.84 62.864
## + drat  1      1.428 158.64 63.229
## - qsec  1     20.225 180.29 63.323
## + cyl   1      0.249 159.82 63.465
## + vs    1      0.249 159.82 63.466
## + gear  1      0.171 159.90 63.481
## - 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
## + hp    1      9.219 160.07 61.515
## + carb  1      8.036 161.25 61.751
## + disp  1      3.276 166.01 62.682
## + cyl   1      1.501 167.78 63.022
## + drat  1      1.400 167.89 63.042
## + gear  1      0.123 169.16 63.284
## + vs    1      0.000 169.29 63.307
## - am    1     26.178 195.46 63.908
## - qsec  1    109.034 278.32 75.217
## - wt    1    183.347 352.63 82.790
```

```
basemodel <- lm(mpg ~ am, data = mtcars)
anova(basemodel, bestmodel)
```

```
## Analysis of Variance Table
```

```
##
```

```
## Model 1: mpg ~ am
```

```
## Model 2: mpg ~ wt + qsec + am
```

```
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
```

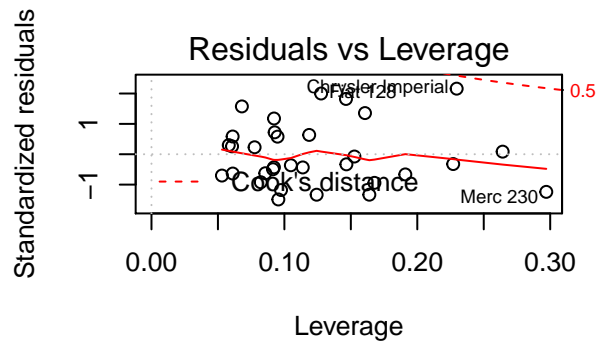
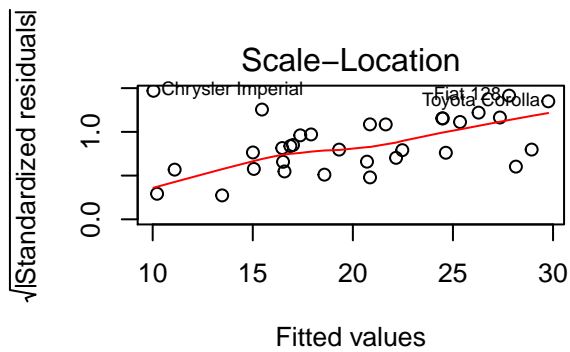
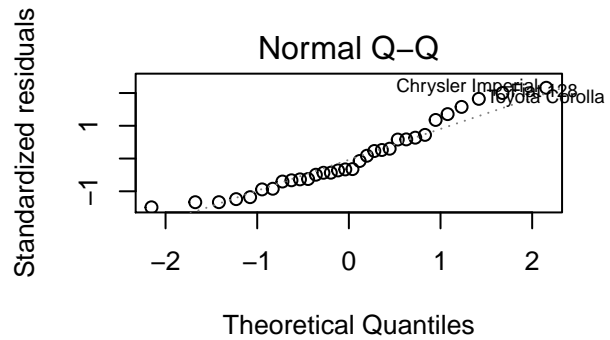
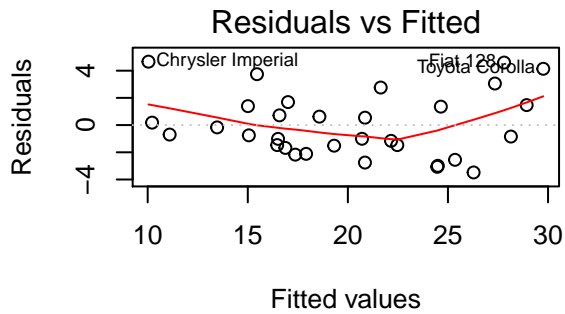
```
## 1      30 720.90
```

```
## 2      28 169.29  2    551.61 45.618 1.55e-09 ***
```

```
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

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



```
t.test(mpg ~ am, data = mtcars)
```

```
##
## 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 1 mean in group 2
## 17.14737 24.39231
```