

Peer-graded Assignment: Regression Models Course Project

anckur khaitan

MOTOR TREND: ROAD TESTS THE EFFECTS OF TRANSMISSION ON MPG

Executive Summary

This is the course project for the Regression Models course offered by Johns Hopkins University through the Coursera specialization. For this assignment, we will analyze the “mtcars” data set in order to explore the relationships between a set of variables and the miles per gallon (MPG).

Our objectives for this research are:

- Is an automatic or manual transmission better for MPG?
- Quantifying how different is the MPG between automatic and manual transmissions?

The results from our analysis were:

- Manual transmission is better for MPG by a factor of 1.8 compared to automatic transmission.
- Means and medians for automatic and manual transmission cars are significantly different.

Set the WD and Load Knitr

Load our working directory and apply Knitr package

```
library(knitr)
```

Data Processing and Transformation

Load the data set, perform the transformations below and factor the necessary variables. We will look at the results of this data in the following section

```
data(mtcars)
mtcars$cyl <- factor(mtcars$cyl)
mtcars$vs <- factor(mtcars$vs)
mtcars$gear <- factor(mtcars$gear)
mtcars$carb <- factor(mtcars$carb)
mtcars$am <- factor(mtcars$am, labels=c('Automatic', 'Manual'))
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 : Factor w/ 3 levels "4","6","8": 2 2 1 2 3 2 3 1 1 2 ...
## $ 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 : Factor w/ 2 levels "0","1": 1 1 2 2 1 2 1 2 2 2 ...
## $ am : Factor w/ 2 levels "Automatic","Manual": 2 2 2 1 1 1 1 1 1 1 ...
## $ gear: Factor w/ 3 levels "3","4","5": 2 2 2 1 1 1 1 2 2 2 ...
## $ carb: Factor w/ 6 levels "1","2","3","4",...: 4 4 1 1 2 1 4 2 2 4 ...
```

Exploratory Data Analysis

Exploring relationships between our data and variables of interest. First, we plot relationships between all variables in the data set (Figure 2 in Appendix). Notice variables: drat, hp, cyl, disp, wt, vs and am have a strong correlation with mpg. We will use linear models 2 to quantify that in the regression analysis.

Use boxplots of the variable mpg when am is Automatic or Manual (Figure 1 in the Appendix). An increase in mpg happens when the transmission is manual.

Regression Analysis

Building linear regression models bason on different variables. We are trying to find the best model fit and compare it with out base model

Model Building and Selection

Based on the pairs plot, where many variables had a high correlation with mpg, we built an initial model with variables as predictors and perform stepwise model selection to select significant predictors for our “best model”. Done with the “step” method with runs “lm” multiple times to build multiple regression models and select the best vars. Use of the AIC algorithm using forward selection and backward elimation.

```
init_model <- lm(mpg ~ ., data = mtcars)
best_model <- step(init_model, 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
```

Best model obtained from the computations and consist of vars: cyl, wt and hp as confounders and “am” as the independant variable.

```
summary(best_model)
```

```
##
## Call:
## lm(formula = mpg ~ cyl + hp + wt + am, data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.9387 -1.2560 -0.4013  1.1253  5.0513
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  33.70832     2.60489   12.940 7.73e-13 ***
## cyl6         -3.03134     1.40728   -2.154  0.04068 *
## cyl8         -2.16368     2.28425   -0.947  0.35225
## hp           -0.03211     0.01369   -2.345  0.02693 *
## wt           -2.49683     0.88559   -2.819  0.00908 **
## amManual      1.80921     1.39630    1.296  0.20646
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.41 on 26 degrees of freedom
## Multiple R-squared:  0.8659,    Adjusted R-squared:  0.8401
## F-statistic: 33.57 on 5 and 26 DF,  p-value: 1.506e-10
```

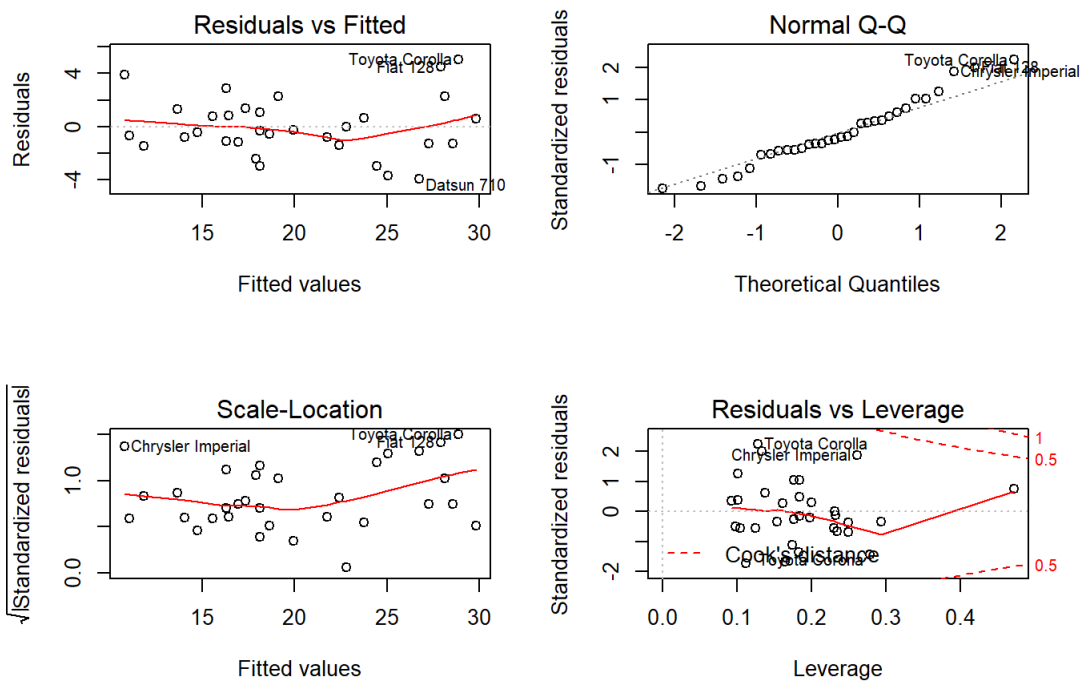
```
base_model <- lm(mpg ~ am, data = mtcars)
anova(base_model, best_model)
```

```
## Analysis of Variance Table
##
## Model 1: mpg ~ am
## Model 2: mpg ~ cyl + hp + wt + am
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      30 720.90
## 2      26 151.03  4    569.87 24.527 1.688e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Residuals and Diagnostics

Residual plots of our regression model and also the computation of the regression diagnostics for our model. Looking for outliers in the data set.

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



Notice:

- Residual vs Fitted plot seems to be randomly scattered on the plot and verify the independent condition
- Normal Q-Q consists of points which fall on the line indicating the residuals are distributed normally
- Scale-Location consists of points scattered in a constant band pattern, indicating constant variance
-
- Some distinct points of interest (outliers) in top right

Compute top 3 points in each case of influence measures

```
leverage <- hatvalues(best_model)
tail(sort(leverage), 3)
```

```
##      Toyota Corona Lincoln Continental      Maserati Bora
##      0.2777872          0.2936819          0.4713671
```

```
influential <- dfbetas(best_model)
tail(sort(influential[,6]), 3)
```

```
## Chrysler Imperial      Fiat 128      Toyota Corona
##      0.3507458      0.4292043      0.7305402
```

Our analysis is correct, the same cars are in the residual plots

Inference

T-test assuming that the transmission data has a normal distribution and clearly see that the manual and automatic transmissions are different

```
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 Automatic      mean in group Manual
##           17.14737           24.39231
```

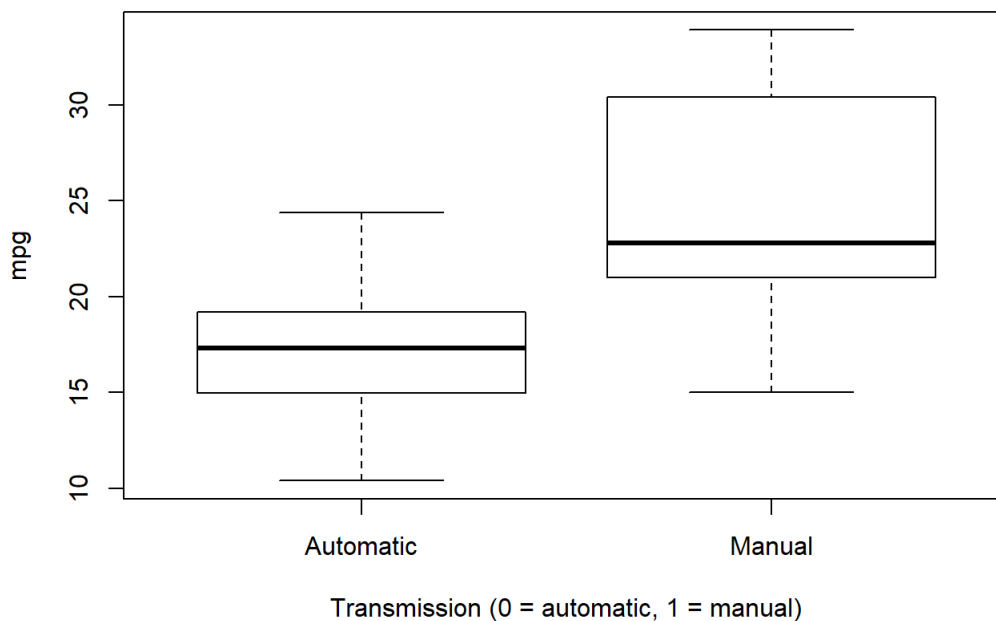
Conclusion

From our best fit model, we can conclude:

1. Cars with "Manual" transmission get more miles per gallon (mpg) vs Automatic transmission
2. MPG will decrease by 2.5 for each 1000lb increase in weight (wt)
3. MPG decreases negligibly with increase of horse power (hp)
4. If number of cylinders, cyl, increases from 4 to 6 to 8 mpg and will decrease by factor of 3 and 2.2 respectively (adjusted by hp, wt and am)

Appendix

```
boxplot(mpg ~ am, data = mtcars, xlab = "Transmission (0 = automatic, 1 = manual)")
```



```
pairs(mtcars)
```

