

Regression Model Project

Mudit

Objective

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”

Executive Summary

For this task we will analyze mtcars data set and explore the relationship between miles per gallon (MPG) and all other variables. Regression Modelling and Exploratory Data Analysis is done to mainly explore how automatic (am = 0) and manual (am = 1) transmissions features affect the MPG feature. The t-test shows that the performance difference between cars with automatic and manual transmission. And it is about 7 MPG more for cars with manual transmission than those with automatic transmission. Then, we fit some linear regression models and select the one with highest Adjusted R-squared value. Based on the final model selected we can say, manual transmission is better for mpg than automatic, by 1.80921 mpg.

Load Data

```
## load the data
data(mtcars)
```

Exploratory Data Analysis

All categorical variables to factors and names were given to am variable for better understanding

```
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 ...
```

```
## convert numeric variables to factors
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'))
attach(mtcars)
summary(mtcars)
```

```
##      mpg      cyl      disp      hp      drat
##  Min.   :10.40   4:11   Min.   : 71.1   Min.   : 52.0   Min.   :2.760
##  1st Qu.:15.43   6: 7   1st Qu.:120.8   1st Qu.: 96.5   1st Qu.:3.080
##  Median :19.20   8:14   Median :196.3   Median :123.0   Median :3.695
##  Mean   :20.09           Mean   :230.7   Mean   :146.7   Mean   :3.597
##  3rd Qu.:22.80           3rd Qu.:326.0   3rd Qu.:180.0   3rd Qu.:3.920
##  Max.   :33.90           Max.   :472.0   Max.   :335.0   Max.   :4.930
##      wt      qsec      vs      am      gear      carb
##  Min.   :1.513   Min.   :14.50   0:18   Automatic:19   3:15   1: 7
##  1st Qu.:2.581   1st Qu.:16.89   1:14   Manual   :13   4:12   2:10
##  Median :3.325   Median :17.71           5: 5   3: 3
##  Mean   :3.217   Mean   :17.85           4:10
##  3rd Qu.:3.610   3rd Qu.:18.90           6: 1
##  Max.   :5.424   Max.   :22.90           8: 1
```

For comparing the means for mpg of automatic and manual boxplot was done and it was observed that Manual (represented by 1) has a higher mean for mpg than automatic (represented by 0). Please check boxplot in appendix

Hypothesis Testing

With the observation made by boxplot, let the null hypothesis be that mpg for automatic is higher than manual.

```
t.test(mpg~am,mtcars,paired=FALSE,var.equal=FALSE)
```

```
##
## 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
```

The p-value is less than 0.05, so we can reject the null hypothesis.

Model Fitting & Selection

First will plot a model against am variable

```
##model against am variable
modelam<-lm(mpg~am,mtcars)
summary(modelam)

##
## 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 ***
## amManual       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
```

From this we can see that automatic has 7.24mpg less than manual. Rsq is 36% and AdjRsqr is 34% which is very low.

Now we will model against all variables and step model using backward step function.

```
modelall<-lm(mpg~.,mtcars)
summary(modelall)

##
## Call:
## lm(formula = mpg ~ ., data = mtcars)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -3.5087 -1.3584 -0.0948  0.7745  4.6251
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  23.87913    20.06582   1.190  0.2525
## cyl6         -2.64870     3.04089  -0.871  0.3975
## cyl8         -0.33616     7.15954  -0.047  0.9632
## disp          0.03555     0.03190   1.114  0.2827
## hp           -0.07051     0.03943  -1.788  0.0939 .
## drat          1.18283     2.48348   0.476  0.6407
## wt           -4.52978     2.53875  -1.784  0.0946 .
## qsec          0.36784     0.93540   0.393  0.6997
```

```
## vs1          1.93085    2.87126    0.672    0.5115
## amManual     1.21212    3.21355    0.377    0.7113
## gear4        1.11435    3.79952    0.293    0.7733
## gear5        2.52840    3.73636    0.677    0.5089
## carb2       -0.97935    2.31797   -0.423    0.6787
## carb3        2.99964    4.29355    0.699    0.4955
## carb4        1.09142    4.44962    0.245    0.8096
## carb6        4.47757    6.38406    0.701    0.4938
## carb8        7.25041    8.36057    0.867    0.3995
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 2.833 on 15 degrees of freedom
## Multiple R-squared:  0.8931, Adjusted R-squared:  0.779
## F-statistic:  7.83 on 16 and 15 DF,  p-value: 0.000124
```

```
fit<-step(modelall, direction="backward")
```

```
## 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
##
## 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
## - wt 1 31.1896 170.21 71.482
##
## 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
## - hp 1 17.7382 157.73 67.044
## - am 1 19.4660 159.46 67.393
## - wt 1 30.7151 170.71 69.574
##
## 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
## - hp 1 18.184 159.42 65.386
## - am 1 18.885 160.12 65.527
## - wt 1 39.645 180.88 69.428
##
## 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
## - am 1 16.443 160.12 63.527
## - hp 1 36.344 180.02 67.275
## - wt 1 41.088 184.77 68.108
##
## 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
## - cyl 2 29.265 180.29 63.323
## - hp 1 31.943 182.97 65.794
## - wt 1 46.173 197.20 68.191

```

```
summary(fit)
```

```
##
## 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
```

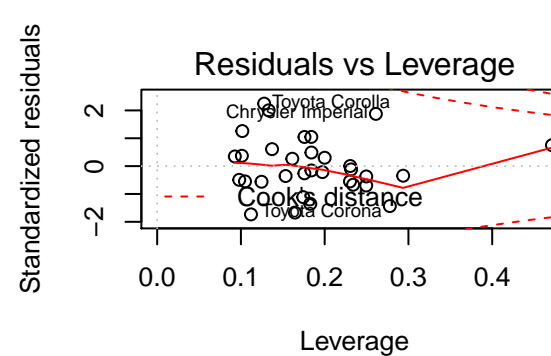
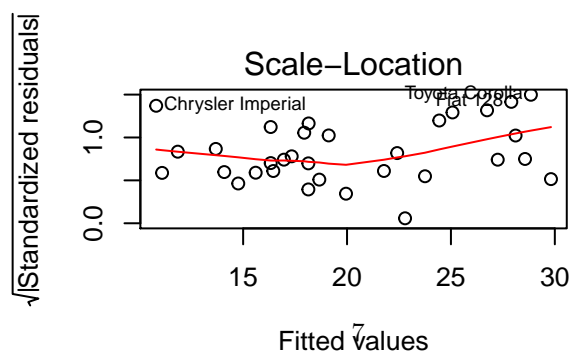
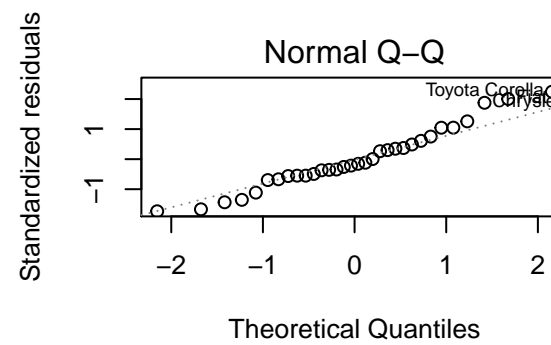
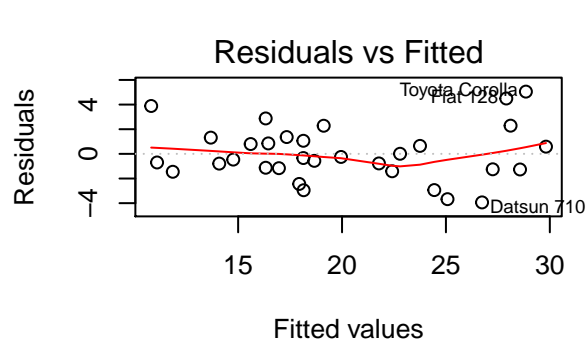
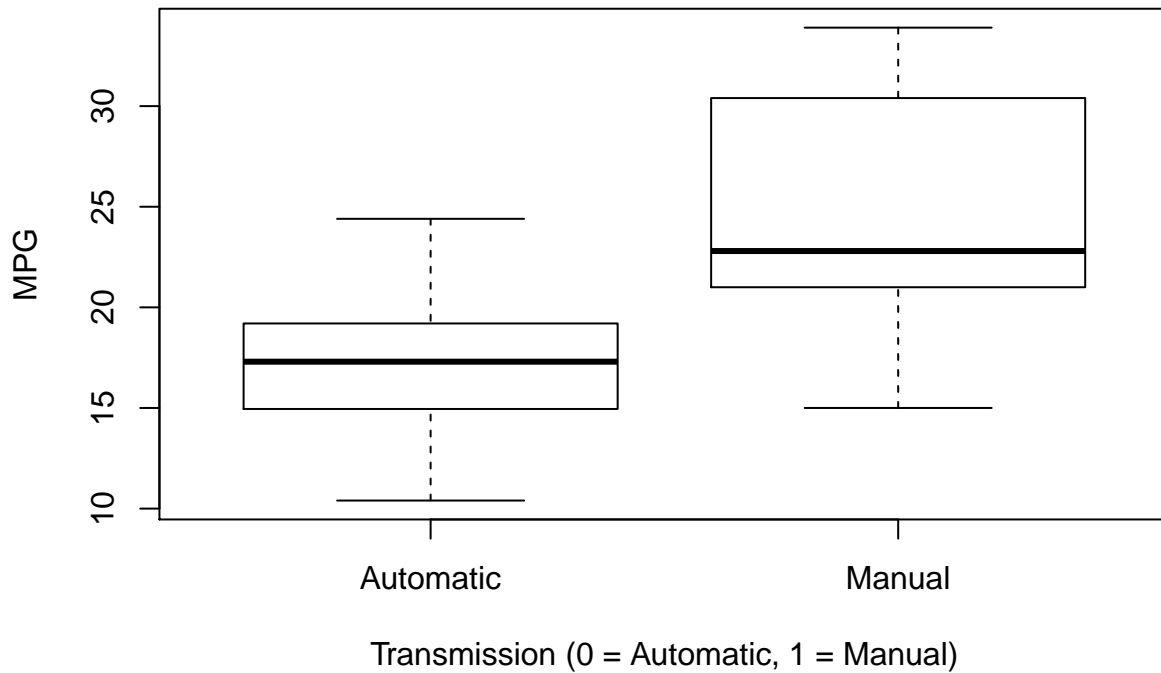
For all variable model we have Rsq 89% and AdjRsq as 78% and for back step model Rsq is 87% and Adjusted Rsq is 84%.

Model with highest AdjRsq is selected i.e. Backward Step Model Take a look at residuals in appendix. Residuals vs Fitted and Scale-Location plots show no pattern, the Normal Q-Q plot indicates that Residuals approximately follow a Normal distributions, and the Residuals vs Leverage plot tell us that there's none particular outlier to be concerned and Scale-Location plot confirms the constant variance assumption.

Appendix

MPG difference between automatic and manual transmissions

Boxplot of MPG vs. Transmission



Conclusion

From this analysis we can conclude that manual transmissions are better than automatic by 1.80921 mpg.