

Regression model Course Project

Tham Nguyen

7/5/2019

Executive Summary:

We estimate the relationship between one variable (i.e. type of transmission (manual or automatic) and other independent variables along with miles per gallon (MPG).

Data Description

We analyze the **'mtcars'** data set through Regression Modelling and exploratory analysis to show how automatic (am = 0) and manual (am = 1) transmissions features affect the MPG feature. The dataset **"mtcars"** is located in the package **"reshape2"** first introduced in the Reshaping Data Course of the same Data Specialization Course. The data was extracted from the 1974 Motor Trend US magazine, which comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973-74 models).

Data processing and transformation

We load package in the data set:

```
library(ggplot2)

## Registered S3 methods overwritten by 'ggplot2':
##   method      from
##   [.quosures   rlang
##   c.quosures   rlang
##   print.quosures rlang

library(datasets)
library(readr)
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
```

Data Preparation

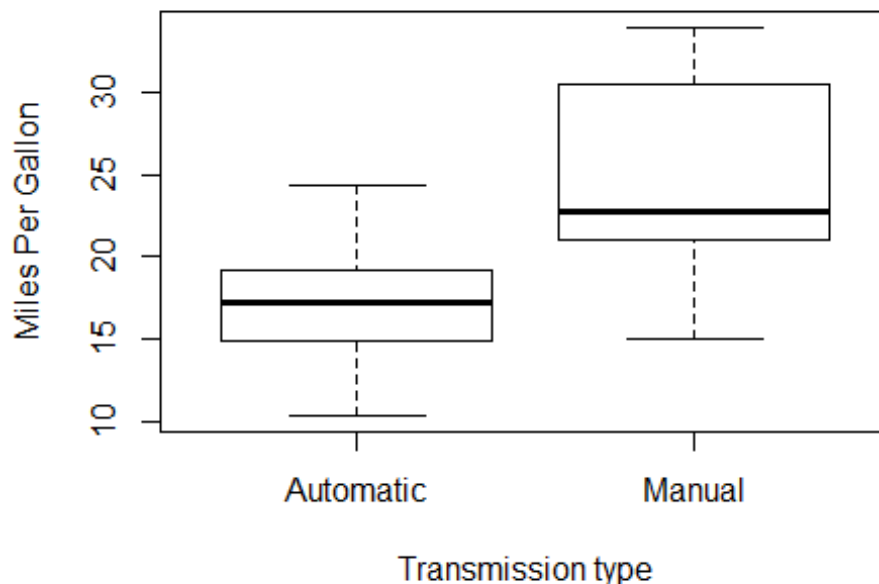
```
data(mtcars)
mtcars <- mtcars %>%
  mutate(vs = as.factor(vs),
         am = as.factor(am),
         cyl = as.ordered(cyl),
         gear = as.ordered(gear),
         carb = as.ordered(carb))
mtcars$am <- factor(mtcars$am, labels=c('Automatic', 'Manual'))
str(mtcars)

## '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 : Ord.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: Ord.factor w/ 3 levels "3"<"4"<"5": 2 2 2 1 1 1 1 2 2 2 ...
## $ carb: Ord.factor w/ 6 levels "1"<"2"<"3"<"4"<...: 4 4 1 1 2 1 4 2 2 4
## ...
```

Exploratory Data Analysis

we plot boxplots of the variable mpg when am is Automatic or Manual (see boxplot in the appendix plot 1). This plot clearly depicts an increase in the **mpg** when the transmission is Manual

```
boxplot(mpg ~ am, data = mtcars, xlab = "Transmission type", ylab = "Miles Per Gallon")
```



Regression Analysis

In this section, we start building linear regression models based on the different variables and try to find out the best model fit and compare it with the base model which we have using **anova**. After model selection, we also perform analysis of residuals.

Model building and selection

First, let start with simple regression model how **am** effect in **mpg**:

```
fit_base <- lm(mpg ~ am, data=mtcars)
summary(fit_base)

##
## 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
```

This shows us that the average MPG for **automatic** is **17.1 MPG**, while **manual** is **7.2 MPG higher**. The **R square value is 0.36** thus telling us this model only explains us **36%** of the variance. As a result, significant effect on MPG otherwise it low bias, so we need more feature to complete the model.

Next: we will try all feature and check how it effect to mpg.

```
fit_all <- lm(mpg ~ ., data = mtcars)
summary(fit_all)

##
## 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)  26.57171   19.56616   1.358   0.1945
## cyl.L        -0.23770    5.06256  -0.047   0.9632
## cyl.Q         2.02541    2.14952   0.942   0.3610
## 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
## gear.L       1.78785    2.64200   0.677   0.5089
## gear.Q       0.12235    2.40896   0.051   0.9602
## carb.L       6.06156    6.72822   0.901   0.3819
## carb.Q       1.78825    2.80043   0.639   0.5327
## carb.C       0.42384    2.57389   0.165   0.8714
## carb^4       0.93317    2.45041   0.381   0.7087
## carb^5      -2.46410    2.90450  -0.848   0.4096
## ---
## 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
```

This result show that fit_all much better than fit_base. But we need to find the best model that fix higher variation. As mentioned, based on the pairs plot where several variables has high correlation with **mpg**, We build an initial model with all the variables as predictors, and stepwise model selection to select significant predictors for the final model which is the best model - fit_final.

```
fit_final <- step(fit_all, 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
```

The final model obtained from the above computations consists of the variables **wt**, **qsec** and **am** the independent variable. Details of the model are depicted below.

```
summary(fit_final)

##
## 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)  31.97665    3.06337   10.438 8.61e-11 ***
## cyl.L        -1.52995    1.61521   -0.947  0.35225
## cyl.Q         1.59177    0.88076    1.807  0.08231 .
## 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
```

From the result above, we observe that the “Adjusted R square value” is **0.83** which is the maximum obtained considering all combinations of variables. Thus, we can conclude that

more than **83%** of the variability is explained by the above model. **Next**, we use anova to compare against our base model and final model that was found through performing **stepwise** selection

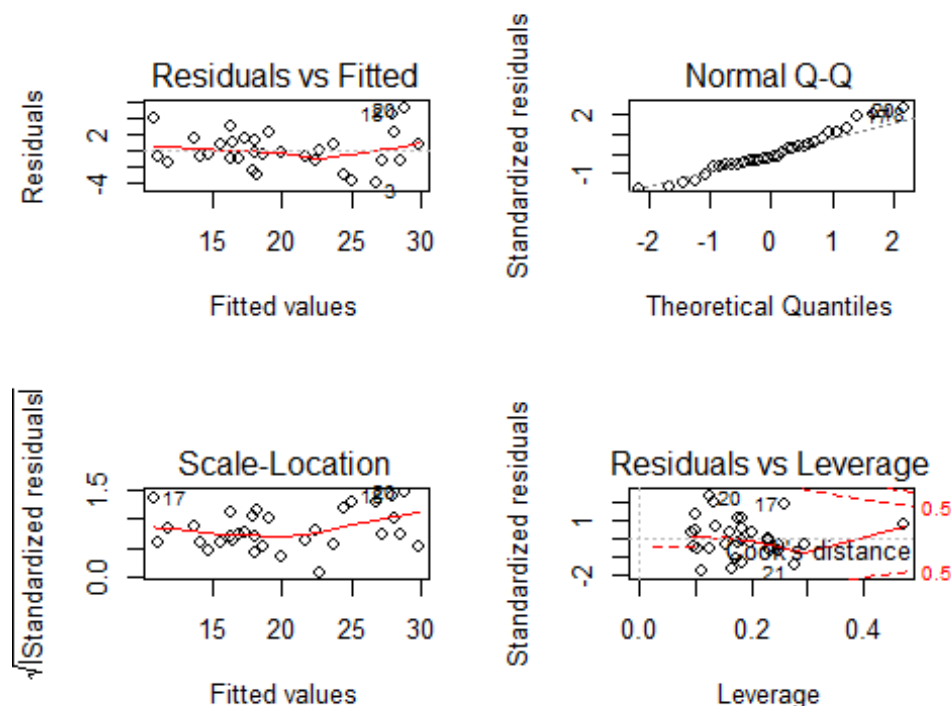
```
anova(fit_base, fit_final)

## 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
```

This results in a **p-value of 1.55e-09**, and we can claim the fit_final model is significantly better than our fit_base simple model. We double-check the residuals for non-normality (Appendix - Plot 3) and can see they are all normally distributed.

Appendix

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



```
pairs(mtcars)
```