

Regression Models Project

V. Boodram

October 22, 2014

Executive Summary

Based on the 1974 statistics for fuel consumption and 10 aspects of automobile design for 32 (1973-4) automobile models, provided by the magazine *Motor Trend* (1974), it would like to be determined if automatic or manual transmissions are better for fuel economy, measured in miles per gallon, and what the quantitative difference in mpg between the two transmission types is.

The appropriate predictive model for determining fuel economy with a set of potential predictors from those provided in the `mtcars` data set was determined to be

$$\text{mpg} = 9.62 - 3.92 \times \text{weight} + 1.23 \times \text{quarter-mile time} + 2.94 \times \text{transmission:manual}$$

With this model, it was determined with 95% confidence, all other variables being equal, that cars with a manual transmission have a fuel economy of 0.04107 mpg to 5.83093 mpg higher than cars with an automatic transmission.

Methodology

The full model (i.e: with every variable) was first evaluated from the skeptical point of view that none of the slopes β_i in the set of potential predictors was in fact a statistically significant predictor of the response variable, mpg; the alternative hypothesis was that at least one predictor is a significant predictor of the response variable

$$H_0 : \beta_i = 0, \forall i$$
$$H_A : \exists i \ni \beta_i \neq 0$$

```
# full model
summary(lm(mpg ~ ., data = mtcars))$fstatistic
```

```
## value numdf dendif
## 13.93 10.00 21.00
```

The small p-value indicates that the model as a whole is significant, so the null hypothesis is rejected.

Having found evidence that there is at least one statistically significant predictor of mpg, an initial model was fit to determine if the transmission type was among these.

```
mtcars$am<-as.factor(mtcars$am)
levels(mtcars$am)<-c("automatic", "manual")
summary(lm(formula = mpg~am, mtcars))$coef
```

```
##           Estimate Std. Error t value Pr(>|t|)
## (Intercept)   17.147      1.125   15.247 1.134e-15
## ammanual       7.245      1.764    4.106 2.850e-04
```

The percentage in variability of fuel economy explained by transmission type, R^2 is 36%, and with a small p-value (0.00029), transmission type was determined to be a significant predictor of fuel economy. The above summary indicates that automatic transmission cars achieve 17.15 mpg, and that all else being held constant, manual cars achieve 7.24 more mpg. The initial model

$$\text{mpg} = 17.17 + 7.24 \times \text{transmission: manual}$$

reflects patterns in the box plots in Appendix Figure 1, composed in the exploratory data analysis.

The plots suggest that higher fuel efficiency is achieved by cars with manual transmissions; however, it may be the case that some other features may have been influencing the fuel consumption. This possibility was evaluated in reference to the bivariate plots in the Appendix Figure 2. These plots show scatter plots between each pair of variables in the data set, as well as the correlation coefficients on the upper half of the matrix. It shows a moderate correlation (0.66) between transmission and fuel efficiency, and significantly stronger correlations between mpg and the number of cylinders, the engine displacement, and the weight of the car, suggesting interactions between the potential set of predictors that may affect the amount of variability in mpg explained by transmission type. The collinearities between transmission type and these other variables are moderately high, with the strongest being between transmission type and weight, and transmission type and rear-axle ratio, which means that some of the information in the latter variables is already being captured by the variable transmission.

The remaining variables were incorporated into the model through a *backwards elimination* process with reference to the p-value. Beginning with the full model (which includes every variable), in each iteration the variable with the highest p-value was discarded, and the model was refit, until all the variables contained in the model were significant, and this was declared to be the parsimonious model. See the Appendix for the function `backwards()`; only the final fit and confidence intervals for the slope are provided here

```
fit<-lm(mpg~., data = mtcars)
final<-backwards(fit); summary(final)$coef
```

| | Estimate | Std. Error | t value | Pr(> t) |
|----------------|----------|------------|---------|-----------|
| ## (Intercept) | 9.618 | 6.9596 | 1.382 | 1.779e-01 |
| ## wt | -3.917 | 0.7112 | -5.507 | 6.953e-06 |
| ## qsec | 1.226 | 0.2887 | 4.247 | 2.162e-04 |
| ## ammanual | 2.936 | 1.4109 | 2.081 | 4.672e-02 |

```
# 95% confidence interval for the slope
2.936+c(-1,1)*qt(p = .975, df = 27)*1.4109
```

```
## [1] 0.04107 5.83093
```

Diagnostics

The computed model is considered valid if the residuals are nearly normal, if they have constant variability, and if they are independent. For a categorical explanatory variable, it makes no sense to seek a linear relationship with the response. The near-normality of the residuals was checked with the Q-Q plot of Figure 3 in the Appendix. Near the tail areas, fairly significant deviations from the mean can be seen, but residual normality is fairly satisfied. Constant variability can be seen to be satisfied by Figure 4, which shows the residuals randomly scattered with a constant width around 0. Finally, it is reasonable to assume independence in the data set, because even though some cars may have the same manufacturer, the specifications of each car is likely to be different. Consequently, the residuals must also be independent, and the model is deemed valid.

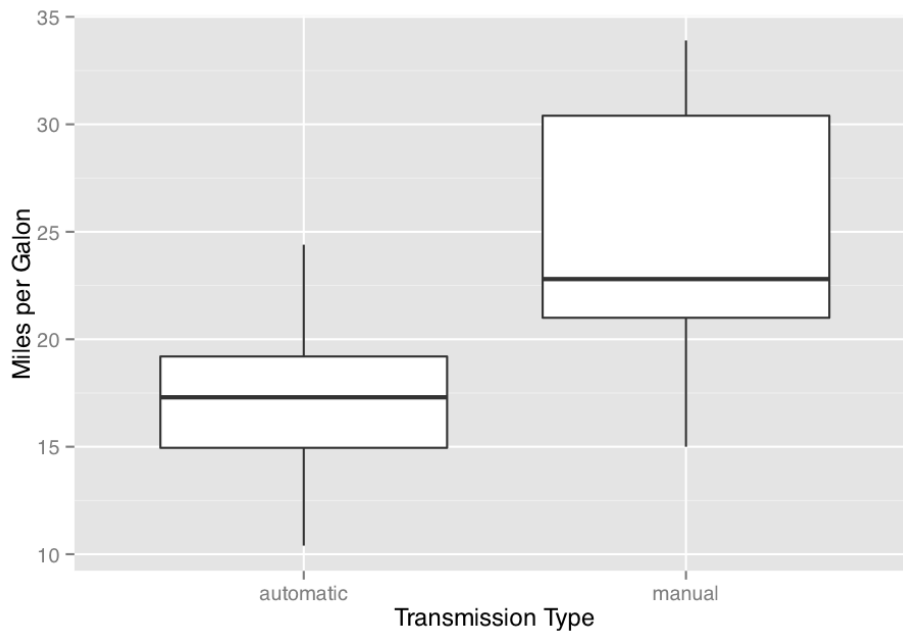
Conclusion

Fuel efficiency is predicted to be between 0.04107 mpg to 5.83093 mpg higher in manual transmission cars than automatic transmission cars with 95% confidence. All else being held equal, a manual transmission car improves the number of miles per gallon by 2.936, compared to a car with an automatic transmission.

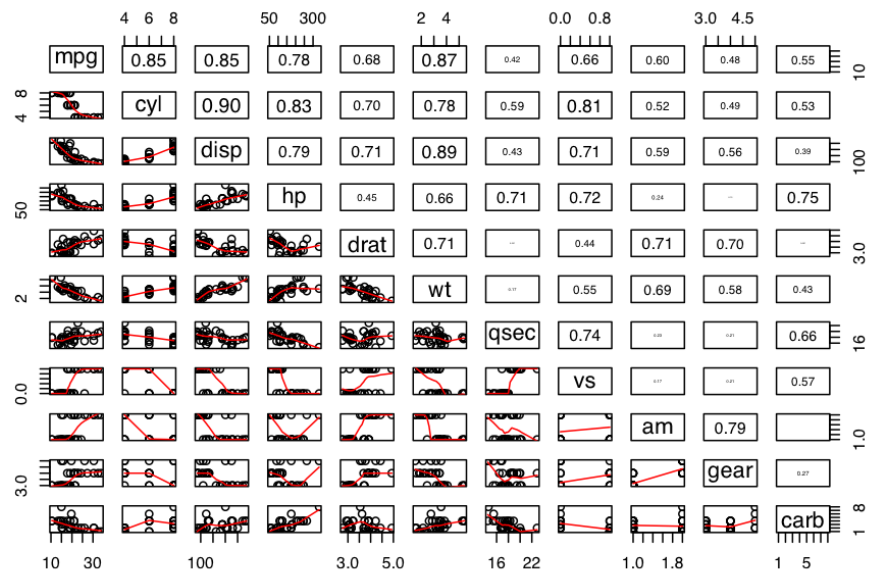
Appendix

Figures

```
library(ggplot2)
options(height=60)
ggplot(mtcars, aes(factor(am), mpg))+geom_boxplot()+labs(x="Transmission Type", y="Miles per Gallon")
```

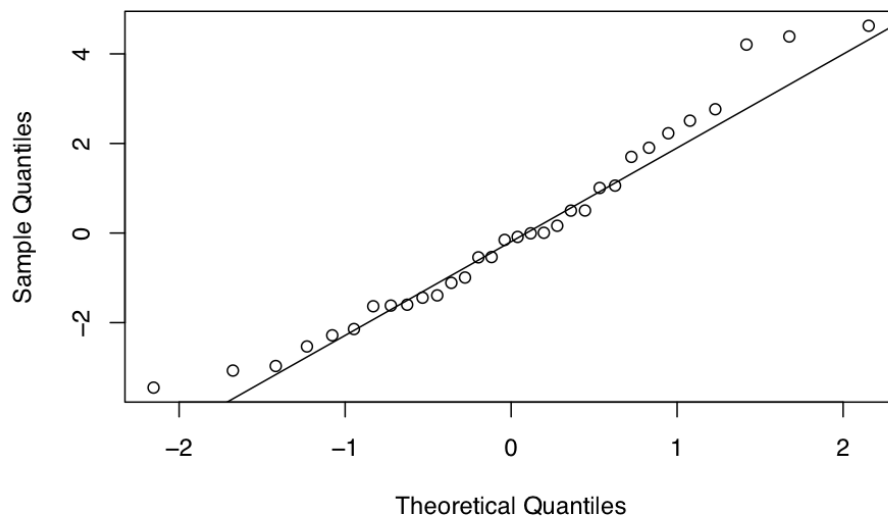


```
panel.cor <- function(x, y, digits = 2, prefix = "", cex.cor, ...){
  usr <- par("usr"); on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y))
  txt <- format(c(r, 0.123456789), digits = digits)[1]
  txt <- paste0(prefix, txt)
  if(missing(cex.cor)) cex.cor <- 0.8/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex.cor * r)}
pairs(mtcars, lower.panel = panel.smooth, upper.panel = panel.cor)
```

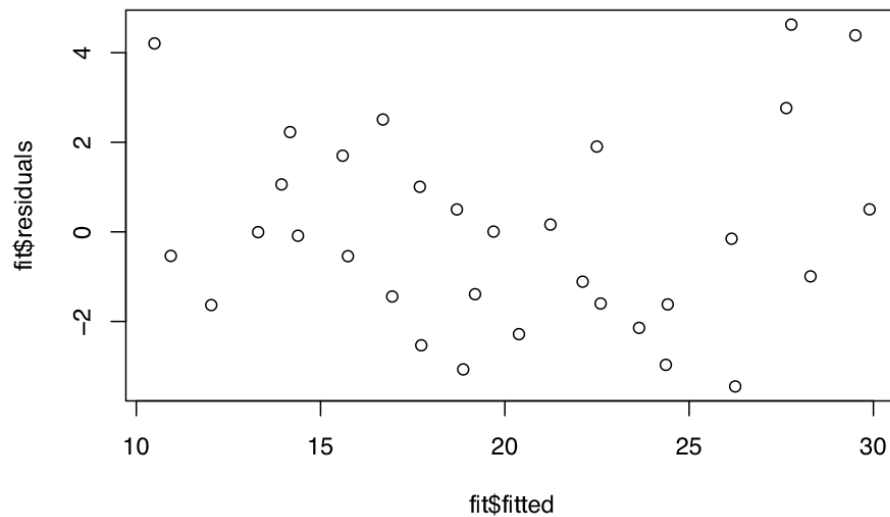


```
qqnorm(fit$residuals)
qqline(fit$residuals)
```

Normal Q-Q Plot



```
plot(fit$residuals~fit$fitted)
```



2. backwards() function

```
backwards<-function(fit){
  while(remove[2]>0.05 ){
    if(remove[1] == "(Intercept)"){
      return(fit)
    }
    # obtain all factors in the current model
    allFactors<-names(summary(fit)$coef[,4])
    # determine the index of the factor to be removed
    removeFactor<-which.is.max(x = summary(fit)$coef[,4])
    # remove the intercept value and the factor computed in the last step
    newFactors<-allFactors[-c(1, removeFactor)]
    # replace ammanual with am in the vector of factors
    restoreFactors<-gsub(pattern = "^.*amm.*$",replacement = "am", x = newFactors )
    # convert the vector to the correct format for use in the formula
    useFactors<-paste(restoreFactors, collapse = "+")
    # fit a new model
    fit<-lm(paste("mpg~",useFactors, sep = ""), data=mtcars)
    # obtain the index of the factor with the largest p-values
    ind<-which.is.max(x = summary(fit)$coef[,4])
    # remove this factor
    remove<-c(rownames(summary(fit)$coef)[ind], summary(fit)$coef[ind, 4])
    fit
  }
}
```