

# Regression Models Course Project

Shadi

5/8/2020

## Summary

The purpose of analysis is to explore which kind of transmission is better for MPG and to quantify the difference between the two kind of transmission. In order to answer these questions exploratory data analysis performed, and we used hypothesis testing and linear regression as methodologies to make inference. We then established both simple and multivariate linear regression analysis. However the result of the multivariable regression model is more promising as it includes the potential effect of other variables on MPG. Finally, we concluded that the manual car are better in terms of MPG as they have a higher MPG compared to automatic car. see the Appendix for full Model Summary and diagnostics.

## Data

The data was extracted from the 1974 Motor Trend US magazine, and comprises fuel consumption and 10 aspects of automobile design and performance for 32 automobiles (1973–74 models). description for variable can be find in appendix

	mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
Mazda RX4	21.0	6	160	110	3.90	2.620	16.46	0	1	4	4
Mazda RX4 Wag	21.0	6	160	110	3.90	2.875	17.02	0	1	4	4
Datsun 710	22.8	4	108	93	3.85	2.320	18.61	1	1	4	1

looking at the data we see 11 variables with all numeric value.

Our variables of interest here are 1. transmission (am) which form the information provided about the data we know that it take value 1 for manual and 0 for automatic. 2. MPG (mile per gallon).

Visualizing MPG by the violin plot over two value of transmission would give us a good idea for our future hypothesis testing. (plot 1, appendix) It appears that automatic cars have a lower miles per gallon, and therefore a lower fuel efficiency, than manual cars do. But it is possible that this apparent pattern happened by random chance- that is, that we just happened to pick a group of automatic cars with low efficiency and a group of manual cars with higher efficiency. So to check whether that's the case, we have to use a statistical test.

In order to determine the relationship between the variables, and to find out which variables should be included in our model, we perform a correlation test and create a correlation heat map (plot 2, appendix).

Finally we take a look at the distribution of mpg to check if it look like normal. (plot 3, appendix)

## Analysis

To address the first question we run a simple t-test checking if the automatic cars have a lower miles per gallon. our null hypothesis is: “there is no difference between the two transmission”

t.statistic	df	p.value	lower.CL	upper.CL	automatic.mean	manual.mean
-3.767123	18.33225	0.0013736	-11.28019	-3.209684	17.14737	24.39231

Here the p-value is 0.001 meaning that we can reject the null.

quantifying the difference between the two transmission, we need to have a proper model, we conduct both simple linear and multivariable linear regression, and select the best model, to have the most accurate prediction

	Model 1	Model 2
(Intercept)	18.90 *** [17.56, 20.23]	17.15 *** [15.30, 18.99]
wt	-3.83 *** [-5.64, -2.03]	
qsec	2.19 ** [0.91, 3.47]	
am	2.94 [-0.32, 6.19]	7.24 ** [3.17, 11.32]
N	32	32
R2	0.85	0.36

All continuous predictors are mean-centered and scaled by 1 standard deviation. Standard errors are heteroskedasticity robust. \*\*\*  $p < 0.001$ ; \*\*  $p < 0.01$ ; \*  $p < 0.05$ .

In table below we can see that for linear regression, both intercept and transmission are significant. and the model shows that mpg for manual car are 7.25 scale higher than automatic, but the value of the R2: 0.35 is low. For the multivariable linear regression, all the variable are significant and the R2 is near 1 which means that it is a better model compared to linear one.

## Model selection

As we see earlier almost all of other variables have high correlation with MPG so in order to have a good predictor in our model using a stepwise selection method which consists of iteratively adding and removing predictors to find the subset of variables in the data set resulting in the best performing model would be best move. following code automatically give us the best model

for selecting the best model one way is to look at the R2 value but to have a better result we conduct a likelihood ratio test and compare the two models.

Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
30	720.8966	NA	NA	NA	NA
28	169.2859	2	551.6107	45.61838	0

Finally, in order to interpret the result we must make sure that the error doesn't have heteroskedasticity and

colinearity this can be find at appendix(Plot 4).

## Conclusion

Using model contruction, we have shown that, adjusted to other strong mpg predictors that we can find in the mtcars dataset, manual transmission is really the best transmission for mpg with  $mmpg=9.61+2.93 \text{ amManual}-3.91 \text{ wt} +1.22 \text{ qsec}$

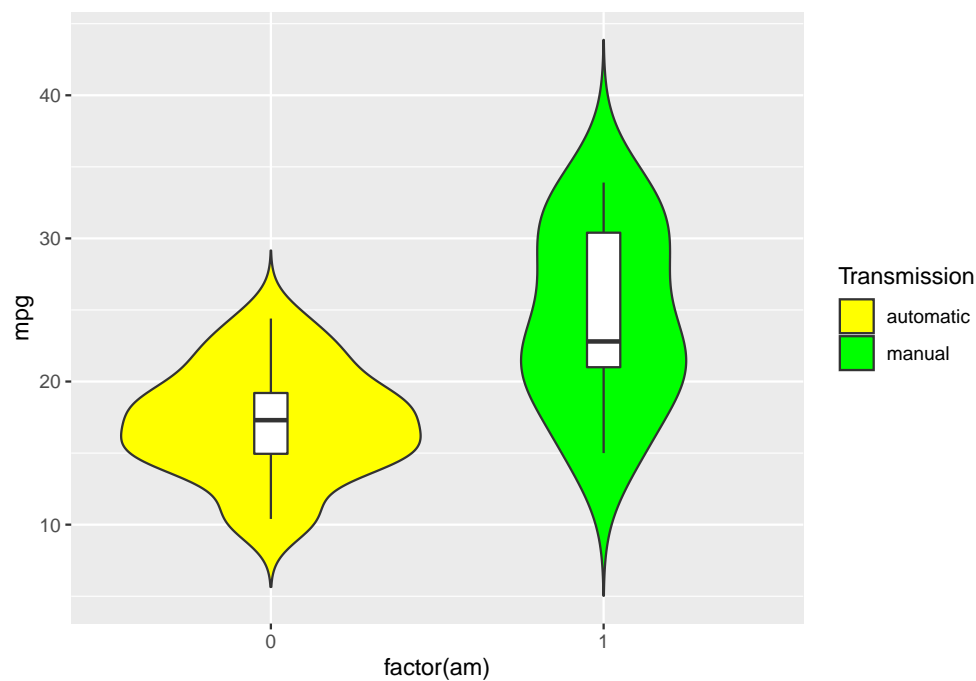
## Appendix

residual plot and some diagnostics

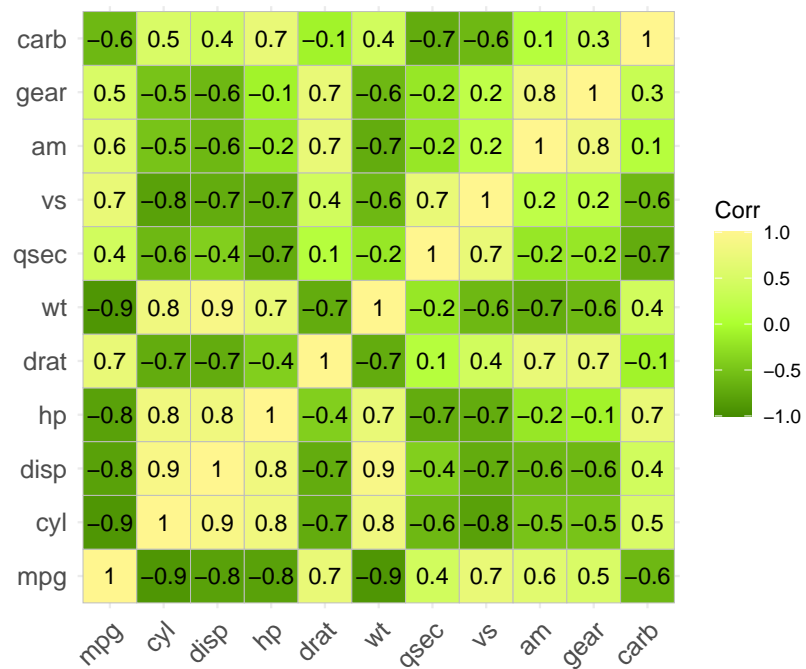
### data description

1-**mpg**: Miles/(US) gallo 2-**cyl**: Number of cylinders 3-**disp**: Displacement (cu.in.) 4-**hp**: Gross horsepower 5-**drat**: Rear axle ratio 6- **wt**: Weight (1000 lbs) 7-**qsec**: 1/4 mile time 8- **vs**: Engine (0 = V-shaped, 1 = straight) 9-**am**:Transmission (0 = automatic, 1 = manual) 10-**gear**: Number of forward gears 11- **carb**: Number of carburetors

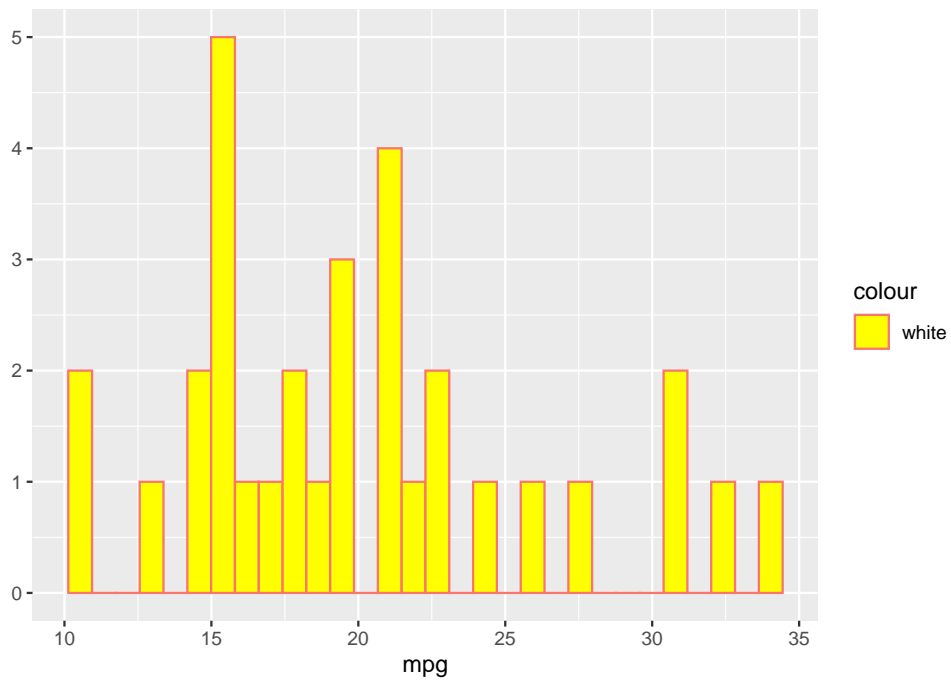
### plot 1



plot 2



plot 3



plot 4

