

Experiment 8:

Write R program to implement linear and multiple regression on 'mtcars' dataset to estimate the value of 'mpg' variable, with best R^2 and plot the original values in 'green' and predicted values in 'red'.

Solution:

The built-in **mtcars** data frame contains information including their weight, fuel efficiency (in miles-per-gallon), speed, etc. of 32 models of automobile from 1973-74 as reported in Motor Trend Magazine. In analyzing the dataset of different collection of cars, we will explore the relationship between a set of eleven variables, and miles per gallon (MPG). Dataset Motor Trend has been used to find out that,

- Is an automatic or manual transmission better for miles per gallon?
- How different is the MPG between automatic and manual transmission?

Read MTCARS dataframe

```
> data ("mtcars")
```

```
> head (mtcars)
```

	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
Hornet 4 Drive	21.4	6	258	110	3.08	3.215	19.44	1	0	3	1
Hornet Sportabout	18.7	8	360	175	3.15	3.440	17.02	0	0	3	2
Valiant	18.1	6	225	105	2.76	3.460	20.22	1	0	3	1

We have to explore the relationship between a set of variables and miles per gallon (mpg), so mpg is our **dependent variable**. Plot dependent variable to check its distribution.

```
x <- mtcars$mpg
```

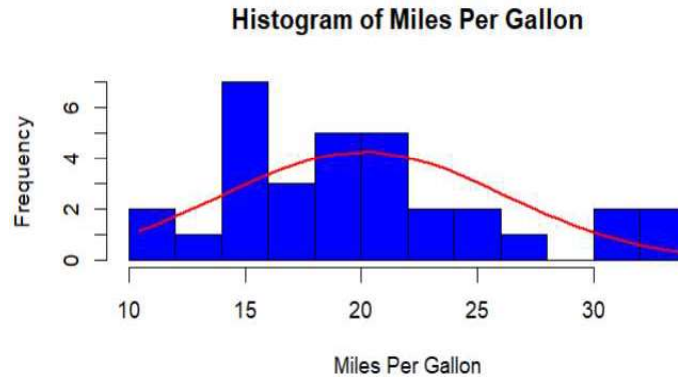
```
h<- hist (x, breaks=10, col="blue", xlab="Miles Per Gallon", main="Histogram of Miles Per Gallon")
```

```
xfit <- seq (min(x),max(x),length=40)
```

```
yfit <- dnorm (xfit, mean=mean(x), sd=sd(x))
```

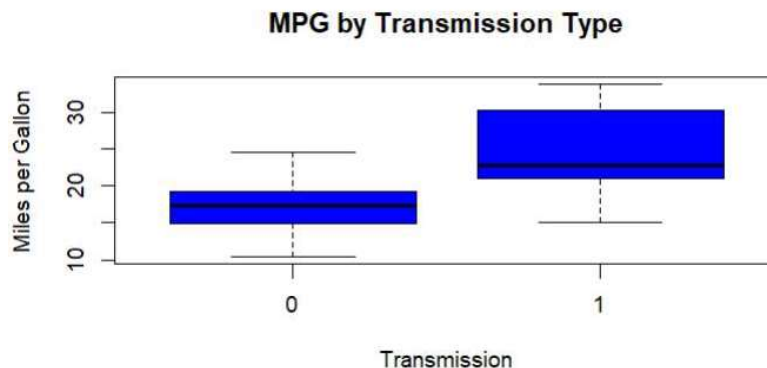
```
yfit <- yfit*diff(h$mids[1:2])*length(x)
```

```
lines (xfit, yfit, type="l", col="red", lwd=2)
```



The distribution of mpg, showing in the graphs is approximately normal and does not contain any outliers. Now we check how mpg can be changed by automatic or manual transmission, by plotting a box plot. From this boxplot, it seems that automatic cars have a lower miles per gallon, and so a lower fuel potency, than manual cars do.

```
> boxplot(mpg~am, data = mtcars, col = c("blue", "blue"), xlab = "Transmission",
  ylab = "Miles per Gallon", main = "MPG by Transmission Type")
```



Correlation analysis

A correlation test is performed to determine the relationship between the variables, and to find out which variables should be included in our model to answer the questions. The correlation matrix is

```
> cor (mtcars, use="complete.obs", method="pearson")
```

	mpg	cy1	disp	hp	drat	wt	qsec	vs	am	gear	carb
mpg	1.0000000	-0.8521620	-0.8475514	-0.7761684	0.6811719	-0.8676594	0.4186840	0.6640389	0.5998324	0.4802848	-0.5509250
cy1	-0.8521620	1.0000000	0.9020329	0.8324475	-0.6999381	0.7824958	-0.5912420	-0.8108118	-0.5226070	-0.4926866	0.5269882
disp	-0.8475514	0.9020329	1.0000000	0.7909486	-0.7102139	0.8879799	-0.4336978	-0.7104159	-0.5912270	-0.5555692	0.3949768
hp	-0.7761684	0.8324475	0.7909486	1.0000000	-0.4487591	0.6587479	-0.7082233	-0.7230967	-0.2432042	-0.1257043	0.7498124
drat	0.6811719	-0.6999381	-0.7102139	-0.4487591	1.0000000	-0.7124406	0.0912047	0.4402785	0.7127113	0.6996101	-0.0907898
wt	-0.8676594	0.7824958	0.8879799	0.6587479	-0.7124406	1.0000000	-0.1747158	-0.5549157	-0.6924952	-0.5832870	0.4276059
qsec	0.4186840	-0.5912421	-0.4336979	-0.7082234	0.0912047	-0.1747159	1.0000000	0.7445354	-0.2298608	-0.2126822	-0.6562492
vs	0.6640389	-0.8108118	-0.7104159	-0.7230967	0.4402784	-0.5549157	0.7445354	1.0000000	0.1683451	0.2060233	-0.5696071
am	0.5998324	-0.5226070	-0.5912270	-0.2432043	0.7127113	-0.6924953	-0.2298608	0.1683451	1.0000000	0.7940588	0.0575343
gear	0.4802848	-0.4926866	-0.5555692	-0.1257043	0.6996101	-0.5832870	-0.2126822	0.2060233	0.7940587	1.0000000	0.2740728
carb	-0.5509251	0.5269883	0.3949769	0.7498125	-0.0907898	0.4276059	-0.6562492	-0.5696071	0.0575343	0.2740728	1.0000000

The values in correlation matrix shows that variables such as *wt*, *cyl*, *disp*, and *hp* are highly correlated with the dependent variable *mpg*. Hence, they should be included in the regression model. From the correlation matrix, it can be also be observed that *cyl* and *disp* are highly correlated with each other. In order to avoid the problem of collinearity only one variable from these two will be included in the model.

Simple Regression model

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

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 ***
am           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 the above summary, there exists a linear relation between the predictor variable MPG and AM. Intercepts and Coefficient can be explained that, on average automatic transmission cars has 17.147 MPG and manual transmission cars has 24.39(17.147 + 7.24). The value of R^2 is 0.3385, which means this model only explain 33.85% of the variance.

Multiple Regression model

In the correlation analysis, it is observed that variables such as *wt*, *cyl*, *am*, and *hp* are highly correlated with the dependent variable *mpg*. So, we apply a multi variant regression for *mpg* on *am*, *wt*, *cyl*, and *hp*.

```
> mfit <- lm (mpg ~ am + cyl + wt + hp, data = mtcars)
> anova (fit, mfit)
```

```

Analysis of Variance Table

Model 1: mpg ~ am
Model 2: mpg ~ am + cyl + wt + hp
  Res.Df  RSS Df Sum of Sq    F    Pr(>F)
1      30 720.9
2      27 170.0  3      550.9 29.166 1.274e-08 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

The p-value of 1.274e-08, it is clear that the multivariate model of regression is different from that of above simple model.

```
> summary(mfit)
```

```

Call:
lm(formula = mpg ~ am + cyl + wt + hp, data = mtcars)

Residuals:
    Min       1Q   Median       3Q      Max
-3.4765 -1.8471 -0.5544  1.2758  5.6608

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  36.14654     3.10478   11.642 4.94e-12 ***
am             1.47805     1.44115    1.026  0.3142
cyl           -0.74516     0.58279   -1.279  0.2119
wt            -2.60648     0.91984   -2.834  0.0086 **
hp            -0.02495     0.01365   -1.828  0.0786 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.509 on 27 degrees of freedom
Multiple R-squared:  0.849,    Adjusted R-squared:  0.8267
F-statistic: 37.96 on 4 and 27 DF,  p-value: 1.025e-10

```

Multivariate regression model explain 84.9% variance. It can be seen that *wt* and up to some extent *hp* confound the relationship between *am* and *mpg*.

Result plots

```
> par(mfrow = c(2,2))
```

```
> plot(mfit, col=3)
```

