

Supervised learning - Task 2 for TSF Internship

Pallabi Mandal

10/08/2020

Table of Contents :-

1. Import Data
2. EDA (Univariate and Bivariate analysis)
3. Correlation check
4. Data splicing
5. Build Linear regression model
6. Check Prediction

#Import libraries

```
library(ggplot2)
library(coefplot)
library(corrplot)
library(psych)
library(tidyverse)
library(MASS)
library(lattice)
library(DataExplorer)
library(ModelMetrics)
library(lmtest)
library(caret)
library(recipes)
library(caTools)
library(Hmisc)
library(GGally)
```

```
#Import data
```

```
setwd("P:/R file")  
getwd()
```

```
## [1] "P:/R file"
```

```
library(readr)  
Score_data <- read.csv("Student_score.csv")  
view(Score_data)  
str(Score_data)
```

```
## 'data.frame': 25 obs. of 2 variables:  
## $ Hours : num 2.5 5.1 3.2 8.5 3.5 1.5 9.2 5.5 8.3 2.7 ...  
## $ Scores: int 21 47 27 75 30 20 88 60 81 25 ...
```

```
summary(Score_data)
```

```
##      Hours      Scores  
## Min.   :1.100   Min.   :17.00  
## 1st Qu.:2.700   1st Qu.:30.00  
## Median :4.800   Median :47.00  
## Mean   :5.012   Mean   :51.48  
## 3rd Qu.:7.400   3rd Qu.:75.00  
## Max.   :9.200   Max.   :95.00
```

```
dim(Score_data)
```

```
## [1] 25 2
```

```
names(Score_data)
```

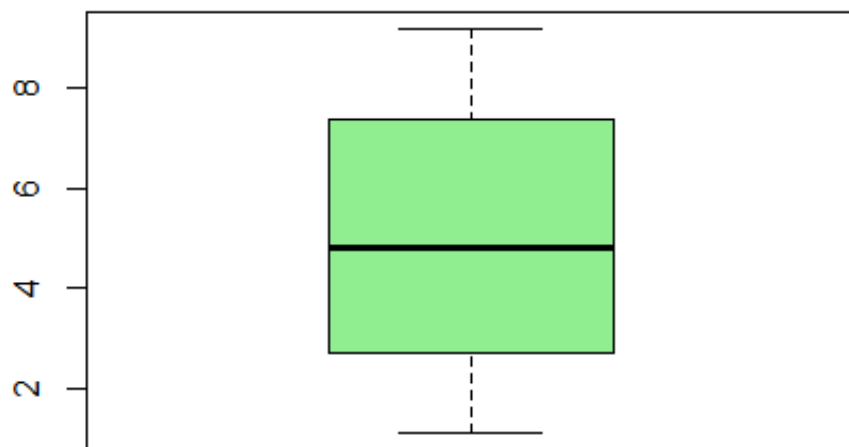
```
## [1] "Hours" "Scores"
```

```
attach(Score_data)
```

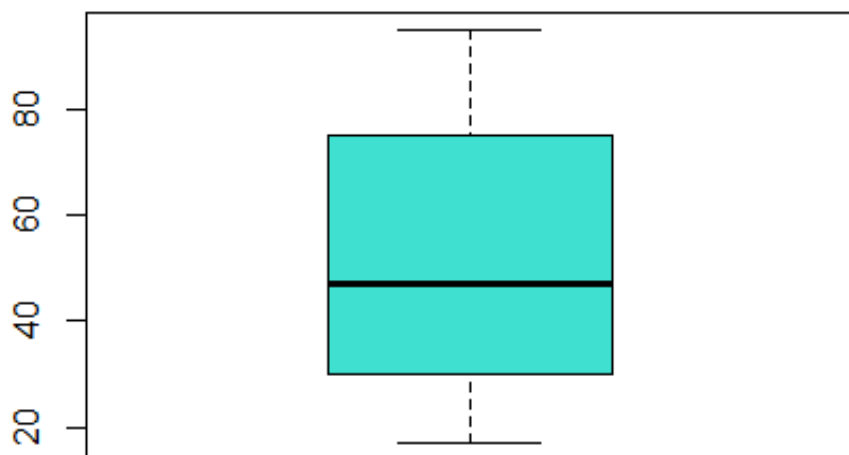
```
#EDA
```

```
#Univariate and Bivariate analysis
```

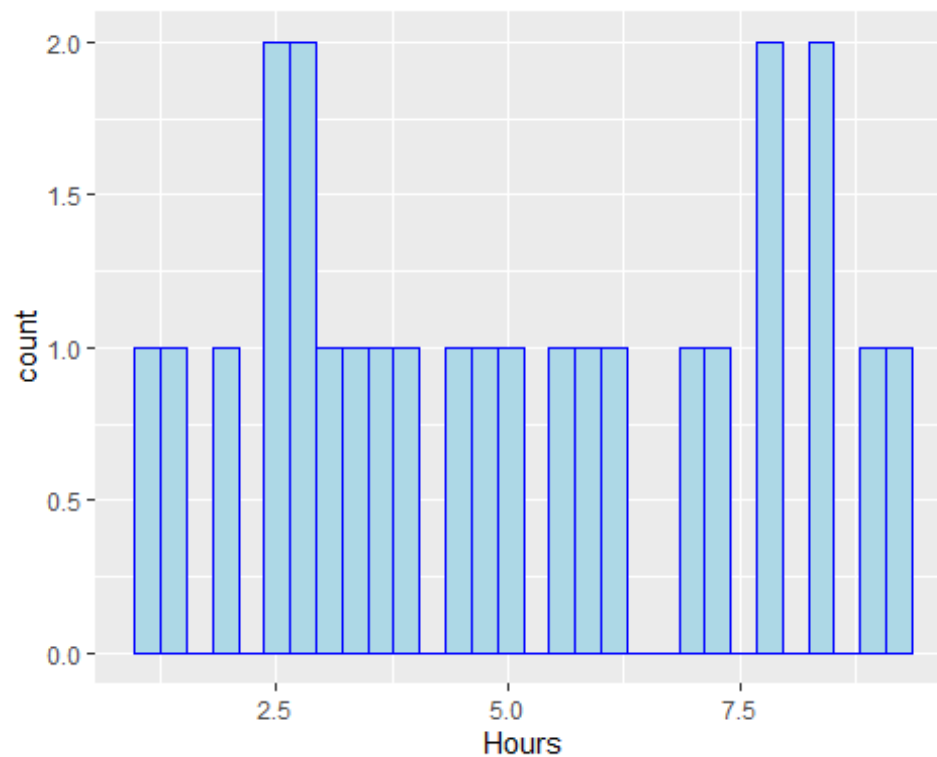
```
boxplot(Hours, col = "lightgreen")
```



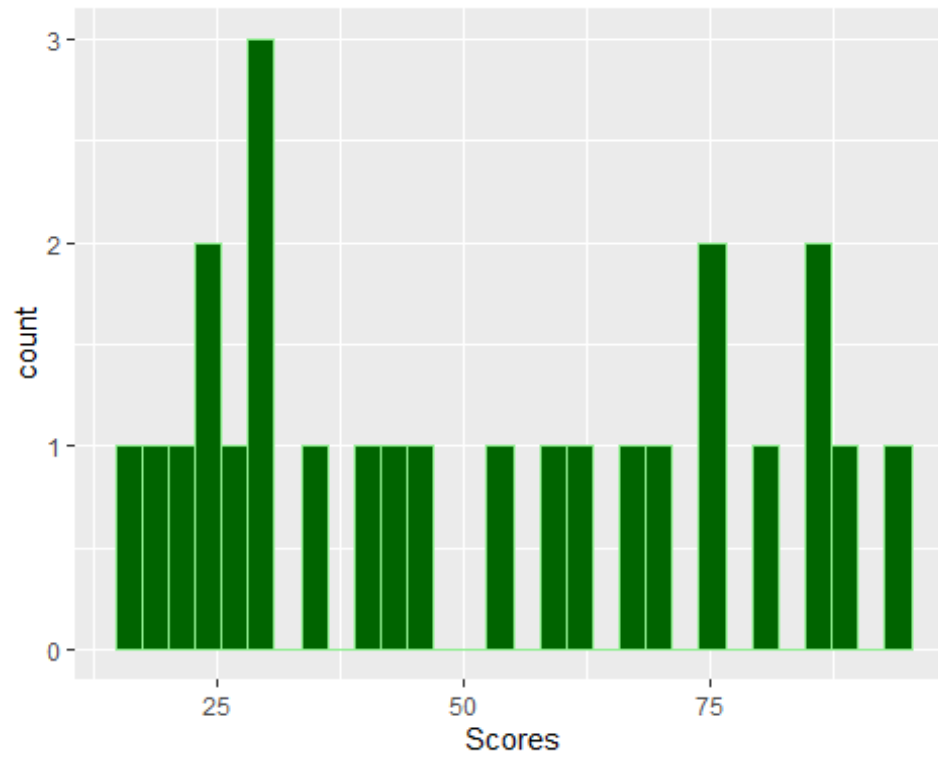
```
boxplot(Scores, col = "turquoise")
```



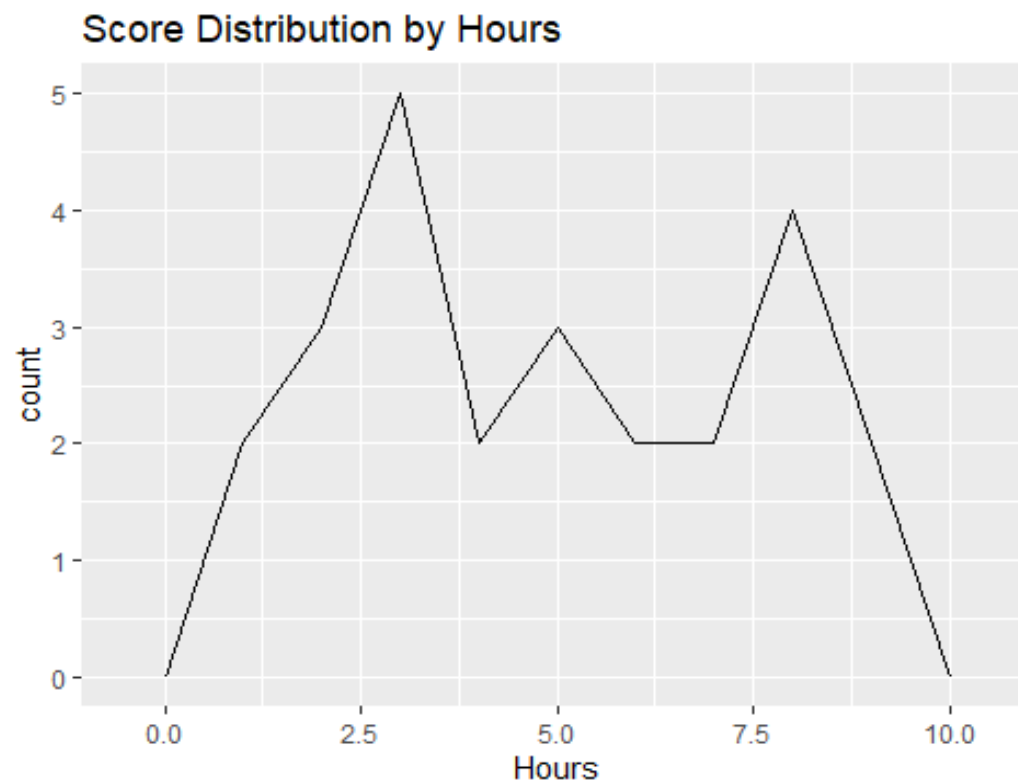
```
ggplot(Score_data, aes(x= Hours)) + geom_histogram(bins = 30, fill = "lightblue", col = "blue")
```



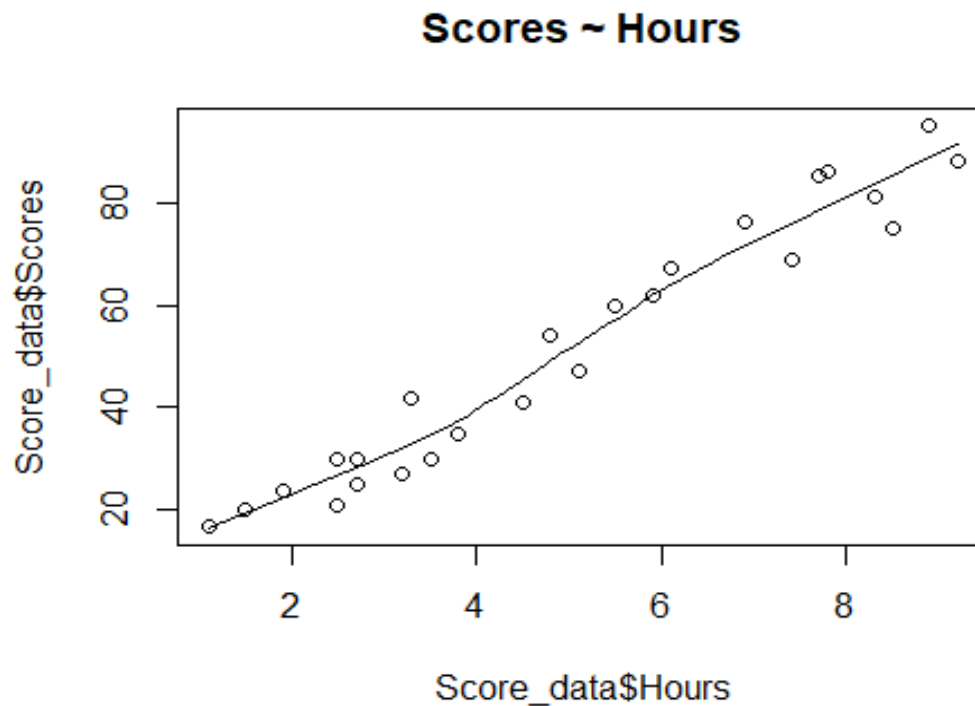
```
ggplot(Score_data, aes(x = Scores)) + geom_histogram(bins = 30, fill = "darkgreen", col = "lightgreen")
```



```
ggplot(Score_data, aes(Hours, colour = Scores)) +  
  geom_freqpoly(binwidth = 1) + labs(title="Score Distribution by Hours")
```



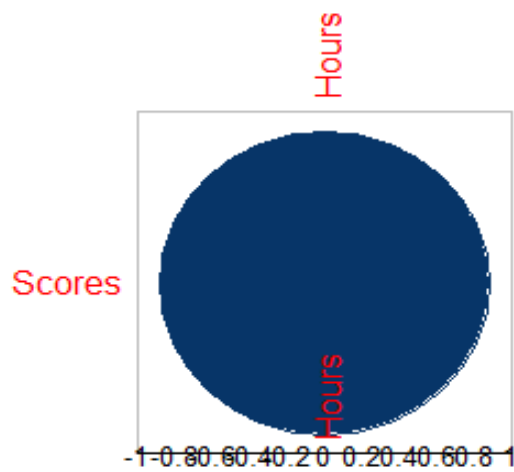
```
scatter.smooth(x=Score_data$Hours, y=Score_data$Scores, main="Scores ~ Hours"
)
```



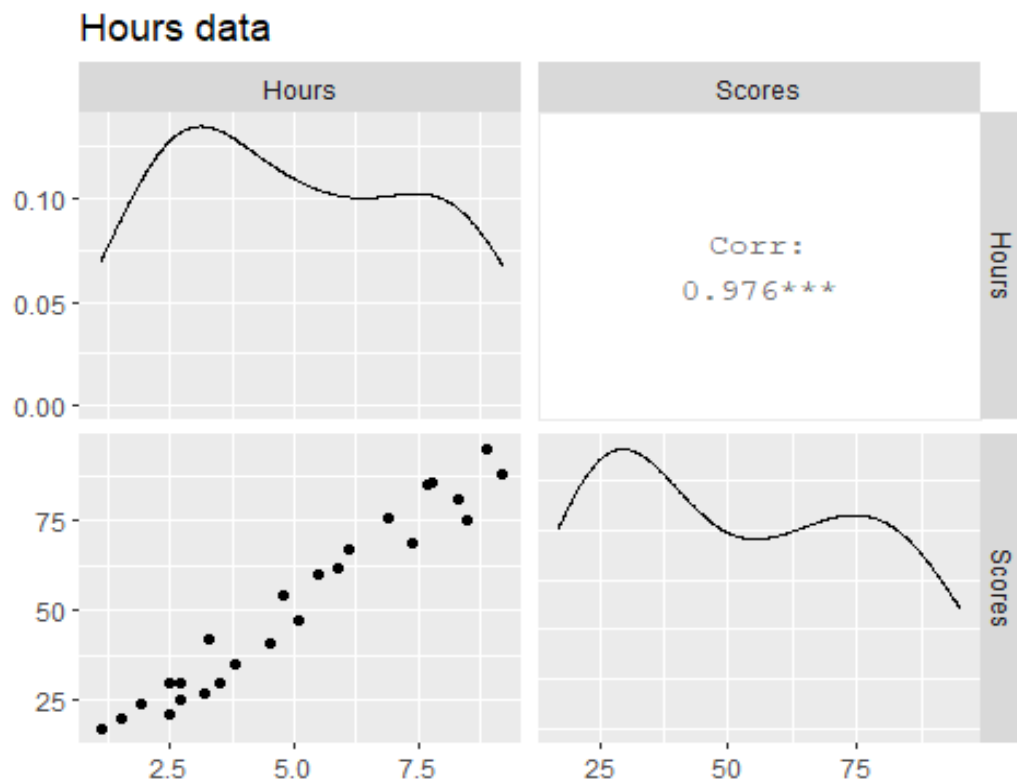
```
#Correlation
Score.cor = cor(Score_data)
Score.cor = cor(Score_data, method = c("pearson"))
library(Hmisc)
Score.rcorr = rcorr(as.matrix(Score_data))
Score.rcorr

##           Hours Scores
## Hours    1.00   0.98
## Scores   0.98   1.00
##
## n= 25
##
##
## P
##           Hours Scores
## Hours          0
## Scores         0

corrplot(Score.cor, type="lower", diag = FALSE)
```



```
ggpairs(data=Score_data, columns=1:2, title="Hours data")
```



```

cor(Score_data$Hours, Score_data$Scores)

## [1] 0.9761907

#Data splicing

#Data is split into 2 part Train and Test with 70:30 ratio.
split <- sample.split(Score_data$Scores, SplitRatio = 0.70)
length(split)

## [1] 25

train<-subset(Score_data, split == TRUE)
test<-subset(Score_data, split == FALSE)

#Model building

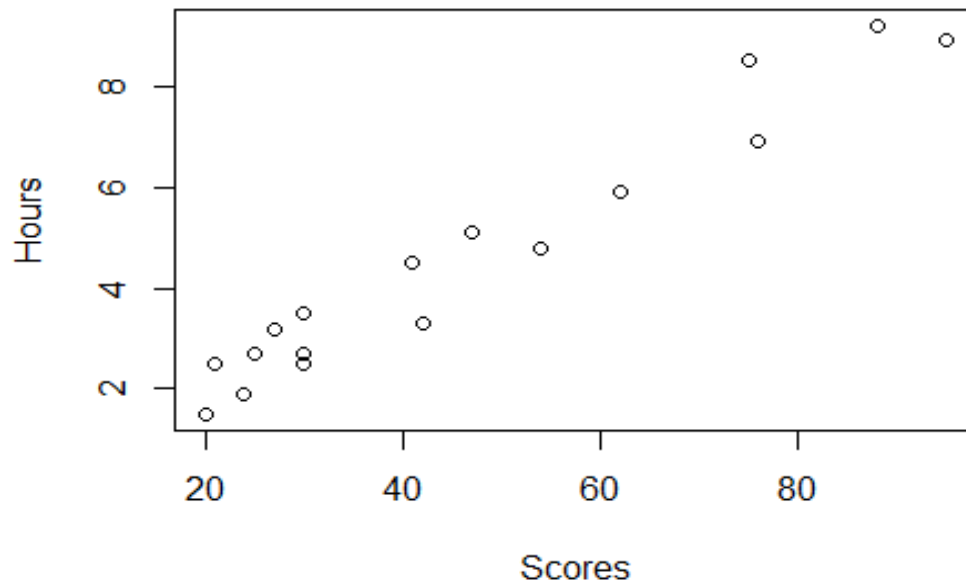
#Simple Linear Regression
set.seed(123)
Mod1 <- lm(Scores ~ ., data = train)
summary(Mod1)

##
## Call:
## lm(formula = Scores ~ ., data = train)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.921 -4.675  1.535  3.448  7.798
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   2.6486     2.9452   0.899   0.383
## Hours         9.5615     0.5698  16.779 3.95e-11 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.696 on 15 degrees of freedom
## Multiple R-squared:  0.9494, Adjusted R-squared:  0.946
## F-statistic: 281.6 on 1 and 15 DF, p-value: 3.949e-11

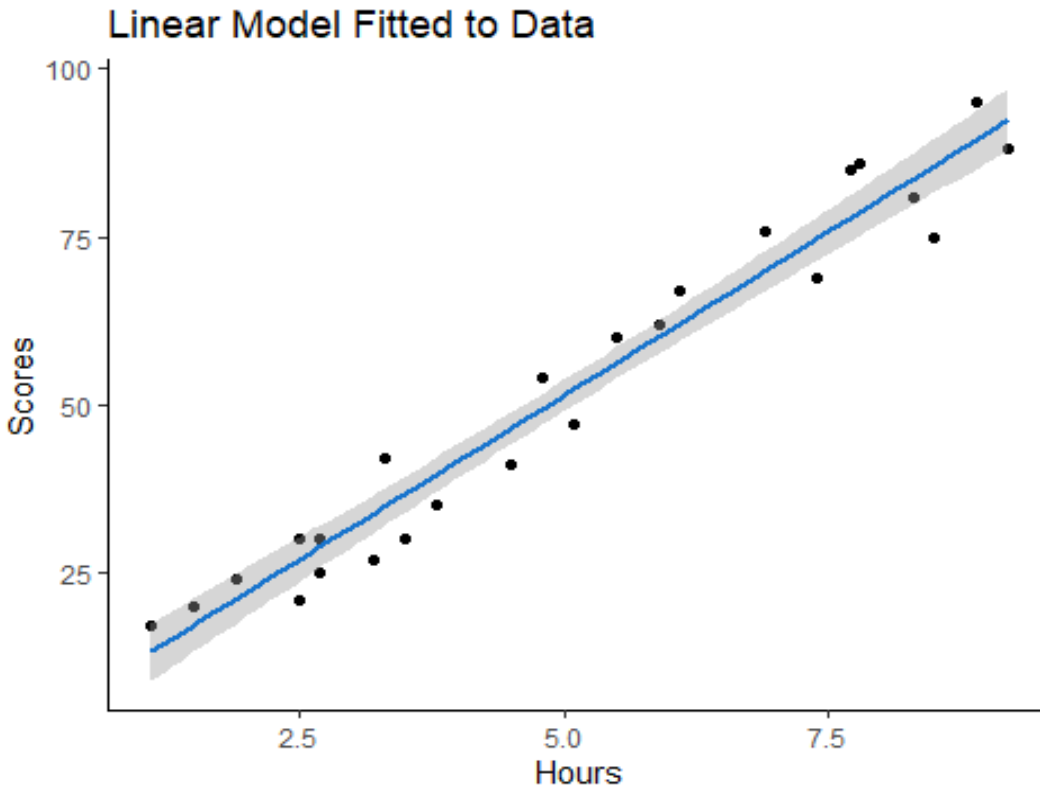
```



```
#Model plots
with(train,plot(Scores, Hours))
  abline(0, 1)
```



```
#Fitted line
ggplot(data = Score_data, aes(x = Hours, y = Scores)) +
  geom_point() +
  stat_smooth(method = "lm", col = "dodgerblue3") +
  theme(panel.background = element_rect(fill = "white"),
        axis.line.x=element_line(),
        axis.line.y=element_line()) +
  ggtitle("Linear Model Fitted to Data")
```



```
#K fold validation
set.seed(123)
train.control <- trainControl(method = "cv", number = 10)
#Train the model
model <- train(Scores~., data = train, method = "lm", trControl = train.contr
ol)
print(model)

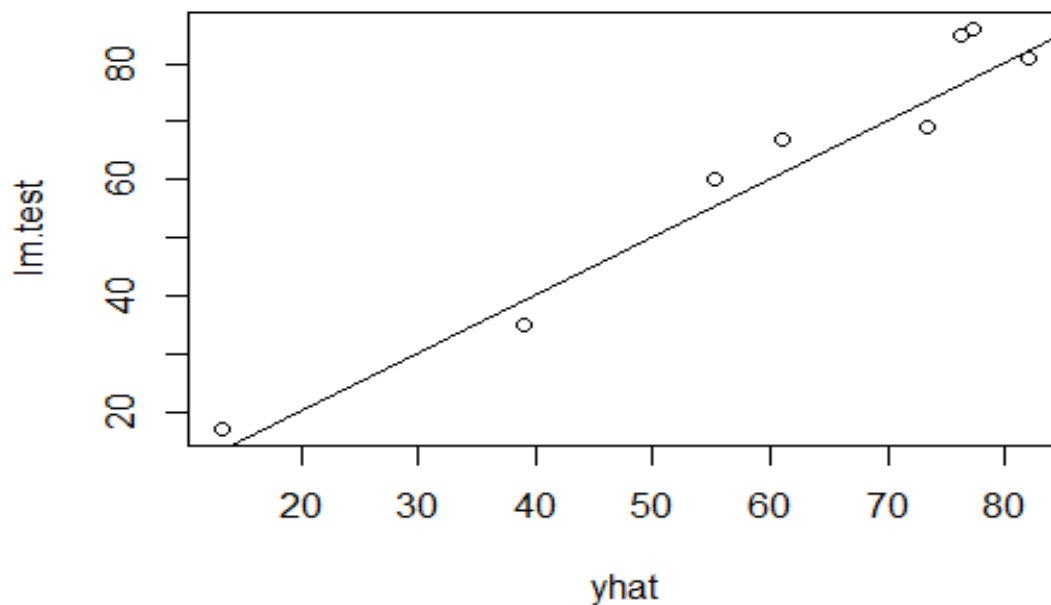
## Linear Regression
##
## 17 samples
## 1 predictor
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 15, 15, 15, 16, 15, 16, ...
## Resampling results:
##
##   RMSE      Rsquared   MAE
## 6.048492    1         5.73039
```

```
##
## Tuning parameter 'intercept' was held constant at a value of TRUE

#Predict MSE
yhat = predict(Mod1,newdata= test)
lm.test=test$Scores # These are the actual values
lm.test # Take a quick Look

## [1] 60 81 85 17 67 69 35 86

plot(yhat,lm.test) # Let's plot predicted vs. actual
abline(0,1) # And draw a line
```



```
mean((yhat-lm.test)^2) # And calculate the MSE

## [1] 32.88531

#Calculate RMSE, MAE score on Test data
predictions <- predict(Mod1, test)
postResample(test$Scores, predictions)

##      RMSE  Rsquared      MAE
## 5.7345714 0.9544756 5.1896996
```