# Business Analytics with Power BI

# Module 3 – Predictive Analytics with

# Power BI and R

*Student Lab Manual – Lab 2 – Creating Machine*

*Learning Models by Using R*

Version 1.0

**Conditions and Terms of Use**

**Microsoft Confidential**

This training package is proprietary and confidential, and is intended only for uses described in the training materials. Content and software is provided to you under a Non-Disclosure Agreement and cannot be distributed. Copying or disclosing all or any portion of the content and/or software included in such packages is strictly prohibited.

The contents of this package are for informational and training purposes only and are provided "as is" without warranty of any kind, whether express or implied, including but not limited to the implied warranties of merchantability, fitness for a particular purpose, and non-infringement.

Training package content, including URLs and other Internet Web site references, is subject to change without notice. Because Microsoft must respond to changing market conditions, the content should not be interpreted to be a commitment on the part of Microsoft, and Microsoft cannot guarantee the accuracy of any information presented after the date of publication. Unless otherwise noted, the companies, organizations, products, domain names, e-mail addresses, logos, people, places, and events depicted herein are fictitious, and no association with any real company, organization, product, domain name, e-mail address, logo, person, place, or event is intended or should be inferred.

**Copyright and Trademarks**

Module 3

Predictive Analytics with Power BI and R

Lab 2 –  Creating Machine Learning Models by Using R

# Lab 2: Creating Machine Learning Models by Using R

## Introduction

In this lab, you will create two machine learning models by using R language. One model will be to predict gas consumption of a car (regression model). The other model will be to classify flowers in three different species using linear discriminant analysis (LDA) and random forests (classification models).

## Objectives

After completing this lab, you will be able to:

- Create machine learning models by using R language.

## Estimated time to complete this lab

45 minutes (depends on experience)

## Resources

| | |
|---|---|
| Virtual machine (VM) Name | **Business Analytics with Power BI - Module 1** |
| Domain | **POWERBI-WIN10** |
| User | **POWERBI-WIN10\LabUser** |
| Password | **P@ssw0rd1!** |
| Lab Files | **E:\Labs\** |
| Asset Files | **E:\Assets\** |

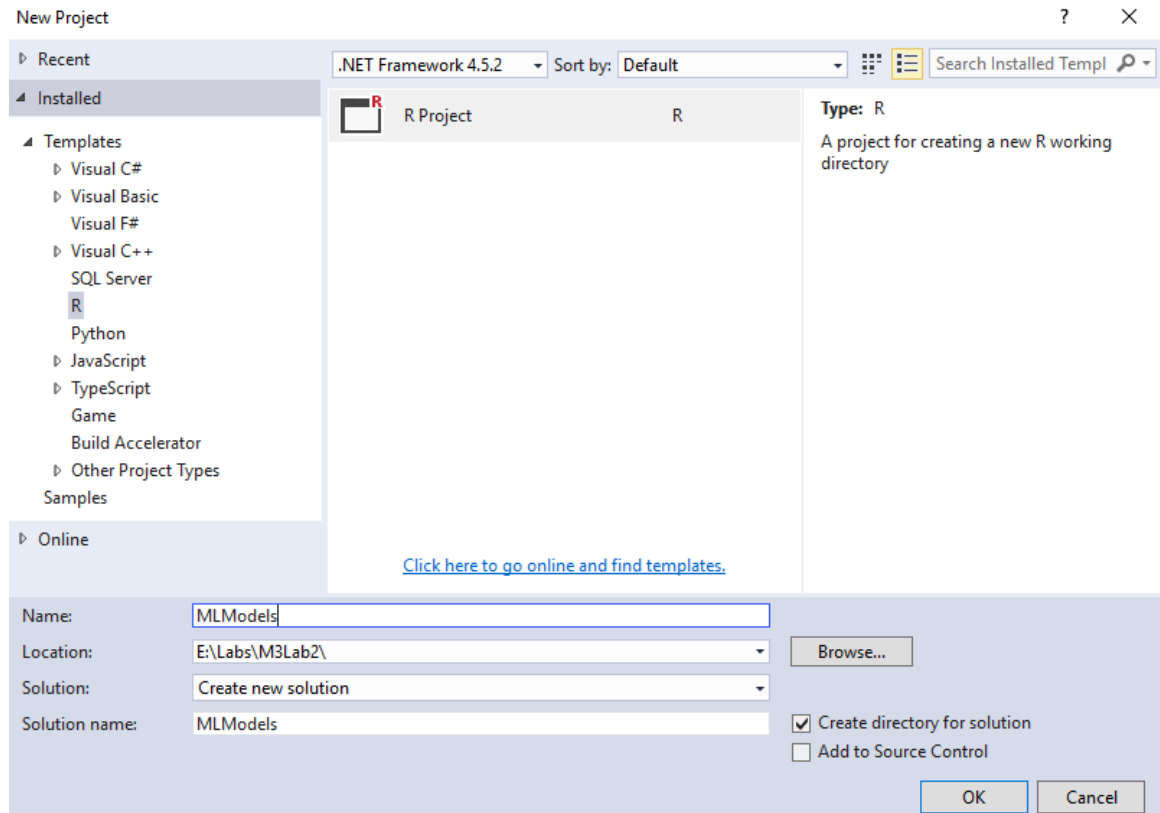# Exercise 1: Creating Machine Learning Models in R

## Creating a Regression Model to Predict Consumption

In this task, you will create a regression model to predict oil consumption based on cars stored in the **mtcars** data frame.

1. If Microsoft Visual Studio is not open, open it. On the taskbar, click the **Visual Studio 2015** shortcut.



2. Click the **File** tab, select **New**, and then select **Project**.

3. In the **New Project** window, click on the left side, **Templates** >> **R**.

4. Select **R Project**.

5. In the **Name** box, type **MLModels**, and in the **Location** box, browse to **E:\Labs\M3Lab2\**.

6. In the **Solution name** box, type **MLModels**.

7. Leave the **Create directory for solution** check box selected. You should see the following screen:
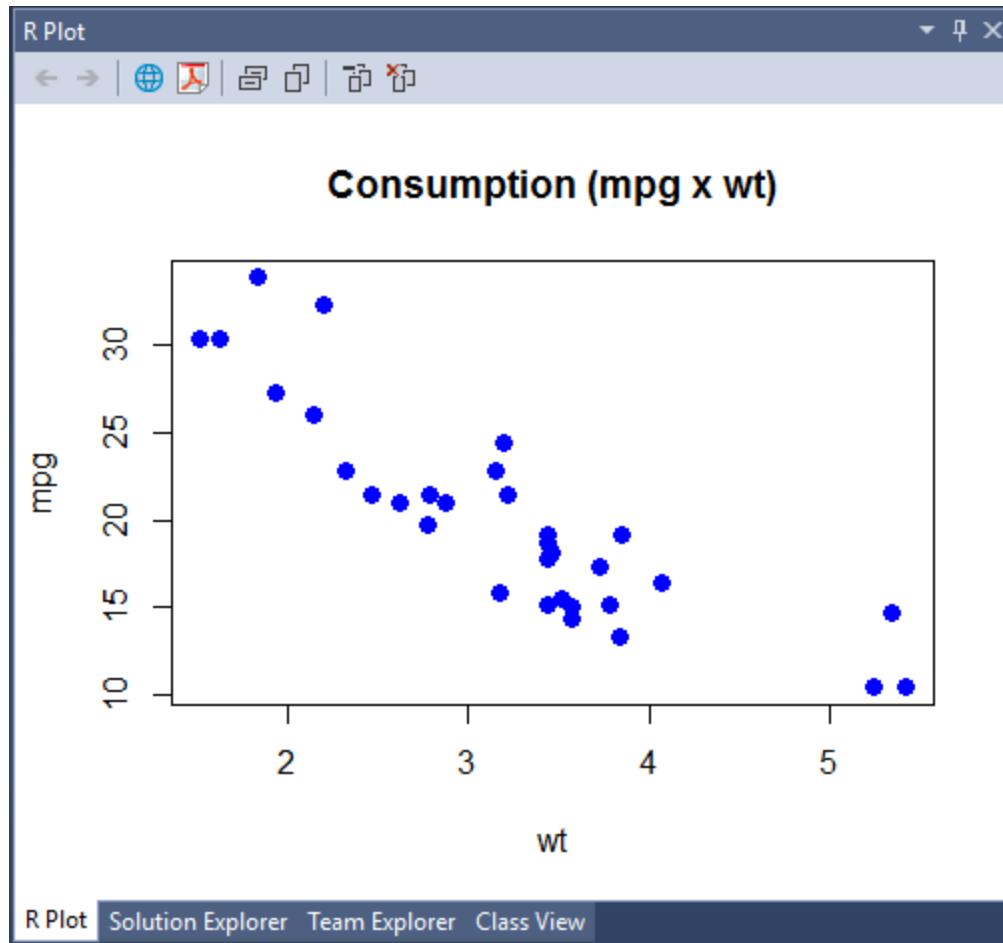
8.  Click **OK**. Make sure you save your progress as you progress through the lab.

9.  In the **Solution Explorer** pane, rename **script.R** to **Consumption.R**.

10. Add the following code to the **Consumption.R** script, and then run the code. Notice that these functions help you to understand your data.

```
tail(mtcars)
str(mtcars)
summary(mtcars)
```

11. Add the following lines and run it.  Notice that there is a correlation between weight (wt) and consumption (mpg).

```
mpg <- mtcars$mpg
wt <- mtcars$wt
plot(wt, mpg, pch = 16, cex = 1.3, col = "blue", main =
"Consumption (mpg x wt)", xlab = "wt", ylab = "mpg")
```
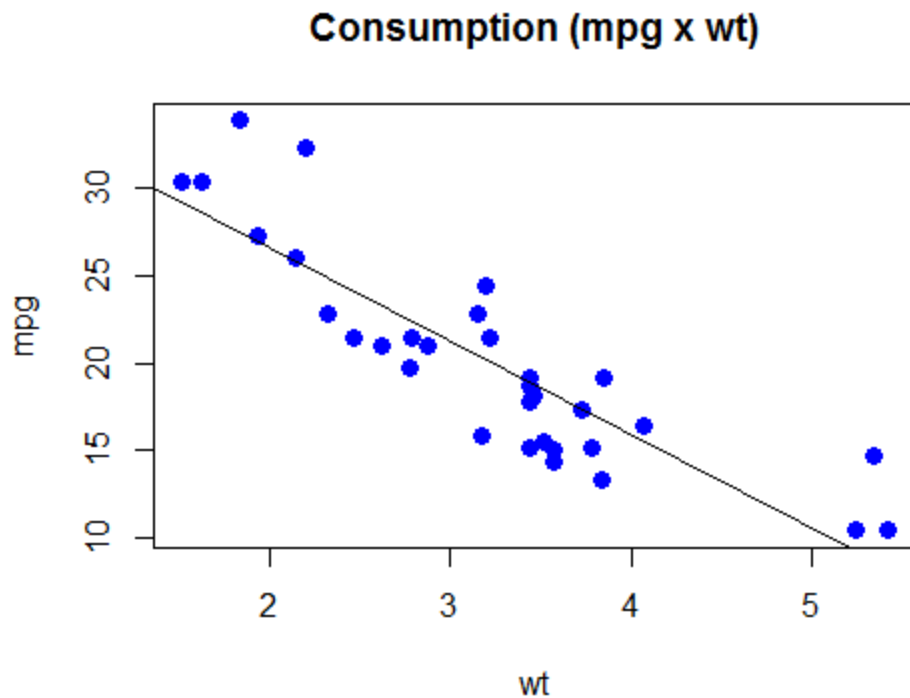
12. You should have the following plot in your **R Plot pane**:

13. Add the following code to the same file, and then run it. This code will create a linear model in which **mpg** is calculated based on **wt** only.

```
m <- lm(mpg ~ wt)
m
```

14. Use the calculated coefficients to plot a line over the last plot you made. Add the following line and then run it. Your plot should look like the following screenshot:

## Consumption (mpg x wt)



15. Test the model you have created. Add the following code to your **Consumption.R** script file and run it.

    ```
    s <- mtcars[1:5,] #get the first 5 lines
    data.frame(actual = s$mpg, predict = predict(m,s))
    ```

16. Your output (**R Interactive** pane) should have similar information. Notice that the predicted values are not exactly the same as the actual ones. This is due to the generalization that is needed in the model. We need to avoid overfitting. Generic models usually work better with real world examples.

    ```
                      actual  predict
    Mazda RX4          21.0 23.28261
    Mazda RX4 Wag      21.0 21.91977
    Datsun 710         22.8 24.88595
    Hornet 4 Drive     21.4 20.10265
    Hornet Sportabout  18.7 18.90014
    ```
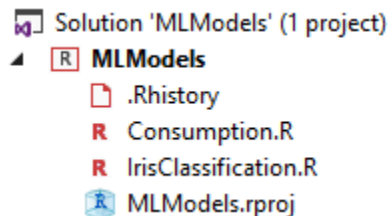
17. Save your script file.

## Creating a Classification Model to Predict a Class of the Iris Plant

In this task, you will create a classification model that will classify one of the three different classes of the Iris plant: Iris Setosa, Iris Versicolour, and Iris Virginica. We will use two different algorithms: Linear Discriminant Analysis (LDA) and Random Forests. Remember that there are several other algorithms that would be helpful in this task as well.

1. If Visual Studio is not open, open it. On the taskbar, click the **Visual Studio 2015** shortcut. Open the project created in the last task (**MLModels**).

2. Add a new script file (right-click the project's name, and then click **Add R Script**). Rename the file name to **IrisClassification.R**. You should have something similar to the following screenshot:
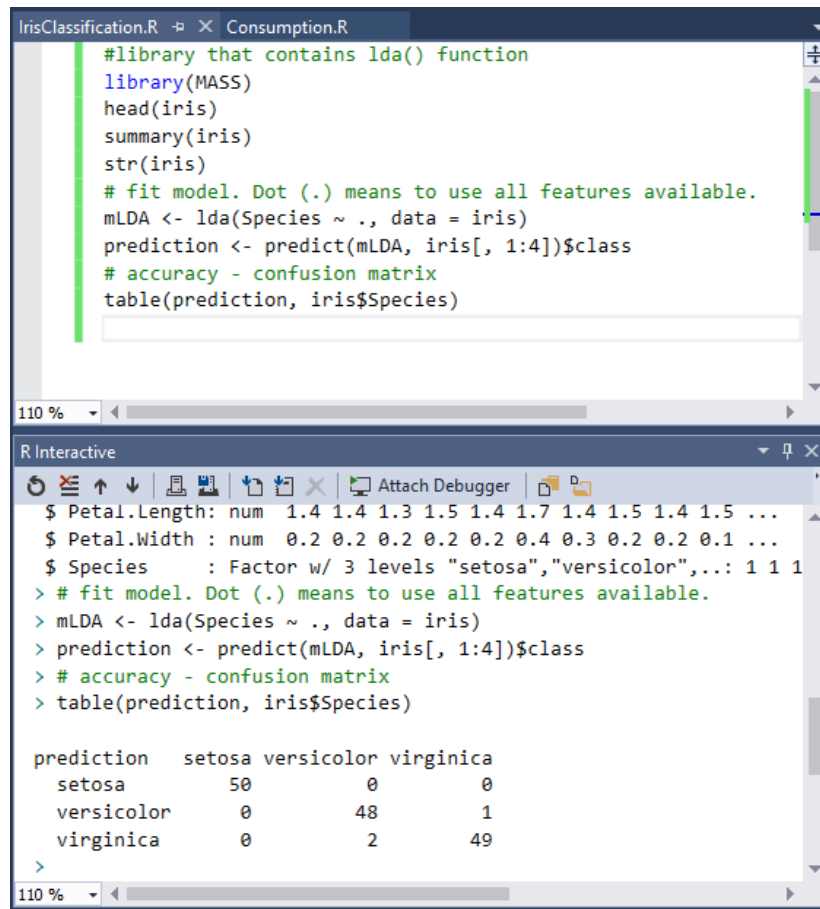


3. Add the following code to your **IrisClassification.R** script. You do not need to copy the comments. Run all the code.

```
#library that contains lda() function
library(MASS)
head(iris)
summary(iris)
str(iris)
# fit model. Dot (.) means to use all features available in iris
# dataframe.
mLDA <- lda(Species ~ ., data = iris)
prediction <- predict(mLDA, iris[, 1:4])$class
# accuracy - confusion matrix
table(prediction, iris$Species)
```

*Note: We usually create at least two datasets from our data: one to train the model and the other to validate it. This example was simplified for the purposes of this lab. We will see a more complex example in the last lab of this module.*

4. You should have the following result. This shows a very good classification, which is not necessarily true, because we used the same dataset to fit the model and to validate it.

In real-world cases, you should always validate a model with elements that were not used to fit/train your model.



5.  Add the following code to your script and then run it. This code simulates a new case to test against our model:

```
new_case <- data.frame(Sepal.Length = 5.1, Sepal.Width = 3.5,
Petal.Length = 1.4, Petal.Width = 0.2, Specie = "setosa")
new_case
data.frame(actual = new_case$Specie,
    predict = predict(mLDA, new_case[, 1:4])$class)
```

6.  Now that you already have a model using LDA, you will create another one using random forests. It is part of a data scientist's job to determine which algorithm and parameters are good enough for the problem.
    Add the following code to your script (**IrisClassification.R**), and then run it. Notice that for this algorithm, we will create two datasets (train and test) to validate our model.

```
install.packages("randomForest")
library(randomForest)

#split in two groups (~70% train | ~30% test)
i <- sample(2,nrow(iris),replace=TRUE,prob=c(0.7,0.3))
train <- iris[i==1,]
test <- iris[i==2,]

#create model
rf <-
randomForest(Species~.,data=train,ntree=100,proximity=TRUE)
print(rf)

#test model against test data
rf.pred <-predict(rf,newdata=test)
table(rf.pred, test$Species)
```

7. See the results created by your script. Save your project and close Visual Studio.