

Installation guide for R and RStudio

Step 1 – Install R

1. Download the R installer from <https://cran.r-project.org/>

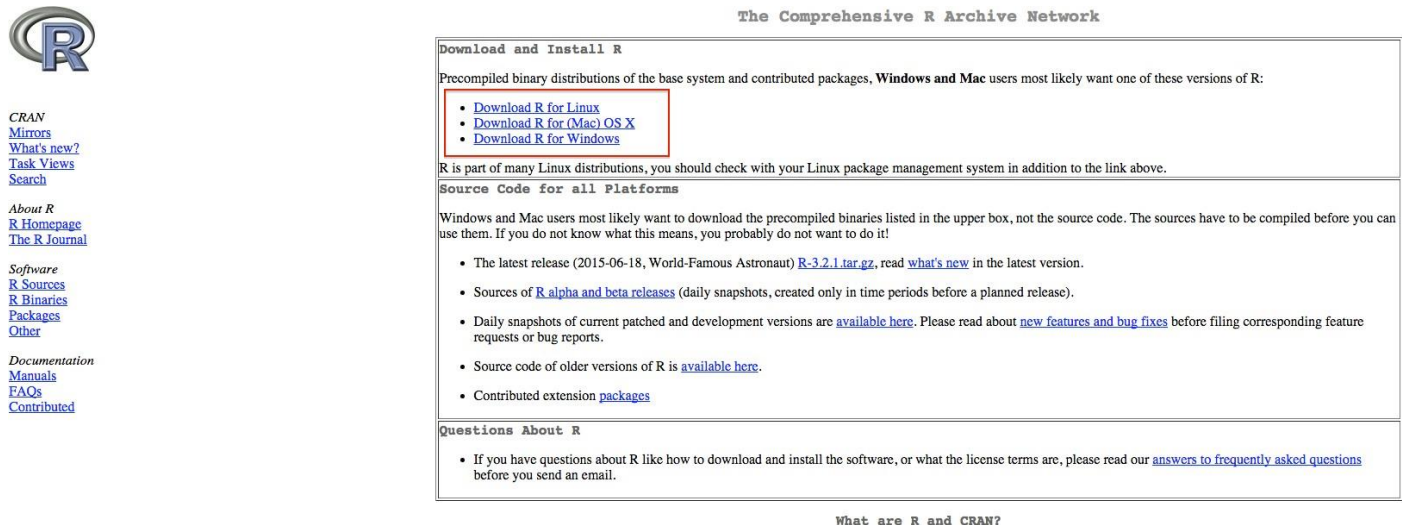


Figure 1. Screenshot of <http://cran.csiro.au/>

2. Run the installer. Default settings are fine. If you do not have admin rights on your laptop, then ask you local IT support. In that case, it is important that you also ask them to give you full permissions to the R directories. Without this, you will not be able to install additional packages later

Step 2 – Install RStudio

1. Download RStudio: <https://www.rstudio.com/products/rstudio/download/>

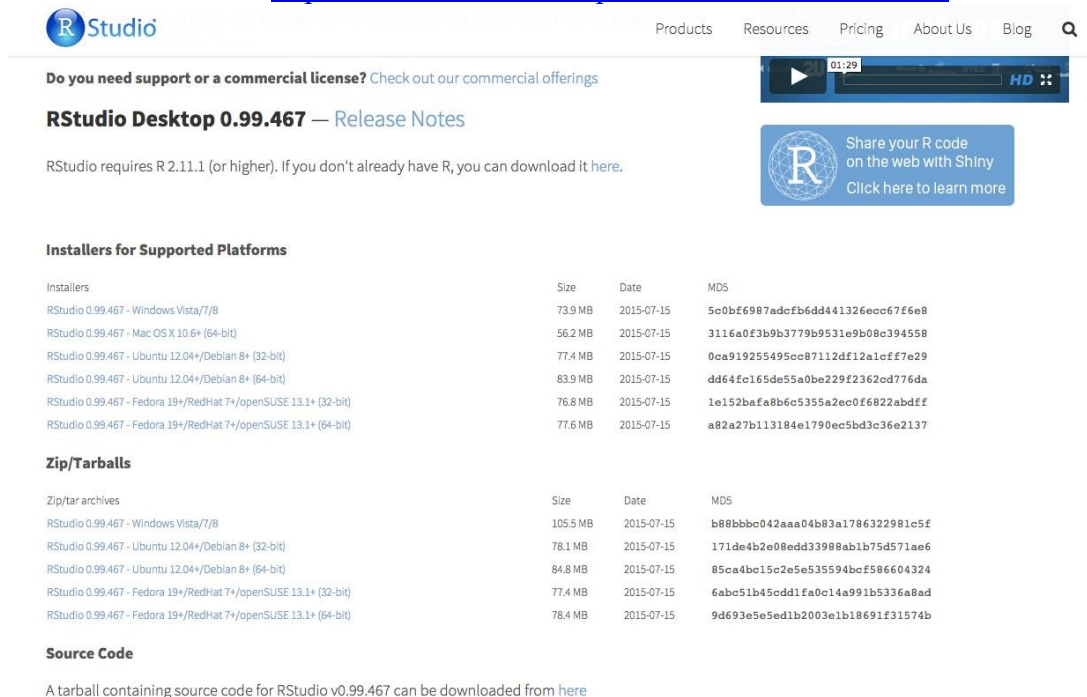


Figure 2. Download RStudio on <https://www.rstudio.com/products/rstudio/download/>

2. Once the installation of R has completed successfully (and not before), run the RStudio installer.

3. If you do not have administrative rights on your laptop, step 2 may fail. Ask your IT Support or download a pre-built zip archive of RStudio which doesn't need installing. The link for this is towards the bottom of the download page, highlighted in Image 2.
 - a. Download the appropriate archive for your system (Windows/Linux only – the Mac version can be installed into your personal “Applications” folder without admin rights).
 - b. Double clicking on the zip archive should automatically unpack it on most Windows machines.

Step 3 – Check that R and RStudio are working

1. Open RStudio. It should open a window that looks similar to image 3 below.
2. In the left hand window, by the '>' sign, type '4+5'(without the quotes) and hit enter. An outputline reading '[1] 9' should appear. This means that R and RStudio are working.
3. If this is not successful, contact us or your local IT support for further advice

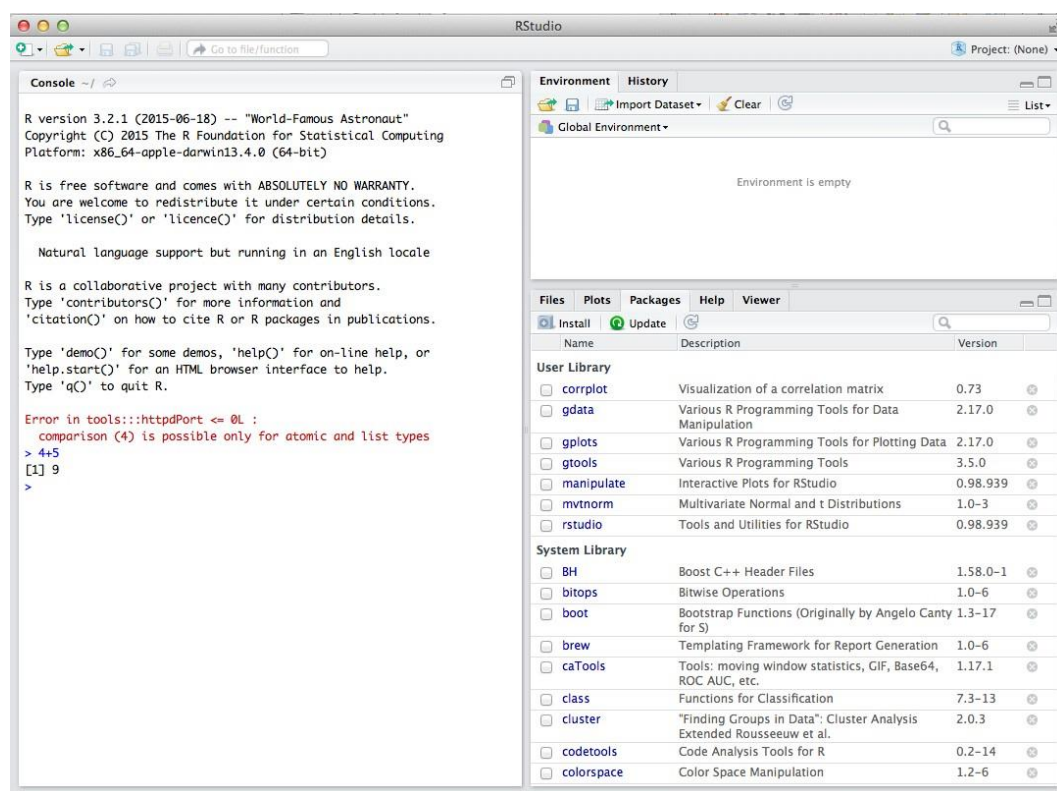


Figure 3. Running R with RStudio

Step 4 – Install R packages required for the workshop

1. Click on the tab ‘ Packages’ then ‘Install’ as shown in Image 4. Or Tools -> Install packages.
2. Install the following packages: mixOmics **version 6.1.0**, mvtnorm, RColorBrewer, corplot, igraph (see Image 4). For apple mac users, if you are unable to install the mixOmics imported library rgl, you will need to install the XQuartz software first <https://www.xquartz.org/>

3. Check that the packages are installed by typing 'library(mixOmics)' (without the quotes) in the prompt and press enter (see Image 5).
4. Then type 'sessionInfo()' and check that mixOmics version 6.1.0 has been installed (image 6).

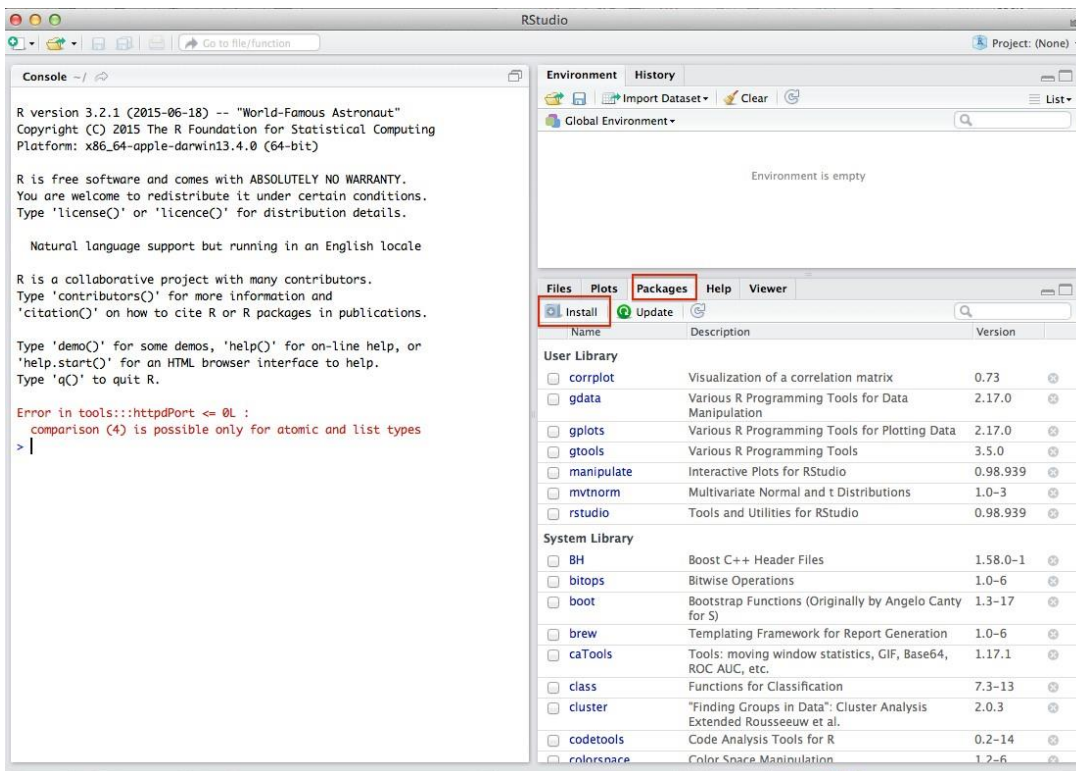


Figure 4. Click on Install to install R packages.

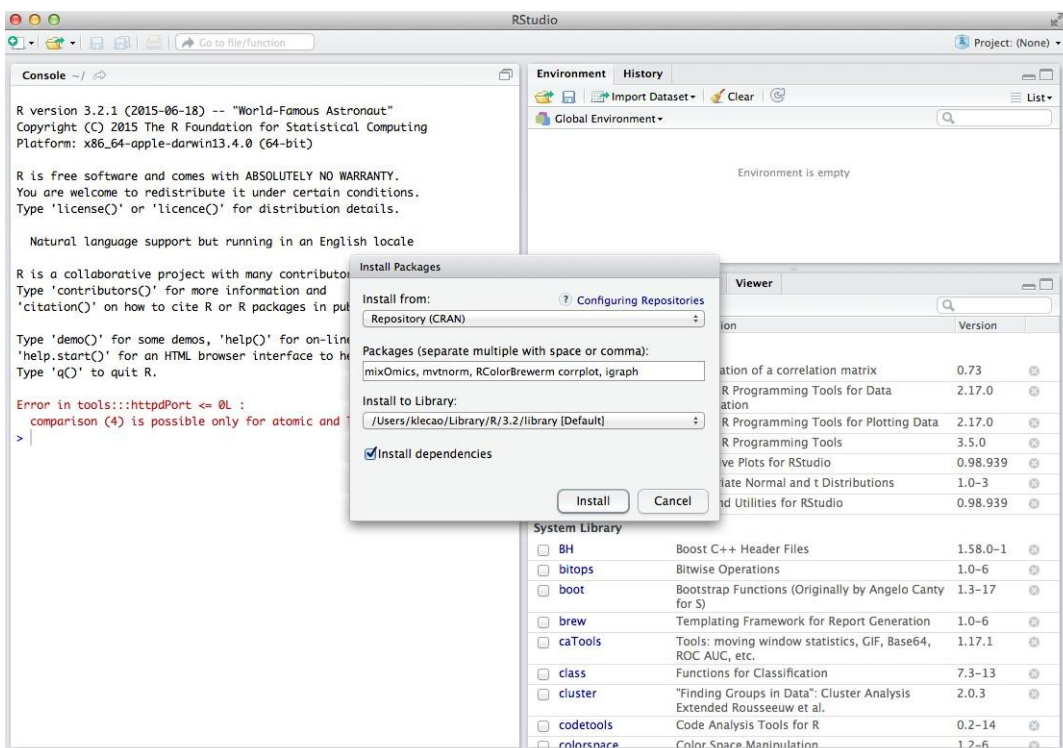


Figure 5. Specify the list of packages to be installed

The screenshot displays the RStudio interface with the following components:

- Console:**

```

> 4+5
[1] 9
> library(mixOmics)
Loading required package: MASS
Attaching package: 'MASS'
The following object is masked by_ '.GlobalEnv':
    genotype
Loading required package: lattice
Loading required package: ggplot2
Loaded mixOmics 6.1.0 ok!
Visit http://www.mixOmics.org for more details about our methods.
Any bug reports or comments? Notify us at mixomics at math.univ-toulouse.fr or https://bitbucket.org/klecao/package-mixomics/issues
Thank you for using mixOmics!
> sessionInfo()
R version 3.3.1 beta (2016-06-11 r70764)
Platform: x86_64-apple-darwin13.4.0 (64-bit)
Running under: OS X 10.9.5 (Mavericks)

locale:
 [1] en_AU.UTF-8/en_AU.UTF-8/en_AU.UTF-8/C/en_AU.UTF-8/en_AU.UTF-8
attached base packages:
[1] stats    graphics  grDevices  utils      datasets  methods   base

other attached packages:
[1] mixOmics_6.1.0  ggplot2_2.1.0  lattice_0.20-33 MASS_7.3-45

loaded via a namespace (and not attached):
 [1] rgl_0.95.1441  Rcpp_0.12.6    tidyr_0.5.0    corpcor_1.6.8
 [5] assertthat_0.1 dplyr_0.5.0    R6_2.1.3       grid_3.3.1
 [9] plyr_1.8.4     DBI_0.5        gtable_0.2.0   magrittr_1.5
[13] ellipse_0.3-8  scales_0.4.0   stringr_1.1.1  reshape2_1.4.1
[17] RColorBrewer_1.1-2 tools_3.3.1    stringr_1.1.0  munsell_0.4.3
[21] igraph_1.0.1   parallel_3.3.1 colorspace_1.2-6 tibble_1.1

```
- Environment:**
 - Data:**
 - data: 32423 obs. of 48 variables
 - data.gene: Large matrix (240000 elements, 2.1 Mb)
 - data.physio: 48 obs. of 10 variables
 - design: num [1:2, 1:2] 0 1 1 0
 - name.gene: 32423 obs. of 82 variables
 - Values:**
 - d: List of 7
 - diablo.res: Large block.splsda (24 elements, 3.6 Mb)
 - genotype: Factor w/ 8 levels "Inedi", "Melod", ...: 1 1 1 1 1 2 2 2 ...
 - k: 48L
 - kee.genes: chr [1:5000] "Heli058698_st" "Heli092737_st" "Heli058195_..."
 - keep.genes: chr [1:5000] "Heli058698_st" "Heli092737_st" "Heli058195_..."
 - keep.name.genes: chr [1:5000] "unknw" "unknw" "unknw" "arginase," "unknw" ...
 - list.data: Large list (2 elements, 2.1 Mb)
- Files Plots Packages Help Viewer:**
 - Install:**
 - ace() and avas() for selecting regression transformations: 1.3-3.3
 - ade4: Analysis of Ecological Data: Exploratory and Euclidean Methods in Environmental Sciences: 1.7-4
 - ALL: A data package: 1.14.0
 - annotate: Annotation for microarrays: 1.50.0
 - AnnotationDbi: Annotation Database Interface: 1.34.4
 - astsa: Applied Statistical Time Series Analysis: 1.4
 - Biobase: Base functions for Bioconductor: 2.32.0
 - BiocGenerics: S4 generic functions for Bioconductor: 0.18.0
 - BiocInstaller: Install/Update Bioconductor, CRAN, and github Packages: 1.22.3
 - BiocParallel: Bioconductor facilities for parallel evaluation: 1.6.2
 - capushe: CALibrating Penalties Using Slope HEuristics: 1.1.1
 - car: Companion to Applied Regression: 2.1-2
 - chron: Chronological Objects which can Handle Dates and Times: 2.3-47
 - cValid: Validation of Clustering Results: 0.6-6

Figure 6. Check that the package mixOmics is installed and has the version 6.1.0.

Exp:7

Implement Linear and Logistic Regressiona

a)Linear regression

```
# Sample data
heights <- c(150, 160, 165, 170, 175, 180, 185)
weights <- c(55, 60, 62, 68, 70, 75, 80)

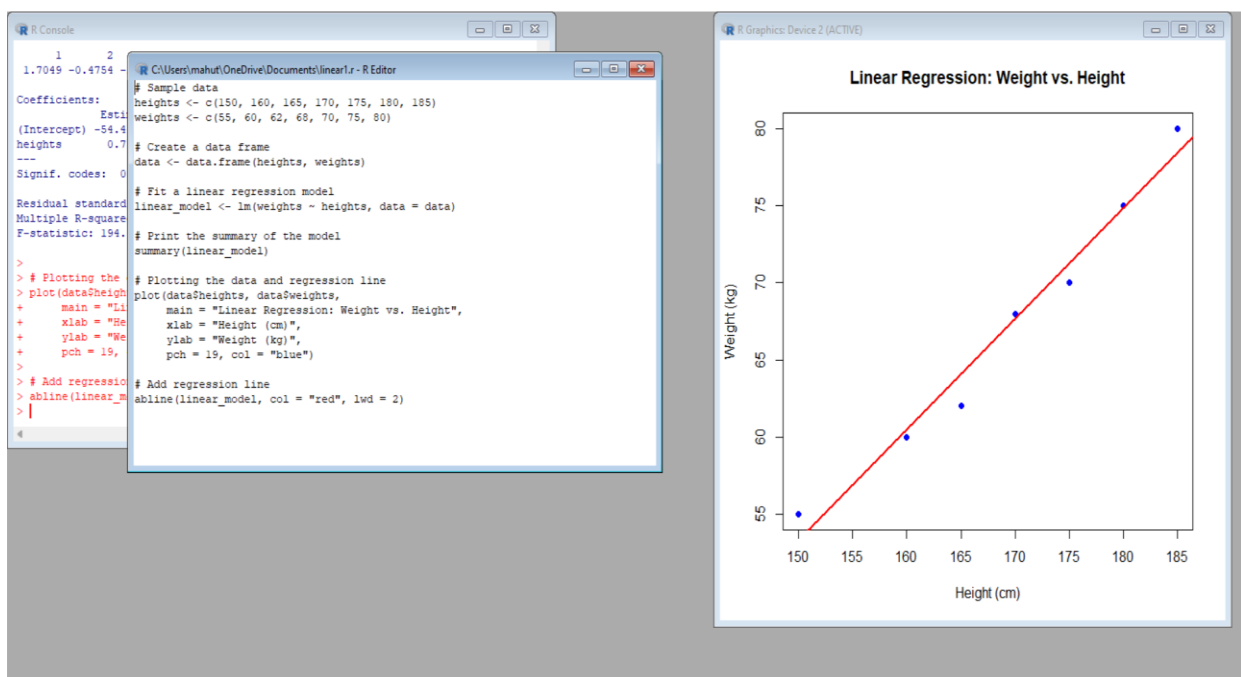
# Create a data frame
data <- data.frame(heights, weights)

# Fit a linear regression model
linear_model <- lm(weights ~ heights, data = data)

# Print the summary of the model
print(summary(linear_model))

# Plotting the data and regression line
plot(data$heights, data$weights,
     main = "Linear Regression: Weight vs. Height",
     xlab = "Height (cm)",
     ylab = "Weight (kg)",
     pch = 19, col = "blue")

# Add regression line
abline(linear_model, col = "red", lwd = 2)
```



b) Logistic regression

```
# Load the dataset
```

```
data(mtcars)
```

```
# Convert 'am' to a factor (categorical variable)
```

```
mtcars$am <- factor(mtcars$am, levels = c(0, 1), labels = c("Automatic", "Manual"))
```

```
# Fit a logistic regression model
```

```
logistic_model <- glm(am ~ mpg, data = mtcars, family = binomial)
```

```
# Print the summary of the model
```

```
print(summary(logistic_model))
```

```
# Predict probabilities for the logistic model
```

```
predicted_probs <- predict(logistic_model, type = "response")
```

```
# Display the predicted probabilities
```

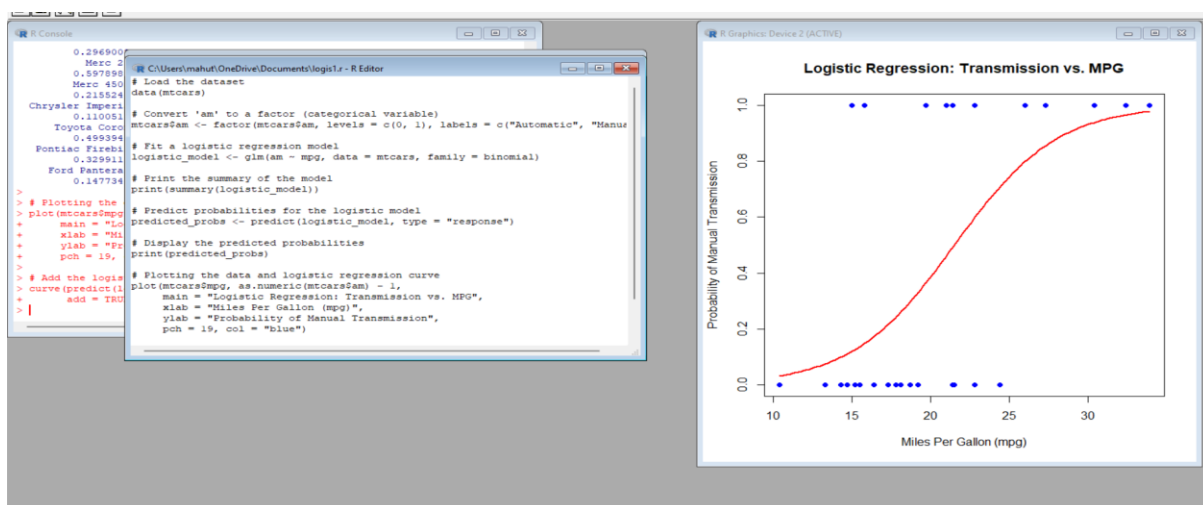
```
print(predicted_probs)
```

```
# Plotting the data and logistic regression curve
```

```
plot(mtcars$mpg, as.numeric(mtcars$am) - 1,  
     main = "Logistic Regression: Transmission vs. MPG",  
     xlab = "Miles Per Gallon (mpg)",  
     ylab = "Probability of Manual Transmission",  
     pch = 19, col = "blue")
```

```
# Add the logistic regression curve
```

```
curve(predict(logistic_model, data.frame(mpg = x), type = "response"),  
      add = TRUE, col = "red", lwd = 2)
```



Exp8:

Implement SVM/Decision tree classification techniques

a) SVM IN R

```
# Install and load the e1071 package (if not already installed)
install.packages("e1071")
library(e1071)

# Load the iris dataset
data(iris)

# Inspect the first few rows of the dataset
head(iris)

# Split the data into training (70%) and testing (30%) sets
set.seed(123) # For reproducibility
sample_indices <- sample(1:nrow(iris), 0.7 * nrow(iris))
train_data <- iris[sample_indices, ]
test_data <- iris[-sample_indices, ]

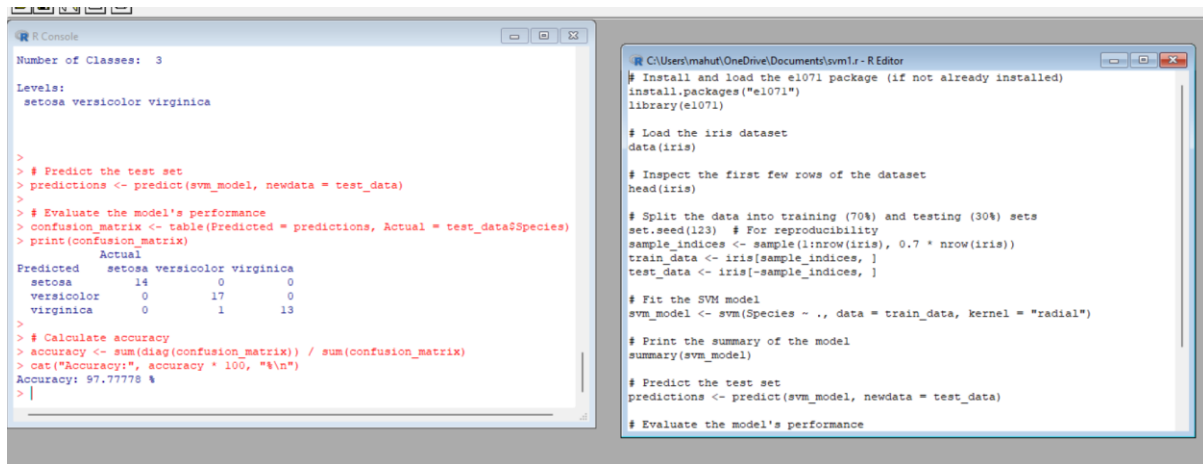
# Fit the SVM model
svm_model <- svm(Species ~ ., data = train_data, kernel = "radial")

# Print the summary of the model
summary(svm_model)

# Predict the test set
predictions <- predict(svm_model, newdata = test_data)

# Evaluate the model's performance
confusion_matrix <- table(Predicted = predictions, Actual = test_data$Species)
print(confusion_matrix)

# Calculate accuracy
accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
cat("Accuracy:", accuracy * 100, "%\n")
```



b) Decision tree in R

Install and load the rpart package (if not already installed)

```
install.packages("rpart")
```

```
library(rpart)
```

Load the iris dataset

```
data(iris)
```

Split the data into training (70%) and testing (30%) sets

```
set.seed(123) # For reproducibility
```

```
sample_indices <- sample(1:nrow(iris), 0.7 * nrow(iris))
```

```
train_data <- iris[sample_indices, ]
```

```
test_data <- iris[-sample_indices, ]
```

Fit the Decision Tree model

```
tree_model <- rpart(Species ~ ., data = train_data, method = "class")
```

Print the summary of the model

```
summary(tree_model)
```

Plot the Decision Tree

```
plot(tree_model)
```

```
text(tree_model, pretty = 0)
```

Predict the test set

```
predictions <- predict(tree_model, newdata = test_data, type = "class")
```

Evaluate the model's performance

```
confusion_matrix <- table(Predicted = predictions, Actual = test_data$Species)
```

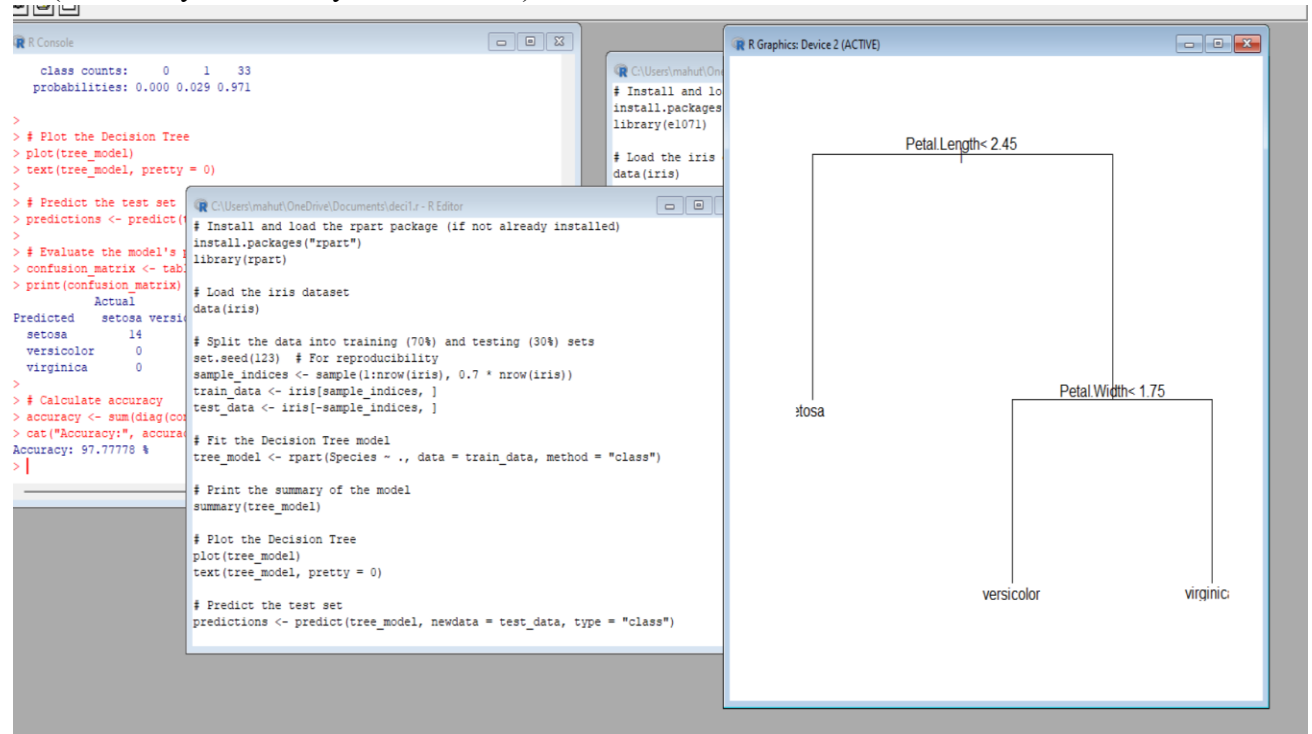
```
print(confusion_matrix)
```



```
# Calculate accuracy
```

```
accuracy <- sum(diag(confusion_matrix)) / sum(confusion_matrix)
```

```
cat("Accuracy:", accuracy * 100, "%\n")
```



Exp:9

Implement clustering techniques – Hierarchical and K-Means

a) HIERARCHIAL CLUSTERING

```
# Load the iris dataset
data(iris)

# Use only the numeric columns for clustering (exclude the Species column)
iris_data <- iris[, -5]

# Standardize the data
iris_scaled <- scale(iris_data)

# Compute the distance matrix
distance_matrix <- dist(iris_scaled, method = "euclidean")

# Perform hierarchical clustering using the "complete" linkage method
hc_complete <- hclust(distance_matrix, method = "complete")

# Plot the dendrogram
plot(hc_complete, main = "Hierarchical Clustering Dendrogram", xlab = "", sub = "", cex = 0.6)

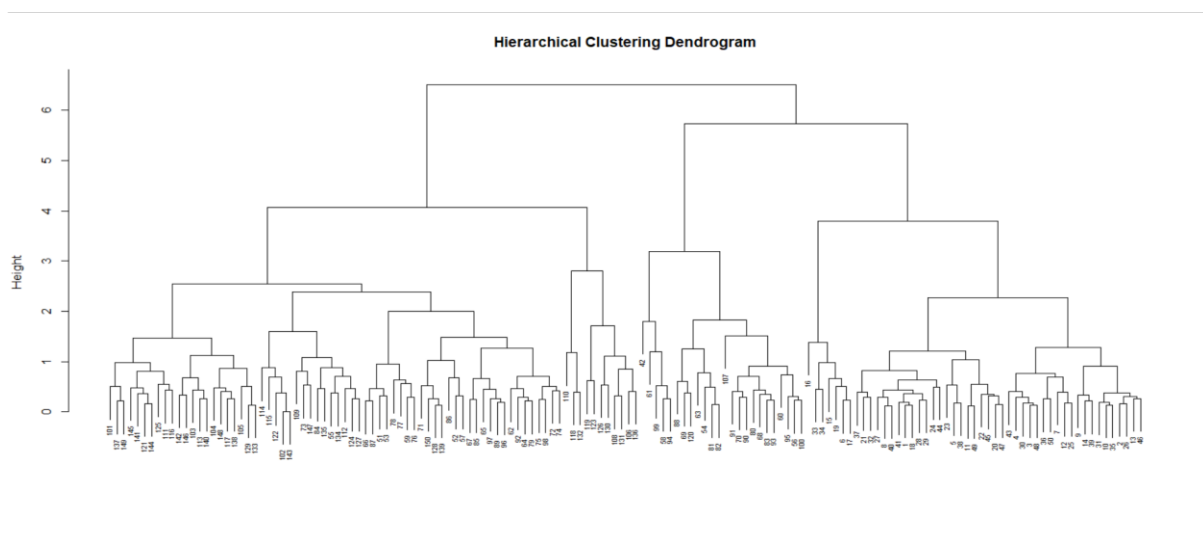
# Cut the tree to form 3 clusters
clusters <- cutree(hc_complete, k = 3)

# Print the cluster memberships
print(clusters)

# Add the clusters to the original dataset
iris$Cluster <- as.factor(clusters)

# Display the first few rows of the updated dataset
head(iris)
```

```
R Console
>
> # Cut the tree to form 3 clusters
> clusters <- cutree(hc_complete, k = 3)
>
> # Print the cluster memberships
> print(clusters)
[1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
[38] 1 1 1 1 2 1 1 1 1 1 1 1 1 3 3 3 2 3 2 3 2 3 2 2 3 2 3 3 3 2 2 2 3 3 3
[75] 3 3 3 3 3 2 2 2 3 3 3 3 2 3 2 2 3 2 2 2 3 3 3 2 2 3 3 3 3 3 3 2 3 3 3
[112] 3 3 3 3 3 3 3 3 2 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3 3
[149] 3 3
>
> # Add the clusters to the original dataset
> iris$Cluster <- as.factor(clusters)
>
> # Display the first few rows of the updated dataset
> head(iris)
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species Cluster
1          5.1         3.5        1.4         0.2   setosa        1
2          4.9         3.0        1.4         0.2   setosa        1
3          4.7         3.2        1.3         0.2   setosa        1
4          4.6         3.1        1.5         0.2   setosa        1
5          5.0         3.6        1.4         0.2   setosa        1
6          5.4         3.9        1.7         0.4   setosa        1
> |
```



b) K-MEANS CLUSTERING

```
# Load the iris dataset
data(iris)
```

```
# Use only the numeric columns for clustering (exclude the Species column)
iris_data <- iris[, -5]
```

```

# Standardize the data
iris_scaled <- scale(iris_data)

# Set the number of clusters
set.seed(123) # For reproducibility
k <- 3 # Number of clusters

# Perform K-Means clustering
kmeans_result <- kmeans(iris_scaled, centers = k, nstart = 25)

# Print the K-Means result
print(kmeans_result)

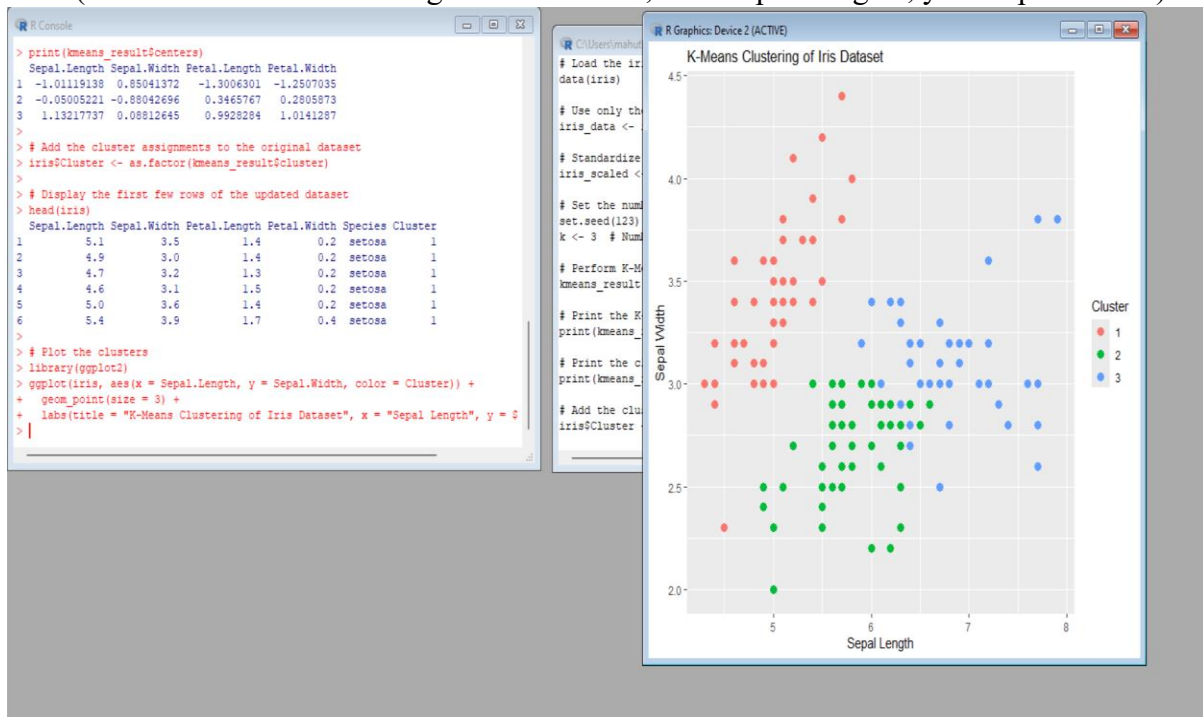
# Print the cluster centers
print(kmeans_result$centers)

# Add the cluster assignments to the original dataset
iris$Cluster <- as.factor(kmeans_result$cluster)

# Display the first few rows of the updated dataset
head(iris)

# Plot the clusters
library(ggplot2)
ggplot(iris, aes(x = Sepal.Length, y = Sepal.Width, color = Cluster)) +
  geom_point(size = 3) +
  labs(title = "K-Means Clustering of Iris Dataset", x = "Sepal Length", y = "Sepal Width")

```



Exp:10

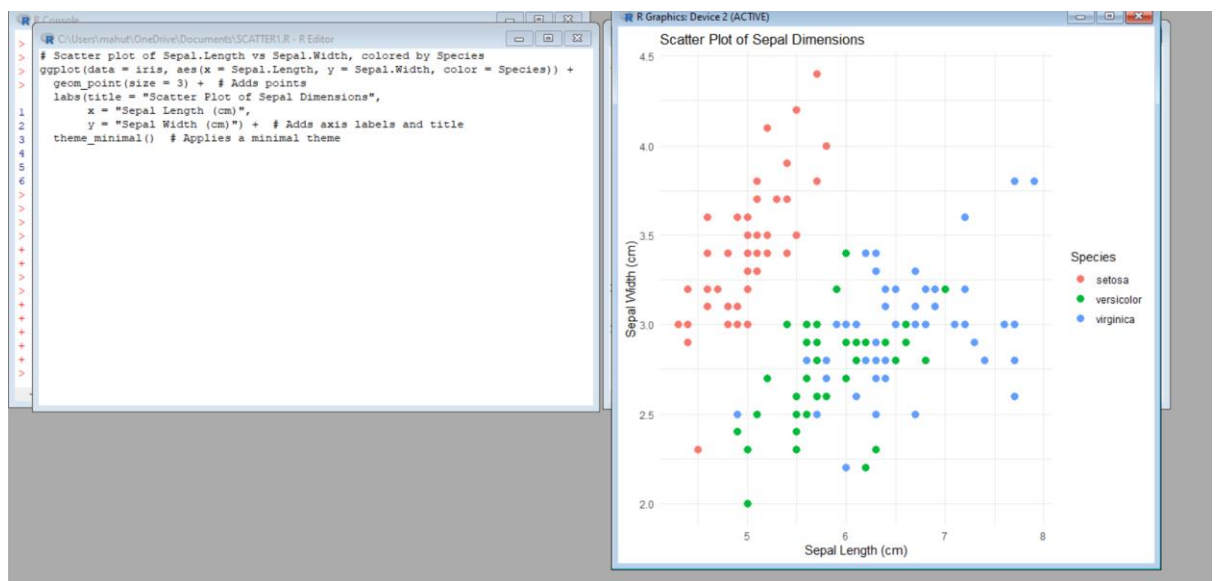
VISUALIZE DATA USING ANY PLOTTING FRAMEWORK

1) SCATTER PLOT

```
# Install ggplot2 (if not already installed)
install.packages("ggplot2")

# Load the ggplot2 package
library(ggplot2)

# Scatter plot of Sepal.Length vs Sepal.Width, colored by Species
ggplot(data = iris, aes(x = Sepal.Length, y = Sepal.Width, color = Species)) +
  geom_point(size = 3) + # Adds points
  labs(title = "Scatter Plot of Sepal Dimensions",
       x = "Sepal Length (cm)",
       y = "Sepal Width (cm)") + # Adds axis labels and title
  theme_minimal() # Applies a minimal theme
```



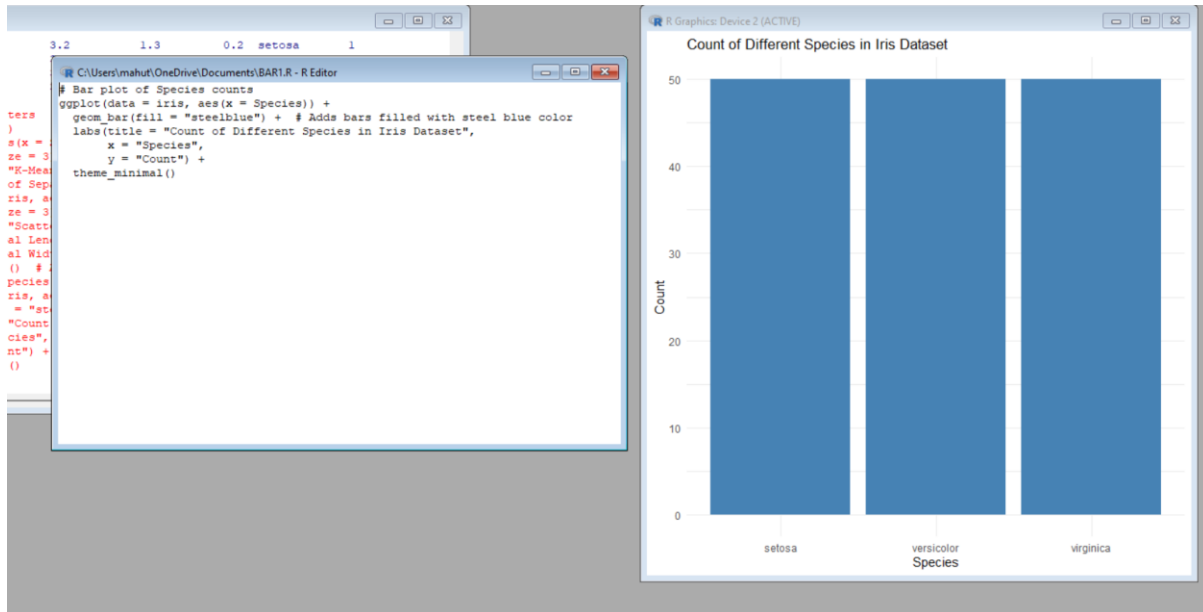
2) BAR CHART

```
# Install ggplot2 (if not already installed)
install.packages("ggplot2")

# Load the ggplot2 package
library(ggplot2)

# Bar plot of Species counts
ggplot(data = iris, aes(x = Species)) +
  geom_bar(fill = "steelblue") + # Adds bars filled with steel blue color
```

```
labs(title = "Count of Different Species in Iris Dataset",
      x = "Species",
      y = "Count") +
theme_minimal()
```

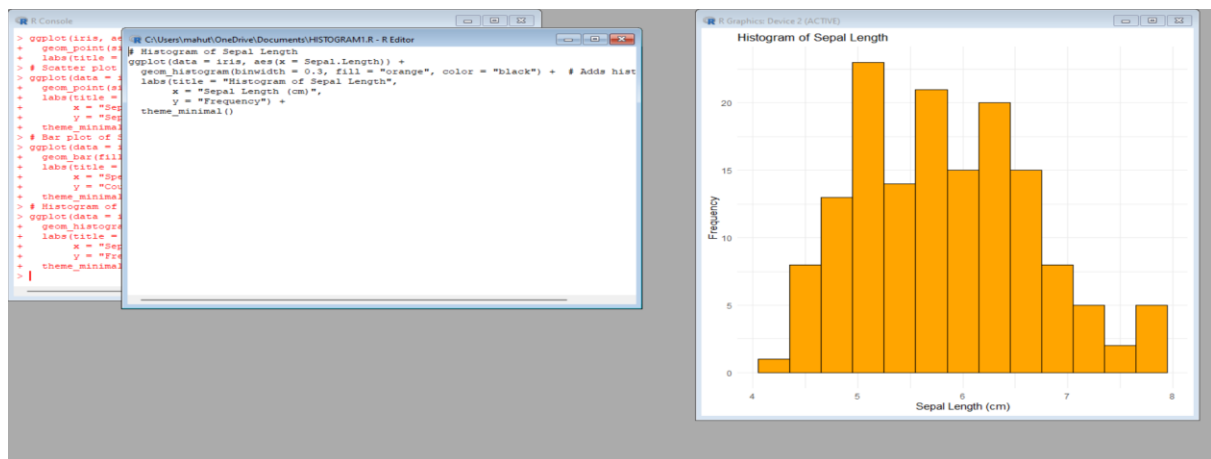


3) HISTOGRAM

```
# Install ggplot2 (if not already installed)
install.packages("ggplot2")
```

```
# Load the ggplot2 package
library(ggplot2)
```

```
# Histogram of Sepal Length
ggplot(data = iris, aes(x = Sepal.Length)) +
  geom_histogram(binwidth = 0.3, fill = "orange", color = "black") + # Adds
  histogram bars
  labs(title = "Histogram of Sepal Length",
        x = "Sepal Length (cm)",
        y = "Frequency") +
  theme_minimal()
```

4)BOX PLOT

Install ggplot2 (if not already installed)
install.packages("ggplot2")

Load the ggplot2 package
library(ggplot2)

Box plot of Sepal Length for each Species
ggplot(data = iris, aes(x = Species, y = Sepal.Length, fill = Species)) +
geom_boxplot() + # Adds box plot
labs(title = "Box Plot of Sepal Length by Species",
x = "Species",
y = "Sepal Length (cm)") +
theme_minimal()

