

Práctica 2: Limpieza y validación de los datos

Hernando Hernández Mariño

29 de mayo, 2021

1. Descripción del dataset. ¿Por qué es importante y qué pregunta/problema pretende responder?

```
# Lectura de datos
# Carpeta de trabajo
setwd("H:/Master_Ciencia_datos")
data <- read.csv("H:/Master_Ciencia_Datos/Práctica_2/iris.csv")
```

```
# Estructura del conjunto de datos
str(data)
```

```
## 'data.frame': 150 obs. of 6 variables:
## $ Id : int 1 2 3 4 5 6 7 8 9 10 ...
## $ SepalLengthCm: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
## $ SepalWidthCm : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
## $ PetalLengthCm: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
## $ PetalWidthCm : num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
## $ Species : chr "Iris-setosa" "Iris-setosa" "Iris-setosa" "Iris-setosa" ...
```

```
# Las primeras 5 filas
head(data,5)
```

##	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
## 1	1	5.1	3.5	1.4	0.2	Iris-setosa
## 2	2	4.9	3.0	1.4	0.2	Iris-setosa
## 3	3	4.7	3.2	1.3	0.2	Iris-setosa
## 4	4	4.6	3.1	1.5	0.2	Iris-setosa
## 5	5	5.0	3.6	1.4	0.2	Iris-setosa

```
# Las últimas 5 filas
tail(data ,5)
```

##	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
## 146	146	6.7	3.0	5.2	2.3	Iris-virginica
## 147	147	6.3	2.5	5.0	1.9	Iris-virginica
## 148	148	6.5	3.0	5.2	2.0	Iris-virginica
## 149	149	6.2	3.4	5.4	2.3	Iris-virginica
## 150	150	5.9	3.0	5.1	1.8	Iris-virginica

```
# Añadir las variables del data frame al entorno global de R.
attach(data)
```

El conjunto de datos objeto de análisis se ha obtenido a partir del enlace en Kaggle, el cual contiene la longitud y la anchura de los pétalos y sépalos y la especie de 150 flores iris. De manera que es un conjunto de datos multivariante comprendido por 5 características (columnas) de 150 flores iris (filas o registros).

El famoso estadístico Sir Ronald. A. Fisher usó este conjunto de datos en su artículo «The Use of Multiple Measurements in Taxonomic Problems» (Annals of Eugenics 7 (1936), pp. 179–188). A veces se llama el conjunto de datos Iris de Anderson porque Edgar Anderson recopiló los datos para cuantificar la variación morfológica de las flores de Iris de tres especies relacionadas. Dos de las tres especies fueron recogidas en la Península de Gaspé “todas del mismo pasto, y recogidas el mismo día y medidas al mismo tiempo por la misma persona con el mismo aparato”.

El conjunto de datos consta de 50 muestras de cada una de las tres especies de Iris (Iris setosa, Iris virginica e Iris versicolor). Se midieron cuatro características de cada muestra: la longitud y la anchura de los sépalos y pétalos, en centímetros. Basándose en la combinación de estas cuatro características, Fisher desarrolló un modelo discriminador lineal para distinguir la especie entre sí.

La idea es realizar con este conjunto de datos un análisis exploratorio o descriptivo que permita resumir, representar y explicar los datos concretos a disposición. Igualmente se pretenden plantear un modelo estadístico que logre predecir o clasificar las tres especies a partir de los 4 atributos enunciados anteriormente, lo cual se convierta en un caso de prueba y aprendizaje para las técnicas de clasificación estadística en el aprendizaje automático.

2. Integración y selección de los datos de interés a analizar.

No se realizaron procesos de integración o fusión de datos tales como añadir nuevos atributos o registros a la base original, pues no se considera necesario, por ahora, al logro de los objetivos planteados.

En cuanto a la selección de los datos se consideran todos los atributos a excepción del primer campo Id, dado que no es un atributo que mida algún tipo de característica relevante que aporte al ejercicio analítico.

```
# eliminar columna Id
data$Id <- NULL
```

Ahora bien al revisar el tipo de atributo o variable del dataset importado se observa que todos los atributos son numéricos a excepción del atributo Species, el cual se ha importado como un vector de palabras: lo indica el chr, de character, en la fila correspondiente del resultado de str. Esta variable Species es de tipo categórico y tiene asociada una descripción, una cadena de caracteres y, al mismo tiempo cuenta con un limitado número de valores posibles. Almacenar estos datos directamente como cadenas de caracteres implica un uso de memoria

innecesario, ya que cada una de las apariciones en la base de datos puede asociarse con un índice numérico sobre el conjunto total de valores posibles, obteniendo una representación mucho más compacta. Para tal fin, el atributo Species se crea como factor, de tal manera que este -el factor- se almacena internamente como un número y las etiquetas asociadas a cada valor se denominan niveles, que en este caso serán tres.

```
# Transformar Species en Factor
data <- data.frame(SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm,
                  Species, stringsAsFactors=TRUE)
```

```
# Verificación
str(data)
```

```
## 'data.frame':   150 obs. of  5 variables:
## $ SepalLengthCm: num  5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
## $ SepalWidthCm : num  3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
## $ PetalLengthCm: num  1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
## $ PetalWidthCm : num  0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
## $ Species      : Factor w/ 3 levels "Iris-setosa",...: 1 1 1 1 1 1 1 1 1 1 ...
```

3. Limpieza de los datos.

3.1. ¿Los datos contienen ceros o elementos vacíos? ¿Cómo gestionarías cada uno de estos casos?

No se identifican valores perdidos o ausentes (NA) en el dataset. En la gestión de datos ausentes se puede optar por eliminarlos o sustituirlos por el valor promedio de la columna o el valor más frecuente e incluso pueden ser reemplazados a partir de un modelo de regresión que predice dicho valor vacío; el camino a seguir dependerá de los datos a disposición y de los objetivos del análisis a realizar.

```
# valores ausentes
anyNA(data)
```

```
## [1] FALSE
```

3.2. Identificación y tratamiento de valores extremos.

Con el fin de identificar valores extremos se presentan diagramas de caja por cada una de las cuatro características y según la especie de flor. Pero antes se presentan estadísticos descriptivos de cada una de las cuatro características y según la especie de flor con el fin de notar diferencias entre las especies.

```
# Estadísticos descriptivos
summary(data)
```

```
## SepalLengthCm    SepalWidthCm    PetalLengthCm    PetalWidthCm
## Min.    :4.300    Min.    :2.000    Min.    :1.000    Min.    :0.100
## 1st Qu.:5.100    1st Qu.:2.800    1st Qu.:1.600    1st Qu.:0.300
## Median :5.800    Median :3.000    Median :4.350    Median :1.300
## Mean   :5.843    Mean   :3.054    Mean   :3.759    Mean   :1.199
## 3rd Qu.:6.400    3rd Qu.:3.300    3rd Qu.:5.100    3rd Qu.:1.800
## Max.   :7.900    Max.   :4.400    Max.   :6.900    Max.   :2.500
##           Species
## Iris-setosa    :50
## Iris-versicolor:50
## Iris-virginica :50
##
##
##
```

```
# Estadísticos descriptivos por especie
tapply(data$SepalLengthCm, data$Species, summary)
```

```
## $`Iris-setosa`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  4.300  4.800   5.000   5.006   5.200   5.800
##
## $`Iris-versicolor`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  4.900  5.600   5.900   5.936   6.300   7.000
##
## $`Iris-virginica`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  4.900  6.225   6.500   6.588   6.900   7.900
```

```
tapply(data$SepalWidthCm, data$Species, summary)
```

```
## $`Iris-setosa`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  2.300  3.125   3.400   3.418   3.675   4.400
##
## $`Iris-versicolor`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  2.000  2.525   2.800   2.770   3.000   3.400
##
## $`Iris-virginica`
##      Min. 1st Qu.  Median      Mean 3rd Qu.      Max.
##  2.200  2.800   3.000   2.974   3.175   3.800
```

```
tapply(data$PetalLengthCm, data$Species, summary)
```

```
## $`Iris-setosa`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.000   1.400   1.500   1.464   1.575   1.900
##
## $`Iris-versicolor`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   3.000   4.000   4.350   4.260   4.600   5.100
##
## $`Iris-virginica`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   4.500   5.100   5.550   5.552   5.875   6.900
```

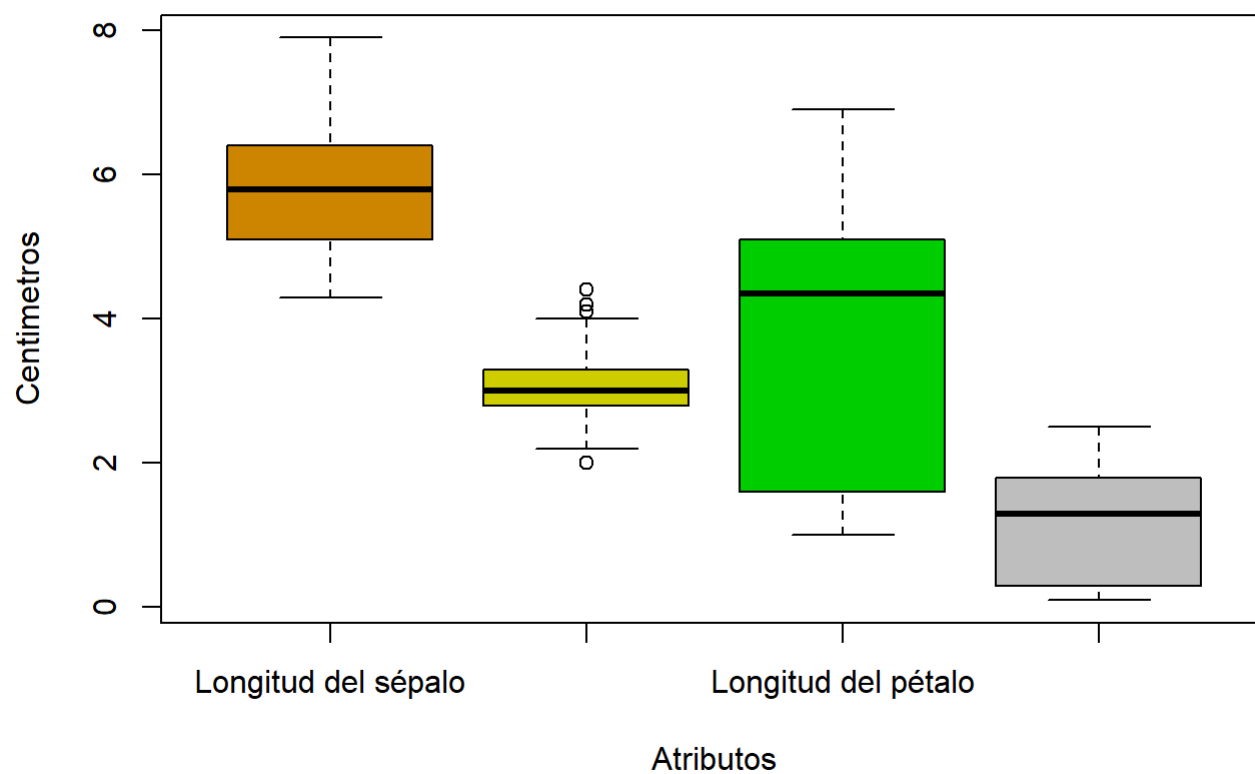
```
tapply(data$PetalWidthCm, data$Species, summary)
```

```
## $`Iris-setosa`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   0.100   0.200   0.200   0.244   0.300   0.600
##
## $`Iris-versicolor`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.000   1.200   1.300   1.326   1.500   1.800
##
## $`Iris-virginica`
##   Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##   1.400   1.800   2.000   2.026   2.300   2.500
```

Los estadísticos de tendencia central de la longitud y del ancho del sépalos entre las especies presentan diferencias marcadas; por ejemplo, la media y mediana de la longitud del sépalos de la especie virginica es mayor que las otras dos especies. En contraste, el ancho del sépalos -su media y mediana- es superior en la especie setosa. En cuanto a la longitud y ancho del sépalos, la especie virginica es mayor frente a las otras dos especies.

```
# Por cada variable en un solo gráfico
boxplot(SepalLengthCm, SepalWidthCm, PetalLengthCm, PetalWidthCm,
        names = c("Longitud del sépalos", "Ancho del sépalos", "Longitud del pétalo",
                  "Ancho del pétalo"), horizontal=FALSE, main="Diagramas de caja Iris",
        col = c("orange3", "yellow3", "green3", "grey"),
        xlab = "Atributos", ylab = "Centímetros")
```

Diagramas de caja Iris



En los diagramas de caja de los cuatro atributos se observan diferencias marcadas en su mediana tanto en las longitudes como en los anchos del sépalo y pétalo de las flores iris. También se identifican algunos valores extremos en el atributo Ancho del sépalo.

```
# Por cada variable según su especie

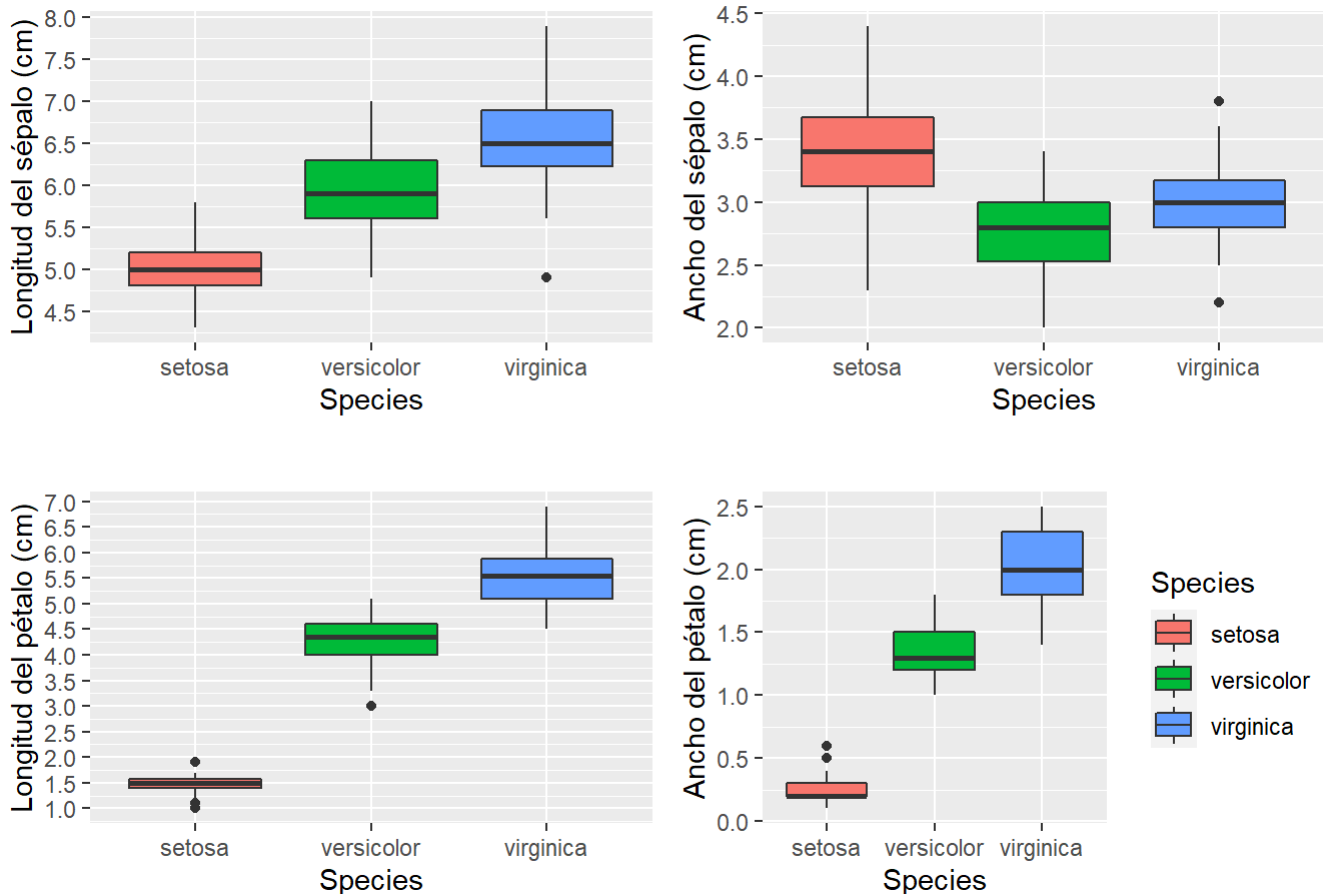
BpSl <- ggplot(iris, aes(Species, SepalLengthCm, fill=Species)) +
  geom_boxplot()+
  scale_y_continuous("Longitud del sépaló (cm)", breaks= seq(0,30, by=.5))+
  theme(legend.position="none")

BpSw <- ggplot(iris, aes(Species, SepalWidthCm, fill=Species)) +
  geom_boxplot()+
  scale_y_continuous("Ancho del sépaló (cm)", breaks= seq(0,30, by=.5))+
  theme(legend.position="none")

BpPl <- ggplot(iris, aes(Species, PetalLengthCm, fill=Species)) +
  geom_boxplot()+
  scale_y_continuous("Longitud del pétalo (cm)", breaks= seq(0,30, by=.5))+
  theme(legend.position="none")

BpPw <- ggplot(iris, aes(Species, PetalWidthCm, fill=Species)) +
  geom_boxplot()+
  scale_y_continuous("Ancho del pétalo (cm)", breaks= seq(0,30, by=.5))+
  labs(title = "Iris Box Plot", x = "Species")

grid.arrange(BpSl + ggtitle(""),
              BpSw + ggtitle(""),
              BpPl + ggtitle(""),
              BpPw + ggtitle(""),
              nrow = 2)
```



Igualmente al realizar los diagramas de caja de los atributos de acuerdo con cada especie de iris se observan diferencias relevantes en las medianas de la longitud y ancho del pétalo, así como la longitud del sépalo; pero en el ancho del sépalo, si bien se presentan diferencias en sus medianas estas son menos marcadas. Por otra parte, también se identifican valores extremos en las características de longitud y ancho del pétalo de la especie setosa y de la especie virginica en las características de longitud y ancho del sépalo. Por ahora, se mantendrán todos los valores extremos identificados en el ejercicio analítico de este dataset.

Una vez realizado sobre el conjunto de datos inicial los procedimientos de integración, validación y limpieza anteriores, procedemos a guardar estos en un nuevo fichero denominado `Automobile_data_clean.csv`:

```
# Exportación de los datos preprocesados
write.csv(data, "data_clean.csv")
```

4. Análisis de los datos.

4.1. Selección de los grupos de datos que se quieren analizar/comparar (planificación de los análisis a aplicar).

Se divide el dataset iris en varios datasets, los cuales contienen cada uno las muestras pertenecientes a una especie de flor que sería interesante analizar y/o comparar; sin embargo, no todos se utilizarían en la realización de pruebas estadísticas posteriores.


```
# Separar en grupos según un factor
setosa <- data[1:50, 1:4]
versicolor <- data[51:100, 1:4]
virginica <- data[101:150, 1:4]
```

```
# verificación
str(setosa)
```

```
## 'data.frame': 50 obs. of 4 variables:
## $ SepalLengthCm: num 5.1 4.9 4.7 4.6 5 5.4 4.6 5 4.4 4.9 ...
## $ SepalWidthCm : num 3.5 3 3.2 3.1 3.6 3.9 3.4 3.4 2.9 3.1 ...
## $ PetalLengthCm: num 1.4 1.4 1.3 1.5 1.4 1.7 1.4 1.5 1.4 1.5 ...
## $ PetalWidthCm : num 0.2 0.2 0.2 0.2 0.2 0.4 0.3 0.2 0.2 0.1 ...
```

```
str(versicolor)
```

```
## 'data.frame': 50 obs. of 4 variables:
## $ SepalLengthCm: num 7 6.4 6.9 5.5 6.5 5.7 6.3 4.9 6.6 5.2 ...
## $ SepalWidthCm : num 3.2 3.2 3.1 2.3 2.8 2.8 3.3 2.4 2.9 2.7 ...
## $ PetalLengthCm: num 4.7 4.5 4.9 4 4.6 4.5 4.7 3.3 4.6 3.9 ...
## $ PetalWidthCm : num 1.4 1.5 1.5 1.3 1.5 1.3 1.6 1 1.3 1.4 ...
```

```
str(virginica)
```

```
## 'data.frame': 50 obs. of 4 variables:
## $ SepalLengthCm: num 6.3 5.8 7.1 6.3 6.5 7.6 4.9 7.3 6.7 7.2 ...
## $ SepalWidthCm : num 3.3 2.7 3 2.9 3 3 2.5 2.9 2.5 3.6 ...
## $ PetalLengthCm: num 6 5.1 5.9 5.6 5.8 6.6 4.5 6.3 5.8 6.1 ...
## $ PetalWidthCm : num 2.5 1.9 2.1 1.8 2.2 2.1 1.7 1.8 1.8 2.5 ...
```

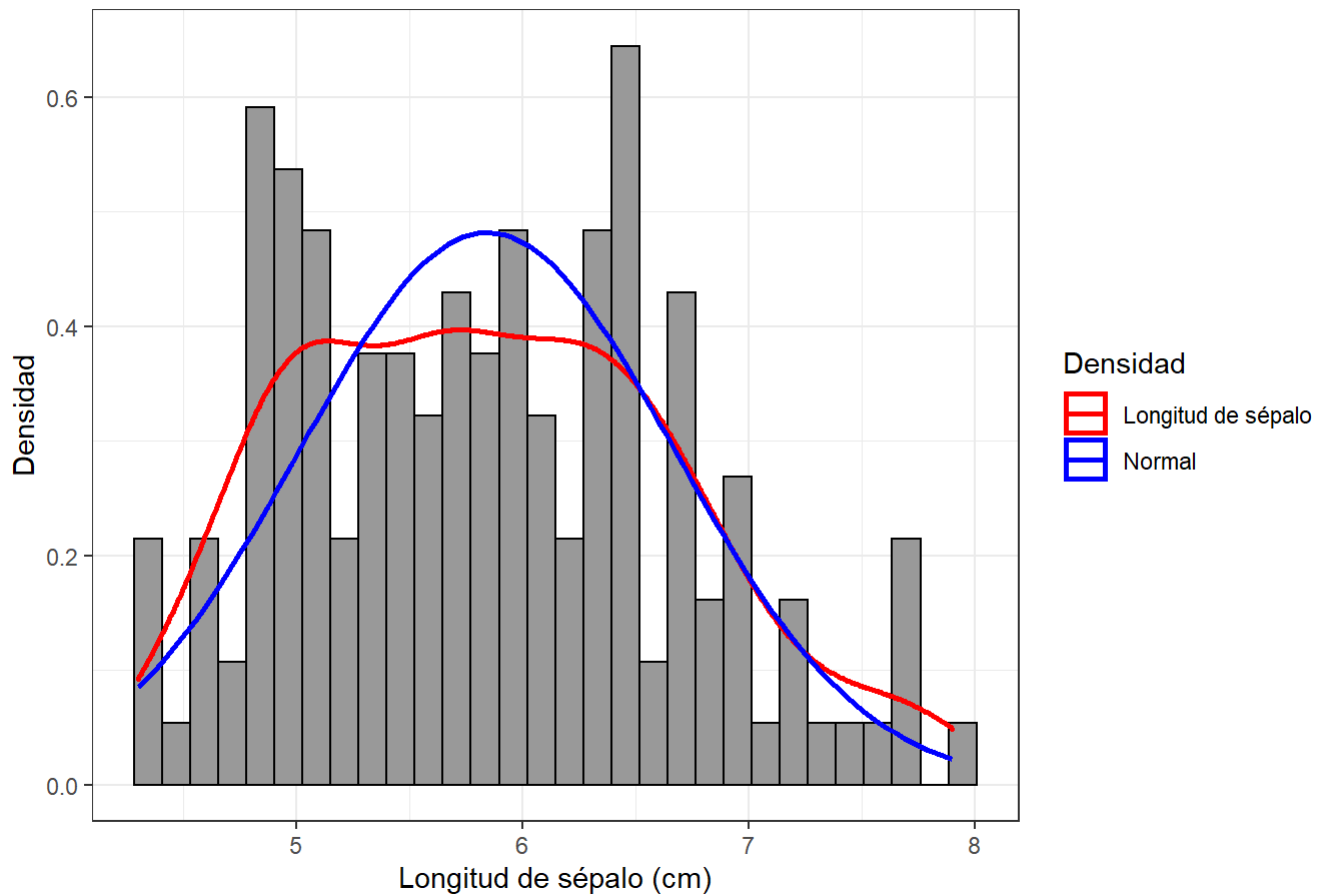
4.2. Comprobación de la normalidad y homogeneidad de la varianza.

Con el fin de comprobar la normalidad de cada uno de los atributos y según su especie de flor se presenta el histograma y curva de densidad, gráfico de cuantiles teóricos (Q-Q plot), así como los test de normalidad Anderson-Darling cuyo nivel de significación se fija en 0.05.

```
# Histograma y curva de densidad
# Longitud del sépalo
ggplot(data = data, aes(x = SepalLengthCm)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "gray60") +
  geom_density(aes(color = "Longitud de sépalo"), lwd = 0.95) +
  stat_function(aes(color = "Normal"), fun = dnorm, lwd = 0.95,
    args = list(mean = mean(data$SepalLengthCm),
      sd = sd(data$SepalLengthCm))) +
  scale_colour_manual("Densidad", values = c("red", "blue")) +
  labs(x = "Longitud de sépalo (cm)", y = "Densidad",
    title = "Distribución de la longitud de sépalo vs curva normal") +
  theme_bw()
```

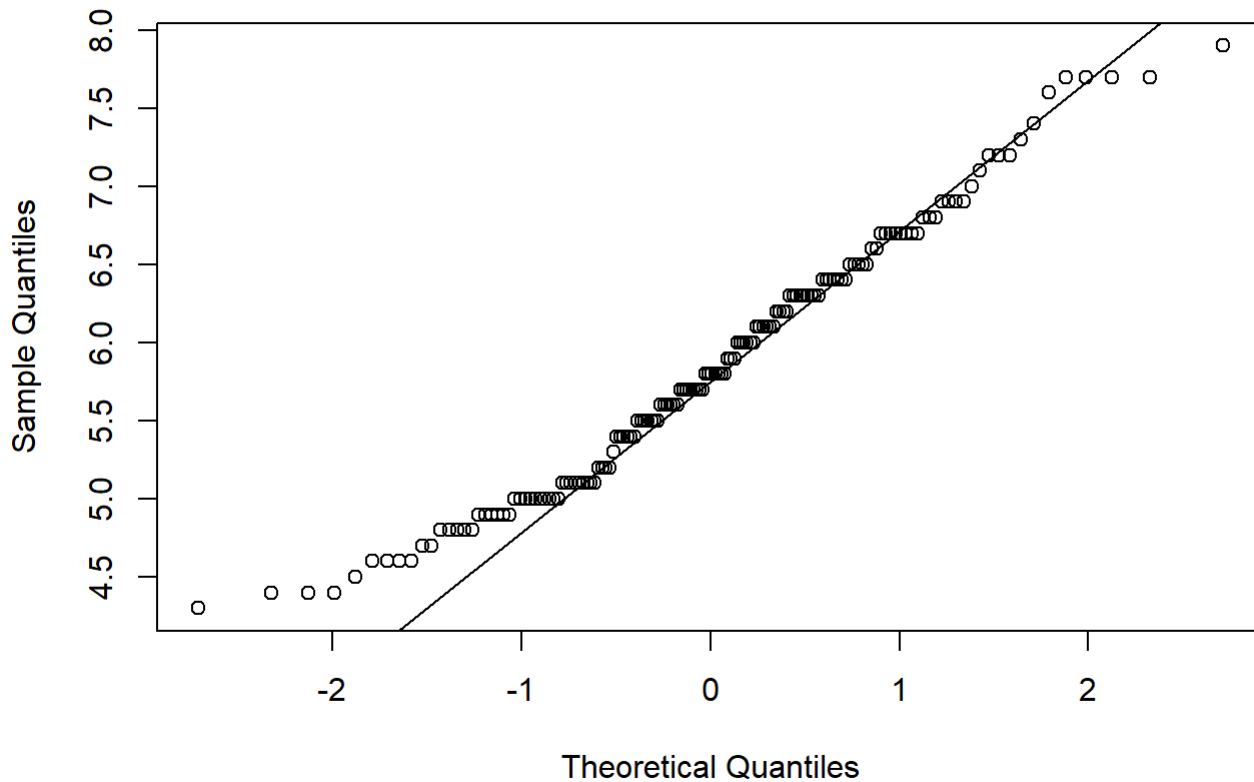
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Distribución de la longitud de sépalo vs curva normal



```
# Longitud del sépalo
# gráfico de cuantiles teóricos (Q-Q plot)
qqnorm(y = data$SepalLengthCm)
qqline(y = data$SepalLengthCm)
```

Normal Q-Q Plot



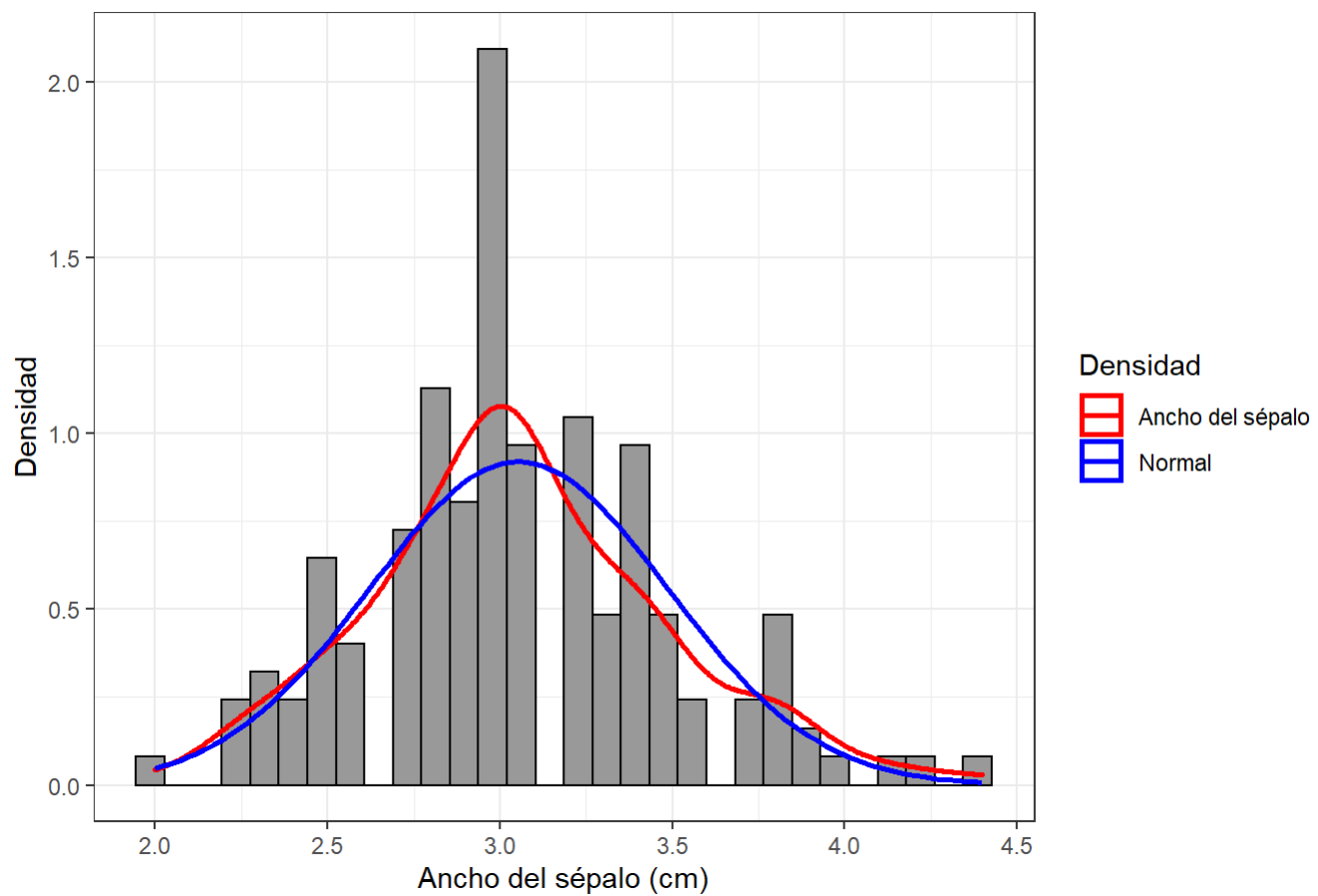
```
# Longitud del sépalo
# Prueba de normalidad
# Anderson-Darling
ad.test(data$SepalLengthCm)
```

```
##
## Anderson-Darling normality test
##
## data: data$SepalLengthCm
## A = 0.8892, p-value = 0.02251
```

```
# Histograma y curva de densidad
# Ancho del sépalo
ggplot(data = data, aes(x = SepalWidthCm)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "gray60") +
  geom_density(aes(color = "Ancho del sépalo"), lwd = 0.95) +
  stat_function(aes(color = "Normal"), fun = dnorm, lwd = 0.95,
    args = list(mean = mean(data$SepalWidthCm),
      sd = sd(data$SepalWidthCm))) +
  scale_colour_manual("Densidad", values = c("red", "blue")) +
  labs(x = "Ancho del sépal (cm)", y = "Densidad",
    title = "Distribución del ancho del sépal vs curva normal") +
  theme_bw()
```

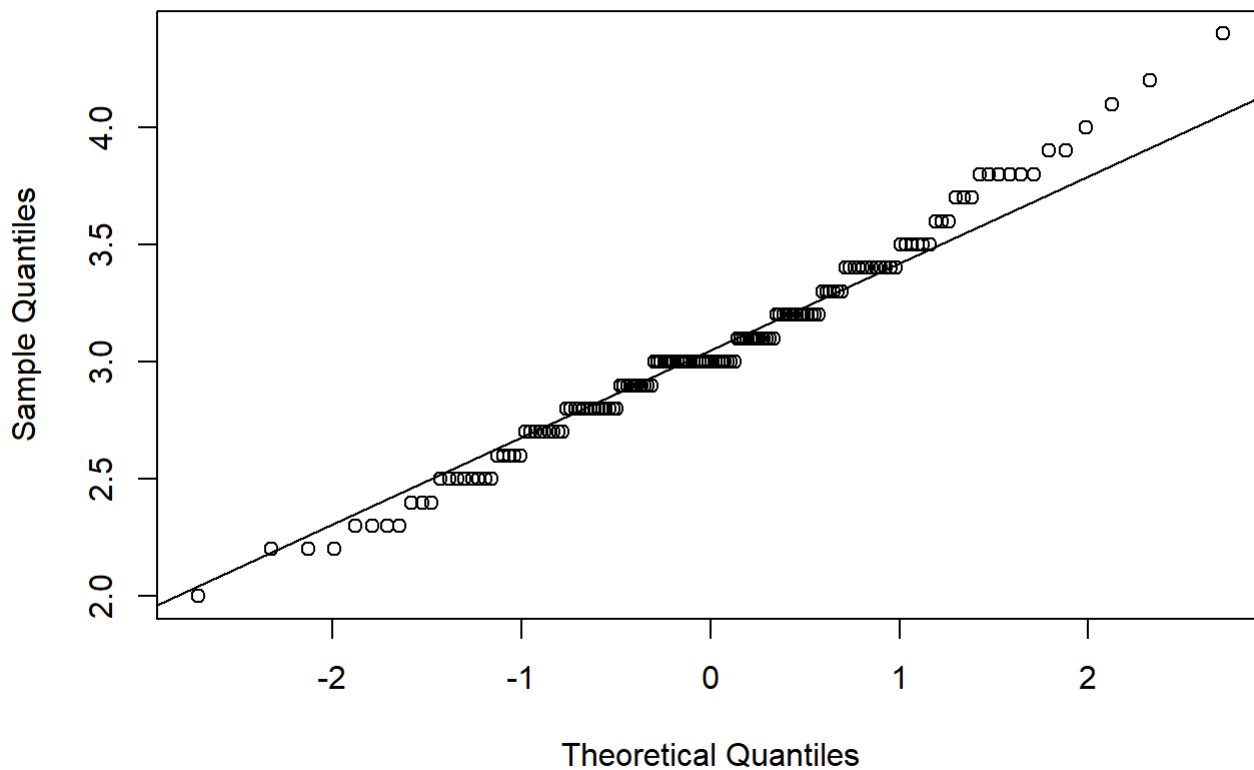
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Distribución del ancho del sépalo vs curva normal



```
# Ancho del sépalo  
# gráfico de cuantiles teóricos (Q-Q plot)  
qqnorm(y = data$SepalWidthCm)  
qqline(y = data$SepalWidthCm)
```

Normal Q-Q Plot



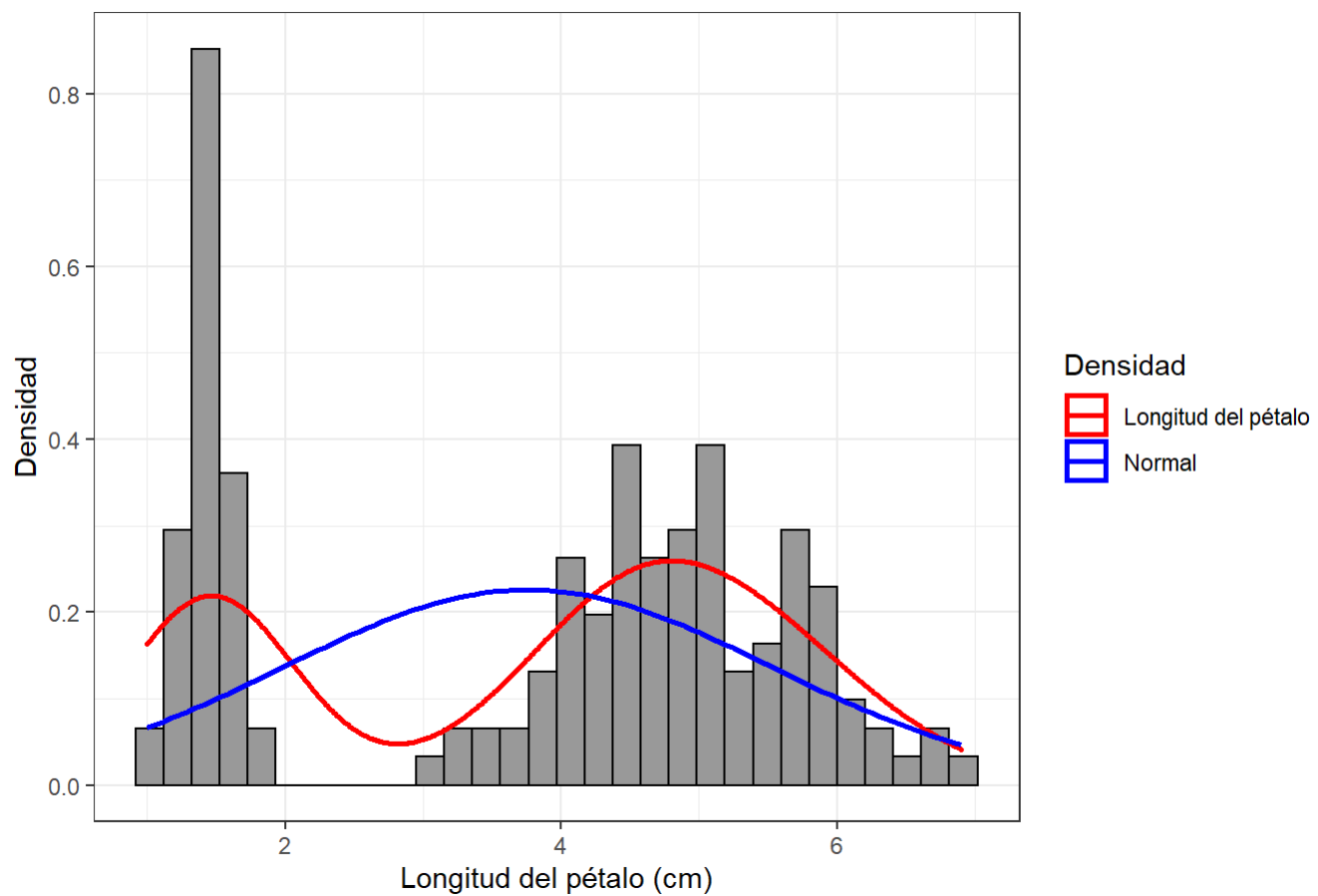
```
# Ancho del sépalo
# Prueba de normalidad
# Anderson-Darling
ad.test(data$SepalWidthCm)
```

```
##
## Anderson-Darling normality test
##
## data: data$SepalWidthCm
## A = 0.96566, p-value = 0.01455
```

```
# Longitud del pétalo
# Histograma y curva de densidad
ggplot(data = data, aes(x = PetalLengthCm)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "gray60") +
  geom_density(aes(color = "Longitud del pétalo"), lwd = 0.95) +
  stat_function(aes(color = "Normal"), fun = dnorm, lwd = 0.95,
    args = list(mean = mean(data$PetalLengthCm),
      sd = sd(data$PetalLengthCm))) +
  scale_colour_manual("Densidad", values = c("red", "blue")) +
  labs(x = "Longitud del pétalo (cm)", y = "Densidad",
    title = "Distribución de la longitud del pétalo vs curva normal") +
  theme_bw()
```

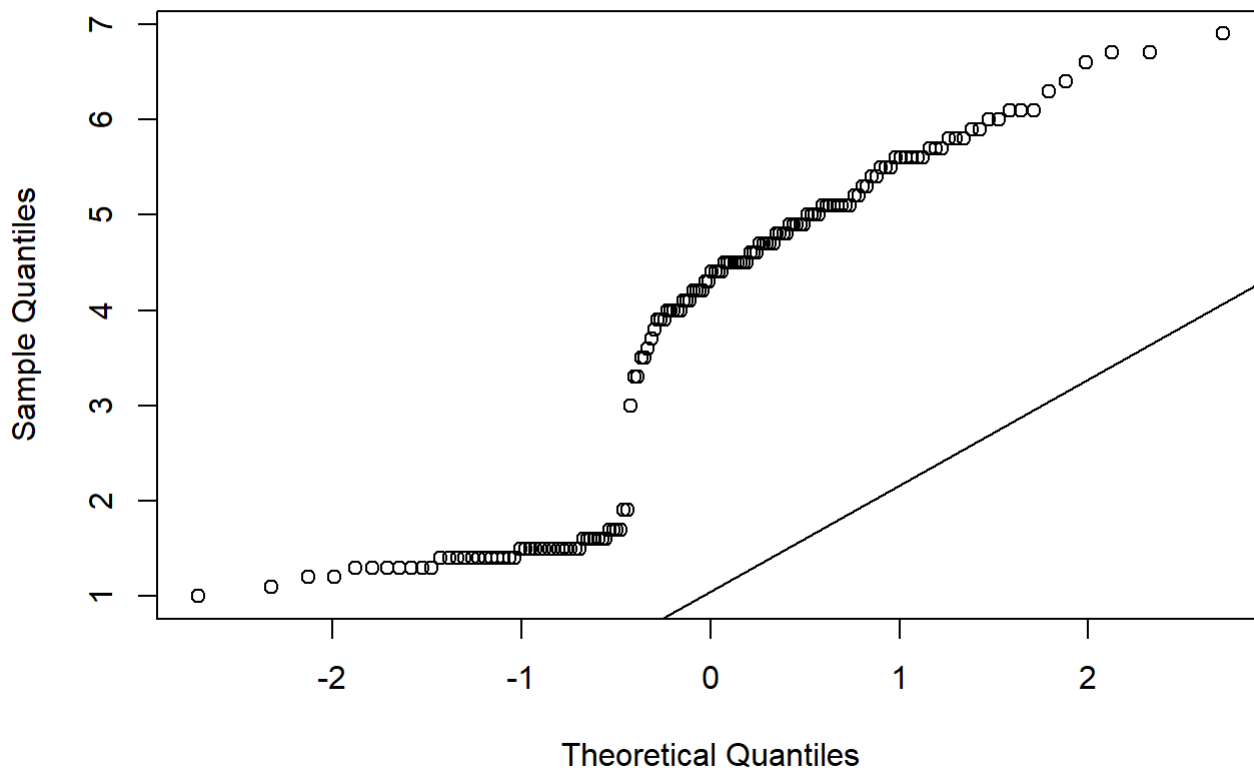
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Distribución de la longitud del pétalo vs curva normal



```
# Longitud del pétalo  
# gráfico de cuantiles teóricos (Q-Q plot)  
qqnorm(y = data$PetalLengthCm)  
qqline(y = data$PetalWidthCm)
```

Normal Q-Q Plot



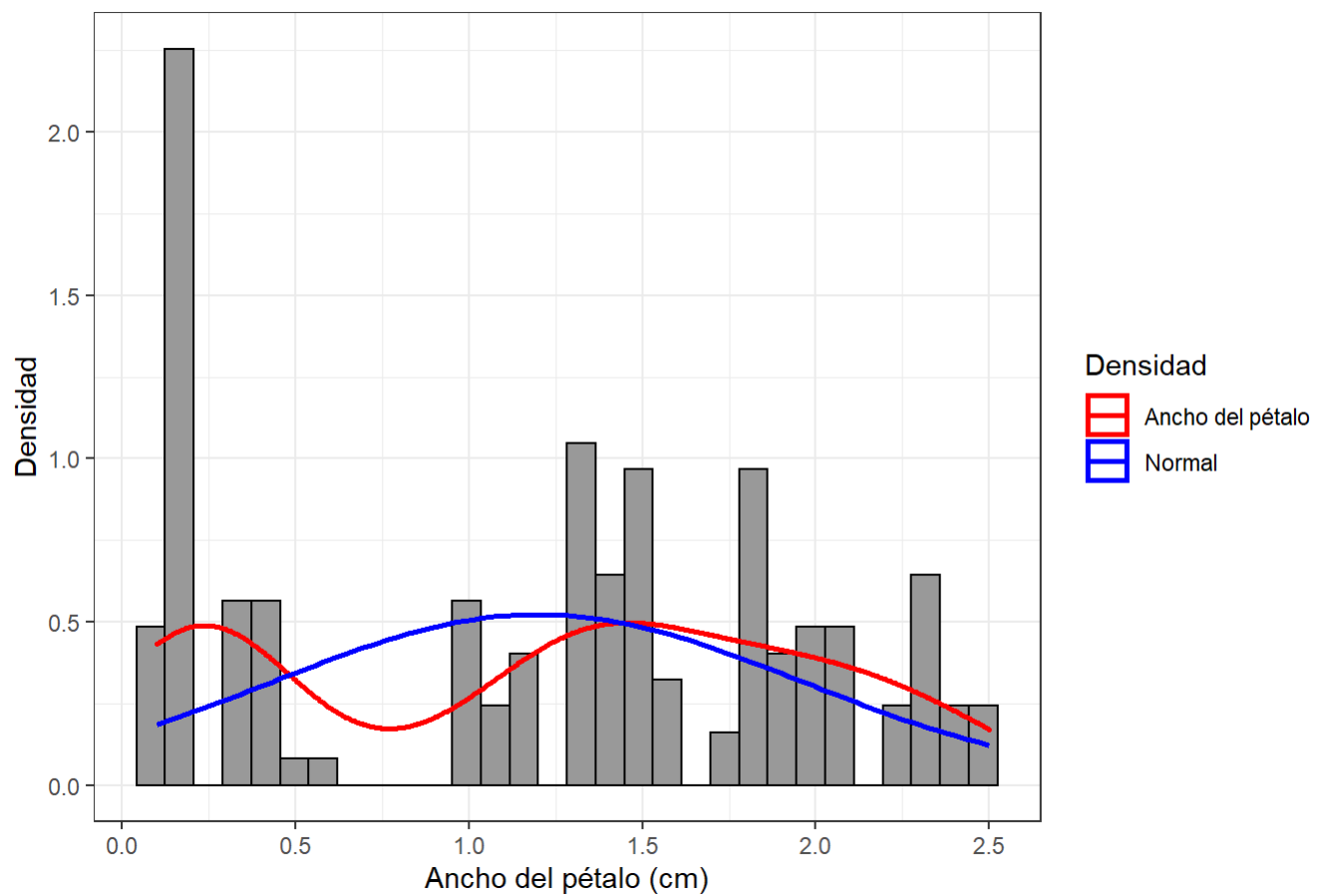
```
# Longitud del pétalo
# Prueba de normalidad
# Anderson-Darling
ad.test(data$PetalLengthCm)
```

```
##
## Anderson-Darling normality test
##
## data: data$PetalLengthCm
## A = 7.6729, p-value < 2.2e-16
```

```
# Ancho del pétalo
# Histograma y curva de densidad
ggplot(data = data, aes(x = PetalWidthCm)) +
  geom_histogram(aes(y = ..density..), color = "black", fill = "gray60") +
  geom_density(aes(color = "Ancho del pétalo"), lwd = 0.95) +
  stat_function(aes(color = "Normal", fun = dnorm, lwd = 0.95,
    args = list(mean = mean(data$PetalWidthCm),
      sd = sd(data$PetalWidthCm)))) +
  scale_colour_manual("Densidad", values = c("red", "blue")) +
  labs(x = "Ancho del pétalo (cm)", y = "Densidad",
    title = "Distribución del ancho del pétalo vs curva normal") +
  theme_bw()
```

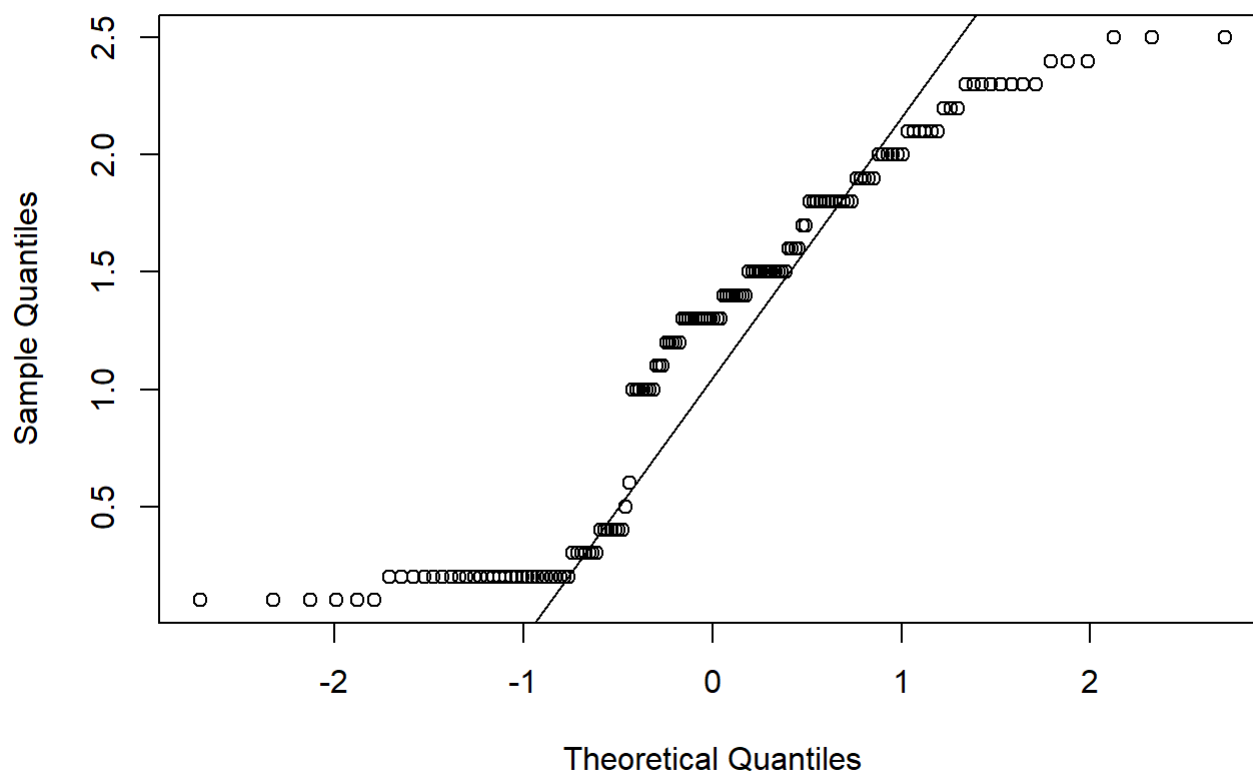
```
## `stat_bin()` using `bins = 30`. Pick better value with `binwidth`.
```

Distribución del ancho del pétalo vs curva normal



```
# Ancho del pétalo  
# gráfico de cuantiles teóricos (Q-Q plot)  
qqnorm(y = data$PetalWidthCm)  
qqline(y = data$PetalWidthCm)
```


Normal Q-Q Plot



```
# Ancho del pétalo
# Prueba de normalidad
# Anderson-Darling
ad.test(data$PetalWidthCm)
```

```
##
## Anderson-Darling normality test
##
## data: data$PetalWidthCm
## A = 5.0628, p-value = 1.427e-12
```

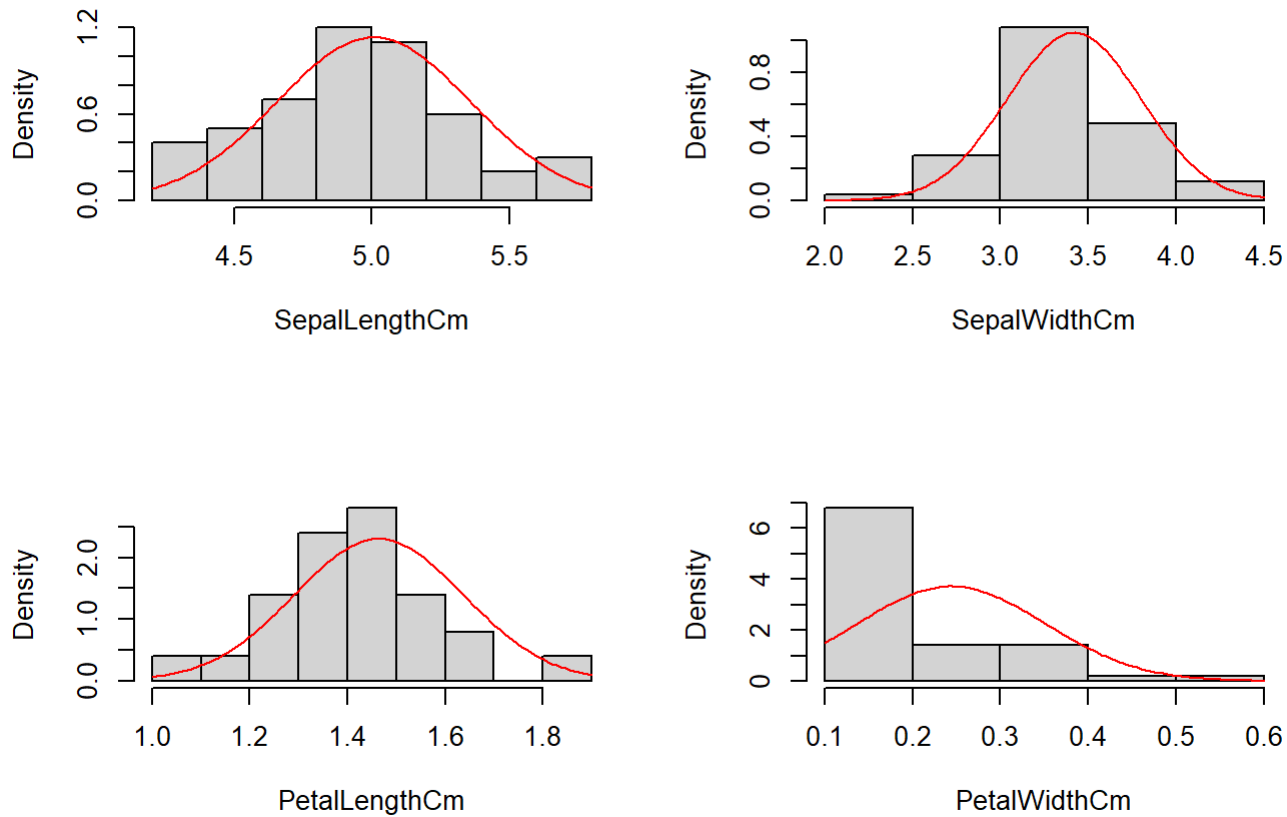
En el caso de los atributos longitud y ancho del sépal, los gráficos indican que su distribución se aleja de la distribución normal, lo cual se verifica en el test aplicado, pues el valor P (0.02251 y 0.01455 respectivamente) es menor que el nivel de significancia (0.05), por tanto, existe evidencia para rechazar hipótesis nula, es decir, que los datos no provienen de una población con distribución normal.

En el caso de los atributos longitud y ancho del pétalo, los gráficos indican que su distribución se aleja de la distribución normal, lo cual se verifica en el test aplicado, pues el valor P (2.2e-16 y 1.427e-12 respectivamente) es menor que el nivel de significancia (0.05), por tanto, existe evidencia para rechazar hipótesis nula, es decir, que los datos no provienen de una población con distribución normal.

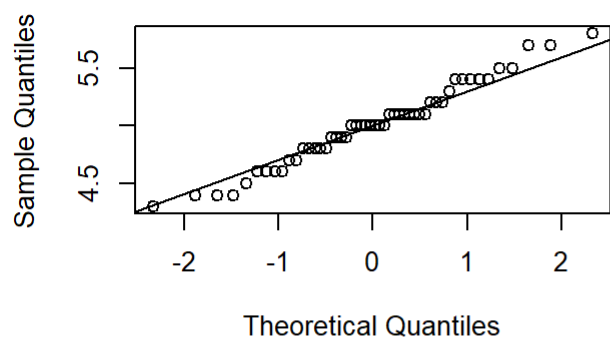
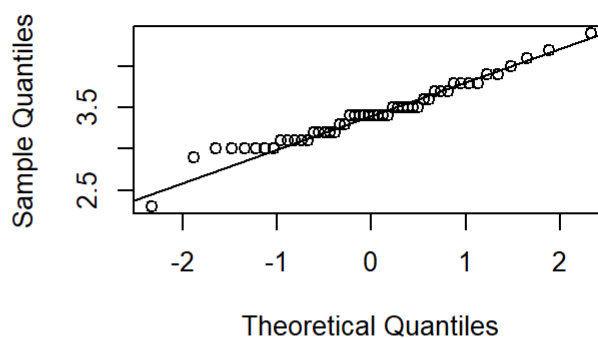
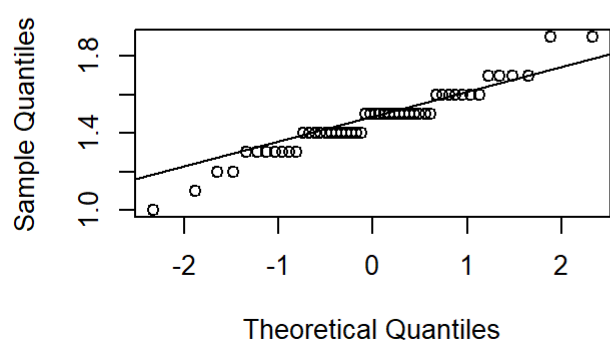
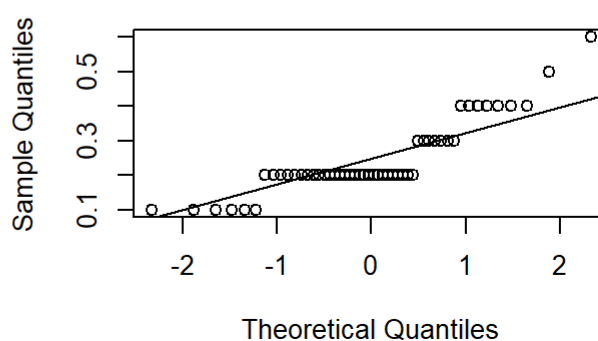
Estos resultados -en particular la distribución de los atributos- también sugieren la presencia de diferentes muestras, es decir, se evidencia la influencia de las tres especies de flores.

Ahora bien, al analizar los atributos de longitud y ancho del sépalo y pétalo según cada una de las tres especies, a partir del histograma y curva de densidad, así como el grafico Q-Q y el test de normalidad Anderson-Darling se observa lo siguiente:

```
# tipo de flor setosa
# Histograma y curva de densidad
result<- mvn(data=setosa, mvnTest="royston", univariatePlot="histogram")
```



```
# tipo de flor setosa
# Gráfico Q-Q
result<-mvn(data=setosa, mvnTest = "royston", univariatePlot = "qqplot")
```

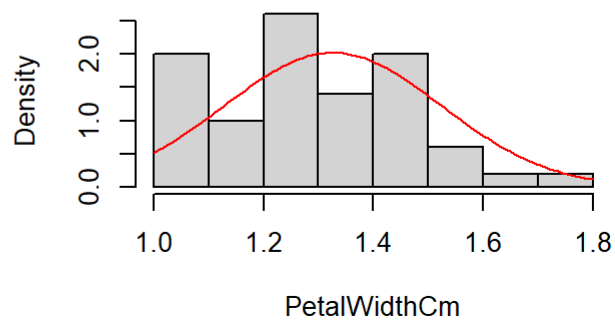
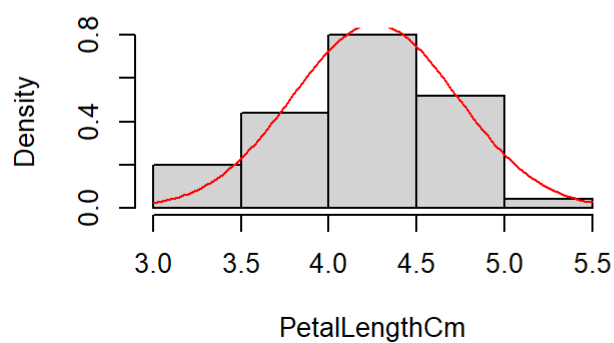
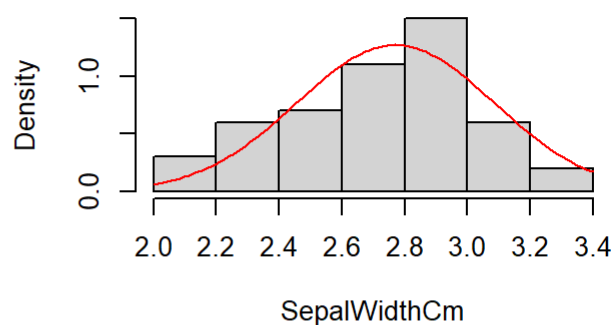
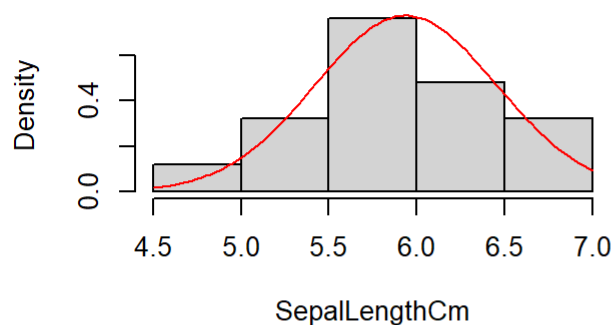
Normal Q-Q Plot (SepalLengthCm)**Normal Q-Q Plot (SepalWidthCm)****Normal Q-Q Plot (PetalLengthCm)****Normal Q-Q Plot (PetalWidthCm)**

```
# tipo de flor setosa  
# Test de Anderson Darling  
result <- mvn(data = setosa, mvnTest = "royston", univariateTest = "AD", desc = TRUE)  
result
```

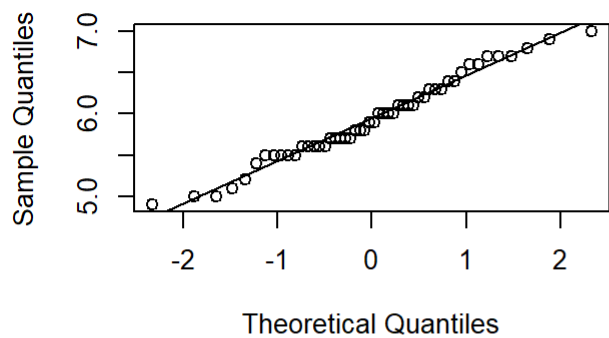
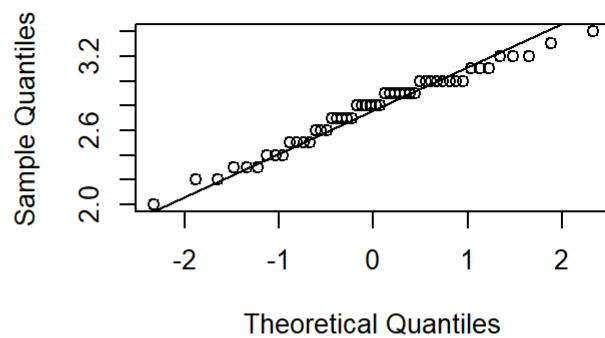
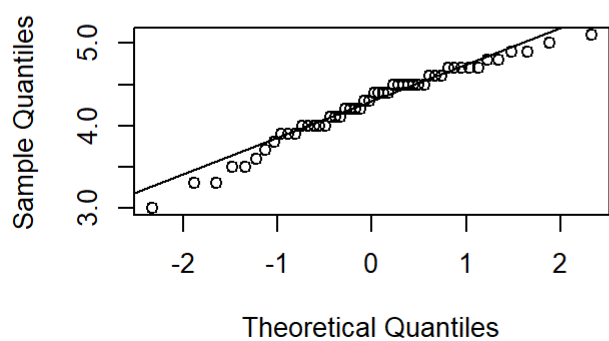
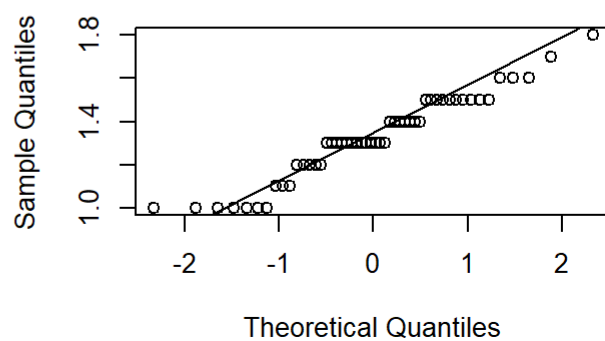
```
## $multivariateNormality
##      Test      H      p value MVN
## 1 Royston 30.37255 3.721164e-06 NO
##
## $univariateNormality
##      Test      Variable Statistic    p value Normality
## 1 Anderson-Darling SepalLengthCm    0.4080 0.3352      YES
## 2 Anderson-Darling SepalWidthCm     0.5635 0.1375      YES
## 3 Anderson-Darling PetalLengthCm    1.0111 0.0106      NO
## 4 Anderson-Darling PetalWidthCm     4.3070 <0.001      NO
##
## $Descriptives
##      n Mean   Std.Dev Median Min Max  25th  75th      Skew
## SepalLengthCm 50 5.006 0.3524897    5.0 4.3 5.8 4.800 5.200 0.11297784
## SepalWidthCm  50 3.418 0.3810244    3.4 2.3 4.4 3.125 3.675 0.10071528
## PetalLengthCm 50 1.464 0.1735112    1.5 1.0 1.9 1.400 1.575 0.06759284
## PetalWidthCm  50 0.244 0.1072095    0.2 0.1 0.6 0.200 0.300 1.12636619
##
##      Kurtosis
## SepalLengthCm -0.4508724
## SepalWidthCm  0.5392028
## PetalLengthCm 0.6626438
## PetalWidthCm  1.1263351
```

En el caso de la especie setosa el ancho del pétalo tiene una distribución sesgada a la derecha mientras que las otras variables tienen distribuciones aproximadamente normales. De acuerdo con el gráfico Q-Q se presentan algunas desviaciones de la línea recta y esto indica posibles desviaciones de una distribución normal, particularmente en el ancho del pétalo. Y según el test de normalidad, las variables longitud y ancho del pétalo no provienen de poblaciones normales, lo cual confirma lo señalado en los gráficos descritos anteriormente.

```
# tipo de flor versicolor
# Histograma y curva de densidad
result<- mvn(data=versicolor, mvnTest="royston", univariatePlot="histogram")
```



```
# tipo de flor versicolor  
# Gráfico Q-Q  
result<-mvn(data=versicolor, mvnTest = "royston", univariatePlot = "qqplot")
```

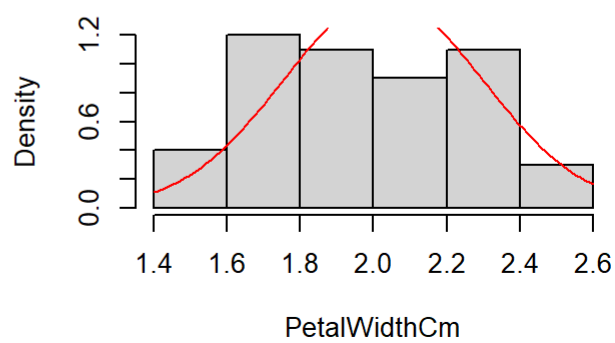
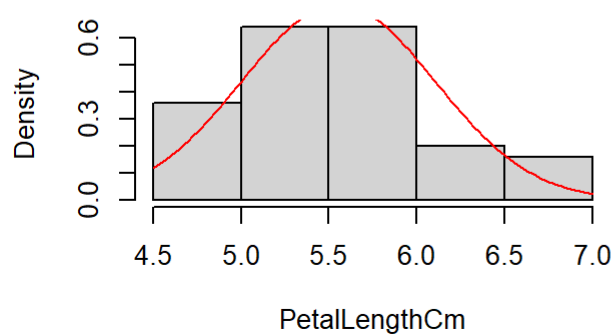
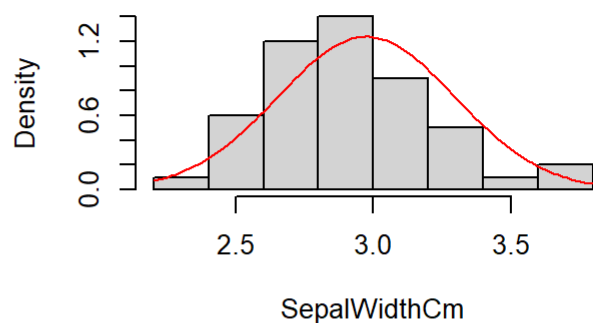
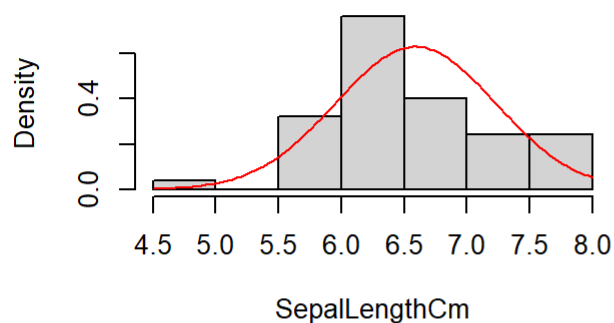
Normal Q-Q Plot (SepalLengthCm)**Normal Q-Q Plot (SepalWidthCm)****Normal Q-Q Plot (PetalLengthCm)****Normal Q-Q Plot (PetalWidthCm)**

```
# tipo de flor versicolor  
# Test de Anderson Darling  
result <- mvn(data = versicolor, mvnTest = "royston", univariateTest = "AD", desc = TRUE)  
result
```

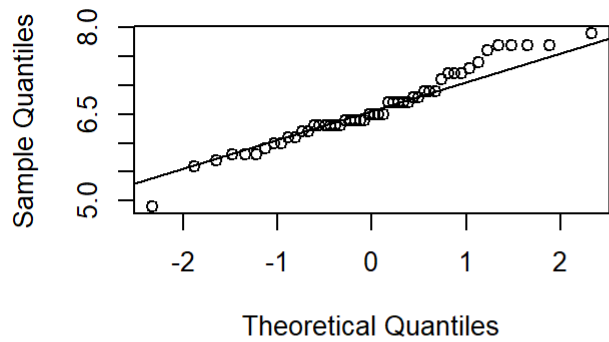
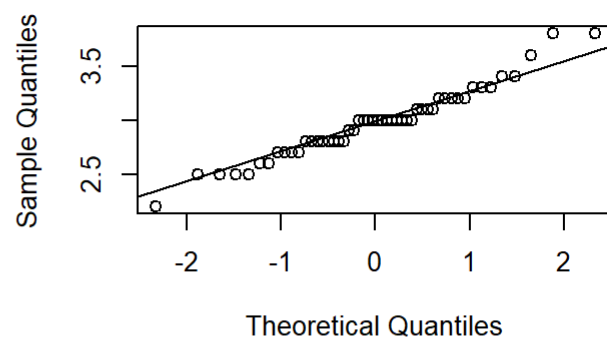
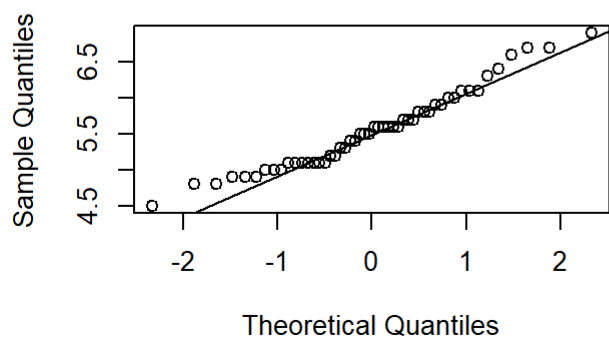
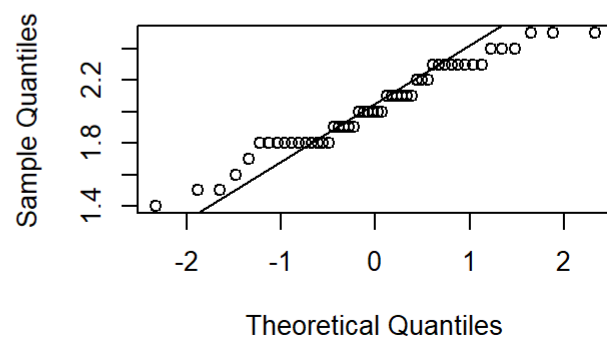
```
## $multivariateNormality
##      Test      H    p value MVN
## 1 Royston 7.85262 0.0847746 YES
##
## $univariateNormality
##      Test      Variable Statistic    p value Normality
## 1 Anderson-Darling SepalLengthCm    0.3608    0.4333    YES
## 2 Anderson-Darling SepalWidthCm     0.5598    0.1406    YES
## 3 Anderson-Darling PetalLengthCm    0.5551    0.1446    YES
## 4 Anderson-Darling PetalWidthCm    0.9569    0.0144    NO
##
## $Descriptives
##      n Mean   Std.Dev Median Min Max  25th 75th      Skew
## SepalLengthCm 50 5.936 0.5161711  5.90 4.9 7.0 5.600  6.3  0.09913926
## SepalWidthCm  50 2.770 0.3137983  2.80 2.0 3.4 2.525  3.0 -0.34136443
## PetalLengthCm 50 4.260 0.4699110  4.35 3.0 5.1 4.000  4.6 -0.57060243
## PetalWidthCm  50 1.326 0.1977527  1.30 1.0 1.8 1.200  1.5 -0.02933377
##
##      Kurtosis
## SepalLengthCm -0.6939138
## SepalWidthCm  -0.5493203
## PetalLengthCm -0.1902555
## PetalWidthCm  -0.5873144
```

En cuanto a la especie versicolor el ancho del pétalo tiene una distribución sesgada a la derecha mientras que las otras variables tienen distribuciones aproximadamente normales. De acuerdo con el gráfico Q-Q se presentan algunas desviaciones de la línea recta y esto indica posibles desviaciones de una distribución normal, particularmente en el ancho del pétalo. Y según el test de normalidad, la variable ancho del pétalo no proviene de poblaciones normales, lo cual confirma lo señalado en los gráficos descritos anteriormente.

```
# tipo de flor virginica
# Histograma y curva de densidad
result<- mvn(data=virginica, mvnTest="royston", univariatePlot="histogram")
```



```
# tipo de flor virginica  
# Gráfico Q-Q  
result<-mvn(data=virginica, mvnTest = "royston", univariatePlot = "qqplot")
```


Normal Q-Q Plot (SepalLengthCm)**Normal Q-Q Plot (SepalWidthCm)****Normal Q-Q Plot (PetalLengthCm)****Normal Q-Q Plot (PetalWidthCm)**

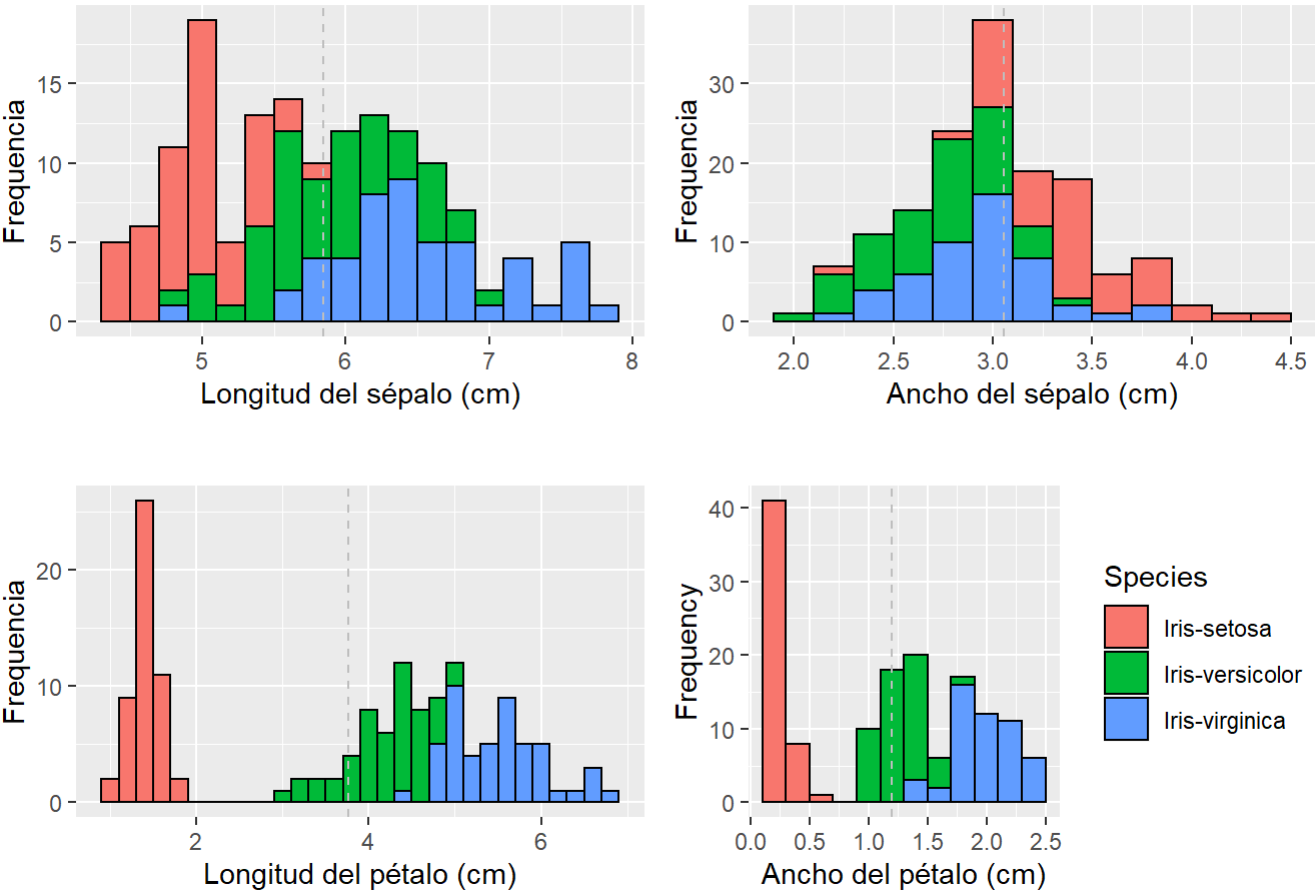
```
# tipo de flor virginica  
# Test de Anderson Darling  
result <- mvn(data = virginica, mvnTest = "royston", univariateTest = "AD", desc = TRUE)  
result
```

```
## $multivariateNormality
##      Test      H      p value MVN
## 1 Royston 8.141444 0.06776605 YES
##
## $univariateNormality
##      Test      Variable Statistic      p value Normality
## 1 Anderson-Darling SepalLengthCm      0.5516      0.1475      YES
## 2 Anderson-Darling SepalWidthCm      0.6182      0.1018      YES
## 3 Anderson-Darling PetalLengthCm      0.6090      0.1074      YES
## 4 Anderson-Darling PetalWidthCm      0.7388      0.0508      YES
##
## $Descriptives
##      n Mean      Std.Dev Median Min Max 25th 75th      Skew
## SepalLengthCm 50 6.588 0.6358796      6.50 4.9 7.9 6.225 6.900 0.1110286
## SepalWidthCm 50 2.974 0.3224966      3.00 2.2 3.8 2.800 3.175 0.3442849
## PetalLengthCm 50 5.552 0.5518947      5.55 4.5 6.9 5.100 5.875 0.5169175
## PetalWidthCm 50 2.026 0.2746501      2.00 1.4 2.5 1.800 2.300 -0.1218119
##
##      Kurtosis
## SepalLengthCm -0.2032597
## SepalWidthCm 0.3803832
## PetalLengthCm -0.3651161
## PetalWidthCm -0.7539586
```

Y finalmente, en la especie virginica, todas las variables tienen distribuciones aproximadamente normales. Si bien se presentan ciertas desviaciones de la línea recta en el gráfico Q-Q no son pronunciadas en las variables en estudio. Y el test de normalidad indica que sin excepción todas las variables provienen de poblaciones normales.

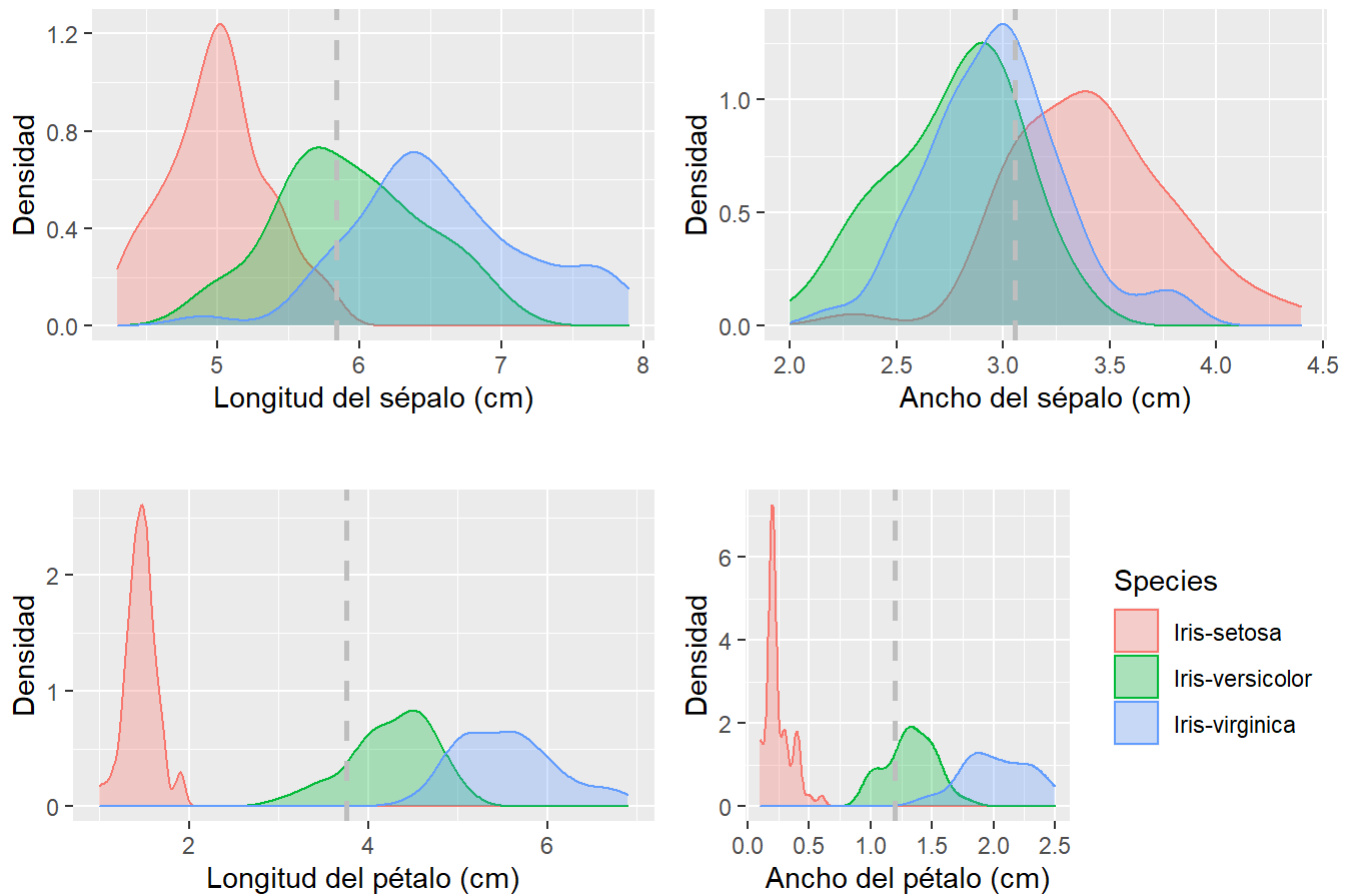
Las gráficas de histograma de frecuencias, curvas de densidad y dispersión -según las especies de flor iris- visualizadas en conjunto permiten identificar la posible separación de las especies y la superposición de valores de cada especie para un atributo en específico.

```
# Análisis de frecuencia con el histograma
# Longitud del sépalo
HisSl <- ggplot(data=data, aes(x=SepalLengthCm))+
  geom_histogram(binwidth=0.2, color="black", aes(fill=Species)) +
  xlab("Longitud del sépalo (cm)") +
  ylab("Frecuencia") +
  theme(legend.position="none")+
  ggtitle("Histograma de la longitud del sépalo")+
  geom_vline(data=data, aes(xintercept = mean(SepalLengthCm)),linetype="dashed",color="grey")
# Ancho del sépalo
HistSw <- ggplot(data=data, aes(x=SepalWidthCm)) +
  geom_histogram(binwidth=0.2, color="black", aes(fill=Species)) +
  xlab("Ancho del sépalo (cm)") +
  ylab("Frecuencia") +
  theme(legend.position="none")+
  ggtitle("Histograma del ancho del sépalo")+
  geom_vline(data=data, aes(xintercept = mean(SepalWidthCm)),linetype="dashed",color="grey")
# Longitud del pétalo
HistPl <- ggplot(data=data, aes(x=PetalLengthCm))+
  geom_histogram(binwidth=0.2, color="black", aes(fill=Species)) +
  xlab("Longitud del pétalo (cm)") +
  ylab("Frecuencia") +
  theme(legend.position="none")+
  ggtitle("Histograma de la longitud del pétalo")+
  geom_vline(data=data, aes(xintercept = mean(PetalLengthCm)),
    linetype="dashed",color="grey")
# Ancho del pétalo
HistPw <- ggplot(data=data, aes(x=PetalWidthCm))+
  geom_histogram(binwidth=0.2, color="black", aes(fill=Species)) +
  xlab("Ancho del pétalo (cm)") +
  ylab("Frequency") +
  theme(legend.position="right" )+
  ggtitle("Histograma del ancho del pétalo")+
  geom_vline(data=data, aes(xintercept = mean(PetalWidthCm)),linetype="dashed",color="grey")
# Visualización en conjunto
grid.arrange(HisSl + ggtitle(""),
  HistSw + ggtitle(""),
  HistPl + ggtitle(""),
  HistPw + ggtitle(""),
  nrow = 2)
```



```
# Análisis de densidad
# Longitud del pétalo
DhistPl <- ggplot(data, aes(x=PetalLengthCm, colour=Species, fill=Species)) +
  geom_density(alpha=.3) +
  geom_vline(aes(xintercept=mean(PetalLengthCm), colour=Species), linetype="dashed", color="grey"
, size=1)+
  xlab("Longitud del pétalo (cm)") +
  ylab("Densidad")+
  theme(legend.position="none")
# Ancho del pétalo
DhistPw <- ggplot(data, aes(x=PetalWidthCm, colour=Species, fill=Species)) +
  geom_density(alpha=.3) +
  geom_vline(aes(xintercept=mean(PetalWidthCm), colour=Species), linetype="dashed", color="grey",
size=1)+
  xlab("Ancho del pétalo (cm)") +
  ylab("Densidad")
# Ancho del sépalos
DhistSw <- ggplot(data, aes(x=SepalWidthCm, colour=Species, fill=Species)) +
  geom_density(alpha=.3) +
  geom_vline(aes(xintercept=mean(SepalWidthCm), colour=Species), linetype="dashed", color="grey"
, size=1)+
  xlab("Ancho del sépalos (cm)") +
  ylab("Densidad")+
  theme(legend.position="none")
# Longitud del sépalos
DhistSl <- ggplot(data, aes(x=SepalLengthCm, colour=Species, fill=Species)) +
  geom_density(alpha=.3) +
  geom_vline(aes(xintercept=mean(SepalLengthCm), colour=Species), linetype="dashed", color="gre
y", size=1)+
  xlab("Longitud del sépalos (cm)") +
  ylab("Densidad")+
  theme(legend.position="none")
# visualización conjunta

grid.arrange(DhistSl + ggtitle(""),
              DhistSw + ggtitle(""),
              DhistPl + ggtitle(""),
              DhistPw + ggtitle(""),
              nrow = 2)
```



Con relación a la homogeneidad de la varianza (la varianza es constante (no varía) en los diferentes niveles de un factor) se utiliza el test de Levene. Este test de Levene se caracteriza, porque en primer lugar se puede comparar 2 o más poblaciones (en este caso son tres muestras) y, en segundo lugar permite elegir entre diferentes estadísticos de centralidad: mediana (por defecto), media, media truncada, lo cual es importante a la hora de contrastar la homocedasticidad, dependiendo de si los grupos se distribuyen de forma normal o no, lo cual como se anotó anteriormente algunas variables no siguen una distribución normal.

```
# Homogeneidad
# Test de Levene
leveneTest(y = data$SepalLengthCm, group = data$Species, center = "median")
```

```
## Levene's Test for Homogeneity of Variance (center = "median")
##      Df F value  Pr(>F)
## group  2  6.3527 0.002259 **
##      147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
leveneTest(y = data$SepalWidthCm, group = data$Species, center = "median")
```

```
## Levene's Test for Homogeneity of Variance (center = "median")
##           Df F value Pr(>F)
## group    2   0.6475 0.5248
##           147
```

```
leveneTest(y = data$PetalLengthCm, group = data$Species, center = "median")
```

```
## Levene's Test for Homogeneity of Variance (center = "median")
##           Df F value    Pr(>F)
## group    2   19.72 2.589e-08 ***
##           147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
leveneTest(y = data$PetalWidthCm, group = data$Species, center = "median")
```

```
## Levene's Test for Homogeneity of Variance (center = "median")
##           Df F value    Pr(>F)
## group    2   19.412 3.302e-08 ***
##           147
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

De acuerdo con el test aplicado se encuentran diferencias entre los tres grupos de especies de iris en todos los atributos a excepción de la característica ancho del sépalo, de lo cual ya se tenía cierto indicio desde el punto de vista gráfico con los diagramas de caja desglasado por especie, pues se indicaba que si bien se presentaban diferencias en las medianas, estas no eran muy marcadas en comparación a los otros atributos según especie de iris.

4.3. Aplicación de pruebas estadísticas para comparar los grupos de datos. En función de los datos y el objetivo del estudio, aplicar pruebas de contraste de hipótesis, correlaciones, regresiones, etc. Aplicar al menos tres métodos de análisis diferentes.

```
# Anova
# Longitud del sépalo
AnovaSL <- aov(SepalLengthCm ~ Species, data = data)
summary(AnovaSL)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## Species      2   63.21   31.606   119.3 <2e-16 ***
## Residuals    147   38.96    0.265
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Anova
# Ancho del sépalo
AnovaSW <- aov(SepalWidthCm ~ Species, data = data)
summary(AnovaSW)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## Species      2   10.98    5.489   47.36 <2e-16 ***
## Residuals    147   17.04    0.116
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Anova
# Longitud del pétalo
AnovaPL <- aov(PetalLengthCm ~ Species, data = data)
summary(AnovaPL)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## Species      2  436.6   218.32  1179 <2e-16 ***
## Residuals    147   27.2    0.19
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
# Anova
# Ancho de pétalo
AnovaPW <- aov(PetalWidthCm ~ Species, data = data)
summary(AnovaPW)
```

```
##              Df Sum Sq Mean Sq F value Pr(>F)
## Species      2   80.60   40.30  959.3 <2e-16 ***
## Residuals    147    6.18    0.04
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

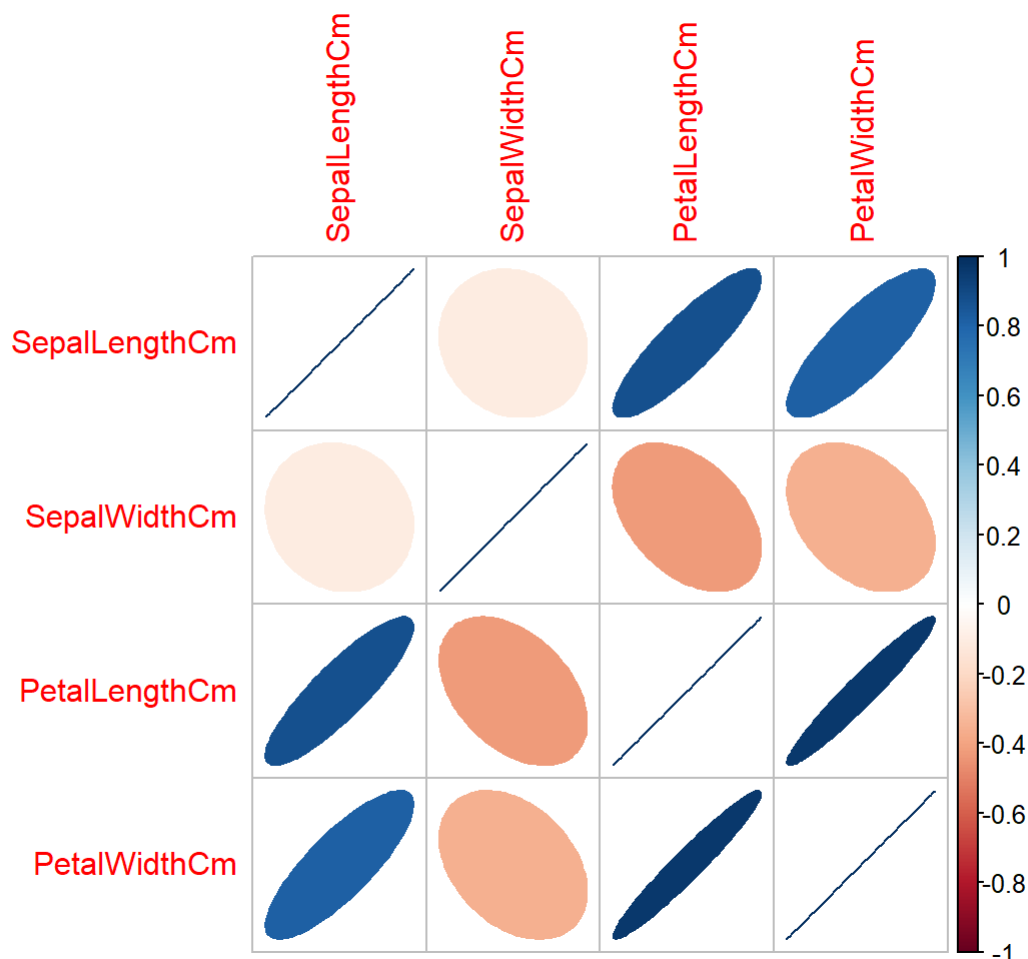
Se realiza un Anova con el fin de comparar las medias de cada uno de los atributos entre los grupos o especies de flores iris. Al establecer el valor alfa en 0.05 y al ver en la tabla que el valor de p es menor a alfa, se rechaza la hipótesis nula de que las medias son iguales, y se concluye que la media de la longitud y ancho del sépalo y pétalo es distinta entre las tres especies en todos los casos.

Al verificar la correlación entre los atributos de longitud y ancho del sépalo y pétalo se observa que entre la longitud del sépalo y la longitud y ancho del pétalo guardan una correlación positiva superior al 80%. Mientras que el ancho del sépalo guarda una correlación negativa con estas mismas variables, pero mucho menor (entre el

35% y 42%). Finalmente entre la longitud y ancho del pétalo su correlación es positiva y es del 96% (muy alta) y, entre la longitud y ancho del sépalo su correlación es negativa y muy baja (10%).

```
# Correlación
# entre todas sin especie
M <- cor(data[,1:4])
```

```
corrplot(M, method = "ellipse")
```



M

```
##          SepalLengthCm SepalWidthCm PetalLengthCm PetalWidthCm
## SepalLengthCm      1.0000000 -0.1093692    0.8717542    0.8179536
## SepalWidthCm      -0.1093692  1.0000000   -0.4205161   -0.3565441
## PetalLengthCm      0.8717542 -0.4205161    1.0000000    0.9627571
## PetalWidthCm      0.8179536 -0.3565441    0.9627571    1.0000000
```

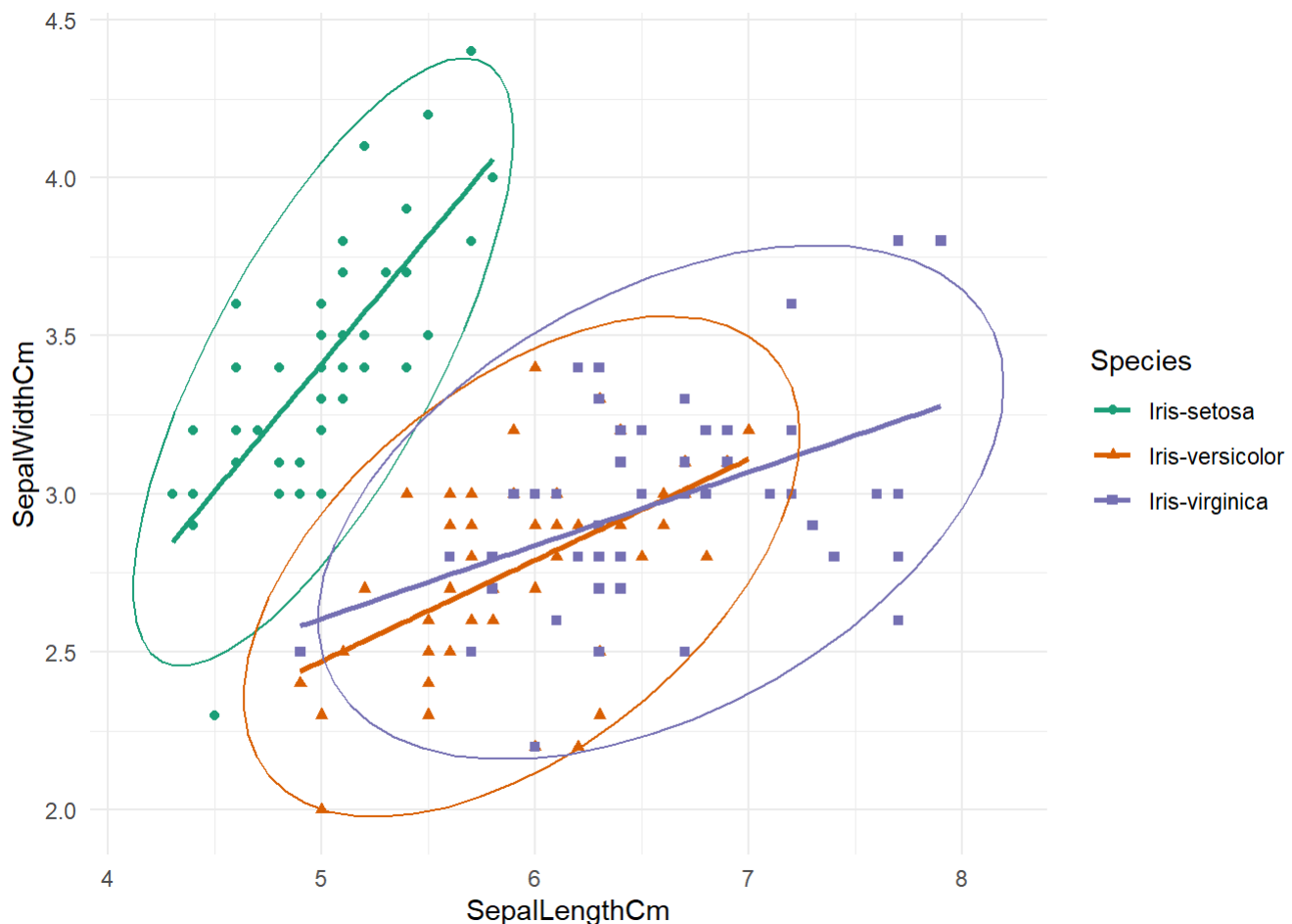
A continuación se presentan las correlaciones y diagramas de dispersión con línea de regresión, por parejas de atributos, pero desglosadas por cada especie de iris.

```
# Longitud del sépalo y ancho del pétalo
lapply(split(data, Species), function(x){cor(x[,1], x[,2])})
```

```
## `$Iris-setosa`
## [1] 0.7467804
##
## `$Iris-versicolor`
## [1] 0.5259107
##
## `$Iris-virginica`
## [1] 0.4572278
```

```
ggplot(data, aes(x=SepalLengthCm, y=SepalWidthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

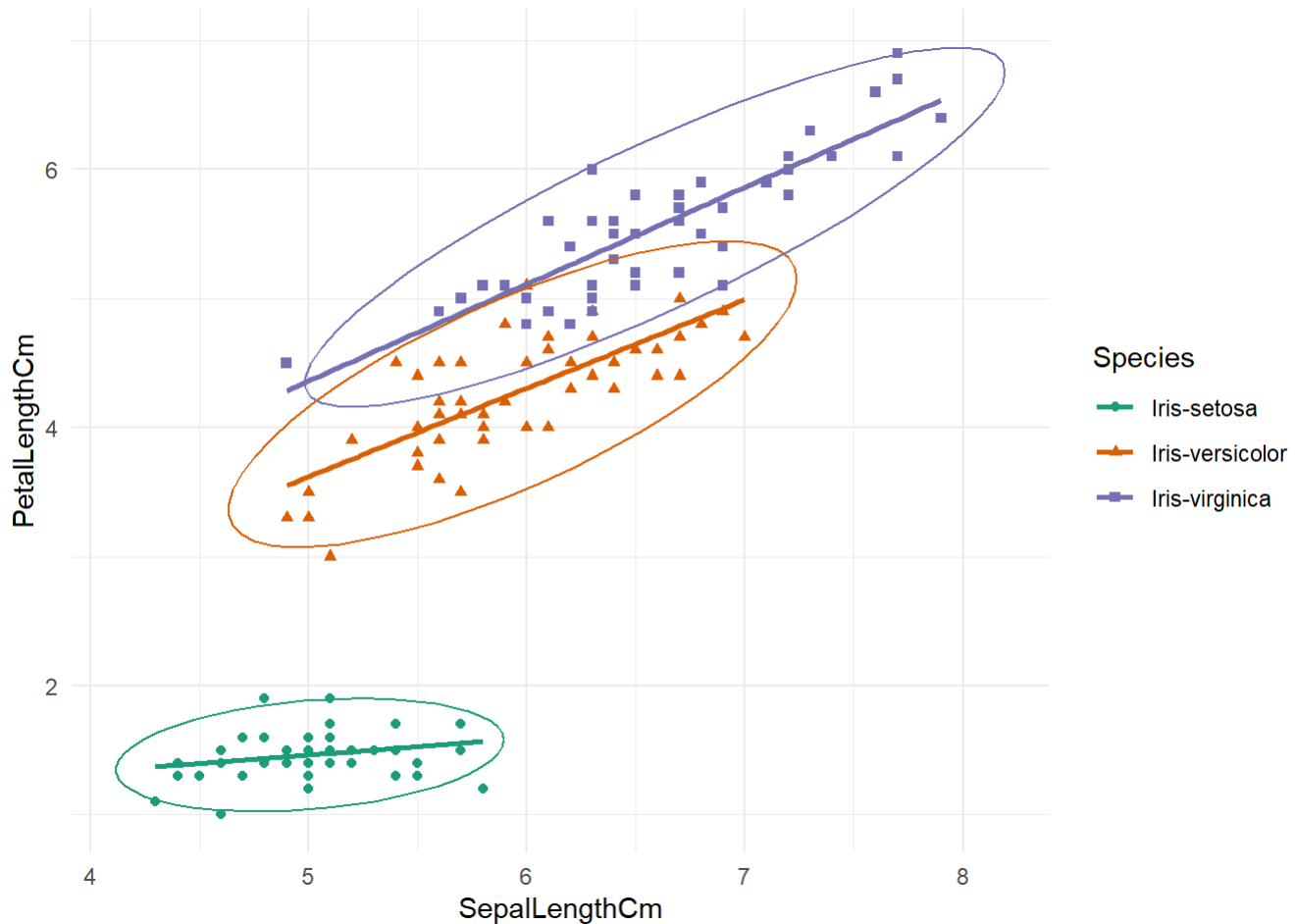


```
# Longitud del sépalo y Longitud del pétalo
lapply(split(data, Species), function(x){cor(x[,1], x[,3])})
```

```
## `$Iris-setosa`
## [1] 0.2638741
##
## `$Iris-versicolor`
## [1] 0.754049
##
## `$Iris-virginica`
## [1] 0.8642247
```

```
ggplot(data, aes(x=SepalLengthCm, y=PetalLengthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

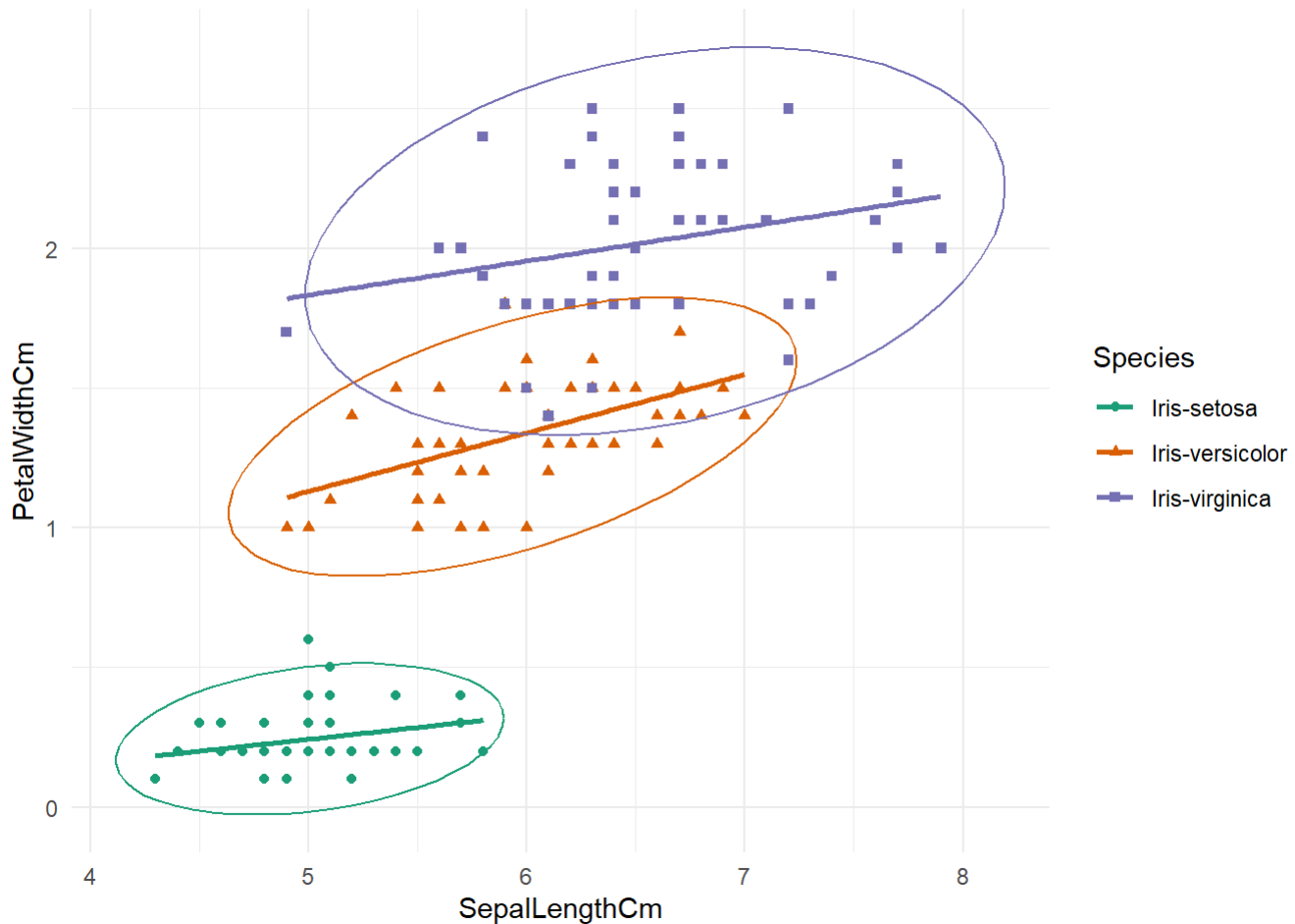


```
# Longitud del sépalo y ancho del pétalo
lapply(split(data, Species), function(x){cor(x[,1], x[,4])})
```

```
## `$Iris-setosa`
## [1] 0.2790916
##
## `$Iris-versicolor`
## [1] 0.5464611
##
## `$Iris-virginica`
## [1] 0.2811077
```

```
ggplot(data, aes(x=SepalLengthCm, y=PetalWidthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

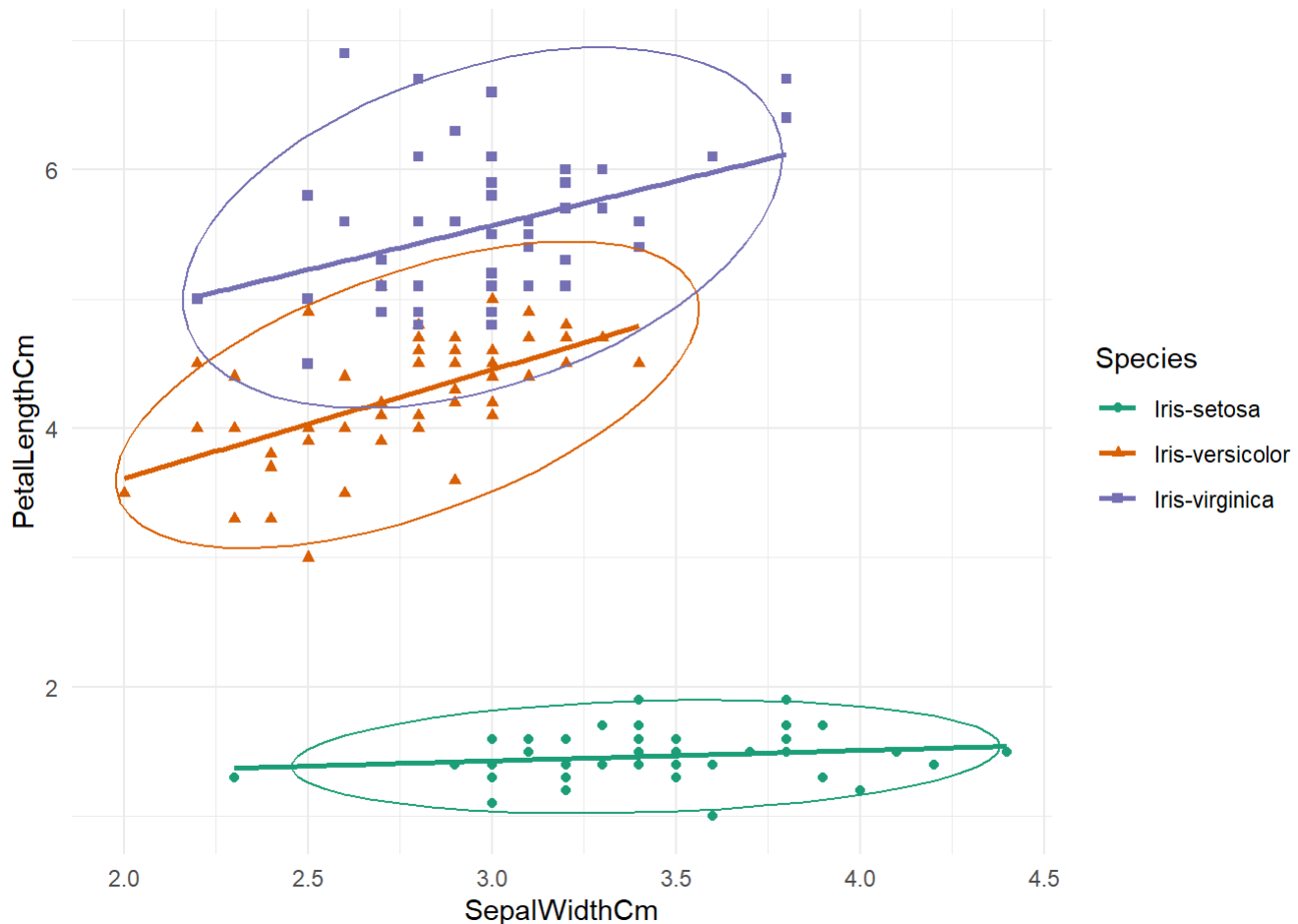


```
# ancho del sépalo y longitud del pétalo
lapply(split(data, Species), function(x){cor(x[,2], x[,3])})
```

```
## `$Iris-setosa`
## [1] 0.1766946
##
## `$Iris-versicolor`
## [1] 0.5605221
##
## `$Iris-virginica`
## [1] 0.4010446
```

```
ggplot(data, aes(x=SepalWidthCm, y=PetalLengthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

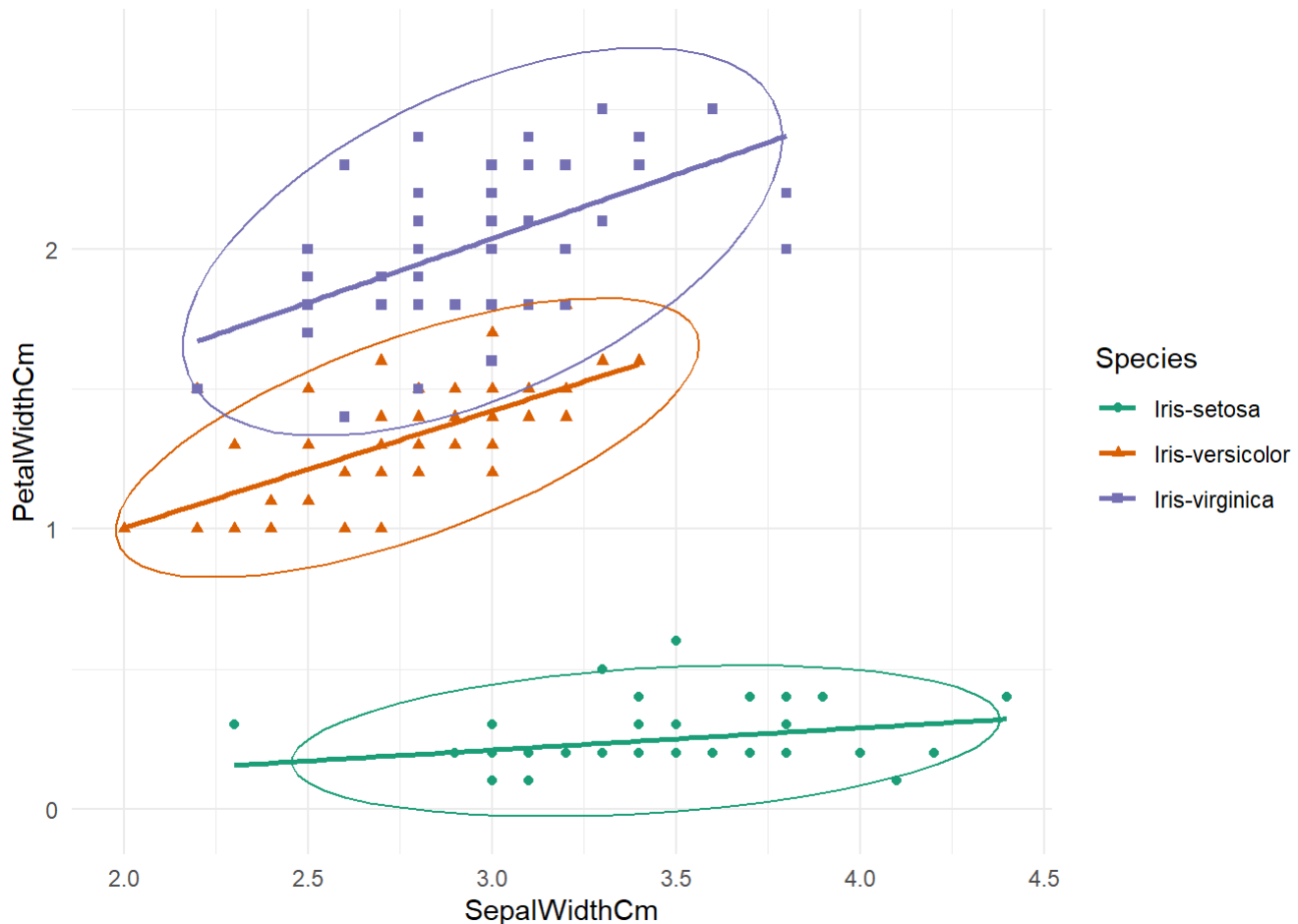


```
# ancho del sépalo y ancho del pétalo
lapply(split(data, Species), function(x){cor(x[,2], x[,4])})
```

```
## `$Iris-setosa`
## [1] 0.2799729
##
## `$Iris-versicolor`
## [1] 0.6639987
##
## `$Iris-virginica`
## [1] 0.537728
```

```
ggplot(data, aes(x=SepalWidthCm, y=PetalWidthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

```
## `geom_smooth()` using formula 'y ~ x'
```

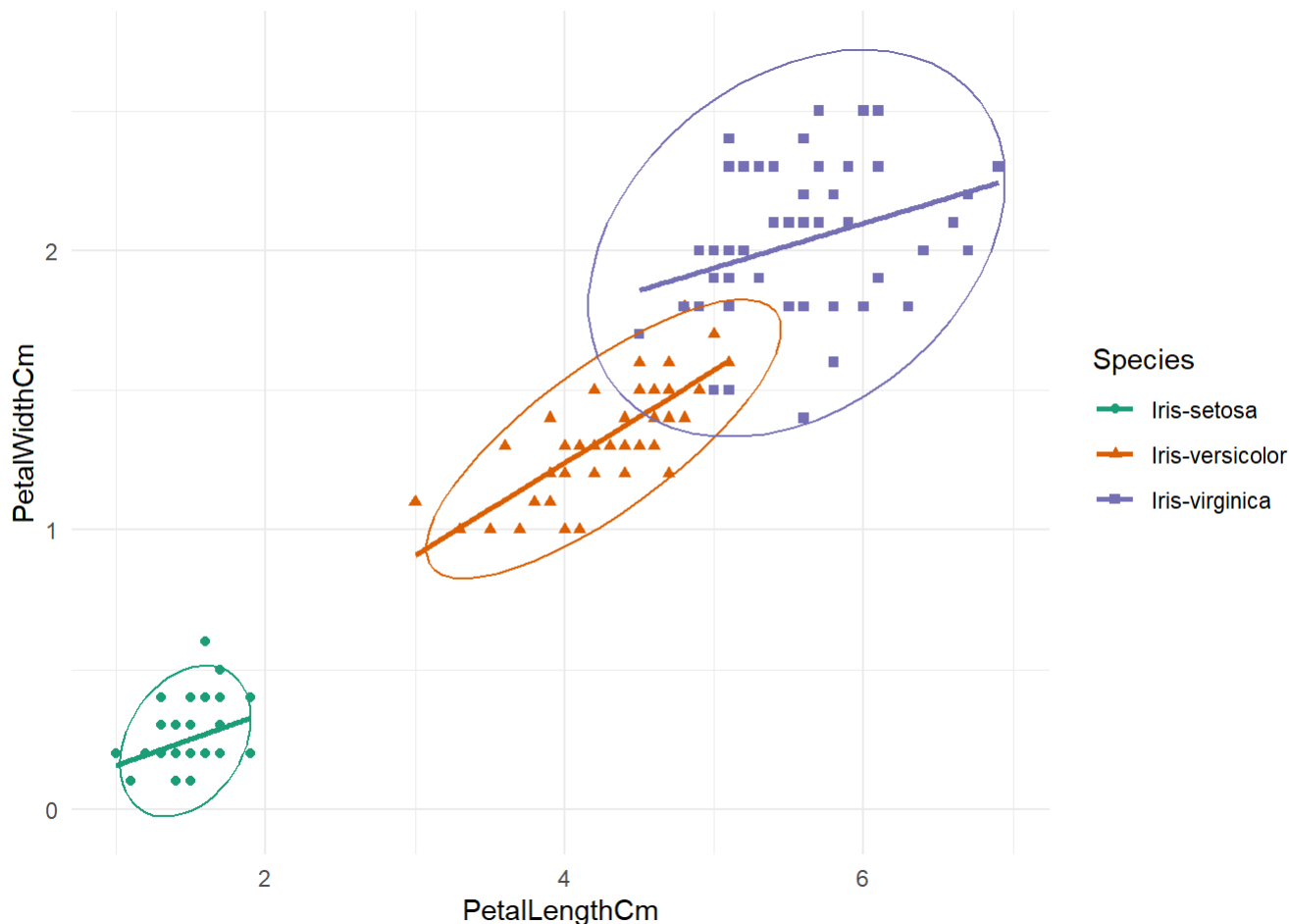


```
# Longitud del pétalo y ancho de pétalo
lapply(split(data, Species), function(x){cor(x[,3], x[,4])})
```

```
## `$Iris-setosa`
## [1] 0.3063082
##
## `$Iris-versicolor`
## [1] 0.7866681
##
## `$Iris-virginica`
## [1] 0.3221082
```

```
ggplot(data, aes(x=PetalLengthCm, y=PetalWidthCm, shape=Species, color=Species))+
  geom_point() +
  geom_smooth(method=lm, se=F, fullrange=F)+
  scale_color_brewer(palette="Dark2")+
  theme_minimal()+
  stat_ellipse(type = "norm")
```

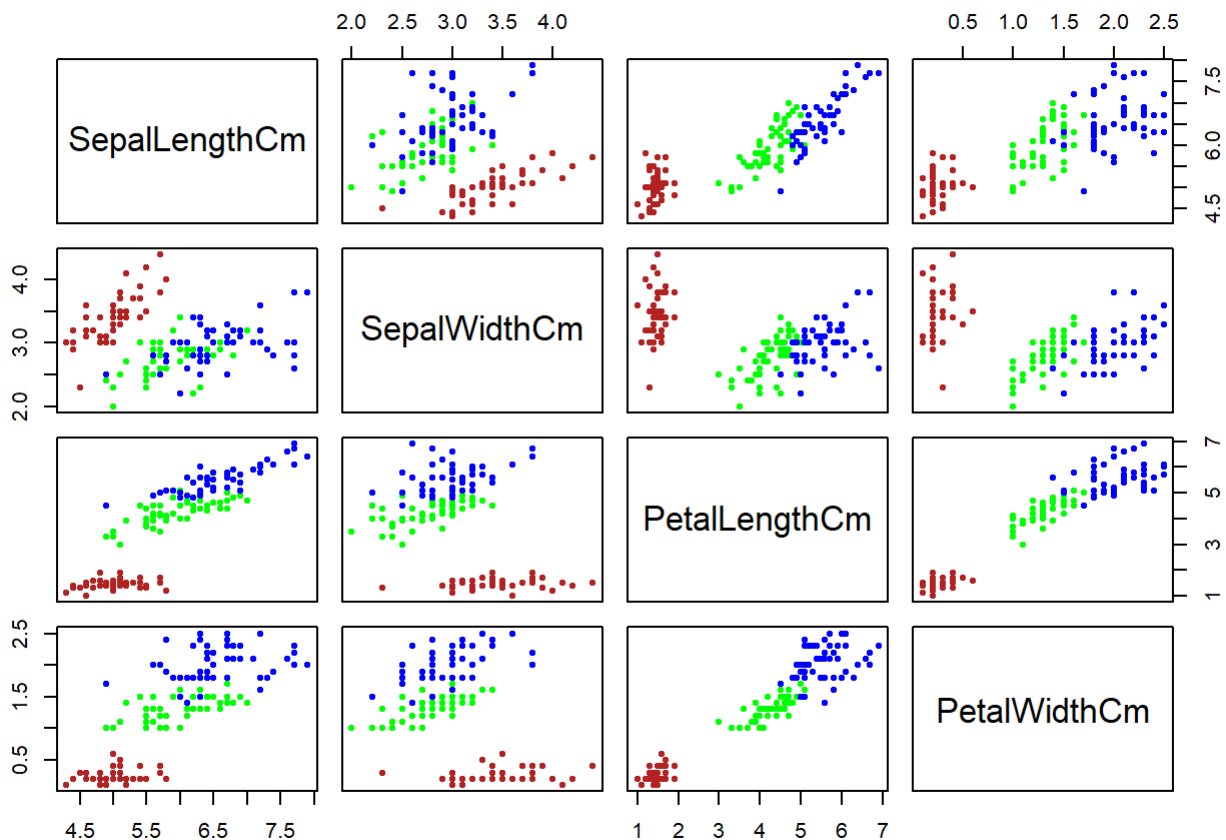
```
## `geom_smooth()` using formula 'y ~ x'
```



En cuanto a la correlación para cada uno de los grupos o especies de flores iris se observa que en el caso de la especie setosa, los atributos con correlaciones relevantes son la longitud y ancho del sépalo (76%). En cuanto a la especie versicolor, las correlaciones superiores al 60% ocurren entre las longitudes del sépalo y pétalo (75%) y entre los anchos del sépalo y pétalo (66%). Y, en la especie virginica, se presenta una única correlación alta entre las longitudes del sépalo y pétalo (86%).

Finalmente, el diagrama de dispersión -según las especies de flor iris- visualizadas en conjunto permiten identificar la posible separación de las especies. Se observa que las variables longitud y ancho del pétalo son las dos variables con más potencial para poder separar entre especies. Sin embargo, como se indicó en párrafos anteriores están altamente correlacionadas, por lo que la información que aportan es en gran medida redundante.

```
# Diagrama de dispersión
pairs(x = data[, -5], col = c("firebrick", "green", "blue")[data$Species],
      pch = 20)
```



```
par(mfrow=c(1,1))
```

Ahora bien con el fin de clasificar las tres especies de iris a partir de sus atributos de longitud y ancho del sépal y pétalo se plantea desarrollar un Análisis Discriminante Lineal o Linear Discriminant Analysis (LDA). La LDA es un método de clasificación de variables cualitativas en el que dos o más grupos son conocidos a priori y nuevas observaciones se clasifican en uno de ellos en función de sus características. Haciendo uso del teorema de Bayes, LDA estima la probabilidad de que una observación, dado un determinado valor de los predictores, pertenezca a cada una de las clases de la variable cualitativa, $P(Y=k|X=x)$. Finalmente se asigna la observación a la clase k para la que la probabilidad predicha es mayor.

Se requieren las siguientes dos condiciones para que el LDA se considere válido:

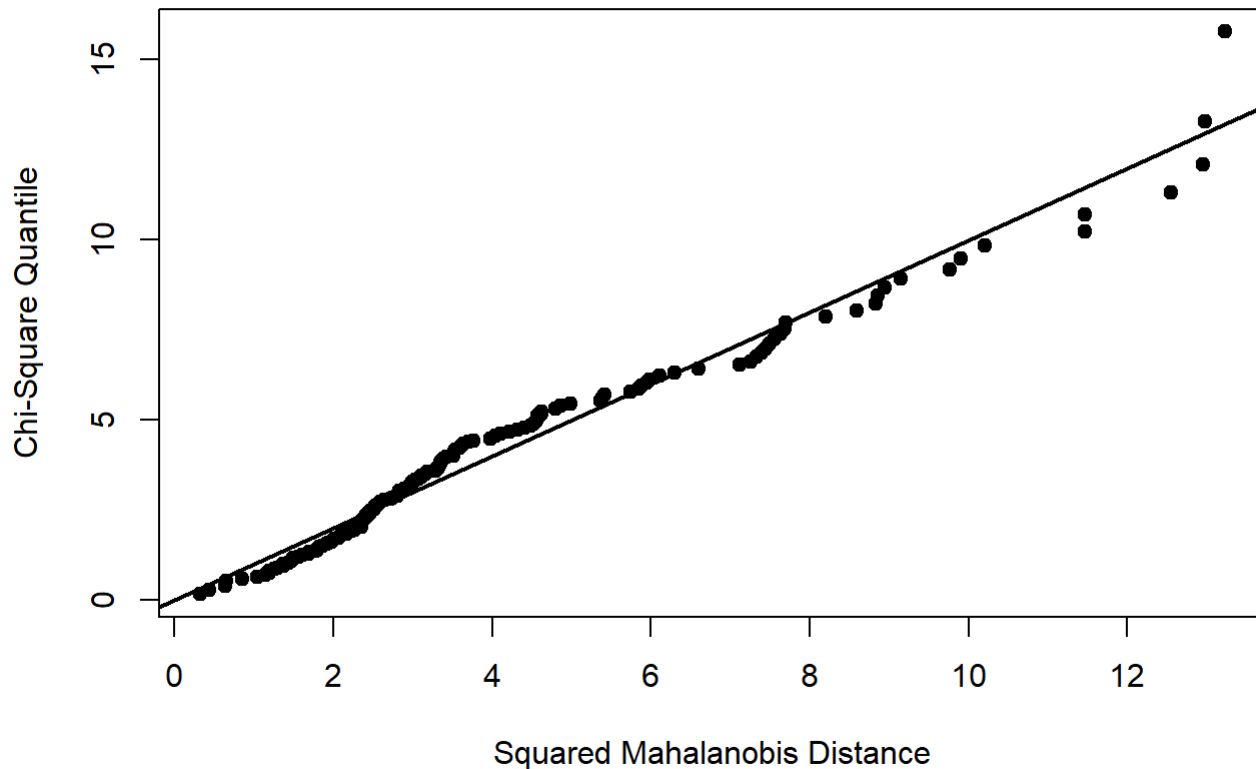
La primera es que cada predictor que forma parte del modelo se distribuye de forma normal en cada una de las clases de la variable respuesta. En un apartado anterior se presentaron los resultados y en general se puede decir que la mayoría de los predictores en cada una de las clases siguen la distribución normal, a excepción de la variable ancho del pétalo, la cual no se distribuye de forma normal en las especies setosa y versicolor.

En el caso de múltiples predictores, las observaciones siguen una distribución normal multivariante en todas las clases.

Con el fin de verificar el cumplimiento de esta condición se aplica el test de normalidad multivariante royston.

```
# LDA
# Verificar normalidad multivariante
# test de royston
royston_test <- mvn(data = data[, -5], mvnTest = "royston", multivariatePlot = "qq")
```

Chi-Square Q-Q Plot



```
royston_test$multivariateNormality
```

```
##      Test      H      p value MVN
## 1 Royston 50.64564 2.754697e-11 NO
```

El test muestra evidencias significativas de falta de normalidad multivariante. El LDA tiene cierta robustez frente a la falta de normalidad multivariante, pero es importante tenerlo en cuenta en la conclusión del análisis.

Y la segunda condición a cumplir es que la varianza del predictor es igual en todas las clases de la variable respuesta. En el caso de múltiples predictores, la matriz de covarianza es igual en todas las clases. Si esto no se cumple se recurre a Análisis Discriminante Cuadrático (QDA).

Con el fin de verificar el cumplimiento de esta condición se aplica el test Box M, el cual se utiliza en el caso multivariante y permite contrastar la igualdad de matrices entre grupos.

```
# ¿La matriz de covarianza es constante en todos los grupos?
boxM(data = data[, -5], grouping = data[, 5])
```

```
##
## Box's M-test for Homogeneity of Covariance Matrices
##
## data: data[, -5]
## Chi-Sq (approx.) = 139.24, df = 20, p-value < 2.2e-16
```

El test Box's M muestra evidencias de que la matriz de covarianza no es constante en todos los grupos, $p(2.2e-16) < \alpha(0.05)$, lo que a priori descartaría el método LDA en favor del QDA. Sin embargo, como el test Box's M es muy sensible a la falta de normalidad multivariante, con frecuencia resulta significativo no porque la matriz de covarianza no sea constante sino por la falta de normalidad, cosa que ocurre para los datos en estudio. Por esta razón se va a asumir que la matriz de covarianza sí es constante y que LDA puede alcanzar una buena precisión en la clasificación. En la evaluación del modelo se verá como de buena es esta aproximación.

Ahora bien, se procede al cálculo de la función discriminante:

```
# Modelo LDA
modelo_lda <- lda(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm +
                  PetalWidthCm, data = data)
```

```
# Modelo LDA
modelo_lda
```

```
## Call:
## lda(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm +
##      PetalWidthCm, data = data)
##
## Prior probabilities of groups:
##      Iris-setosa Iris-versicolor Iris-virginica
##      0.3333333    0.3333333    0.3333333
##
## Group means:
##              SepalWidthCm SepalLengthCm PetalLengthCm PetalWidthCm
## Iris-setosa             3.418           5.006           1.464           0.244
## Iris-versicolor         2.770           5.936           4.260           1.326
## Iris-virginica          2.974           6.588           5.552           2.026
##
## Coefficients of linear discriminants:
##              LD1           LD2
## SepalWidthCm  1.5478732  2.15471106
## SepalLengthCm 0.8192685  0.03285975
## PetalLengthCm -2.1849406 -0.93024679
## PetalWidthCm  -2.8538500  2.80600460
##
## Proportion of trace:
##      LD1      LD2
## 0.9915 0.0085
```

5. Representación de los resultados a partir de tablas y gráficas.

Se procede a realizar la predicción con el mismo dataset (original)

```
# Realiza la predicción con el modelo LDA
prediccion<-predict(modelo_lda,data[-5])
```

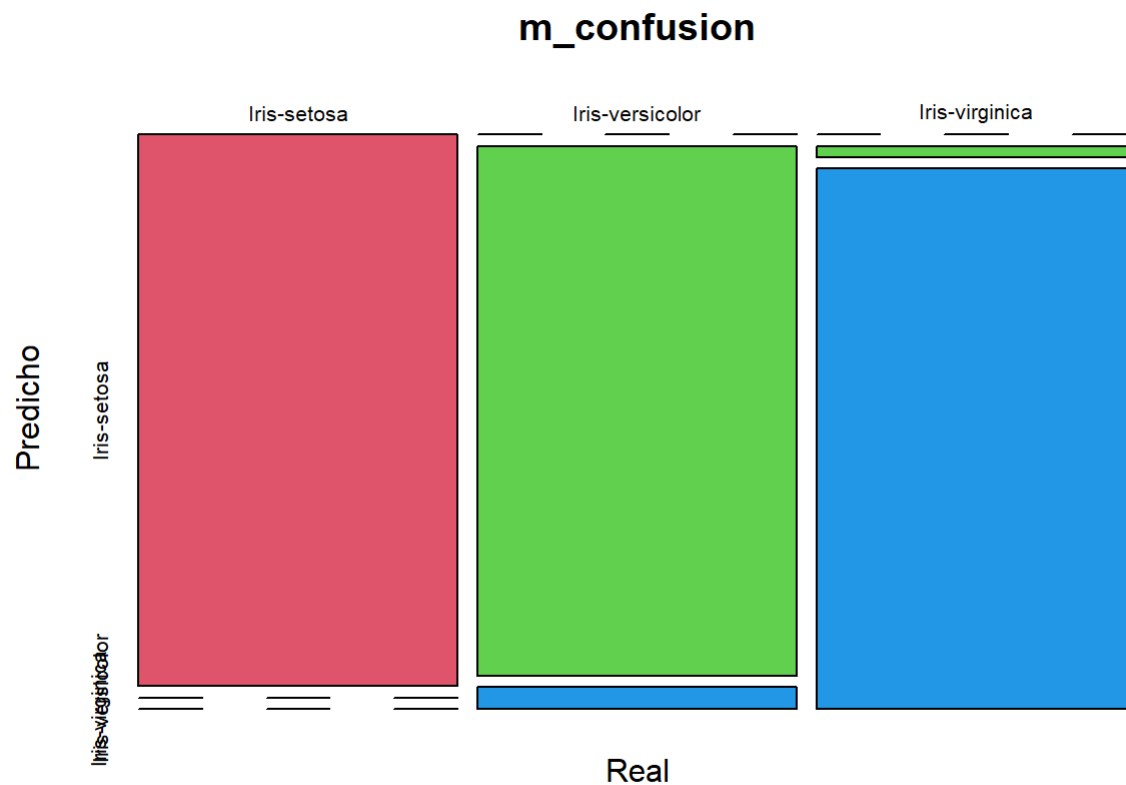
Una vez que las normas de clasificación se han establecido, se tiene que evaluar como de buena es la clasificación resultante. En otras palabras, evaluar el porcentaje de aciertos en las clasificaciones. Para tal fin se presenta la matriz de confusión, la cual presenta el número de verdaderos positivos, verdaderos negativos, falsos positivos y falsos negativos.

```
# Matriz de confusión
m_confusion<-table(data$Species,prediccion$class,
                    dnn=c("Real","Predicho"))
```

```
# Matriz de confusión
m_confusion
```

```
##              Predicho
## Real          Iris-setosa Iris-versicolor Iris-virginica
##  Iris-setosa           50             0             0
##  Iris-versicolor        0             48             2
##  Iris-virginica         0             1             49
```

```
# Matriz de confusión
mosaicplot(m_confusion,col=2:4)
```



Solo 3 de las 150 predicciones que ha realizado el modelo han sido erróneas.

Ahora bien, para evaluar el error de clasificación se emplean las mismas observaciones con las que se ha creado el modelo, obteniendo así lo que se denomina el training error. Si bien esta es una forma sencilla de estimar la precisión en la clasificación, tiende a ser muy optimista. Es más adecuado evaluar el modelo empleando observaciones nuevas que el modelo no ha visto, obteniendo así el test error.

```
# Presición (training error)
precision=mean(data$Species==prediccion$class)
```

```
# Presición (training error)
precision
```

```
## [1] 0.98
```

```
# Presición (training error)
error= (1-precision)*100
```

```
error
```

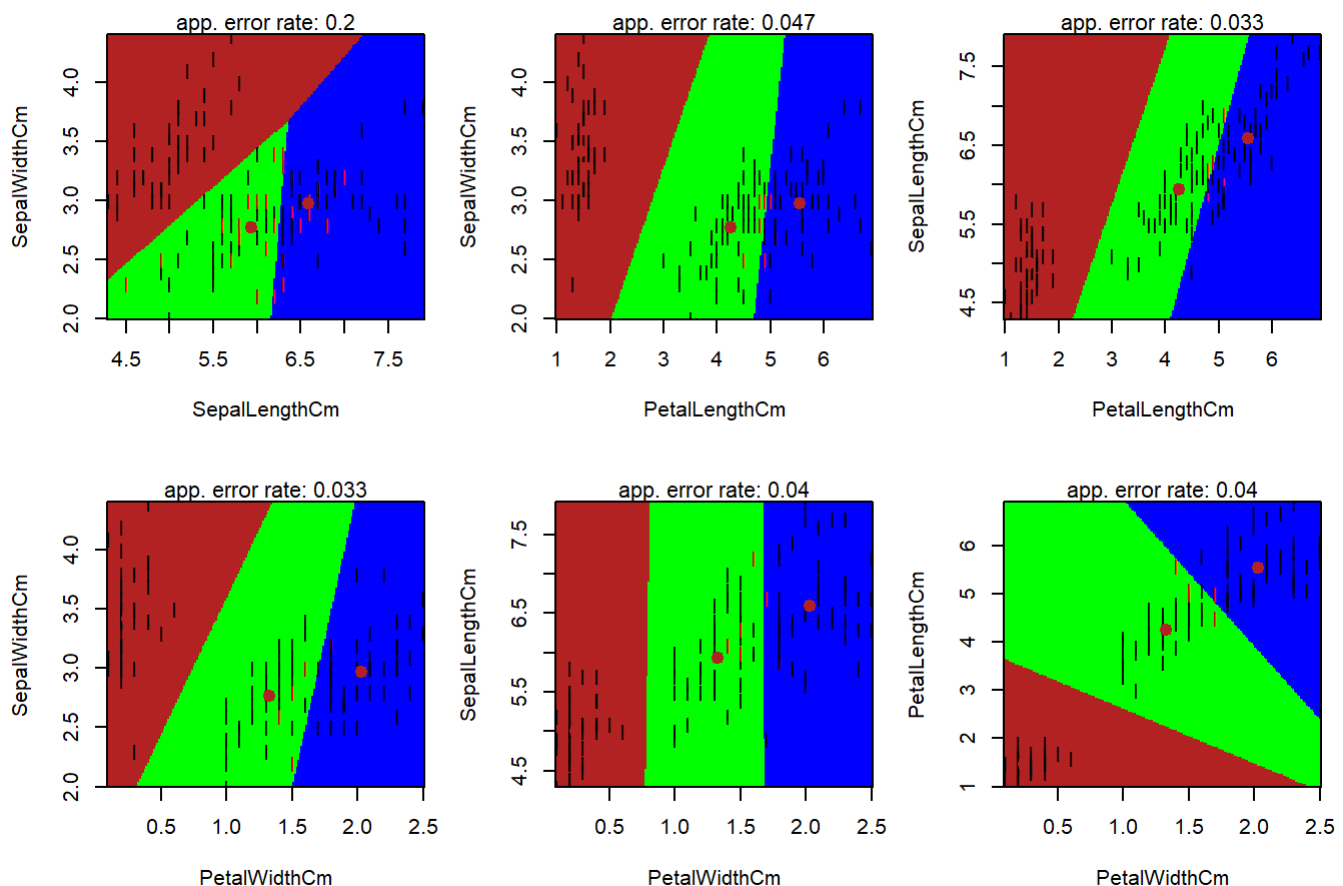
```
## [1] 2
```

El trainig error es muy bajo (2%), lo que apunta a que el modelo es bueno. Sin embargo, para validarlo es necesario un nuevo set de datos con el que calcular el test error o recurrir a validación cruzada.

Se presenta a continuación una visualización que representa los límites de clasificación de un modelo discriminante lineal para cada par de predictores. Cada color representa una región de clasificación acorde al modelo, se muestra el centroide de cada región y el valor real de las observaciones.

```
# Visualización de las clasificaciones
partimat(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm + PetalWidthCm,
  data = data, method = "lda", prec = 200,
  image.colors = c("firebrick", "green", "blue"),
  col.mean = "firebrick")
```

Partition Plot



Creación de un sets de entrenamiento y prueba

Se crea un set de entrenamiento para generar un modelo predictivo, y un set de prueba, para comprobar la eficacia de este modelo para hacer predicciones correctas.

Se obtiene un subconjunto del dataset original, que consiste en 70% del total de ellos. Y se obtiene el subconjunto de datos complementario al de entrenamiento para el set de prueba, esto es, el 30% restante.

```
# Creación de un sets de entrenamiento y prueba
set.seed(1234)
data_entrenamiento <- sample_frac(data, .7)
data_prueba <- setdiff(data, data_entrenamiento)
```

```
# Verificación del set de entrenamiento y prueba  
str(data_entrenamiento)
```

```
## 'data.frame':    105 obs. of  5 variables:  
## $ SepalLengthCm: num  5.2 5.7 6.3 6.5 6.3 6.4 6.8 7.9 6.2 7.1 ...  
## $ SepalWidthCm : num  3.5 2.6 3.3 3.2 3.4 2.8 3.2 3.8 2.9 3 ...  
## $ PetalLengthCm: num  1.5 3.5 6 5.1 5.6 5.6 5.9 6.4 4.3 5.9 ...  
## $ PetalWidthCm : num  0.2 1 2.5 2 2.4 2.2 2.3 2 1.3 2.1 ...  
## $ Species      : Factor w/ 3 levels "Iris-setosa",...: 1 2 3 3 3 3 3 3 2 3 ...
```

```
str(data_prueba)
```

```
## 'data.frame':    45 obs. of  5 variables:  
## $ SepalLengthCm: num  5.1 4.6 4.8 5.8 5.1 4.6 5.1 5.2 4.8 5.2 ...  
## $ SepalWidthCm : num  3.5 3.4 3.4 4 3.5 3.6 3.3 3.4 3.1 4.1 ...  
## $ PetalLengthCm: num  1.4 1.4 1.6 1.2 1.4 1 1.7 1.4 1.6 1.5 ...  
## $ PetalWidthCm : num  0.2 0.3 0.2 0.2 0.3 0.2 0.5 0.2 0.2 0.1 ...  
## $ Species      : Factor w/ 3 levels "Iris-setosa",...: 1 1 1 1 1 1 1 1 1 1 ...
```

Se procede a calcular la función discriminante con el set de entrenamiento:

```
# LDA con el set de entrenamiento  
modelo_lda <- lda(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm +  
                  PetalWidthCm, data = data_entrenamiento)
```

```
modelo_lda
```

```
## Call:
## lda(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm +
##     PetalWidthCm, data = data_entrenamiento)
##
## Prior probabilities of groups:
##      Iris-setosa Iris-versicolor  Iris-virginica
##      0.3238095    0.3238095    0.3523810
##
## Group means:
##              SepalWidthCm SepalLengthCm PetalLengthCm PetalWidthCm
## Iris-setosa      3.370588    4.991176    1.473529    0.2323529
## Iris-versicolor  2.767647    5.920588    4.255882    1.3500000
## Iris-virginica   3.024324    6.640541    5.654054    2.0351351
##
## Coefficients of linear discriminants:
##              LD1      LD2
## SepalWidthCm  1.252587 -2.69650836
## SepalLengthCm 1.059991  0.20893018
## PetalLengthCm -2.475107  0.08323739
## PetalWidthCm  -2.740318 -1.03024775
##
## Proportion of trace:
##      LD1      LD2
## 0.9931 0.0069
```

Se realiza la predicción con el modelo LDA sobre el set de prueba:

```
# Predicción con el set de prueba
prediccion<-predict(modelo_lda,data_prueba[-5])
```

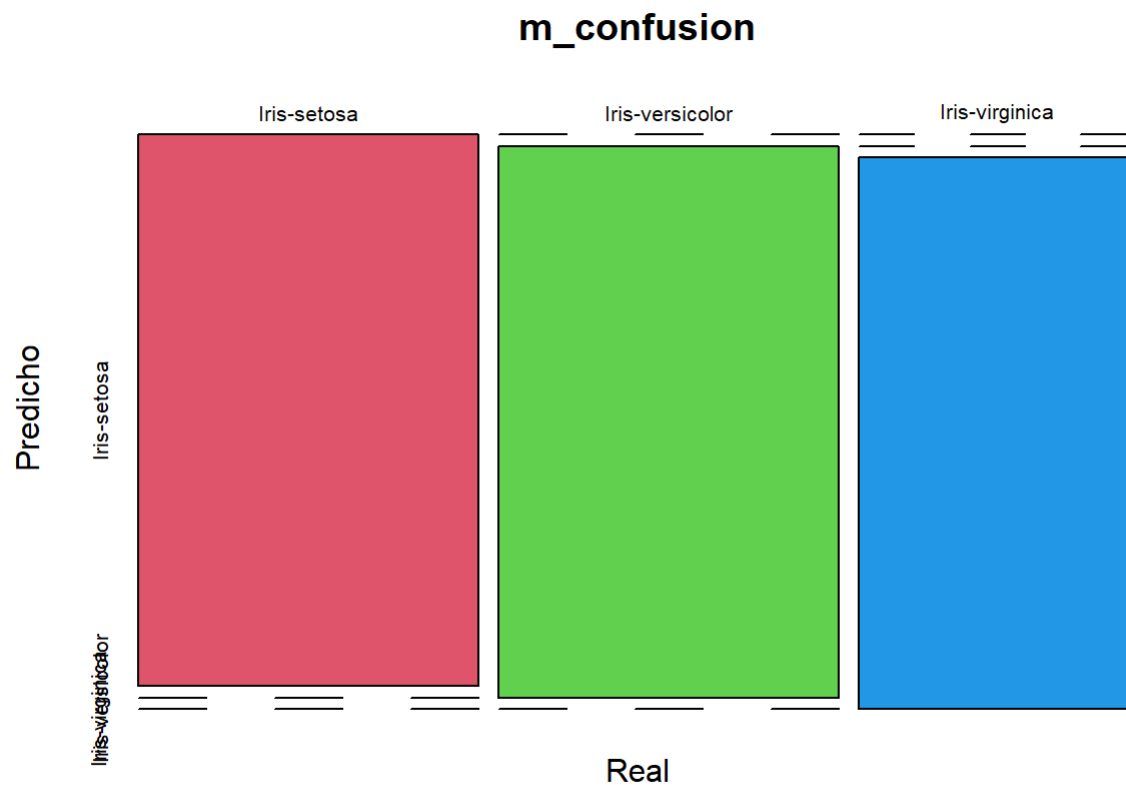
Se procede a evaluar la clasificación del modelo LDA:

```
# Matriz de confusión
m_confusion<-table(data_prueba$Species,prediccion$class,
                    dnn=c("Real","Predicho"))
```

```
m_confusion
```

```
##              Predicho
## Real      Iris-setosa Iris-versicolor Iris-virginica
## Iris-setosa      16           0           0
## Iris-versicolor   0          16           0
## Iris-virginica    0           0          13
```

```
mosaicplot(m_confusion,col=2:4)
```



Se presenta la precisión del modelo LDA:

```
# Presición (test error)
precision=mean(data_prueba$Species==prediccion$class)
```

```
precision
```

```
## [1] 1
```

```
error=(1-precision)*100
```

```
error
```

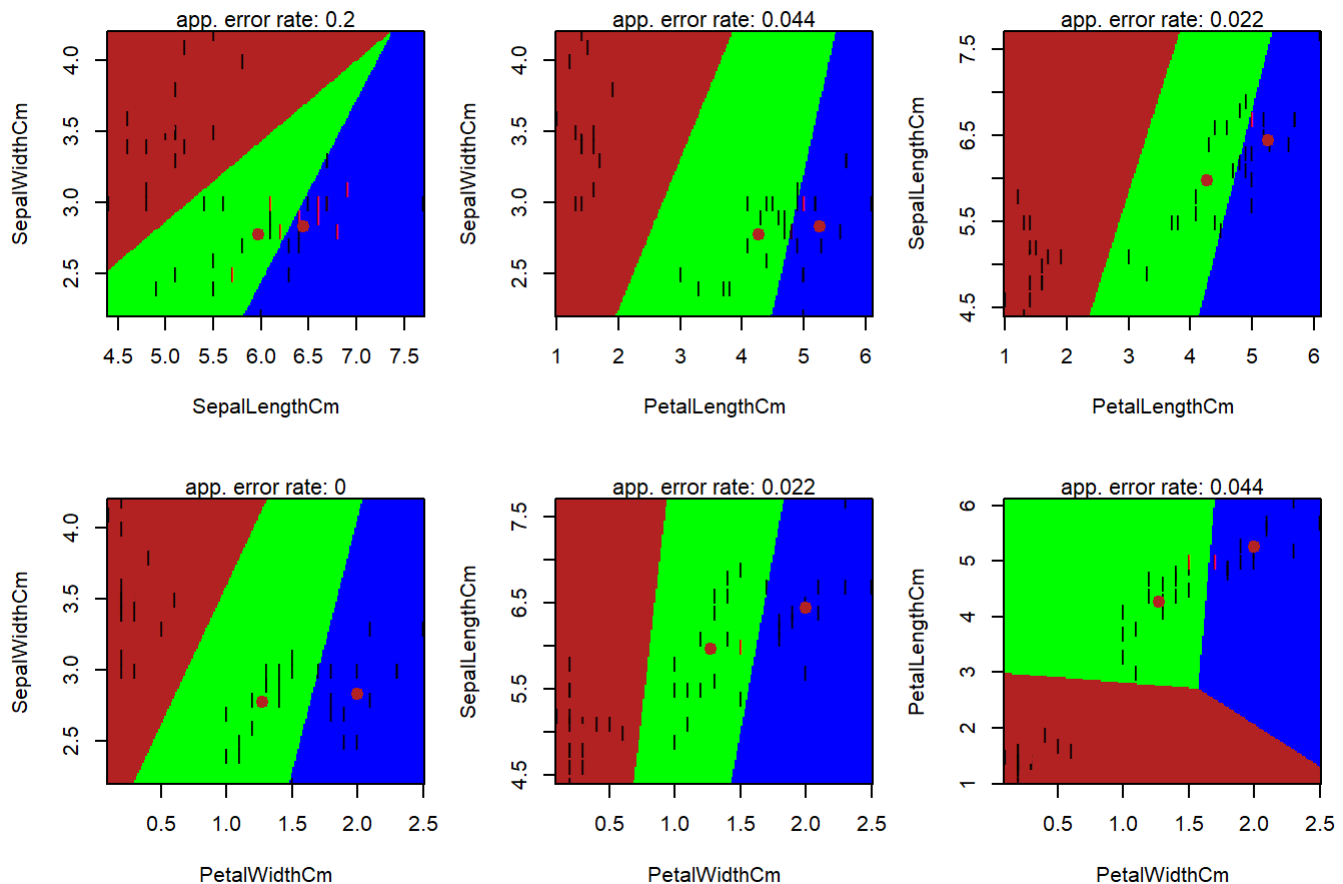
```
## [1] 0
```

De acuerdo con el modelo ninguna de las 45 predicciones realizadas ha sido incorrecta; el test de error es del 0%.

Se presenta la visualización del resultado del modelo LDA:


```
partimat(Species ~ SepalWidthCm + SepalLengthCm + PetalLengthCm + PetalWidthCm,
  data = data_prueba, method = "lda", prec = 200,
  image.colors = c("firebrick", "green", "blue"),
  col.mean = "firebrick")
```

Partition Plot



6. Resolución del problema. A partir de los resultados obtenidos, ¿cuáles son las conclusiones? ¿Los resultados permiten responder al problema?

El modelo de clasificación LDA desarrollado presentan muy buenos resultados en la tarea de predecir o clasificar las especies de flores iris a partir de los atributos de longitud y ancho del sépalos y pétalo, tanto en los datos originales como con datos de prueba.