

# Tema 8 - Estadística descriptiva con datos cuantitativos

Juan Gabriel Gomila & María Santos

## Descripción de datos cuantitativos

## Datos cuantitativos

Los datos cuantitativos son los que expresan cantidades que se representan mediante números. Éstos se suelen clasificar en continuos y discretos.

- ▶ Los datos continuos son los que, si existiese la posibilidad de medirlos con precisión infinita, en principio podrían tomar todos los valores de un intervalo de la recta real. A modo de ejemplo, el peso, la altura, el tiempo... son datos de este tipo.
- ▶ Por su parte, los datos discretos son los que pueden tomar un solo conjunto contable de valores. El número de colores de un gato, el número de individuos que conforman una población son algunos ejemplos de este tipo de datos.

Conviene tener en cuenta que esta división es solo teórica. Es decir, en la práctica, todos estos datos son discretos puesto que la precisión infinita no existe. Sin embargo, es necesario de vez en cuando suponer los datos de tipo continuo para así poder utilizar técnicas específicas en su análisis.

# Datos cuantitativos

A la hora de estudiar variables cuantitativas, podemos utilizar las frecuencias que hemos visto hasta el momento: absoluta, relativa, acumulada y relativa acumulada. Esto se debe a que podemos ordenar los datos cuantitativos en el orden natural de los números reales.

En este caso, disponemos de muchas otras técnicas descriptivas aparte de las frecuencias, puesto que estamos trabajando con números reales y podemos operar con ellos.

Los datos cuantitativos admiten dos tipos de tratamiento según trabajemos con los raw data (datos brutos u originales) o bien los agrupemos en clases o intervalos.

En esta lección trabajaremos sobre la primera situación. En la siguiente, estudiaremos la descripción de datos cuantitativos agrupados.

Frecuencias

# Frecuencias de datos cuantitativos

El tratamiento de las frecuencias de datos cuantitativos es similar al de los datos ordinales. La cosa cambia ligeramente debido a que no se tienen en cuenta todos los niveles posibles, sino únicamente los observados.

## Ejemplo 1

### Ejemplo 1

Se han pedido las edades a 20 clientes de un museo. Las respuestas obtenidas han sido las siguientes:

```
edad = c(15, 18, 25, 40, 30, 29, 56, 40, 13, 27, 42, 23, 11, 26, 25, 32, 30)
```

Recordemos que solamente nos interesan las frecuencias de las edades observadas. Es decir, solamente nos interesan

```
table(edad)
```

edad

|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 11 | 13 | 15 | 18 | 23 | 25 | 26 | 27 | 29 | 30 | 32 | 33 | 40 | 42 | 56 |
| 1  | 1  | 1  | 1  | 1  | 2  | 1  | 1  | 2  | 2  | 1  | 1  | 3  | 1  | 1  |

## Ejemplo 1

Calculemos el resto de frecuencias como ya sabemos

```
round(prop.table(table(edad)),3)
```

edad

|      |      |      |      |      |      |      |      |      |      |      |      |
|------|------|------|------|------|------|------|------|------|------|------|------|
| 11   | 13   | 15   | 18   | 23   | 25   | 26   | 27   | 29   | 30   | 32   | 33   |
| 0.05 | 0.05 | 0.05 | 0.05 | 0.05 | 0.10 | 0.05 | 0.05 | 0.10 | 0.10 | 0.05 | 0.05 |

```
cumsum(table(edad))
```

|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 11 | 13 | 15 | 18 | 23 | 25 | 26 | 27 | 29 | 30 | 32 | 33 | 40 | 42 | 56 |
| 1  | 2  | 3  | 4  | 5  | 7  | 8  | 9  | 11 | 13 | 14 | 15 | 18 | 19 | 20 |



## Ejemplo 1

```
round(cumsum(prop.table(table(edad))),3)
```

|      |      |      |      |      |      |      |      |      |      |      |      |
|------|------|------|------|------|------|------|------|------|------|------|------|
| 11   | 13   | 15   | 18   | 23   | 25   | 26   | 27   | 29   | 30   | 32   | 33   |
| 0.05 | 0.10 | 0.15 | 0.20 | 0.25 | 0.35 | 0.40 | 0.45 | 0.55 | 0.65 | 0.70 | 0.75 |

## Frecuencias de datos cuantitativos

En general, supongamos que tenemos  $n$  observaciones de una propiedad que se mide con un número real y obtenemos la variable cuantitativa formada por los datos

$$x_1, \dots, x_n$$

Sean ahora  $x_1, \dots, x_k$  los valores distintos que aparecen en esta lista de datos y considerémoslos ordenados

$$x_1 < x_2 < \dots < x_k$$

# Frecuencias de datos cuantitativos

Entonces, en esta variable cuantitativa

- ▶ La frecuencia absoluta de  $X_i$  es el número  $n_i$  de elementos que son iguales a  $X_i$
- ▶ La frecuencia relativa de  $X_i$  es  $f_i = \frac{n_i}{n}$
- ▶ La frecuencia absoluta acumulada de  $X_i$  es  $N_i = \sum_{j=1}^i n_j$
- ▶ La frecuencia relativa acumulada de  $X_i$  es  $F_i = \frac{N_i}{n}$

## Ejemplo 2

### Ejemplo 2

Lanzamos 25 veces un dado de 6 caras y anotamos las puntuaciones obtenidas en cada tirada.

En este caso,  $n = 25$  y, los distintos valores observados son

$$X_1 = 1, X_2 = 2, X_3 = 3, X_4 = 4, X_5 = 5, X_6 = 6$$

Nos interesa ahora calcular las frecuencias de este experimento. Además, las organizaremos en un data frame para observarlas de forma más clara y sencilla en una tabla.

```
set.seed(162017)
dados = sample(1:6,25,replace = TRUE)
dados
```

```
[1] 1 1 5 5 5 5 1 6 5 4 1 3 1 3 2 2 1 1 1 4 2 1 6 3 1
```

```
set.seed(NULL)
```

## Ejemplo 2

```
table(dados)
```

dados

| 1  | 2 | 3 | 4 | 5 | 6 |
|----|---|---|---|---|---|
| 10 | 3 | 3 | 2 | 5 | 2 |

```
round(prop.table(table(dados)),2)
```

dados

| 1    | 2    | 3    | 4    | 5    | 6    |
|------|------|------|------|------|------|
| 0.40 | 0.12 | 0.12 | 0.08 | 0.20 | 0.08 |

```
cumsum(table(dados))
```

| 1  | 2  | 3  | 4  | 5  | 6  |
|----|----|----|----|----|----|
| 10 | 13 | 16 | 18 | 23 | 25 |



## Ejemplo 2

```
datos.df
```

|   | Puntuacion | Fr.abs | Fr.rel | Fr.acu | Fr.racu |
|---|------------|--------|--------|--------|---------|
| 1 | 1          | 10     | 0.40   | 10     | 0.40    |
| 2 | 2          | 3      | 0.12   | 13     | 0.52    |
| 3 | 3          | 3      | 0.12   | 16     | 0.64    |
| 4 | 4          | 2      | 0.08   | 18     | 0.72    |
| 5 | 5          | 5      | 0.20   | 23     | 0.92    |
| 6 | 6          | 2      | 0.08   | 25     | 1.00    |

¡OJO! Para entrar una tabla unidimensional como una variable en un data frame, es conveniente transformarla en vector con `as.vector`. Si no, cada `table` y cada `prop.table` añadirían una columna extra con los nombres de los niveles.

## Medidas de tendencia central



# Medidas de tendencia central

Las medidas de tendencia central son las que dan un valor representativo a todas las observaciones. Algunas de las más importantes son:

- ▶ La media aritmética o valor medio

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{\sum_{j=1}^k n_j X_j}{n} = \sum_{j=1}^k f_j X_j$$

- ▶ La mediana, que representa el valor central en la lista ordenada de observaciones.
- ▶ La moda es el valor (o valores) de máxima frecuencia (absoluta o relativa, el resultado será el mismo).

# La mediana

La definición formal de la mediana es la siguiente. Denotando por

$$x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)}$$

los datos de la variable cuantitativa ordenados de menor a mayor, la mediana es

- ▶ Si  $n$  par, la medio de los dos datos centrales

$$\frac{x_{(\frac{n}{2})} + x_{(\frac{n}{2}+1)}}{2}$$

- ▶ Si  $n$  impar, el dato central  $x_{(\frac{n+1}{2})}$

## Ejemplo 1

Recordemos el ejemplo de las edades.

```
sort(edad) #Ordenamos los datos por su orden natural
```

```
[1] 11 13 15 18 23 25 25 26 27 29 29 30 30 32 33 40 40 40
```

```
table(edad)
```

edad

|    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 11 | 13 | 15 | 18 | 23 | 25 | 26 | 27 | 29 | 30 | 32 | 33 | 40 | 42 | 56 |
| 1  | 1  | 1  | 1  | 1  | 2  | 1  | 1  | 2  | 2  | 1  | 1  | 3  | 1  | 1  |

En este caso, la moda es 40, la mediana es  $\frac{29+29}{2} = 29$  y la media aritmética es

$$\frac{11+13+15+18+23+25+25+26+27+29+29+30+30+32+33+40+40+40+42+56}{20} = 29.2$$

## Ejemplo 2

Recordemos el ejemplo de los dados.

```
dados.df
```

|   | Puntuacion | Fr.abs | Fr.rel | Fr.acu | Fr.racu |
|---|------------|--------|--------|--------|---------|
| 1 | 1          | 10     | 0.40   | 10     | 0.40    |
| 2 | 2          | 3      | 0.12   | 13     | 0.52    |
| 3 | 3          | 3      | 0.12   | 16     | 0.64    |
| 4 | 4          | 2      | 0.08   | 18     | 0.72    |
| 5 | 5          | 5      | 0.20   | 23     | 0.92    |
| 6 | 6          | 2      | 0.08   | 25     | 1.00    |

En este caso, la moda son dos valores: el 2 y el 3. La mediana es  $x_{(13)} = 3$  y la media aritmética es 2.8

## Medidas de tendencia central en R

Vamos a calcular la media aritmética, mediana y moda de los dos ejemplos anteriores con instrucciones de R.

```
mean(edad) #La media aritmética
```

```
[1] 29.2
```

```
mean(dados)
```

```
[1] 2.8
```

```
median(edad) #La mediana
```

```
[1] 29
```

## Medidas de tendencia central en R

```
median(dados)
```

```
[1] 2
```

```
as.numeric(names(which(table(edad)==max(table(edad))))) #L
```

```
[1] 40
```

```
as.numeric(names(which(table(dados)==max(table(dados)))))
```

```
[1] 1
```

Cuando trabajamos con datos cuantitativos, es conveniente que el resultado lo demos como un número. De ahí que hayamos aplicado la función `as.numeric`.

## Medidas de posición

# Medidas de posición

Las medidas de posición estiman qué valores dividen las observaciones en unas determinadas proporciones.

Los valores que determinan estas posiciones son conocidos como los cuantiles.

Pensándolo de este modo, la mediana puede interpretarse como una medida de posición, debido a que divide la variable cuantitativa en dos mitades.



## Medidas de posición

Dada una proporción  $p \in (0, 1)$ , el cuantil de orden  $p$  de una variable cuantitativa,  $Q_p$ , es el valor más pequeño tal que su frecuencia relativa acumulada es mayor o igual a  $p$ .

Dicho de otro modo, si tenemos un conjunto de observaciones  $x_1, \dots, x_n$  y los ordenamos de menor a mayor, entonces  $Q_p$  será el número más pequeño que deja a su izquierda (incluyéndose a sí mismo) como mínimo a la fracción  $p$  de los datos. Es decir,  $p \cdot n$ .

Así, ahora es más claro ver que la mediana vendría a ser  $Q_{0.5}$ , el cuantil de orden 0.5.

### Ejemplo 3

Consideremos un experimento en el que lanzamos 50 veces un dado de 4 caras y obtenemos los siguientes resultados

[1] 50

```
dado = sort(dado) #Los ordenamos de menor a mayor
dato
```

```
[1] 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2
[39] 4 4 4 4 4 4 4 4 4 4 4 4
```

### Ejemplo 3

```
df.dado = data.frame(Puntuacion = 1:4,  
                      Fr.abs = as.vector(table(dado)),  
                      Fr.rel = as.vector(round(prop.table(t  
                      Fr.acu = as.vector(cumsum(table(dado)  
                      Fr.racu = as.vector(round(cumsum(prop  
  
df.dado
```

|   | Puntuacion | Fr.abs | Fr.rel | Fr.acu | Fr.racu |
|---|------------|--------|--------|--------|---------|
| 1 | 1          | 16     | 0.32   | 16     | 0.32    |
| 2 | 2          | 15     | 0.30   | 31     | 0.62    |
| 3 | 3          | 5      | 0.10   | 36     | 0.72    |
| 4 | 4          | 14     | 0.28   | 50     | 1.00    |

Si nos piden el cuantil  $Q_{0.3}$ , sabemos que este es el primer elemento de la lista cuya frecuencia relativa acumulada es mayor o igual a 0.3. Este se corresponde con la puntuación 1.

## Ejemplo 3

También podríamos hallarlo de otro modo: fijándonos en la lista ordenada de puntuaciones, el cuantil  $Q_{0.3}$  sería el primer elemento de dicha lista tal que fuera mayor o igual que, como mínimo, el 30% de los datos. Si calculamos el 30% de 50, obtenemos que es 15. Esto lo que nos dice es que el cuantil que buscamos es el número que se encuentra en la quinceava posición de la lista ordenada.

```
dato[15]
```

```
[1] 1
```

# Cuantiles

Algunos cuantiles tienen nombre propio:

- ▶ Los cuartiles son los cuantiles  $Q_{0.25}$ ,  $Q_{0.5}$  y  $Q_{0.75}$ .  
Respectivamente, son llamados primer, segundo y tercer cuartil. El primer cuartil,  $Q_{0.25}$ , será el menor valor que es mayor o igual a una cuarta parte de las observaciones y  $Q_{0.75}$ , el menor valor que es mayor o igual a tres cuartas partes de los datos observados.
- ▶ El cuantil  $Q_{0.5}$  es la mediana
- ▶ Los deciles son los cuantiles  $Q_p$  con  $p$  un múltiplo de 0.1.
- ▶ Los percentiles son los cuantiles  $Q_p$  con  $p$  un múltiplo de 0.01.

## Cuantiles

La definición de cuantil anteriormente dada es orientativa. La realidad es que, exceptuando el caso de la mediana, no hay consenso sobre cómo deben calcularse los cuantiles. En verdad, existen diferentes métodos que pueden dar lugar a soluciones distintas.

Al fin y al cabo, nuestro objetivo no es el de encontrar el primer valor de una muestra cuya frecuencia relativa acumulada en la variable sea mayor o igual a  $p$ , sino estimar el valor de esta cantidad para el total de la población.

Para calcular los cuantiles de orden  $p$  de una variable cualitativa  $x$  con R, se utiliza la instrucción `quantile(x,p)`, la cual dispone de 9 métodos diferentes que se especifican con el parámetro `type`. El valor por defecto es `type = 7` y no hace falta especificarlo, como veremos en el siguiente ejemplo. Para más información sobre todos los valores posibles de este parámetro, haced click en el enlace a [Wikipedia](#)

## Ejemplo 4

```
set.seed(0)
dados2 = sample(1:6,15, replace = TRUE)
dados2
```

```
[1] 6 1 4 1 2 5 3 6 2 3 3 1 5 5 2
```

```
set.seed(NULL)
quantile(dados2,0.25) #Primer cuartil
```

```
25%
  2
```

```
quantile(dados2,0.8)
```

```
80%
  5
```

## Medidas de dispersión



## Medidas de dispersión

Las medidas de dispersión evalúan lo dispersos que están los datos. Algunas de las más importantes son:

- ▶ El rango o recorrido, que es la diferencia entre el máximo y el mínimo de las observaciones.
- ▶ El rango intercuartílico, que es la diferencia entre el tercer y primer cuartil,  $Q_{0.75} - Q_{0.25}$ .
- ▶ La varianza, a la que denotaremos por  $s^2$ , es la media aritmética de las diferencias al cuadrado entre los datos  $x_i$  y la media aritmética de las observaciones,  $\bar{x}$ .

$$s^2 = \frac{\sum_{j=1}^n (x_j - \bar{x})^2}{n} = \frac{\sum_{j=1}^k n_j (X_j - \bar{x})^2}{n} = \sum_{j=1}^k f_j (X_j - \bar{x})^2$$

## Medidas de dispersión

- ▶ La desviación típica es la raíz cuadrada positiva de la varianza,  $s = \sqrt{s^2}$ .
- ▶ La varianza muestral es la corrección de la varianza. La denotamos por  $\tilde{s}^2$  y se corresponde con

$$\tilde{s}^2 = \frac{n}{n-1} s^2 = \frac{\sum_{j=1}^n (x_i - \bar{x})^2}{n-1}$$

- ▶ La desviación típica muestral, que es la raíz cuadrada positiva de la varianza muestral,  $\tilde{s} = \sqrt{\tilde{s}^2}$

# Propiedades de la varianza

## Propiedades de la varianza.

- ▶  $s^2 \geq 0$ . Esto se debe a que, por definición, es una suma de cuadrados de números reales.
- ▶  $s^2 = 0 \implies x_j - \bar{x} = 0 \ \forall j = 1, \dots, n$ . En consecuencia, si  $s^2 = 0$ , entonces todos los datos son iguales.
- ▶  $s^2 = \frac{\sum_{j=1}^n x_j^2}{n} - \bar{x}^2$ . Es decir, la varianza es la media de los cuadrados de los datos menos el cuadrado de la media aritmética de estos.

## Varianza y varianza muestral

La diferencia entre ambas definiciones viene por la interrelación entre la estadística descriptiva y la inferencial.

Por un lado, es normal medir cómo varían los datos cuantitativos mediante su varianza definida como la media aritmética de las distancias al cuadrado de los datos a su valor medio. No obstante, por otro lado, el conjunto de nuestras observaciones, por lo normal, será una muestra de una población mucho mayor y nos interesará estimar entre otras muchas cosas su variabilidad.

La varianza de una muestra suele dar valores más pequeños que la varianza de la población, mientras que la varianza muestral tiende a dar valores alrededor de la varianza de la población.

## Varianza y varianza muestral

Esta corrección, para el caso de una muestra grande no es notable. Dividir  $n$  entre  $n - 1$  en el caso de  $n$  ser grande no significa una gran diferencia y aún menos si tenemos en cuenta que lo que tratamos es de estimar la varianza de la población, no de calcularla de forma exacta.

En cambio, si la muestra es relativamente pequeña (digamos  $n < 30$ ), entonces la varianza muestral de la muestra aproxima significativamente mejor la varianza de la población que la varianza.

La diferencia entre desviación típica y desviación típica muestral es análoga.

Con R, calcularemos la varianza y la desviación típica **muestrales**. Con lo cual, si queremos calcular las que no son muestrales, tendremos que multiplicarlas por  $\frac{n-1}{n}$ , donde  $n$  es el tamaño de la muestra. Lo veremos a continuación.

## Varianza y desviación típica

Nótese que tanto la varianza como la desviación típica dan una información equivalente. Entonces, es comprensible preguntarse por qué se definen ambas medidas si con una basta. Pues bien, las unidades de la varianza (metros, litros, años. . . ), ya sea muestral o no, están al cuadrado, mientras que las de la desviación típica no.

## Medidas de dispersión con R

| Medida de dispersión       | Instrucción                                      |
|----------------------------|--|
| Valores mínimo y máximo    | <code>range(x)</code>                            |
| Rango                      | <code>diff(range(x))</code>                      |
| Rango intercuartílico      | <code>IQR(x, type = ...)</code>                  |
| Varianza muestral          | <code>var(x)</code>                              |
| Desviación típica muestral | <code>sd(x)</code>                               |
| Varianza                   | <code>var(x)*(length(x)-1)/length(x)</code>      |
| Desviación típica          | <code>sd(x)*sqrt((length(x)-1)/length(x))</code> |

## Ejemplo 4

```
datos2
```

```
[1] 6 1 4 1 2 5 3 6 2 3 3 1 5 5 2
```

```
diff(range(datos2))
```

```
[1] 5
```

```
IQR(datos2)
```

```
[1] 3
```

```
var(datos2)
```

```
[1] 3.209524
```



## Ejemplo 4

```
sd(dados2)
```

```
[1] 1.791514
```

```
n = length(dados2)
var(dados2)*(n-1)/n
```

```
[1] 2.995556
```

```
sd(dados2)*sqrt((n-1)/n)
```

```
[1] 1.730767
```

## Función `summary()`

La función `summary` aplicada a un vector numérico o a una variable cuantitativa nos devuelve un resumen estadístico con los valores mínimo y máximo del vector, sus tres cuartiles y su media.

Al aplicar esta función a un data frame, esta se aplica a todas sus variables de forma simultánea. De este modo, podemos observar rápidamente si hay diferencias notables entre sus variables numéricas.

## Ejemplo 5

```
cangrejos = read.table("../data/datacrab.txt", header = TRUE)
cangrejos = cangrejos[-1] #Eliminamos la primera columna
summary(cangrejos) #Aplicamos la función summary
```

| color         | spine         | width        | satell     |
|---------------|---------------|--------------|------------|
| Min. :2.000   | Min. :1.000   | Min. :21.0   | Min. : 0   |
| 1st Qu.:3.000 | 1st Qu.:2.000 | 1st Qu.:24.9 | 1st Qu.: 0 |
| Median :3.000 | Median :3.000 | Median :26.1 | Median : 2 |
| Mean :3.439   | Mean :2.486   | Mean :26.3   | Mean : 2   |
| 3rd Qu.:4.000 | 3rd Qu.:3.000 | 3rd Qu.:27.7 | 3rd Qu.: 5 |
| Max. :5.000   | Max. :3.000   | Max. :33.5   | Max. :15   |

## Ejemplo 5

Si nos interesase comparar numéricamente los pesos y las anchuras de los cangrejos con 3 colores con los que tienen 5 colores, utilizaríamos las siguientes instrucciones:

```
summary(subset(cangrejos, color == 3, c("weight", "width")))
```

| weight       | width        |
|--------------|--------------|
| Min. :1300   | Min. :22.5   |
| 1st Qu.:2100 | 1st Qu.:25.1 |
| Median :2500 | Median :26.5 |
| Mean :2538   | Mean :26.7   |
| 3rd Qu.:3000 | 3rd Qu.:28.2 |
| Max. :5200   | Max. :33.5   |

## Ejemplo 5

```
summary(subset(cangrejos, color == 5, c("weight", "width")))
```

| weight       | width         |
|--------------|---------------|
| Min. :1300   | Min. :21.00   |
| 1st Qu.:1900 | 1st Qu.:23.90 |
| Median :2125 | Median :25.50 |
| Mean :2174   | Mean :25.28   |
| 3rd Qu.:2400 | 3rd Qu.:26.57 |
| Max. :3225   | Max. :29.30   |

Y deducimos así que los cangrejos con 5 colores pesan ligeramente menos y tienen menos anchura que los que tienen 3 colores.

## La función `by()`

La función `by()` se utiliza para aplicar una determinada función a algunas columnas de un data frame segmentándolas según los niveles de un factor.

La sintaxis de esta función es `by(columnas, factor, FUN = función)`.

Con lo cual, haciendo uso de la función `by` y especificando `FUN = summary`, podremos calcular el resumen estadístico anteriormente comentado a subpoblaciones definidas por los niveles de un factor.

## Ejemplo 6

### Ejemplo 6

Para este ejemplo, haremos uso del famoso dataset iris.

Si nos interesase calcular de forma rápida y sencilla las longitudes de sépalos y pétalos en función de la especie, necesitaríamos hacer uso de la instrucción mostrada a continuación.

Por motivos de espacio, no se muestran los resultados proporcionados por R.

```
by(iris[,c(1,3)], iris$Species, FUN = summary)
```

## Función aggregate()

Tanto la función `by` como la función `aggregate` son equivalentes. No obstante, los resultados se muestran de forma diferente en función de cual utilicemos.

En el caso del ejemplo anterior, convenía más hacer uso de la función `by`.

Podéis comprobarlo introduciendo por consola la siguiente instrucción:

```
aggregate(cbind(Sepal.Length,Petal.Length)~Species, data=iris)
```



# NA

La mayoría de las funciones vistas a lo largo de este tema no funcionan bien con valores NA.

Para no tenerlos en cuenta a la hora de aplicar estas funciones, hay que especificar el parámetro `na.rm = TRUE` en el argumento de la función.

## Ejemplo 7

```
datosNA = c(dados2, NA)  
datosNA
```

```
[1] 6 1 4 1 2 5 3 6 2 3 3 1 5 5 2 NA
```

```
mean(datosNA)
```

```
[1] NA
```

```
mean(datosNA, na.rm = TRUE)
```

```
[1] 3.266667
```

## Diagramas de caja

## Diagramas de caja

El conocido diagrama de caja o box plot es un tipo de gráfico que básicamente, remarca 5 valores estadísticos:

- ▶ La mediana, representada por la línea gruesa que divide la caja
- ▶ El primer y tercer cuartil, que son los lados inferior y superior, respectivamente. De este modo, la altura de la caja es el rango intercuantílico
- ▶ Los extremos, los valores  $b_{inf}$ ,  $b_{sup}$ , son los bigotes (whiskers) del gráfico. Si  $m$  y  $M$  son el mínimo y máximo de la variable cuantitativa, entonces los extremos se calculan del siguiente modo:

$$b_{inf} = \max\{m, Q_{0.25} - 1.5(Q_{0.75} - Q_{0.25})\}$$

$$b_{sup} = \min\{M, Q_{0.75} + 1.5(Q_{0.75} - Q_{0.25})\}$$

- ▶ Valores atípicos o outliers, que son los que están más allá de los bigotes. Se marcan como puntos aislados.

## Más sobre los bigotes

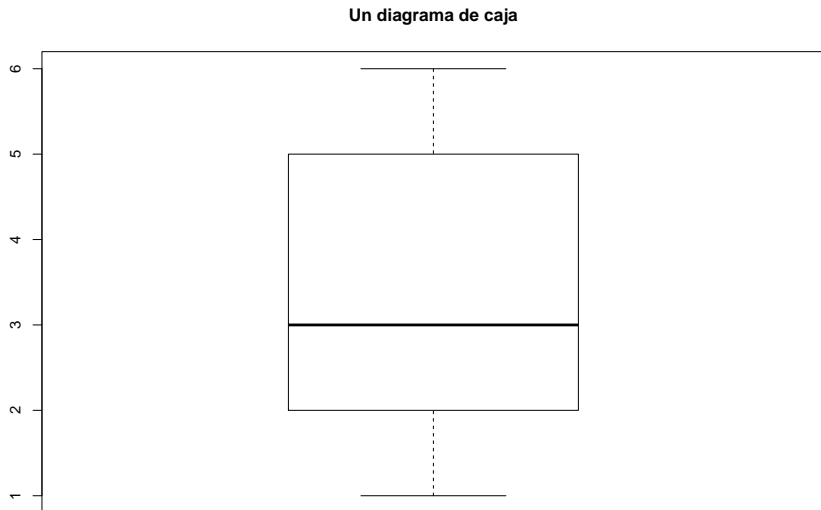
Por su definición, concluimos que los bigotes marcan el mínimo y máximo de la variable cuantitativa, a no ser que haya datos muy alejados de la caja intercuantílica.

En tal caso, el bigote inferior marca el valor 1.5 veces el rango intercuantílico por debajo de  $Q_{0.25}$ , mientras que el superior marca el valor 1.5 veces el rango intercuantílico por encima de  $Q_{0.75}$

## La función boxplot

La instrucción `boxplot()` dibuja diagramas de caja en R.

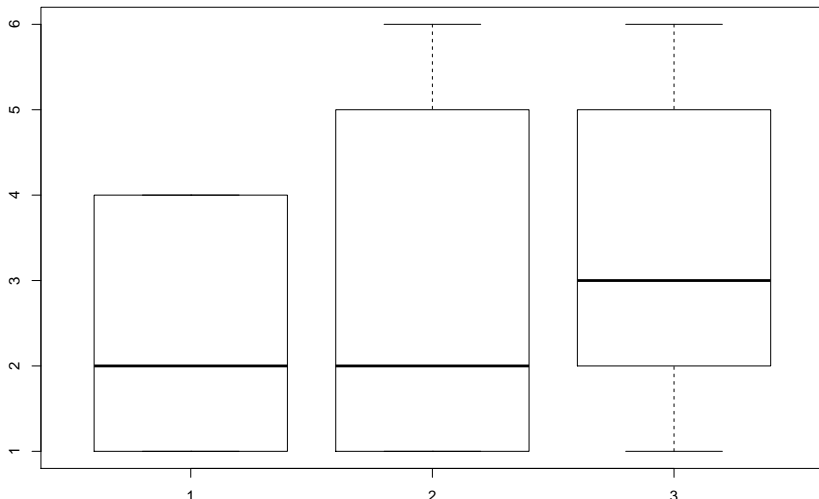
```
boxplot(dados2, main = "Un diagrama de caja")
```



## La función `boxplot`

También podemos dibujar diversos diagramas de caja en un mismo gráfico. De este modo, se pueden comparar con mayor facilidad:

```
boxplot(dado, dados, dados2)
```



## La función `boxplot`

Además, podemos dibujar el diagrama de caja de todas las variables de un data frame en un solo paso aplicando la instrucción `boxplot(data.frame)`.

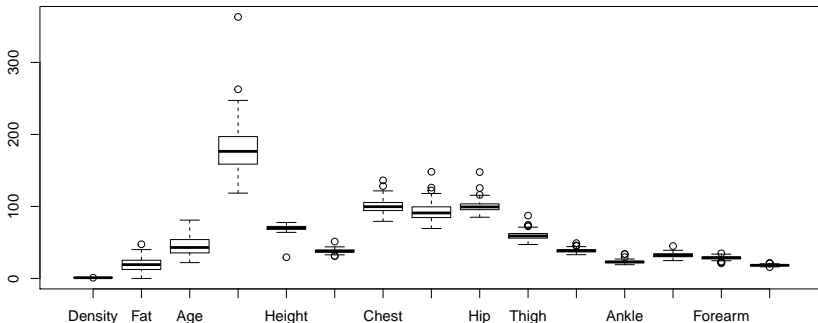
La mayoría de veces, dicho gráfico no será del todo satisfactorio. Dibujar diagramas de factores no tiene sentido alguno. Estos gráficos se pueden manipular incluyendo solo las variables de interés, cambiando los nombres...

Veamos un ejemplo:



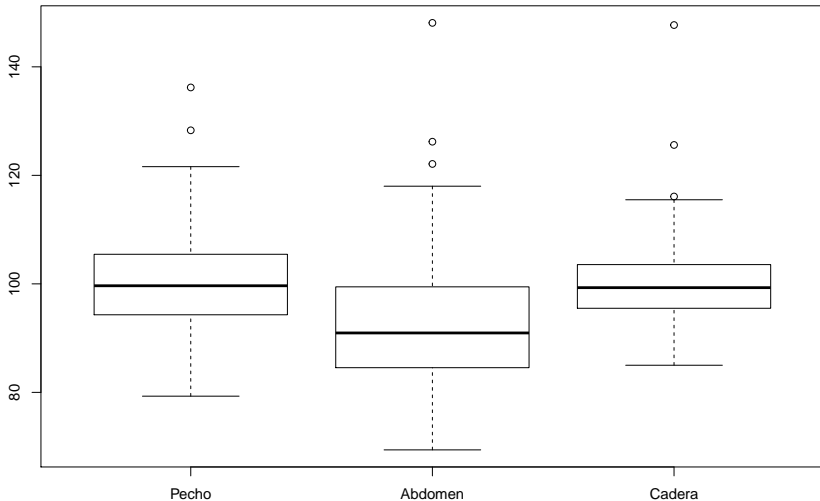
## Ejemplo 8

```
body = read.table("../data/bodyfat.txt", header = TRUE)  
boxplot(body)
```



## Ejemplo 8

```
boxplot(body[,7:9], names = c("Pecho", "Abdomen", "Cadera"))
```



## La función `boxplot`

Agrupar varios diagramas de caja en un solo gráfico tiene por objetivo poder compararlos visualmente, lo cual tiene sentido cuando las variables tienen significados parecidos o cuando comparamos una misma variable de poblaciones distintas.

La mayoría de las veces, queremos comparar diagramas de cajas de una misma variable cuantitativa segmentada por los niveles de un factor.

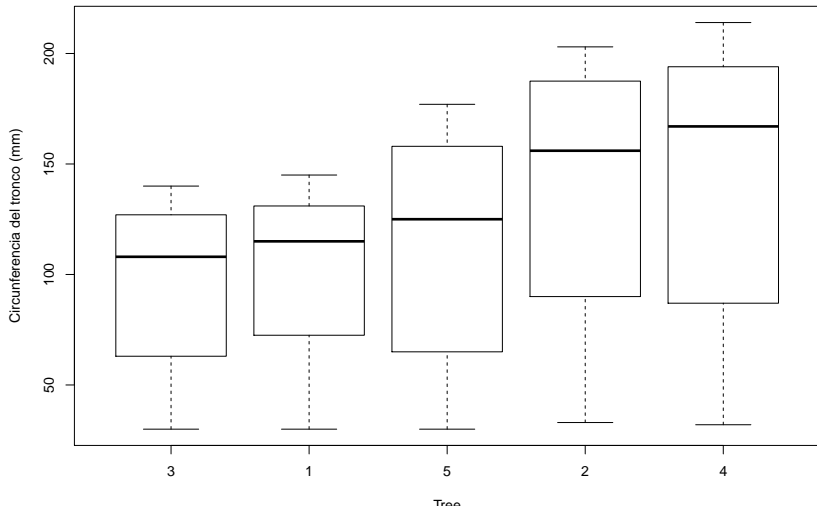
La sintaxis de la instrucción para dibujar en un único gráfico los diagramas de caja de una variable numérica de un data frame en función de los niveles de un factor del mismo data frame es

```
boxplot(var.numérica~factor, data = data frame)
```

## Ejemplo 9

```
boxplot(circumference~Tree, data = Orange, ylab = "Circunfe  
main = "Boxplot de los naranjos en función del tipo
```

Boxplot de los naranjos en función del tipo de árbol



## Parámetros de la función `boxplot`

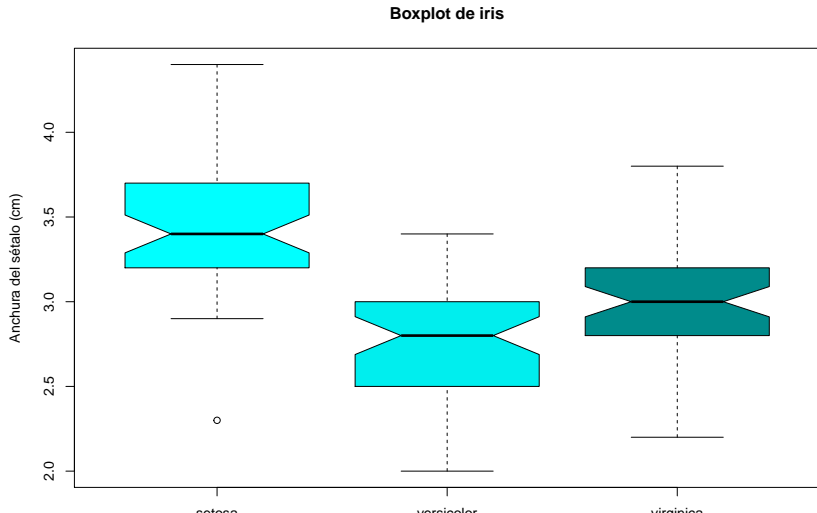
Todos los parámetros de la función `plot()` que tengan sentido pueden ser utilizados en los argumentos de la función `boxplot()`.

Aparte, la función `boxplot()` dispone de algunos parámetros específicos, de los cuales mencionaremos:

- ▶ `notch` igualado a `TRUE` añade una muesca en la mediana de la caja. Si se da el caso en que las muescas de dos diagramas de cajas no se solapan, entonces con alto grado de confianza, concluimos que las medianas de las poblaciones correspondientes son diferentes.

## Ejemplo 10

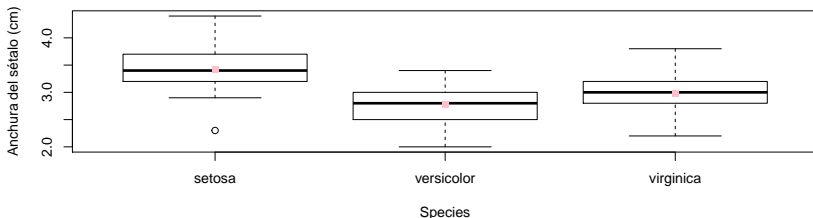
```
boxplot(Sepal.Width~Species, data = iris, ylab = "Anchura c  
notch = TRUE, col = c("cyan","cyan2","cyan4"),  
main = "Boxplot de iris")
```



## Ejemplo 10

Si quisiéramos marcar de alguna forma en un diagrama de caja, cosa que puede ser muy útil en ocasiones, la media aritmética de la variable correspondiente, podríamos hacerlo mediante la función `points`:

```
boxplot(Sepal.Width~Species, data = iris, ylab = "Anchura de  
medias = aggregate(Sepal.Width~Species, data = iris, FUN =  
points(medias, col = "pink", pch = 15)
```



## Ejemplo 10

La primera instrucción del chunk anterior genera el diagrama de cajas de las anchuras de los sépalos en función de la especie. Por su parte, la segunda instrucción lo que hace es calcular las medias aritméticas de las anchuras según la especie. Finalmente, la tercera instrucción lo que hace es añadir al diagrama un punto cuadrado a cada caja en la ordenada correspondiente a su media aritmética.



## La estructura interna de boxplot

Como ya sabemos, podemos estudiar la función interna de algunos objetos con la función `str`.

Dicha función aplicada a un boxplot, nos produce una list. Podéis ver esta list si introducís por consola la siguiente instrucción:

```
str(boxplot(circumference~Tree, data = Orange))
```

Destacaremos dos de sus componenetes aquí:

- ▶ `stats` nos devuelve los valores  $b_{inf}$ ,  $Q_{0.25}$ ,  $Q_{0.5}$ ,  $Q_{0.75}$ ,  $b_{sup}$
- ▶ `out` nos retorna los valores atípicos. En caso de haber diversos diagramas en un plot, la componente `group` nos indica a qué diagramas pertenecen estos outliers.