

ANÁLISIS DE DATOS EN INVESTIGACIÓN BIOMÉDICA **MÉTODOS ESTADÍSTICOS**

Dra. Mariana Díaz-Almirón



Investigación Biomédica

1. Recordar definición de variables
2. Depuración de los datos
3. Toma de contacto con la consola de R
4. Estadística descriptiva

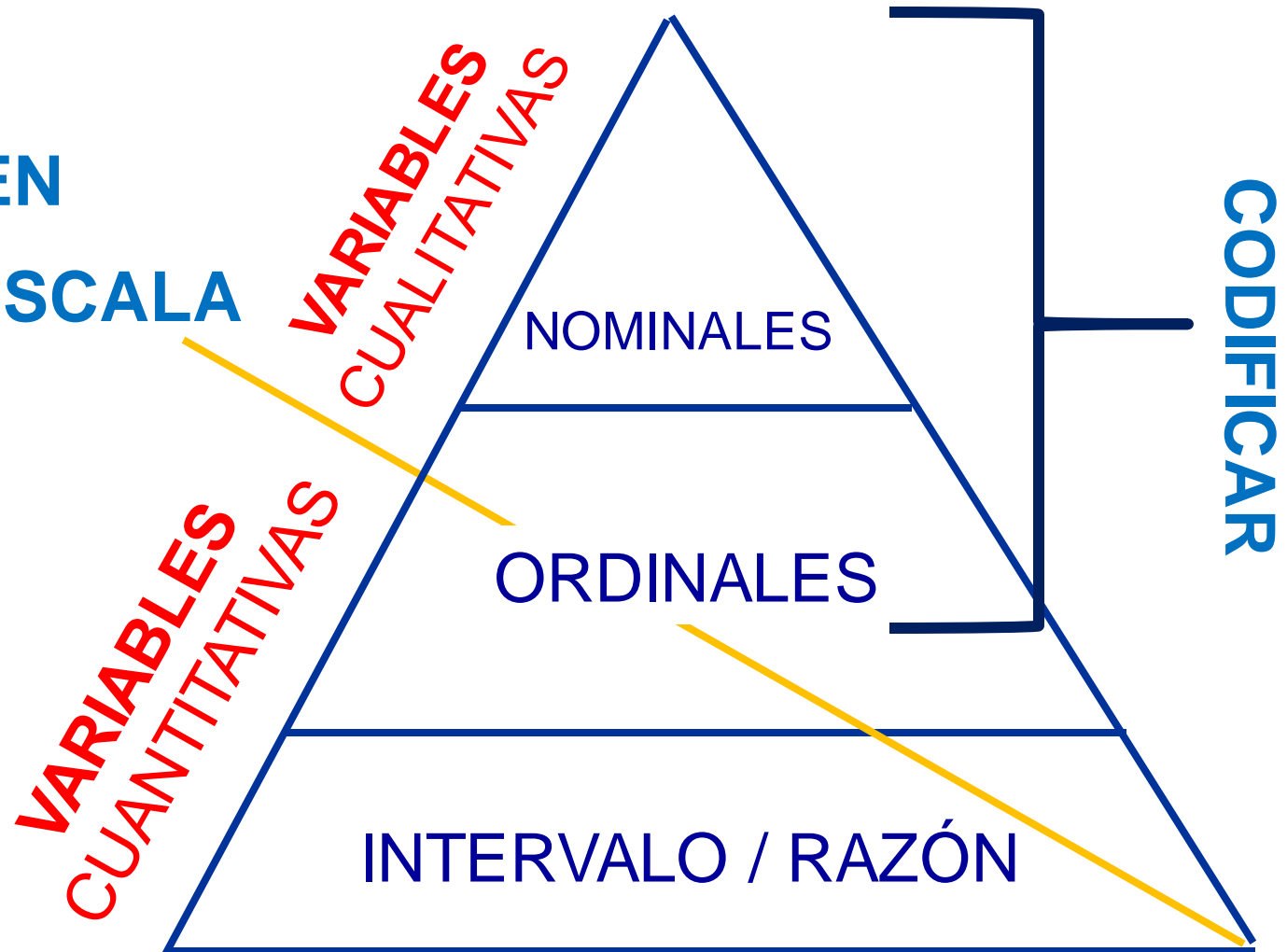
Investigación Biomédica

1. Recordar definición de variables
2. Depuración de los datos
3. Toma de contacto con la consola de R
4. Estadística descriptiva

Recordar definición de variables

RESUMEN

TIPO ~ ESCALA



Investigación Biomédica

1. Recordar definición de variables
2. Depuración de los datos
 - Detección de errores de transcripción
 - Detección de “outliers”,...
3. Toma de contacto con la consola de R
4. Estadística descriptiva

Investigación Biomédica

1. Recordar definición de variables
2. Depuración de los datos
3. Toma de contacto con la consola de R
4. Estadística descriptiva



- CREAR UN CONJUNTO DE DATOS
- DEFINIR LAS VARIABLES
- CREAR (CALCULAR) UNA VARIABLE
- CREAR UNA VARIABLE A PARTIR DE OTRA
- SELECCIONAR CASOS

VARIABLES

Variable	Significado	Códigos	Escala
DATOS PERSONALES			
ID	Nº de Identificación		
SEXO	Sexo	0=Mujer 1=Hombre	Nominal
EDAD	Edad		Continua
TRATAMIENTO	Tratamiento aleatorizado	0=Referencia 1=Nuevo	Nominal
TALLA	Talla basal		Continua
ACTIVIDAD	actividad física	1 'leve' 2 'moderada' 3 'severa'	Ordinal
MOTIVACION	Score de Motivación (1-12, siendo 12 lo mejor)		Ordinal
ANTECEDENTES PERSONALES			
HTA	Hipertensión	0=No 1=Sí	Nominal
DISLIPEMIA	Dislipemia	0=No 1=Sí	Nominal
GLUCEMIA	Alteración de la Glucemia	0=No 1=Sí	Nominal
DE	Disfunción Endotelial	0=No 1=Sí	Nominal
LDLox	Altos niveles de LDL Oxidada	0=No 1=Sí	Nominal
FUMA	Fuma	0=No 1=Sí	Nominal
CCI	Índice Cintura-Cadera Elevado	0=No 1=Sí	Nominal
DATOS BASALES			
P0	Peso basal		Continua
TG0	Triglicéridos basal		Continua
ALBUMINA0	Albumina basal		Continua
MGRASA0	Masa Grasa basal (kg)		Continua
GLUCOSA0	Glucosa basal		Continua
INSULINA0	Insulina basal		Continua
AcURICO0	Ácido Úrico basal		Continua

DATOS EN SUCESIVAS CONSULTAS			
P1	Peso en la segunda consulta		Continua
P2	Peso final		Continua
SEGUIMIENTO	Tiempo desde P0 a P2		Continua
TCAMBIO	Tiempo hasta perder al menos 10K		Continua
TG2	Triglicéridos final		Continua
CCI2	Índice Cintura-Cadera Elevado final	0=No 1=Sí	Nominal
ALBUMINA2	Albumina final		Continua
MGRASA2	Masa Grasa final (kg)		Continua
GLUCOSA2	Glucosa final		Continua
INSULINA2	Insulina final		Continua
AcURICO2	Ácido Úrico final		Continua
SATISFFINAL	Satisfacción final		Ordinal



Usamos R (R Gui / Rstudio / TinnR) o un compilador de R que ejecute el código:

<https://rdr.io/snippets/>

Para más información del software, un manual de R se puede descargar aquí:

- <https://www.dropbox.com/s/293cw299i0lb32x/R%20para%20principiantes.pdf?dl=0>
- <https://javieralvarezliebana.es/curso-intro-R/>



- **CREAR UN CONJUNTO DE DATOS**
- CREAR (CALCULAR) UNA VARIABLE
- DEFINIR LAS VARIABLES
- CREAR UNA VARIABLE A PARTIR DE OTRA
- SELECCIONAR CASOS

DATOS

ESTRUCTURA

```
datos = data.frame(  
  ID=c(1, 2, 3, 4, ...),  
  F_INI=c(42917, 42989, ...),  
  F_FIN=c(43049, 43116, ...),  
  F_10K=c(NA, 43116, NA...),  
  SEXO=c(0, 1, 0, ...),  
  EDAD=c(49, 48, 22, ...),  
  TALLA=c(165.83, 190.94, ...),  
  ACTIVIDAD=c(1, 1, 1...),  
  TRATAMIENTO=c(1, 0, 0, ...),  
  MOTIVACION=c(4, 4, 4, 9...),  
  HTA=c(1, 0, 0...),  
  DISLIPEMIA=c(1, 1, ...),  
  GLUCEMIA=c(0, 1, 1...),  
  DE=c(0, 1, 1...0),  
  LDLox=c(1, 1, 1, ...),  
  FUMA=c(1, 1, 1,...),  
  CCI=c(1, 1, 1, ...),  
  TG0=c(82.61, 76.41, 77.81, ...),  
  TG2=c(76.52, 72.59, 72.63, ...),  
  ALBUMINA0=c(4.06, 4.75, 4.74,...),  
  ALBUMINA2=c(2.06, 4.17, 4.3, ...),  
  MGRASA0=c(38.31, 36.52, 27.96, ...),  
  MGRASA2=c(26.84, 34.69, 25.25, ...),  
  GLUCOSA0=c(54.91, 95.52, 98.97,...),  
  GLUCOSA2=c(45.93, 93.06, 95.9, ...),  
  INSULINA0=c(33.33, 7.53, 20.66, ...),  
  INSULINA2=c(22.71, 3.82, 19.54,...),  
  AcURICO0=c(14.86, 3.68, ...),  
  AcURICO2=c(1.89, 1.02, 5.63, ...),  
  P0=c(97.53, 103.81, 103.91, ...),  
  P1=c(94.905, 89.31,...),  
  P2=c(92.28, ...),  
  SATISFFINAL=c(5, 8, 5...)  
)  
head(datos)
```

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

[illegible]



- CREAR UN CONJUNTO DE DATOS
- **CREAR (CALCULAR) UNA VARIABLE**
- CREAR UNA VARIABLE A PARTIR DE OTRA
- DEFINIR LAS VARIABLES
- SELECCIONAR CASOS



CREAR (CALCULAR) UNA VARIABLE

nº de antecedentes personales totales (NAT)

$$\text{NAT} = \{\text{HTA} + \text{DISLIPEMIA} + \text{GLUCEMIA} + \text{DE} + \text{LDLox} + \text{FUMA} + \text{CCI}\}$$

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

```
# Creamos: NAT= {HTA + DISLIPEMIA + GLUCEMIA + DE + LDLox + FUMA + CCI}
```

```
datos$NAT= datos$HTA + datos$DISLIPEMIA + datos$GLUCEMIA +  
datos$DE + datos$LDLox + datos$FUMA + datos$CCI
```

```
head(datos)
```



- CREAR UN CONJUNTO DE DATOS
- CREAR (CALCULAR) UNA VARIABLE
- **CREAR UNA VARIABLE A PARTIR DE OTRA**
- DEFINIR LAS VARIABLES
- SELECCIONAR CASOS

CREAR UNA VARIABLE A PARTIR DE OTRA

Codifica la variable NAT en otra variable que sea NAT2c, de tal manera que:

$$\text{NAT2c} = \{1='0-2' \text{ y } 2='=>3'\}$$

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

```
# Creamos: NAT2c = {1='0-2', 2='=>3'}
```

```
datos$NAT2c= ifelse( datos$NAT<=2, 1, 2)
```

```
head(datos)
```



- CREAR UN CONJUNTO DE DATOS
- CREAR (CALCULAR) UNA VARIABLE
- CREAR UNA VARIABLE A PARTIR DE OTRA
- **DEFINIR LAS VARIABLES**
- SELECCIONAR CASOS

DATOS

Con la ayuda del archivo que contiene la definición de variables, vamos a definir nuestra BD

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

```
# Convertimos en formato fecha los valores numéricos obtenidos del excel  
datos$F_INI = as.Date(datos$F_INI, origin = "1899-12-30")
```

```
# Convertimos las variables nominales y añadimos las etiquetas  
datos$SEXO = factor(datos$SEXO, levels = c(0,1), labels = c('Mujer', 'Hombre'),  
ordered = FALSE)
```

Fíjate bien que los valores de las etiquetas llevan comillas

```
# Las variables continuas no hay que definirlas
```



- CREAR UN CONJUNTO DE DATOS
- CREAR (CALCULAR) UNA VARIABLE
- CREAR UNA VARIABLE A PARTIR DE OTRA
- DEFINIR LAS VARIABLES
- **SELECCIONAR CASOS**

SELECCIONAR CASOS

Seleccionamos los pacientes tratados con el nuevo tratamiento.

- Si no hemos puesto las etiquetas, la codificación es 1 (valor numérico), en otro caso, la etiqueta es 'Referencia' (texto).

El resultado será la nueva base solo con esos pacientes.

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos con las etiquetas:

```
# Creamos la sub-base datos2 cuyos pacientes solo recibieron el tratamiento de referencia.
```

```
datos2= datos[datos$TRATAMIENTO=='Referencia', ]  
head(datos2)
```

```
#datos2= datos[which(datos$TRATAMIENTO=='Referencia'), ]  
#head(datos2)
```

fíjate bien que Referencia lleva comillas porque es una etiqueta, no un número.

Copiamos datos sin las etiquetas:

```
# Creamos la sub-base datos2 cuyos pacientes solo recibieron el tratamiento de referencia.
```

```
datos2= datos[datos$TRATAMIENTO==1,]  
head(datos2)
```

```
#datos2= datos[which(datos$TRATAMIENTO==1),]  
#head(datos2)
```

Investigación Biomédica

1. Recordar definición de variables
2. Depuración de los datos
3. Toma de contacto con la consola de R
4. Estadística descriptiva

Investigación Biomédica

Estadística descriptiva

Métodos estadísticos **DESCRIPTIVOS**:

Sirven para **describir y resaltar numéricamente** aquello que es esencial en los resultados de un estudio usando los métodos apropiados

Métodos estadísticos **INFERENCIALES**:

Conjunto de **métodos** que permiten **obtener conclusiones**, de una manera objetiva, sobre los datos que se están investigando.

Investigación Biomédica

Recordatorio. Medidas descriptivas

TIPO DE VARIABLE	GRÁFICA	MEDIDAS CENTRALES	MEDIDAS DE DISPERSIÓN
NOMINAL	<ul style="list-style-type: none"> •DIAG. BARRAS •DIAG. SECTORES 	MODA (%)	
ORDINAL	<ul style="list-style-type: none"> •DIAG. BARRAS •BOX-PLOT 	MEDIANA	RANGO $P_{75}-P_{25}$
INTERVALO / RAZÓN	<ul style="list-style-type: none"> •BOX-PLOT •HISTOGRAMAS 	MEDIA	DESV. TÍPICA

Investigación Biomédica

Estadística descriptiva

Variable categórica

- ✓ *Frecuencias absolutas*
- ✓ *Frecuencias relativas (Porcentajes)*

Variable continua

- ✓ *Medidas de centralización y orden*
- ✓ *Medidas de dispersión*

Investigación Biomédica

Estadística descriptiva

Variable categórica

- ✓ *Frecuencias absolutas*
- ✓ *Frecuencias relativas (Porcentajes)*

Variable continua

- ✓ *Medidas de centralización y orden*
- ✓ *Medidas de dispersión*

Investigación Biomédica

Estadística descriptiva

Descripción numérica - cualitativas

- Frecuencias absolutas

Contabilizan el número de individuos de cada atributo

- Frecuencias relativas (%)

Ídem, pero se divide por el total

Investigación Biomédica

Estadística descriptiva

Descripción gráfica - cualitativas

- Diagramas de barras

Alturas proporcionales a las frecuencias (absolutas o relativas).

Se pueden aplicar también a variables discretas.

- Diagramas de sectores (tartas, polares)

No usarlo con variables ordinales.

El área de cada sector es proporcional a su frecuencia (abs. o rel.).

DESCRIPCIÓN NUMÉRICA

Describir la variable ÉXITO (pacientes que perdieron > 10 Kg, PCAMBIO).

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

Creamos las frecuencias absolutas y relativas en %

```
frec = table(datos$PCAMBIO)
```

```
frec
```

```
frecRel = 100*table(datos$PCAMBIO)/length(datos$PCAMBIO)
```

```
frecRel
```

No	Si
10	30

→ Hace referencia a la frecuencia absoluta

No	Si
25	75

→ Hace referencia al %

DESCRIPCIÓN GRÁFICA

Describir la variable ÉXITO (pacientes que perdieron > 10 Kg, PCAMBIO).

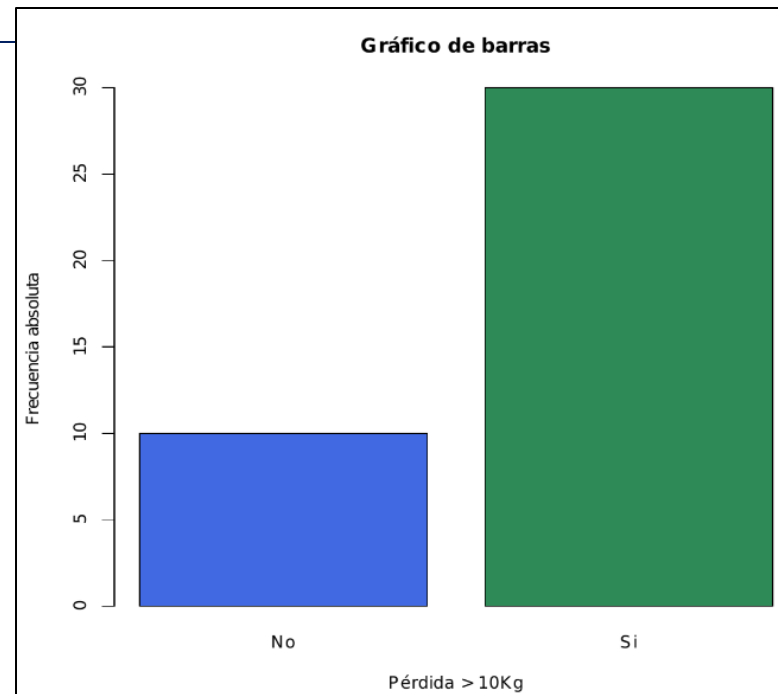
DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

Creamos el gráfico de barras

```
plot(x = datos$PCAMBIO, main = "Gráfico de barras",  
     xlab = "Pérdida > 10Kg", ylab = "Frecuencia absoluta",  
     col = c("royalblue", "seagreen"))
```



DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

```
# Creamos el gráfico de sectores
```

```
frecRel = 100*table(datos$PCAMBIO)/length(datos$PCAMBIO)
```

```
pie(x = frecRel, main = "Gráfico de sectores",
```

```
  xlab = "Pérdida >10Kg", ylab = "%",
```

```
  col = c("royalblue", "seagreen"))
```



Investigación Biomédica

Estadística descriptiva

Variable categórica

- ✓ *Frecuencias absolutas*
- ✓ *Frecuencias relativas (Porcentajes)*

Variable continua

- ✓ *Medidas de centralización y orden*
- ✓ *Medidas de dispersión*

Descripción - cuantitativas

Orden

Dividen un conjunto ordenado de datos en grupos con la misma cantidad de individuos.

Cuantiles, Percentiles, Cuartiles, Deciles,...

Centralización

Indican valores con respecto a los que los datos parecen agruparse.

Media, Mediana y Moda

Dispersión

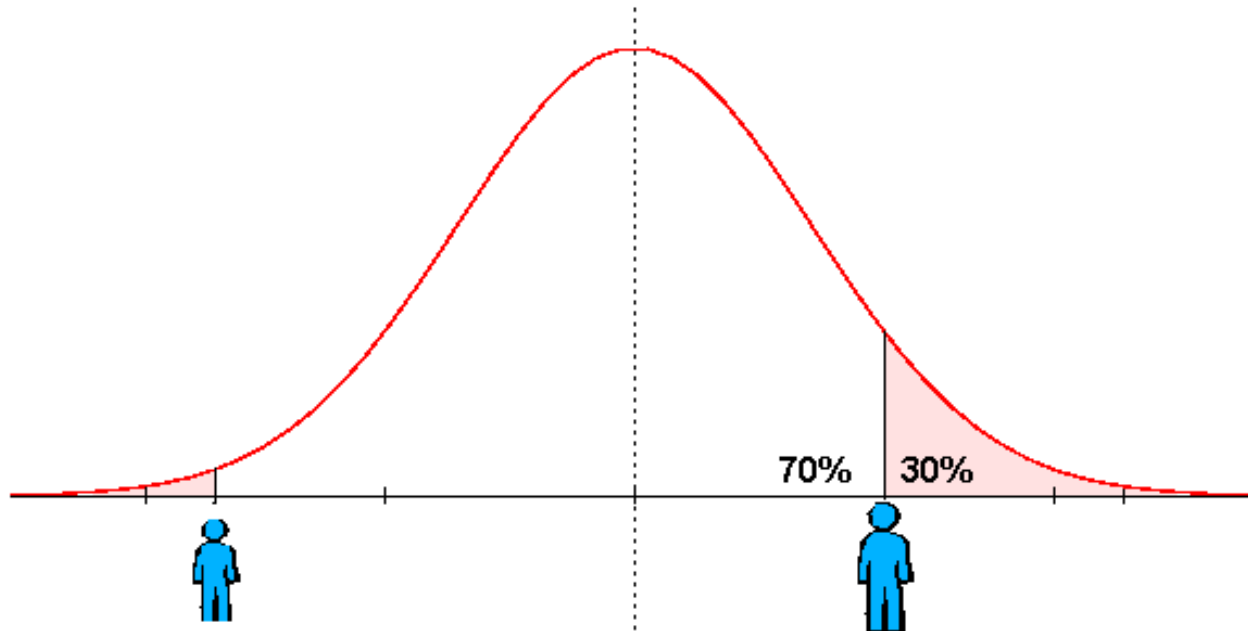
Indican la mayor o menor concentración de los datos con respecto a las medidas de centralización.

Varianza, Desviación típica (st), Coeficiente de Variación (CV), Rango, Rango Inter cuartílico (RI)

Descripción - cuantitativas

Orden - Cuantil

Se define el cuantil de orden α como un valor de la variable por debajo del cual se encuentra una frecuencia acumulada α



Descripción - cuantitativas

Orden - Cuantil

✓ **Cuartiles**, Q_k con $k=1, 2$ ó 3

Dividen a la muestra en 4 grupos con frecuencias similares.

$$Q_K : \frac{K * (N + 1)}{4}, K = 1, 2, 3$$

- Q_1 = Cuantil 0,25 (Percentil 25, P_{25})
- Q_2 = Cuantil 0,5 = mediana (Percentil 50, P_{50})
- Q_3 = Cuantil 0,75 (Percentil 75, P_{75})

Descripción - cuantitativas

Centralización

Son medidas que buscan posiciones (valores) con respecto a los cuales los datos muestran tendencia a agruparse.

Media

Promedio de los valores de una variable. Suma de los valores dividido por el tamaño muestral.

Conveniente cuando los datos se concentran simétricamente con respecto a ese valor.

Muy sensible a valores extremos.

$$\bar{x} = \frac{\sum_i x_i}{n}$$

Mediana

Valor que divide a las observaciones en dos grupos con el mismo número de individuos (cuartil 2). Si el número de datos es par, se elige la media de los dos datos centrales.

Es conveniente cuando los datos son asimétricos.

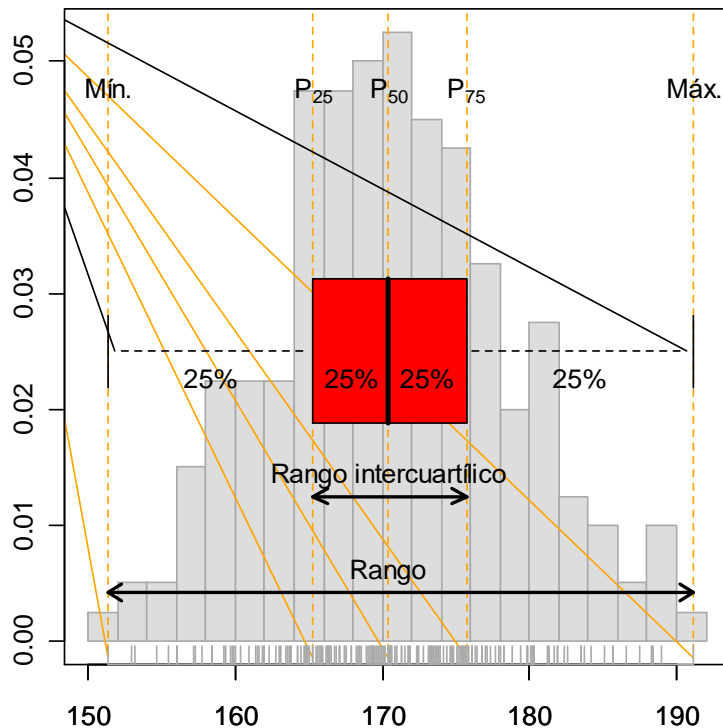
No es sensible a valores extremos.

Moda Valor/es donde la distribución de frecuencia alcanza un máximo.

Descripción - cuantitativas

Dispersión

Miden el grado de dispersión (variabilidad) de los datos, independientemente de su causa.



Amplitud o Rango

Máximo - Mínimo

Es muy sensible a los valores extremos.

Rango intercuartílico

$Q_3 - Q_1$

No es tan sensible a valores extremos.

Descripción - cuantitativas

Dispersión

Varianza muestral S^2

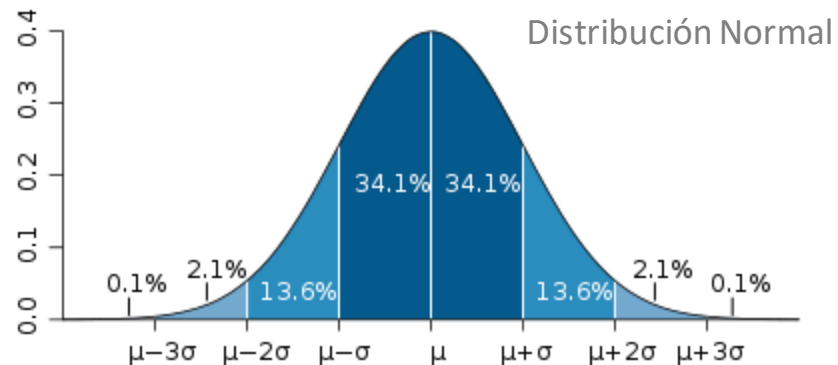
Mide el promedio de las desviaciones (al cuadrado) de las observaciones con respecto a la media.

Es sensible a valores extremos (alejados de la media).

$$S^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$$

Desviación típica S

$$S = \sqrt{S^2}$$



Centrado en la media (μ) y a 1 dt (σ) de distancia hay aproximadamente el 68% de las observaciones. A dos desviaciones típicas tenemos el 95% (aprox.)

Descripción - cuantitativas

Dispersión

Coeficiente de Variación CV

Es la razón entre la desviación típica y la media.

- Es frecuente mostrarla en porcentajes.

Si la media es 80 y la desviación típica 20 entonces $CV=20/80=0,25=25\%$ (variabilidad relativa).

- Interesante para comparar la variabilidad de diferentes variables.

Si el peso tiene $CV=30\%$ y la altura tiene $CV=10\%$, los individuos presentan más dispersión en peso que en altura.

- No debe usarse cuando la variable presenta valores negativos.

$$CV = \frac{S}{\bar{x}}$$

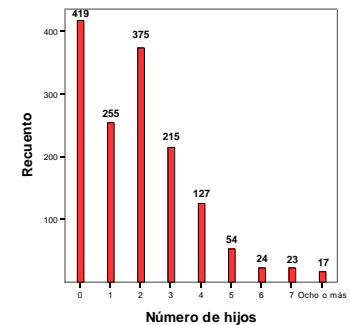
Descripción gráfica - cuantitativas

■ Histograma

Alturas proporcionales a las frecuencias (abs. o rel.)

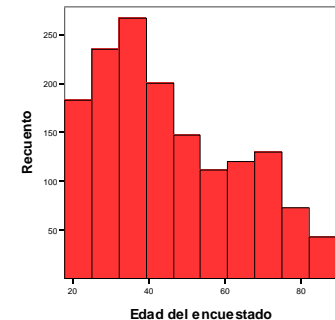
Diagramas barras para v. discretas

Se deja un hueco entre barras para indicar los valores que no son posibles.



Histogramas para v. continuas

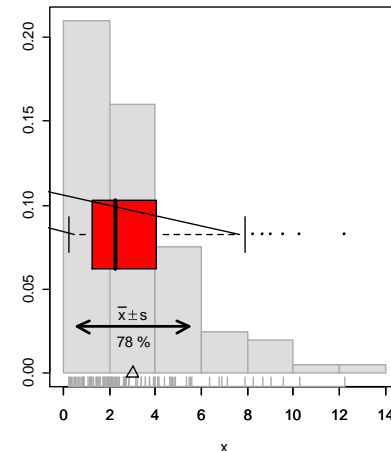
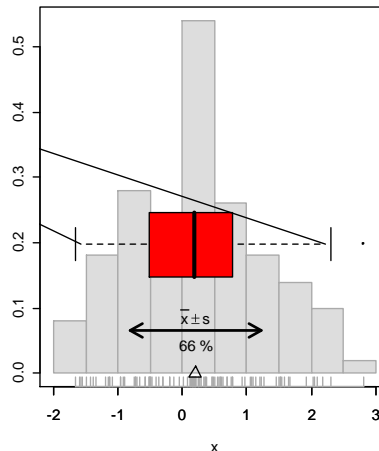
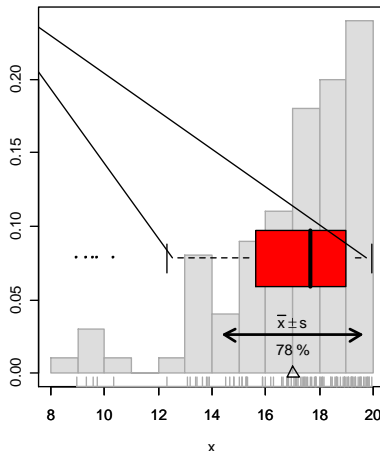
El área que hay bajo el histograma entre dos puntos cualesquiera indica la cantidad (% o frecuencia) de individuos en el intervalo.



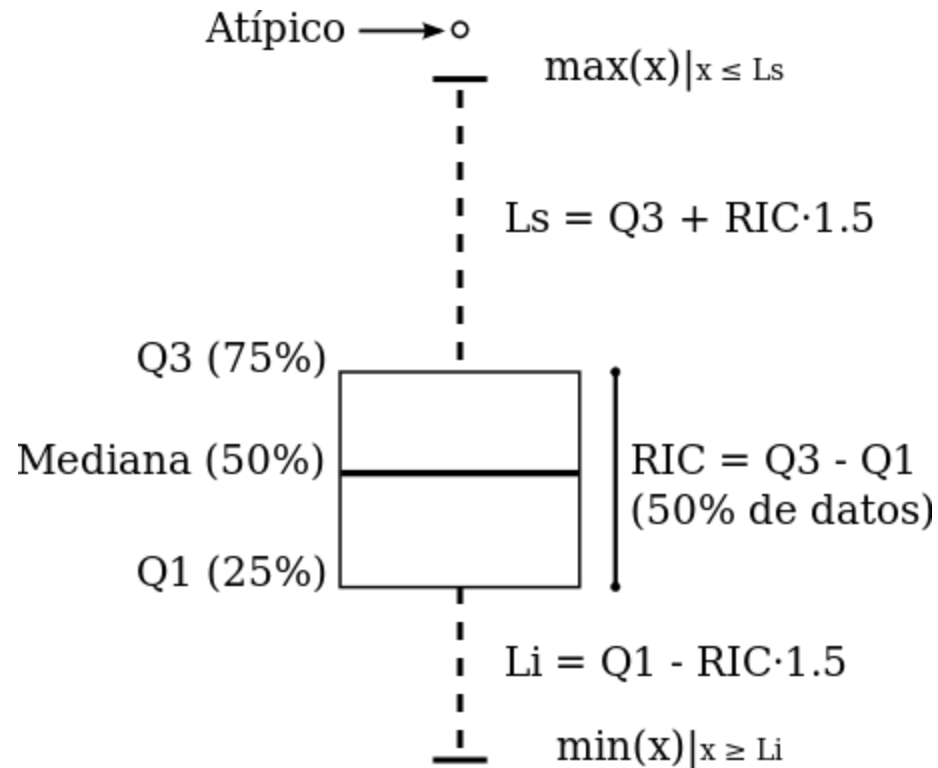
Descripción gráfica - cuantitativas

■ Box-Plot (diagrama de cajas y bigotes)

- Buena aproximación de la distribución.
- La zona central, 'caja', contiene al 50% central de las observaciones. Su tamaño se llama 'rango intercuartílico' (RI)
- Es costumbre que 'los bigotes', no lleguen hasta los extremos, sino hasta las observaciones que se separan de la caja en no más de 1,5 RI.
- Más allá de esa distancia se consideran anómalas.



- Box-Plot (diagrama de cajas y bigotes)



DESCRIPCIÓN NUMÉRICA

Describir la variable P0 (peso basal).

DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

#media

mean(datos\$P0)

#desviación típica

sd(datos\$P0)

#varianza

var(datos\$P0)

#CV

100*sd(datos\$P0)/mean(datos\$P0)

#mediana

median(datos\$P0)

rango intercuartílico

IQR(datos\$P0)

rango

max(datos\$P0)-min(datos\$P0)

Media	$\bar{x} = \frac{\sum_i x_i}{n}$	106,78
D.t	$S^2 = \frac{1}{n} \sum_i (x_i - \bar{x})^2$	6,86
Varianza	$S = \sqrt{S^2}$	47,10
CV	$CV = \frac{S}{\bar{x}}$	6,42
Mediana	Q_2	107,54
RI	$Q_3 - Q_1$	7,02
Rango	Max – Min	31,80

DESCRIPCIÓN GRÁFICA

Describir la variable P0 (peso basal).

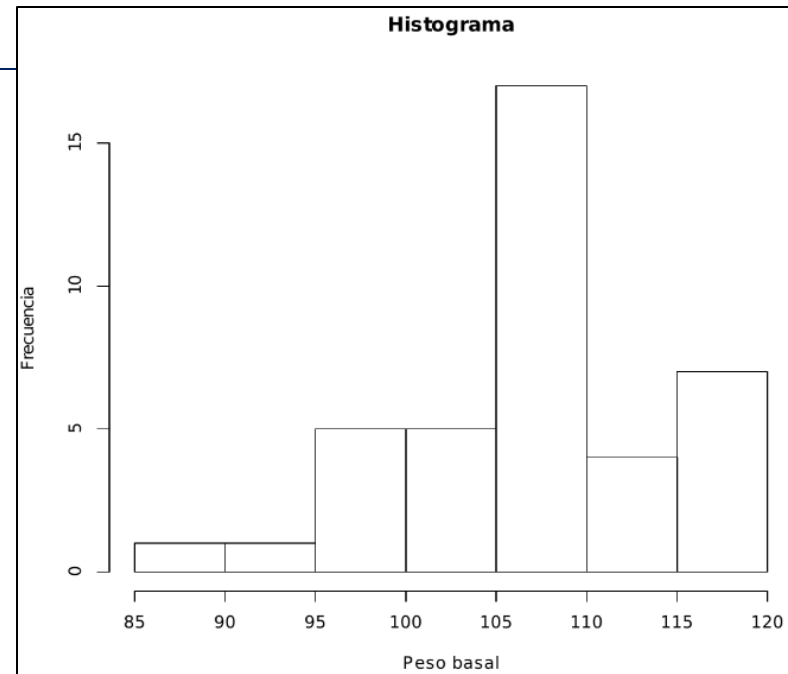
DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

Creamos el histograma

```
hist(x = datos$P0, main = "Histograma",  
     xlab = "Peso basal", ylab = "Frecuencia")
```

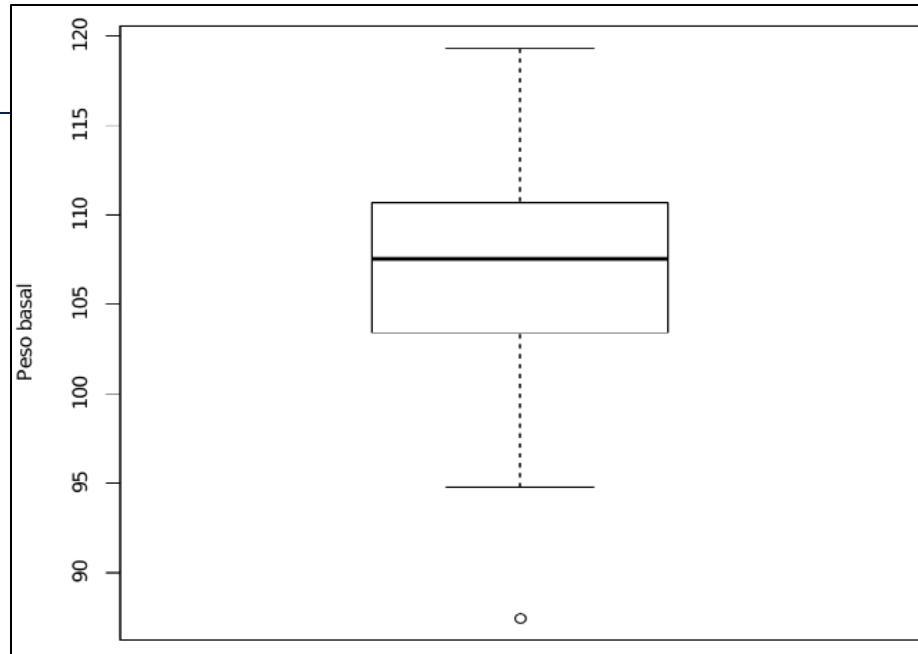


DATOS

Sintaxis (copiar el texto entero y pegar en el compilador):

Copiamos datos

```
# Creamos el diagrama de cajas - boxplot  
boxplot(datos$P0, ylab = "Peso basal")
```





THANK YOU !

mariana.diaz@salud.madrid.org

BIOESTADÍSTICA