# ESTADÍSTICA DESCRIPTIVA

# 1 Conceptos básicos:

**Población:** Conjunto de individuos que presentan una o varias características en común. Esto es, el conjunto objeto de estudio. Por ejemplo:

- Todos los habitantes de una determinada ciudad.
- Las piezas fabricadas por una misma máquina.
- Todos los enfermos de una misma enfermedad.

#### ¿Cómo podemos estudiar las poblaciones?

- mediante un **censo o estudio exhaustivo** que consiste en observar todos y cada uno de los individuos que integran la población.
- por **muestreo** que consiste en estudiar un subconjunto representativo de la población que se llama **muestra**. Se suele considar una muestra porque no siempre es posible estudiar exhaustivamente toda la población por motivos de tiempo, coste económico u otro tipo de dificultad.

<u>Característica</u>: Propiedad que deseamos observar entre los elementos de la población. Los diferentes estados o valores que presenta una característica se suelen llamar modalidades de la característica.

Atendiendo a la característica que se estudia, ésta puede clasificarse en:

- característica cualitativa, categórica o atributo: representa una cualidad del individuo.
- característica cuantitativa: aquellas característica que toman valores numéricos. A las características cuantitativas también se les llaman variables estadísticas y se dividen en:
  - variables estadísticas discretas: aquellas que toman un número finito o infinito numerable de valores.
  - variables estadísticas continuas: aquellas que toman un número no numerable de valores.

**Ejemplo:** Para los habitantes de una determinada ciudad se pueden estudiar las características: sexo, estado civil, profesión, edad, estatura, nivel de estudios,....

Una vez que hemos clasificado la característica que estudiamos, el siguiente paso es ordenar y presentar los datos en tablas y gráficos con el fin de resumir la información que contienen.



Techos Allend Mai

# 2 Distribución de frecuencias:

Supongamos que tenemos una muestra de tamaño n (que puede tomar m ( $m \le n$ ) modalidades o valores):

$$x_1, x_2, ..., x_n$$

Frecuencia absoluta de la modalidad  $x_i$ ,  $n_i$ : número de veces que aparece  $x_i$  en la muestra. Se verifica que:

$$n_1 + n_2 + \dots + n_m = n$$

Frecuencia relativa de la modalidad  $x_i$ ,  $f_i$ : proporción de veces que aparece  $x_i$  en la muestra, esto es

$$f_i = \frac{n_i}{n}$$

Trivialmente se verifica que:

$$f_1 + f_2 + \dots + f_m = 1$$

Se suele presentar en porcentaje sin más que multiplicar por 100, esto es  $100 \times f_i\%$ .

Los valores que toma la característica, junto con las frecuencias de dichos valores se suelen presentar en una tabla que se llama distribución de frecuencias o tabla de frecuencias:

valor, $x_i$	frecuencia absoluta, $n_i$	frecuencia relativa, $f_i$
$x_1$	$n_1$	$f_1$
$x_2$	$n_2$	$f_2$
÷		
$x_m$	$n_m$	$f_m$
totales	n	1

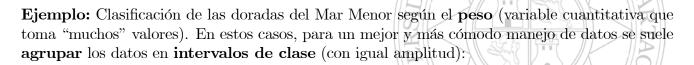
Ejemplo: Clasificación de las doradas del Mar Menor según su sexo (atributo)

sexo	frecuencia absoluta, $n_i$	frecuencia relativa, $f_i$
$\overline{m}$	10	$10/27 \simeq 0.37$
h	17	$17/27 \simeq 0.63$
totales	27	1

**Ejemplo:** Clasificación de las doradas del Mar Menor según la **longitud** (variable cuantitativa que toma "pocos" valores)

longitud	fr. absoluta, $n_i$	fr. relativa, $f_i$
1	6	0.22
2	5	0.19
3	6	0.22
4	4	0.15
5	6	0.22
totales	27	1





- ullet Número de intevalos de clase k
  - El número entero más próximo por exceso a  $\sqrt{n}$ . Además k debe verificar que:  $5 \le k \le 20$ .
  - Regla de Sturges:

$$k = 1 + log_2 n$$

- Amplitud =  $h \equiv \frac{x_{\text{max.}} x_{\text{min.}}}{k}$
- ullet Los datos deben de clasificarse sin ambigüedad en única clase  $\Longrightarrow$

$$[e_0,e_1], (e_1,e_2],....,(e_{k-1},e_k]$$

o bien,

$$[e_0,e_1), [e_1,e_2), ...., [e_{k-1},e_k]$$

Conceptos que aparecen con el tratamiento de los datos en clases:

Clase i-ésima:  $(e_{i-1}, e_i]$ 

Frecuencia absoluta de la clase i-ésima,  $n_i$ : número de observaciones que caen dentro de  $I_i$ .

Frecuencia relativa de la clase i-ésima,  $f_i \equiv \frac{n_i}{n}$ .

Frecuencia absoluta acumulada de la clase i-ésima  $N_i :\equiv n_1 + n_2 + ... + n_i$ 

Frecuencia relativa acumulada la clase i-ésima,  $F_i :\equiv f_1 + f_2 + ... + f_i = \frac{N_i}{n}$ .

Marca de la clase i-ésima:  $m_i \equiv \frac{e_{i-1} + e_i}{2}$ , como representante de todas las observaciones de la clase i-ésima.

Y estos valores también se incorporan en la tabla de frecuencias:



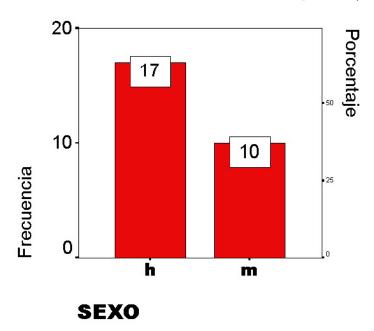
peso	marcas de	frecuencia	fr. absoluta	frecuencia fr. relativa
	clase, $m_i$	absoluta, $n_i$	acumulada, $N_i$	relativa, $f_i$ acumulada, $F_i$
[0.5, 1.92]	1.21	6	6	0.22 0.22
(1.92, 3.34]	2.63	5	11   函	0.19
(3.34, 4.76]	4.05	2	13	0.07
(4.76, 6.18]	5.47	4	17	0.15
(6.18, 7.6]	6.89	4	21	0.15 0.78
(7.6, 9.02]	8.31	6	27	0.22
totales	_	27	_	Siller Si

# 3 Representaciones gráficas

Las represntaciones gráficas proporcinan una síntesis visual de la distribución de frecuencias. Las gráficas más utilizadas son las siguientes:

#### Características cualitativas:

• Diagrama de Pareto: En el eje de abscisas se asocia a cada modalidad un rectángulo de base constante y de altura proporcional a la frecuencia correspondiente (las modalidades se suelen disponer en orden decreciente según su frecuencia de aparición).



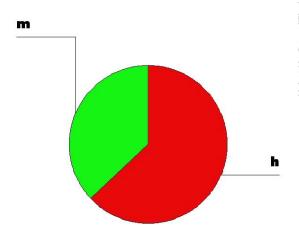
• Diagrama de sectores: A cada modalidad se le asigna un sector circular de amplitud  $w_i$  proporcional a su frecuencia relativa $\Rightarrow$ 

$$w_i = 2\pi \times f_i$$
, para cada modalidad i.



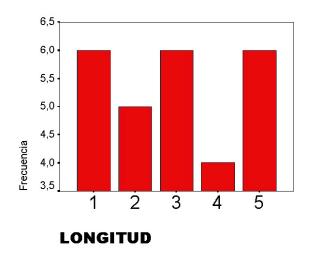
El área de los sectores, para la variable sexo, sería como sigue:

$$\mathbf{h} = 2\pi \times 0.63 \qquad \mathbf{m} = 2\pi \times 0.37$$



## Variables estadísticas discretas o continuas no agrupadas en intervalos de clase:

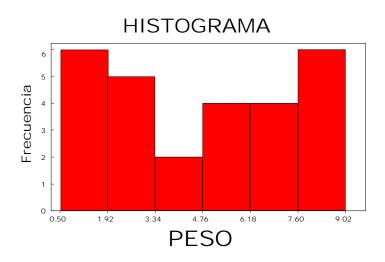
• Diagrama de barras: En el eje de abscisas se representan los valores que tome la variables y sobre cada uno de ellos se dibuja una barra de altura igual o proporcional a la frecuencia absoluta (o relativa) correspondiente.



Variables estadísticas continuas agrupadas en intervalos de clase:



• **Histograma:** En el eje de abscisas se representan los extremos de los intervalos de clase de la variable y sobre cada uno se construye un rectángulo de área proporcional a la frecuencia absoluta de cada clase. Cuando todos los intervalos son de la misma amplitud se suele tomar la altura de cada rectángulo igual a la frecuencia absoluta de dicha clase.



# 4 Características que debemos identificar ante un conjunto de datos

Nota: A partir de ahora, salvo que se diga lo contrario, nos centraremos en datos cuantitativos no agrupados; esto es,

muestra de tamaño  $\mathbf{n}: x_1, x_2, ..., x_n$ 

A continuación vamos a definir medidas numéricas que describen los aspectos más relevantes de la distribución de frecuencias. Estas características se clasifican según la información que tratan de resumir en:

- Medidas de posición o localización: describen cómo se comportan globalmente los datos y localizan la distribución de frecuencias.
- Medidas de dispersión: miden la variabilidad de los datos entre sí o respecto de una medida de centralización.



• Medidas de forma: informan sobre la asimetría de la distribución (medidas de asimetría) y sobre la concentración de las observaciones en torno a la zona central (medidas de apuntamiento o kurtosis).

#### 4.1 Medidas de Posición

#### 4.1.1 Medidas de Posición Central

Indican dónde se sitúa la zona central de la distribución de frecuencias.

• Media aritmética:

$$\overline{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

Si los datos están agrupados en intervalos de clase, se toman como observaciones las marcas de clase, esto es:

$$\overline{\overline{x}} = \frac{1}{n} \sum_{i=1}^{k} m_i \times n_i$$

donde  $m_i$  es la marca de clase del intervalo i y  $n_i$  es la frecuencia absoluta del intervalo i, con i = 1, ..., k.

Es la medida de centralización más utilizada ya que es el centro de gravedad del conjunto de datos, sin embargo, es muy sensible a valores extremos lo que la hace poco representativa.

• **Mediana**,  $M_e$ : Es aquel valor que divide en dos partes iguales la distribución de frecuencias. Esto es,

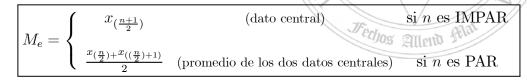
$$\underbrace{x_{(1)}, x_{(2)}, \ldots, M_e}_{50\% \text{ de los datos}}, \underbrace{\dots, x_{(n)}}_{50\% \text{ de los datos}}$$

#### Forma de calcularla para datos no agrupados:

- Ordenamos los datos de menor a mayor:

$$x_{(1)}, x_{(2)}, ..., x_{(n)}$$

- Entonces,



Si los datos están agrupados en intervalos de clase, hablamos de **intervalo mediano** como el primer intervalo de clase cuya frecuencia absoluta acumulada es  $\geq \frac{n}{2}$ . Si necesitamos un valor numérico para la mediana, podemos tomar la marca de clase del intervalo mediano.

• Moda,  $M_o$ : Valor de la variable que presenta mayor frecuencia. Existe para datos cualitativos.

Si los datos están agrupados en intervalos de clase, hablamos de **intervalo modal** como aquel que presenta mayor frecuencia. Si necesitamos un valor numérico para la moda, podemos tomar la marca de clase del intervalo modal.

Notar que puede existir más de una moda, así como no existir.

Puede no situarse en el centro de la distribución de frecuencias.

#### 4.1.2 Medidas de posición no central

Proporcionan información sobre la estructura interna de los datos.

• Cuantiles: Se define el cuantil de valor  $\alpha$  (0 <  $\alpha$  < 1) de una distibución de frecuencias como el valor  $C_{\alpha}$  que deja a su izquierda el  $100\alpha\%$  de las observaciones y a la derecha el  $100(1-\alpha)\%$  restantes de las observaciones.

Casos particulares de cuantiles para valores concretos de  $\alpha$ :

- **Percentiles**:  $P_1, P_2, ..., P_{99}$ , para  $\alpha = 0.01, 0.02, ..., 0.99$ , respectivamente.
- **Deciles**:  $D_1, D_2, ..., D_9$ , para  $\alpha = 0.1, 0.2, ..., 0.9$ , respectivamente.



• Cuartiles:  $Q_1, Q_2, Q_3$ , para  $\alpha = 0.25, 0.50, 0.75$ , respectivemente.



- Ordenamos los datos de menor a mayor:

$$x_{(1)}, x_{(2)}, ..., x_{(n)}$$

y determinamos la  $M_e$ .

– Entonces  $Q_1$  es la mediana del conjunto de datos que hay a la izquierda de la  $M_e$  (excluida la  $M_e$ ) y  $Q_3$  es la mediana del conjunto de datos que hay a la derecha de la  $M_e$  (excluida la  $M_e$ ).

Veámos cómo se calcula la  $M_e$  y los cuartiles con unos ejemplos:

a) Consideremos una muestra de tamaño 11 que toma los valores siguientes:

b) Consideremos una muestra de tamaño 10 que toma los valores siguientes:

$$5 \quad 5.3 \quad \boxed{6.1} \quad 7 \quad \boxed{7.2} \quad 7.5 \quad 7.8 \quad \boxed{8.1} \quad 8.6 \quad 8.9$$

## 4.2 Medidas de Dispersión

## 4.2.1 Medidas de Dispersión Absoluta

• Rango o recorrido: Amplitud del intevalo donde se encuentran distribuidas todas las observaciones.

$$R = x_{\text{max}} - x_{\text{min}}$$

Es muy sensible a valores extremos.

• Rango Intercuartílico: Amplitud del intevalo donde se encuentran distribuidas el 50% de las observaciones.

$$RIQ = Q_3 - Q_1$$

• Varianza: Medida de dispersión asociada la media aritmética y se define por:

$$s_X^2 = Var(X) = \frac{1}{n-1} \sum_{i=1}^n (x_i - \overline{x})^2 = \frac{n}{n-1} (\overline{x^2} - \overline{x}^2)$$

Si los datos están agrupados en intervalos de clase, se toman como observaciones las marcas de clase, esto es:

$$s_X^2 = Var(X) = \frac{1}{n-1} \sum_{i=1}^k (m_i - \overline{x})^2 \times n_i = \frac{n}{n-1} (\overline{x^2} - \overline{x}^2)$$

donde  $m_i$  es la marca de clase del intervalo i y  $n_i$  es la frecuencia absoluta del intervalo i, con i = 1, ..., k.

A la raíz cuadarada positiva de la varianza se le denomina **desviación típica** o **estándar** y se denota por s o  $s_X$ .

La varianza viene dada en unidades al cuadrado, mientras que la desviación típica viene en las mismas unidades físicas de los datos.

#### Importante:

- Como la media, es muy sensible a valores extremos.
- No es sensible a traslaciones.
- Es sensible a cambios de escala.

#### 4.2.2 Medidas de Dispersión Relativa

Medida adimensional, esto es, no tiene unidades.

## • Coeficiente de variación de Pearson:

$$CV = \frac{s}{\overline{x}}$$

Interpretación: Mide la representatividad de la media como medida que resume toda la información de la variable cuando comparamos dos o más distribuciones de frecuencias. Cuanto menor sea el valor de este coeficiente mayor representatividad de la media, ya que significa que los datos están más agrupados entorno a su valor medio.

#### Importante:

- Es sensible a traslaciones.
- No es sensible a cambios de escala.



## 4.3 Medidas de Forma

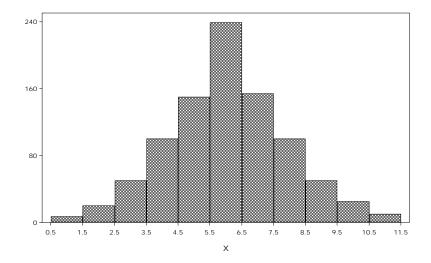
Nos proporcionan información sobre el perfil de la distribución de frecuencias.

Los atributos relacionados con la forma los vamos a establecer de manera aproximada observando la correspondiente representación gráfica de los datos y éstos son:

• Asimetría (skew): Coeficiente de asimetría de Fisher:

$$CA_F = \frac{\frac{1}{n-1} \sum_{i=1}^{n} (x_i - \overline{x})^3}{s^3}$$

Una distribución de frecuencias es simétrica si su correspondiente representación gráfica (diagrama de barra o histograma) es simétrica respecto de un eje vertical  $(CA_F \simeq 0)$ .



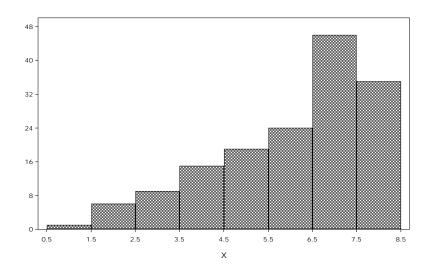


Wechos Allend Mar

Distribución asimétrica a la derecha  $(CA_F>0)$ 

# Distribución asimétrica a la izquierda $(CA_F < 0)$ .

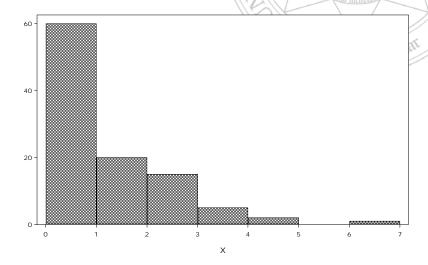
1.5



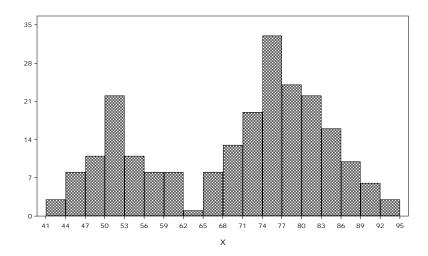


• Existencia o no de datos atípicos o anómalos que son datos que a simple vista se encuentran muy alejados del resto de los datos (en el apartado siguiente veremos una forma para identificarlos). Estos valores pueden ser debidos a un error de medida o de transcripción de los datos o corresponde a un verdadero valor de la variable.

Distribución asimétrica con existencia de valores extremos



• Unimodal, bimodal o multimodal. Cuando existe más de una moda, en ocasiones es posible identificar aproximadamente la existencia de subgrupos de datos. En estos casos, las medidas resumen globalmente pueden llegar a ser engañosas, por lo que siempre que sea posible, conviene explorar las características en dichos subgrupos de datos. Por ejempo los ingresos familiares, se espera globalmente dos grupos de datos según si una o dos personas de la unidad familiar trabajan.







Es un resumen gráfico que permite visualizar, para un conjunto de datos, la tendencia central, la dispersión y la presencia de valores extremos. Otra característica de este tipo de gráfico es que nos da información sobre la asimetría del conjunto de datos.

La mayor utilidad de los diagramas de caja y bigotes es para comparar dos o más conjuntos de datos.

Un diagrama de caja y bigotes se construye de la siguiente manera:

- 1. Ordenamos los datos de la muestra de menor a mayor y obtenemos los tres cuartiles.
- 2. Dibujamos un rectángulo cuyos extremos son  $Q_1$  y  $Q_3$  y dividimos el rectángulo por un segmento central a la altura de la  $M_e$
- 3. Calculamos los límites superior e inferior admisibles que nos servirán para identificar los datos atípicos. Éstos son:

$$L_{SUPERIOR} = Q_3 + 1.5 \times RIQ$$

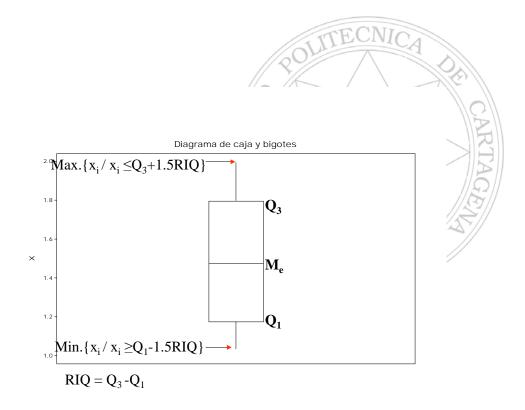
$$L_{INFERIOR} = Q_1 - 1.5 \times RIQ$$

Clasificamos como datos atípicos aquellas observaciones situadas fuera del intervalo  ${f C}$ 

$$[Q_1 - 1.5 \times RIQ, Q_3 + 1.5 \times RIQ]$$

4. Los segmentos  $1.5 \times RIQ$  (llamados **bigotes**) se acortan hasta: el dato del conjunto inmediatamente superior a  $Q_1 - 1.5 \times RIQ$  para el bigote inferior, esto es, hasta el  $min.x_i$ tal que  $x_i \geq Q_1 - 1.5 \times RIQ$ ; y el dato inmediatamente anterior a  $Q_3 + 1.5 \times RIQ$  para el bigote superior, esto es, hasta el  $max.x_i$  tal que  $x_i \leq Q_3 + 1.5 \times RIQ$ .





Cuando existen datos atípicos en el conjunto, representamos los datos atípicos como puntos

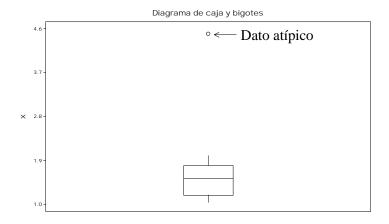
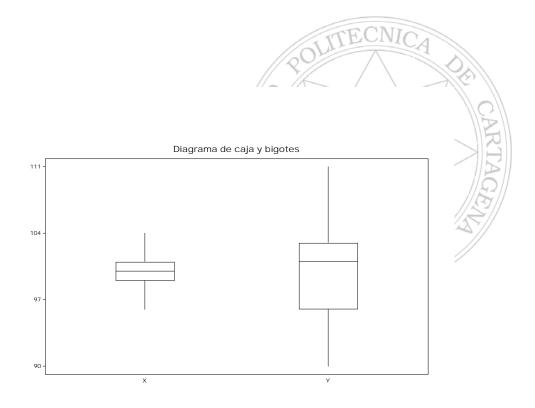
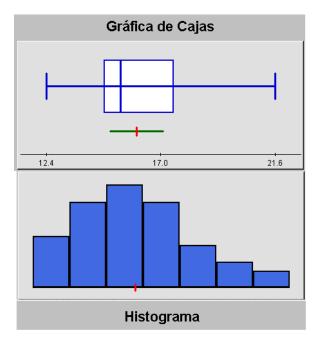


Diagrama de caja y bigote múltiple, cuando lo utilizamos para comparar dos o más conjuntos de datos:

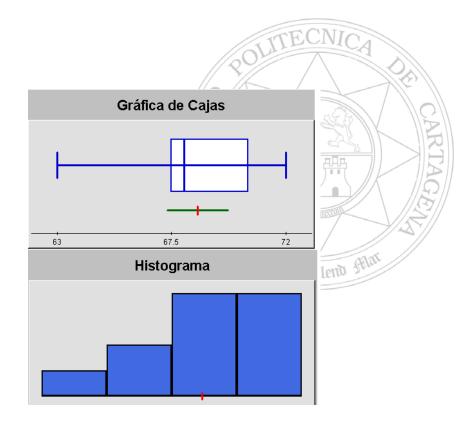




La información que nos proporciona el histograma en poblaciones unimodales también es proporcionado por el diagrama de cajas:







Esto no es cierto en poblaciones bimodales.

