

# Poblaciones y Muestras

Población:	Observaciones	Ejemplos
Conjunto de Individuos	personas o no	Los habitantes de Madrid Los coches de Madrid
que comparten una o más características	Lo que comparten es el hecho de tener la característica, no necesariamente sus valores.  Si todos los individuos no comparten los valores, a la característica la llamamos <b>variable</b> .  Si todos los individuos si comparten el valor de la característica, a la característica la llamamos <b>constante</b> .	Los habitantes de Madrid comparten la característica: Tener un color de ojos, pero cada uno los tiene de un color. El color de Ojos es una variable.  Los coches de Madrid comparten la característica: Tener un número de ruedas. Además todos comparten el valor: 4 ruedas. El número de ruedas de los coches de Madrid es una constante.
objeto de observación y estudio.	Para estudiar la característica necesitamos conocer los distintos valores que puede tomar. A estos valores los llamamos <b>Modalidades o Clases</b> .  El <b>proceso de medición</b> se encarga de definir cómo vamos a apuntar las diferentes modalidades que puede tomar una característica	La variable color de ojos puede presentar las siguientes modalidades: Azul, verde, marrón, ...

Cuando lo que estudiamos no es una muestra completa, sino una parte (subconjunto) de la misma, decimos que lo que estamos estudiando es una **Muestra**.

El número de individuos de una muestra o una población lo denominamos: **Tamaño de la muestra o de la población** y lo representaremos por la letra **n**.

# Proceso de Medición

Utilizaremos distintas escalas a la hora de anotar los valores (también llamados observaciones y puntuaciones directas) que toman las variables para cada individuo de la población o muestra.

Escala de Medición				Variables
Nominal	Otorga <b>nombres</b> a las diferentes modalidades que puede tomar la variable.	Permiten ver si dos individuos tienen <b>igual o diferente</b> modalidad.		Cualitativas
		Permiten <b>clasificar</b> los individuos en función del valor o modalidad.	Hombre, Mujer  Azul, Marrón, Negro, Verde	
Ordinal	Otorga <b>nombres</b> a las diferentes modalidades que puede tomar la variable, pero con la peculiaridad de que existe una <b>ordenación</b> entre ellos	Las modalidades son <b>comparables</b> en grado.	Poco, algo, bastante, mucho	Cualitativas
		Los valores se pueden <b>ordenar</b> entre sí.	Alto, Medio, Bajo	
Intervalo	Otorgan <b>medidas</b> a las diferentes modalidades que puede tomar la variable.  Tienen un 0 no absoluto, o sea, que si la variable toma valor 0 sigue teniendo sentido.	Aportan <b>medidas precisas</b> .  Las <b>medidas se pueden sumar o restar entre sí</b> .	Temperatura en grados Centígrados  (Una temperatura de 0 Grados centígrados tiene sentido)	Cuantitativas
Razón	Requieren de una <b>unidad de medida</b> . Nota: No tener en cuenta sólo las unidades de medida tradicionales: kg, litro, metro. Por ejemplo: <ul style="list-style-type: none"> <li>La nota de un examen se mide en puntos.</li> <li>El número de hijos que tiene una familia se mide en hijos.</li> </ul>		Número de alumnos en un curso  (Si hay 0 alumnos en un curso, es que no hay curso... no tiene sentido)	

# Estadística

Ciencia dedicada al estudio de **Poblaciones**

- **Descriptiva:**

**Describe** una población o una muestra mediante gráficos e índices (como la media, la mediana, la desviación típica...).

- **Inferencial:**

Trata de inferir (**averiguar**) los índices de una población si se conocen los índices de una muestra de la misma.

Estos índices recibirán nombres diferentes dependiendo de si se obtienen al estudiar una población o una muestra:

- Si los índices se calculan de una **población** se denominan **Parámetros**
- Si los índices se calculan de una **muestra** se denominan **Estadísticos**.

# Distribuciones de frecuencias

Para estudiar más fácilmente cómo se distribuyen los valores que toma una variable dentro de una muestra o población, utilizamos las tablas de frecuencias. En ellas se indican para cada modalidad o rango de modalidades que admite la variable los siguientes datos:

<b>Frecuencia Absoluta</b> Se representa por la letra $n_i$	El <b>número total de veces que se repite la modalidad</b> entre los individuos de una muestra o población.	Útil <b>cuando interesa conocer el número total</b> de casos que presentan una modalidad. La <b>suma de todas las frecuencias absolutas</b> debe ser <b>igual al tamaño de la población</b>
<b>Frec. Relativa</b> Se representa por la letra $f_r$ $n_r = n_i / N$	El <b>número total de veces que se repite la modalidad</b> entre los individuos de una muestra o población, dividido entre el tamaño de la población o de la muestra.	Útil <b>cuando interesa comparar</b> datos de dos poblaciones con tamaños diferentes. La <b>suma de todas las frecuencias relativas</b> debe ser <b>igual a 1</b> .
<b>Porcentaje</b> Se representa por la el símbolo % $\% = n_r \times 100$	El porcentaje de individuos que presentan una determinada modalidad.	Útil <b>cuando interesa comparar</b> datos de dos poblaciones con tamaños diferentes, y cuando interesa solamente saber las <b>diferencias en proporción</b> entre las diferentes modalidades. La <b>suma de todos los porcentajes</b> debe ser <b>igual a 100</b> .
<b>Frec. Acumulada</b> Se representa por la letra $N_i$	El <b>número total de veces que se repite la modalidad</b> , sumado al número total de veces que se repiten las modalidades mostradas más arriba en la tabla.	La frecuencia acumulada de la <b>última modalidad</b> de la tabla de frecuencias debe ser <b>igual al tamaño de la población</b> .
<b>Frec. Acum. Relativa</b> Se representa por la letra $N_r$ $N_r = N_i / N$	La frecuencia acumulada relativa de la <b>última modalidad</b> de la tabla de frecuencias debe ser <b>igual a 1</b> .	
<b>% Acumulado</b> Se representa por la el símbolo $\%_A$ $\%_A = N_r \times 100$	El porcentaje acumulado de la <b>última modalidad</b> de la tabla de frecuencias debe ser <b>igual a 100</b> .	

# Medidas de tendencia central

Dan una idea de la **magnitud general** de la variable en una muestra o población.

Se **miden en la misma unidad** que la variable.

<b>Moda</b>  M	<p>La moda es el <b>valor de la variable</b> (modalidad) <b>que más se repite</b>.</p> <p>Puede haber <b>varias modas</b>.</p>	<p>Para <b>cualquier tipo de variable</b></p>	<p>Corresponde con el valor de la variable que tiene mayor frecuencia.</p>
<b>Mediana</b>  Med	<p>La mediana es el <b>valor que divide a la población</b> en 2 partes iguales (50%).</p> <p>Para calcularlo:</p> <ol style="list-style-type: none"> <li><b>Ordenamos</b> los valores de las variables</li> <li>Buscamos el término central y tomamos el valor de la variable:             <ul style="list-style-type: none"> <li><b>Impar:</b> término <math>(N+1)/2</math></li> <li><b>Par:</b> Hay 2 términos, el <math>N/2</math> y el siguiente. Hacemos la media de los valores de la variable para ambos términos</li> </ul> </li> </ol> <p>Sólo usa los <b>valores de los términos centrales</b>.</p>	<p>Requiere <b>ordenar</b> los valores:</p> <ul style="list-style-type: none"> <li>Sólo para variables <b>cuasicuantitativas (ordinales) y cuantitativas</b>.</li> </ul> <p>Estadístico Resistente / Robusto:</p> <ul style="list-style-type: none"> <li>No se ve afectada por valores extremos / atípicos.</li> <li>Lo usamos para <b>distribuciones asimétricas</b>.</li> </ul>	<p>Divide a la superficie bajo la línea de frecuencias en 2 partes iguales. Está entre la media y la moda.</p>
<b>Media</b>  $\bar{X}$	<p>Es la media aritmética de <b>todos los valores</b> de la variable.</p>	<p><b>Suma</b> valores:</p> <ul style="list-style-type: none"> <li>Sólo para variables <b>cuantitativas</b>.</li> </ul> <p>Estadístico no Resistente / no Robusto:</p> <ul style="list-style-type: none"> <li>Si se ve afectada por valores extremos / atípicos</li> <li>Lo usamos para <b>distribuciones simétricas</b>.</li> </ul>	<p>Se coloca de forma que deje a la mediana entre la media y la moda.</p>

# Medidas de Posición:

## Percentiles, deciles, cuartiles

Son **los valores de la variable** que dejan por debajo una **determinada porción de la población**. Se miden en las **mismas unidades que la variable**.

Tipos	Símbolo K	Partes en que dividen a la población: P	¿Cuántos son?	Tamaño de la población entre 2 partes consecutivas
Mediana	Med	2	1	50%
Cuartiles	Q	4	3	25%
Deciles	D	10	9	10%
Percentiles	P	100	99	1%

Entre 2 percentiles, cuartiles, etc, se deja una determinada parte de la población, pero no podemos decir nada de las medidas que contienen.

A priori no hay ninguna relación entre las medidas de ellos. Sólo habría relación si la distribución es simétrica. Entonces ocurre por ejemplo que:

$$Q_1 = Q_3$$

$$D_1 = D_9; D_3 = D_7$$

$$P_{20} = P_{80}; P_{30} = P_{70}$$

## Calculo de las medidas de posición:

1. Se ordenan los términos
2. Se calcula el término central:
  - a. Si el tamaño de la población es impar, tomamos el valor de la variable para el término central. El término central lo calculamos con la siguiente formula:

$$\text{Término para } K_j: \frac{j}{P} \cdot (N + 1)$$

- b. Si el tamaño de la población es par, tomamos el valor de la variable para cada uno de los 2 términos centrales y hacemos su media. El segundo término es el que hay a continuación del primero. El primero lo calculamos con la siguiente formula:

$$\text{Término para } K_j: \frac{j}{P} \cdot (N)$$

# Medidas de Variabilidad

Miden como de dispersos están los valores de una distribución: Si difieren mucho o no de la magnitud general de la distribución.

- Valores **altos** de la variabilidad implican:

Valores de la variable **muy dispersos**: Distribución **heterogénea**.

- Valores **bajos** de la variabilidad implican:

Valores de la variable **poco dispersos**: Distribución **homogénea**.

Son **siempre iguales o mayores a 0**

**No tienen nada que ver con la simetría o asimetría**

No robustos / No resistentes	Robustos / Resistentes
<ul style="list-style-type: none"> <li>Muy afectados por valores extremos / atípicos</li> <li>Los usamos para distribuciones simétricas.</li> <li>Se basan en la media.</li> <li>Sólo variables cuantitativas</li> </ul>	<ul style="list-style-type: none"> <li>Poco afectados por valores extremos / atípicos</li> <li>Los usamos para distribuciones asimétricas.</li> <li>Se basan en la mediana y los cuartiles.</li> <li>Variables ordinales y cuantitativas</li> </ul>
<b>Rango:</b> $R = X_{\text{Max}} - X_{\text{Min}}$	<b>Rango Intercuartil:</b> $R_{\text{IC}} = Q_3 - Q_1$
<b>Varianza:</b> $S^2 = \frac{\Sigma(X_i - \bar{X})^2}{n}$	
<b>Desviación Típica:</b> $S = \sqrt{S^2} = \sqrt{\frac{\Sigma(X_i - \bar{X})^2}{n}}$	<b>Rango Semi-Intercuartil:</b> $R_{\text{SIC}} = \frac{Q_3 - Q_1}{2}$
<b>Coef. de variación:</b> $C_V = \frac{s}{\bar{X}}$ <ul style="list-style-type: none"> <li>Permite comparar la variabilidad de 2 distribuciones con medias y unidades diferentes.</li> <li>Sólo se puede calcular para distribuciones con media positiva</li> </ul>	<b>Coef. de variación Intercuartílico:</b> $C_{\text{VI}}$ <ul style="list-style-type: none"> <li>Permite comparar la variabilidad de 2 distribuciones con medias y unidades diferentes.</li> <li>Sólo se puede calcular para distribuciones con media positiva</li> </ul>



# Puntuaciones

Llamamos puntuaciones a cada una de las mediciones que hacemos de una variable para cada uno de los individuos de una población o una muestra.

<b>Directa</b>  $X_i$	La medida que anotamos de una observación.	Da el valor absoluto de la variable medido para un individuo concreto.	Puede ser positiva o negativa en función de la variable que estemos midiendo
<b>Diferencial</b>  $P_D = X_i - \bar{X}$	Su diferencia a la medida (por encima o por debajo).	Permite saber si la medición de una variable para un individuo está por encima o por debajo de la media.	Puede ser positiva o negativa en función de si la medida se encuentra por encima o por debajo de la media.  Si alguna de estas puntuaciones vale 0, significa que la medida directa se encuentra es igual a la media.
<b>Típica</b>  $P_T = \frac{X_i - \bar{X}}{S}$	Indica el número de desviaciones típicas que una medida se acerca o aleja de la media.	Permite comparar las posiciones relativas de las medidas de: <ul style="list-style-type: none"> <li>• 1 variables en 2 grupos</li> <li>• 2 variables en 1 grupo</li> <li>• 2 variables en 2 grupos</li> </ul> Cuanto mayor sea, mejor es la puntuación relativa de un individuo dentro de un grupo.	

Cada Puntuación directa tiene una puntuación típica y una diferencial.

2 puntuaciones directas iguales (2  $X_i$  iguales) no implica que sus puntuaciones típicas ni diferenciales sean iguales. Dependerá siempre de la media de su distribución, y en el caso de la puntuación típica también depende de la desviación típica.

# Correlaciones

Dan una idea de la **relación que presentan 2 variables entre sí**. Es muy importante tener en cuenta que **correlación no implica causalidad**.

Que dos variables estén correlacionadas significa que existe una **tendencia general** que se manifiesta entre ambas variables, pero que no tiene que cumplirse para la totalidad de los casos (solamente se mostraría para la totalidad de los datos si estuviéramos hablando de una relación perfecta).

Al hacer un estudio de correlación tenemos en cuenta 2 factores:

1. El **sentido** de la correlación:

- **Directo:**

- Individuos con valores altos en una variable tienden a tener valores altos en la otra.
- Individuos con valores bajos en una variable tienden a tener valores bajos en la otra.

- **Inversa:** Individuos con valores altos en una variable tienden a tener valores bajos en la otra.

2. La **intensidad** de la correlación: Ósea, si la tendencia general se presenta de una forma muy marcada o poco marcada.

Los **diagramas de dispersión** representan los valores (puntuaciones directas) que toman 2 variables medidas en el mismo conjunto de individuos, y son muy útiles para ayudarnos a identificar qué tipo de relación presentan dos variables entre sí (en caso de presentar alguna).

Dos variables pueden presentar una **relación lineal** (donde al dibujar un diagrama de dispersión de ambas variables los datos se distribuyen de forma similar a una línea recta); pero también pueden presentar otro tipo de relación (donde los datos se distribuirían de acuerdo a otro tipo de patrones que no sean una línea recta).

Existen ciertos **índices** que nos ayudan a identificar **relaciones lineales** entre variables. Si éstos índices valen 0, significará que **no hay relación lineal**, pero podría existir otro tipo de relación entre las mismas.

Lo que miden cualquiera de estos índices es lo cerca de una línea recta que están de **media** todos los puntos dibujados en un diagrama de dispersión. No hay que olvidar que hablamos de la tendencia general (media), es decir, en unos tramos los puntos pueden estar más cerca y en otros más alejados de la línea recta, incluso podría haber puntos muy alejados, también llamados atípicos. Por eso, **el valor de un índice de correlación medido en un tramo no tiene nada que ver con el valor del índice medido en otro tramo**.

### Cualitativas (Nominales):

En cuanto una de las variables es nominal se usan estos índices:

$$\chi^2 = \sum \frac{(f_e \cdot f_t)}{f_t}$$

$$\chi^2 \geq 0$$

- $f_e$  es la frecuencia empírica (la que medimos)
- $f_t$  es la frecuencia teórica (calculada como el producto de las frecuencias marginales entre el total de casos)

Este índice no tiene límite superior. No da información ni de intensidad ni de sentido. Usamos mejor el **índice de Contingencia**:

$$C = \sqrt{\frac{\chi^2}{n + \chi^2}} \quad 0 \leq C < 1$$

Su valor depende del número de filas y columnas de la tabla de contingencia. No es comparable entre variables que presenten diferente número de modalidades.

Para tablas cuadradas (con k filas y columnas), su valor máximo está acotado por la expresión:

$$C_{max} = \sqrt{\frac{k-1}{k}}$$

Sólo da información de intensidad. Para obtener el sentido de la correlación hay que comparar  $f_e$  con  $f_t$ .

### Cuasi-Cuantitativas (Ordinales):

Se usa cuando analizamos:

- 2 variables ordinales
- 1 ordinal y otra cuantitativa:

**Índice de Spearman:**

$$R_s = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

$$-1 \leq R_s \leq 1$$

donde:  $d_i$  es la distancia (resta) entre las posiciones (orden) que ocupa cada individuo dentro de su grupo para cada variable a correlacionar.

El signo del índice (tanto del de Spearman, como del de Pearson) da información del sentido de la correlación:

- (+) Relación Directa entre las variables
- (-) Relación Inversa entre las variables

El valor absoluto del índice (sin signo) da información de la intensidad de la relación:

Valor 0: No hay relación lineal entre las variables

Valor 1: Hay una relación lineal perfecta

### Cuantitativas (Pearson):

Se usa cuando analizamos 2 variables cuantitativas.

**Índice de Pearson:**

$$\begin{aligned} R_{xy} &= \frac{\sum \left( \frac{x_i - \bar{x}}{S_x} \right) \left( \frac{y_i - \bar{y}}{S_y} \right)}{n} \\ &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\frac{n}{S_x S_y}} \\ &= \frac{S_{xy}}{S_x \cdot S_y} \end{aligned}$$

$$-1 \leq R_{xy} \leq 1$$

Para su cálculo se usan las desviaciones típicas y la covarianza:

$$S_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n} = \frac{\sum x_i y_i}{n} - \bar{x} \bar{y}$$

