

# Estadística descriptiva

PFG-JLF

UAM

Estadística I, 2018-2019

# Datos/muestra

Punto de partida: disponemos

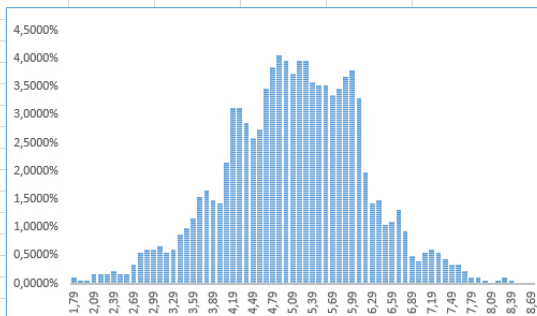
- de un conjunto de valores/datos (**muestra**)
- de una cierta característica/**variable**  $X$
- en una **población** específica.

Notación para la muestra:

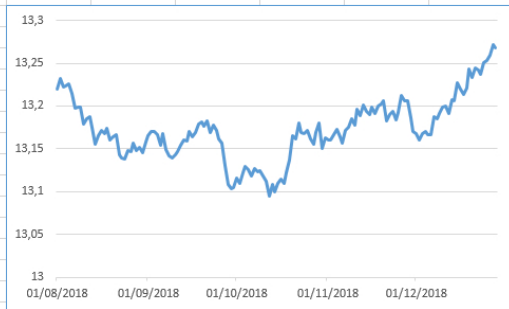
$$x_1, x_2, \dots, x_n$$

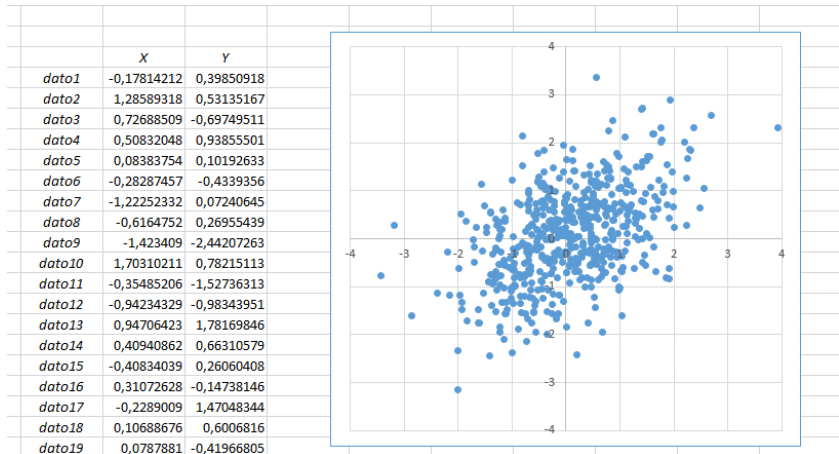
- El orden de los datos es irrelevante.
- El tamaño de la muestra es  $n$ .

<i>dato1</i>	1,793
<i>dato2</i>	4,910
<i>dato3</i>	4,181
<i>dato4</i>	3,444
<i>dato5</i>	3,932
<i>dato6</i>	5,573
<i>dato7</i>	7,242
<i>dato8</i>	5,654
<i>dato9</i>	6,089
<i>dato10</i>	4,450
<i>dato11</i>	5,490
<i>dato12</i>	5,957
<i>dato13</i>	5,828
<i>dato14</i>	3,992
<i>dato15</i>	4,816
<i>dato16</i>	3,592
<i>dato17</i>	5,480
<i>dato18</i>	5,939



01/08/2018	13,22
02/08/2018	13,23
03/08/2018	13,22
04/08/2018	13,22
05/08/2018	13,23
06/08/2018	13,21
07/08/2018	13,20
08/08/2018	13,20
09/08/2018	13,20
10/08/2018	13,18
11/08/2018	13,19
12/08/2018	13,19
13/08/2018	13,18
14/08/2018	13,16
15/08/2018	13,17
16/08/2018	13,17
17/08/2018	13,17
18/08/2018	13,17
19/08/2018	13,16





# Representación de los datos

- Se determinan una serie de **clases**  $C_1, \dots, C_k$  (generalmente intervalos),
- de manera que todos los datos caigan en alguna de las clases;
- se cuenta el número de datos en cada clase:  $n_1, \dots, n_k$  (números que suman  $n$ );
- o mejor, la frecuencia relativa en cada clase:  $f_1, \dots, f_k$ , donde  $f_j = n_j/n$  (estos números  $f_j$  suman 1);
- y se representan gráficamente, generalmente con un diagrama de barras.

# Medidas descriptivas de la muestra

Media muestral:

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

**Mediana:** el “valor” que deja tantos datos a su izquierda como a su derecha. Se ordenan los datos de menor a mayor:

$$x_{i_1} \leq x_{i_2} \leq \cdots \leq x_{i_n}.$$

- Si  $n$  es impar,  $n = 2r + 1$ ,  $\text{MED}_x = x_{i_{r+1}}$ .
- Si  $n$  es par,  $n = 2r$ , se toma (habitualmente)  
$$\text{MED}_x = \frac{x_{i_r} + x_{i_{r+1}}}{2}.$$

# Medidas descriptivas de la muestra

## Varianza muestral:

Definición:

$$V_x = \frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2.$$

Regla habitual de cálculo:

$$V_x = \overline{x^2} - \bar{x}^2,$$

donde

$$\overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2.$$



Cuasivarianza:

$$s_x^2 = \frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2.$$

Desviación típica:

$$\sqrt{V_x} = \sqrt{\frac{1}{n} \sum_{j=1}^n (x_j - \bar{x})^2} = \sqrt{\overline{x^2} - \bar{x}^2}$$

Cuasidesviación típica:

$$s_x = \sqrt{\frac{1}{n-1} \sum_{j=1}^n (x_j - \bar{x})^2}.$$

## Observaciones:

- $(n - 1) s_x^2 = n V_x$ .
- Las unidades de la desviación típica son las “correctas”.
- Varianzas/desviaciones típicas muestrales grandes indican que los datos están bastante “dispersos” con respecto a la media muestral.
- $V_x = 0$  si y solo si datos constantes.

# Cambios de escala/tipificación

**Cambio de escala:** para  $a, b \in \mathbb{R}$ , con  $b \neq 0$ ,

$$x_1, \dots, x_n \longrightarrow z_1, \dots, z_n,$$

donde

$$z_i = a + bx_i \quad \text{para cada } i = 1, \dots, n.$$

Se tiene que

$$\bar{z} = a + b\bar{x}, \quad V_z = b^2 V_x \quad \text{y} \quad s_z^2 = b^2 s_x^2.$$

Tipificación. El cambio

$$x_i \mapsto z_i = \frac{x_i - \bar{x}}{\sqrt{V_x}}$$

transforma los datos  $x_1, \dots, x_n$  en datos  $z_1, \dots, z_n$  tales que

$$\bar{z} = 0 \quad V_z = 1.$$

# Otras medidas

**Cuartiles.** Para datos ya ordenados

$$x_1 \leq x_2 \leq \cdots \leq x_n$$

se definen los cuartiles

$$Q_0, Q_1, Q_2, Q_3, Q_4$$

como sigue:

- $Q_0 = \min(x_1, \dots, x_n)$ ,  $Q_4 = \max(x_1, \dots, x_n)$ .
- $Q_2 = \text{MED}_x$ .
- $Q_3$  es la mediana de los datos que están entre  $Q_2$  y  $Q_4$ , y  $Q_1$  es la mediana de los que están entre  $Q_0$  y  $Q_2$ .

- rango de la muestra:  $Q_4 - Q_0$ ;
- rango intercuartílico:  $RIC = Q_3 - Q_1$ ;
- valores atípicos:
  - ▶ mayores que  $Q_3 + 1.5 \times RIC$ ,
  - ▶ menores que  $Q_1 - 1.5 \times RIC$ .

Coeficiente de asimetría muestral:

$$ASIM_x = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{V_x^{3/2}}$$

# Datos bidimensionales

Dos magnitudes,  $X$  e  $Y$ , medidas sobre los mismos individuos.

Muestra:

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

Se representan con un diagrama de dispersión.

# Medidas de dependencia lineal

## Covarianza muestral

$$\text{cov}_{x,y} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \overline{xy} - \bar{x} \cdot \bar{y},$$

donde

$$\overline{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i.$$

- $\text{cov}_{x,y} > 0$  indica dependencia positiva (relación directa entre las variables  $X$  e  $Y$ );
- $\text{cov}_{x,y} < 0$  indica dependencia negativa (relación inversa entre las variables  $X$  e  $Y$ ).



Se tiene (desigualdad de Cauchy–Schwarz) que

$$|\text{cov}_{x,y}| \leq \sqrt{V_x} \sqrt{V_y}.$$

Con igualdad si y sólo si los datos  $(x_j, y_j)$  están *todos sobre una misma recta*.

# Coeficiente de correlación

Si  $V_x \neq 0$ ,  $V_y \neq 0$ ,

$$\rho_{x,y} = \frac{\text{COV}_{x,y}}{\sqrt{V_x} \sqrt{V_y}}.$$

Propiedades:

- su signo tiene el mismo significado que el de la covarianza;
- $-1 \leq \rho_{x,y} \leq 1$ ;
- invariante bajo cambios de escala;
- $\rho_{x,y} = \pm 1$  si y sólo si datos sobre una recta.

# La recta de regresión

Tenemos una muestra del par de variables  $(X, Y)$ :

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n).$$

De entre todas las rectas  $y = a + bx$ , ¿cuál es la que “mejor” aproxima/explica la muestra?

¿Qué significa “mejor ajuste”? Para  $a, b \in \mathbb{R}$ , definimos el error cuadrático medio

$$E(a, b) = \frac{1}{n} \sum_{i=1}^n (y_i - (a + b x_i))^2$$

(errores “verticales”). Buscamos  $a, b$  que **minimicen** esta cantidad.

Escribimos

$$\begin{aligned}
 E(a, b) &= \frac{1}{n} \sum_{i=1}^n (y_i - (a + b x_i))^2 \\
 &= a^2 + \overline{x^2} b^2 - 2\overline{y} a + 2\overline{x} ab - 2\overline{xy} b + \overline{y^2} \\
 &= \underbrace{(a, b) \begin{pmatrix} 1 & \overline{x} \\ \overline{x} & \overline{x^2} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix}}_{\text{forma cuadrática}} \underbrace{-2\overline{y} a - 2\overline{xy} b}_{\text{términos lineales}} + \underbrace{\overline{y^2}}_{\text{constante}}
 \end{aligned}$$

La forma cuadrática es definida positiva (el determinante es justamente  $V_x$ ), así que tiene al menos un mínimo.

Cálculo del mínimo:

$$\begin{cases} \frac{\partial E(a, b)}{\partial a} = 2a + 2\bar{x}b - 2\bar{y} = 0, \\ \frac{\partial E(a, b)}{\partial b} = 2\bar{x}^2b + 2\bar{x}a - 2\overline{xy} = 0. \end{cases}$$

Es decir,

$$\begin{cases} \bar{y} = a + b\bar{x}, \\ \overline{xy} = a\bar{x} + b\bar{x}^2. \end{cases}$$

Solución:

$$\hat{b} = \frac{\text{COV}_{x,y}}{V_x}$$

$$\hat{a} = \bar{y} - \hat{b}\bar{x} = \bar{y} - \left(\frac{\text{COV}_{x,y}}{V_x}\right)\bar{x}$$

La recta de ecuación

$$y = \hat{a} + \hat{b}x$$

que da el mínimo error cuadrático medio, es la **recta de regresión de Y sobre X**.

Escrituras alternativas:

$$y - \bar{y} = \hat{b}(x - \bar{x}) \quad \text{o bien} \quad \frac{y - \bar{y}}{\sqrt{V_y}} = \rho_{x,y} \frac{x - \bar{x}}{\sqrt{V_x}}.$$

## Bondad de ajuste

El (mínimo) error cuadrático es  $E(\hat{a}, \hat{b})$ , que se puede escribir como

$$E(\hat{a}, \hat{b}) = \sqrt{V_y} \sqrt{1 - \rho_{x,y}^2}.$$

La cantidad  $\rho_{x,y}^2$  es conocida como el **coeficiente de determinación**  $R^2$ .

Valores de  $R^2$  próximos a 1 indican buen ajuste de la recta de regresión. Por ejemplo,  $R^2$  del orden de 0.8 o 0.9 son “buenos” ajustes.

En el análisis de la recta de regresión se suelen dar

- los valores de  $\hat{a}$  y  $\hat{b}$ ,
- y el valor de  $R^2$ .

# Transformación de datos

**Ajuste logarítmico.** Se quiere ajustar una curva del tipo

$$y = B \ln(x) + A$$

a unos datos  $(x_i, y_i)$ , con  $x_i > 0$ .

Procedimiento:

- nueva variable  $Z = \ln(X)$ ,
- transformamos los datos de la muestra: definimos  $z_i = \ln(x_i)$ ,
- ajustamos recta de regresión a los  $(z_i, y_i)$ ,

$$y = \hat{a} + \hat{b}z.$$

- el ajuste a los datos originales será

$$y = \hat{b} \ln(x) + \hat{a},$$

es decir,  $A = \hat{a}$  y  $B = \hat{b} = \text{cov}_{\ln(x), y} / V_{\ln(x)}$ .



Ajuste exponencial. Se quiere ajustar una curva del tipo

$$y = C e^{Dx}$$

a unos datos  $(x_i, y_i)$ , con  $y_i > 0$ .

Procedimiento:

- nueva variable  $W = \ln(Y)$ ; equivalentemente,  $Y = e^W$ ,
- transformamos los datos de la muestra: definimos  $w_i = \ln(y_i)$ ,
- ajustamos recta de regresión a los  $(x_i, w_i)$ ,

$$w = \hat{a} + \hat{b}x,$$

- el ajuste a los datos originales será

$$y = e^w = e^{\hat{a}} e^{\hat{b}x},$$

es decir,  $C = e^{\hat{a}}$  y  $D = \hat{b} = \text{cov}_{x, \ln(y)} / V_x$ .