

# 1. Estadística descriptiva

$X, Y$  variables con muestras  $x_1, \dots, x_n$ , e  $y_1, \dots, y_n$ , respectivamente.

- La **media**  $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$
- La **varianza**  $V_x = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \overline{x^2} - \bar{x}^2$  y la desviación típica  $\sqrt{V_x}$
- La **asimetría**  $\text{asim}_x = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3}{V_x^{3/2}}$
- La **covarianza**, cumple  $|\text{cov}_{x,y}| < \sqrt{V_x} \sqrt{V_y}$

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

- El **coeficiente de correlación** (adimensional, tipificado)

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

- La **recta de regresión** de  $Y$  sobre  $X$ ,  $y = \hat{b}x + \hat{a}$ , es

$$\hat{b} = \frac{\text{cov}_{x,y}}{V_x}, \quad \hat{a} = -\hat{b}\bar{x} + \bar{y}$$

$$\frac{y - \bar{y}}{\sqrt{V_y}} = \rho_{x,y} \frac{x - \bar{x}}{\sqrt{V_x}}$$

- La **bondad del ajuste**  $\sqrt{E(a,b)} = \sqrt{V_y} \sqrt{1 - \rho_{x,y}^2}$  donde  $\rho_{x,y}^2 = R^2$ ,  $E(a,b) = \frac{1}{n} \sum_{i=1}^n (y_i - (\hat{a} + \hat{b}x_i))^2$
- El ajuste logarítmico  $y = B \ln(x) + A$ : cambio de variable  $Z = \ln(X)$ .
- El ajuste exponencial  $y = Ce^{Dx}$ : cambio de variable  $W = \ln(Y)$  y tomar  $C = e^{\hat{a}}$ ,  $D = \hat{b}$ .
- El ajuste potencial  $y = CX^H$ : cambios de variable  $W = \ln(Y)$  y  $Z = \ln(X)$ .
- Ajuste por norma euclídea (sin asumir  $Y = f(X)$ ): tomar  $E(a,b) = \frac{1}{n} \sum_{i=1}^n \|y_i, (\hat{a} + \hat{b}x_i)\|^2$

# 2. Variables aleatorias

- $X$  **variable aleatoria** es una lista de valores  $x_1, \dots, x_n$  con sus probabilidades  $p_1, \dots, p_n$  con  $p_i = P(X = x_i)$ .
- Función de masa o de densidad**  $f_X(x)$  da las probabilidades de cada dato. Cumple  $f_X \geq 0$ ,  $\sum_{i=1}^n f_X(x_i) = \int_{\mathbb{R}} f_X = 1$  en el caso discreto y continuo, resp.
- Función de distribución**  $F_X(x)$  cumple  $0 \leq F_X \leq 1$ ,  $F_X$  creciente,  $F'_X = f_X$

$$F_X(x) = P(X \leq x) = \sum_{j|x_j \leq x} p_j = \int_{-\infty}^x f_X(t) dt$$

- La **esperanza**  $E(X) = \sum_{i=1}^n x_i P(X = x_i) = \int_{\mathbb{R}} x f_X(x) dx$ 
  - es lineal:  $E(aX + b) = aE(X) + b$ ,  $E(X + Y) = E(X) + E(Y)$
  - Cauchy-Schwarz:  $|E(XY)|^2 \leq E(X^2)E(Y^2)$
- La **covarianza**  $\text{cov}(X, Y) = E[(X - E(X))(Y - E(Y))] = E(XY) - E(X)E(Y)$  es un producto escalar (**bilineal!**):

- $\text{cov}(X_1 + X_2, X_3) = \text{cov}(X_1, X_3) + \text{cov}(X_2, X_3)$
- $\text{cov}(\lambda X_1, X_2) = \lambda \text{cov}(X_1, X_2)$

- el **coeficiente de correlación**  $\rho_{X,Y} = \frac{\text{cov}_{X,Y}}{\sqrt{V(X)} \sqrt{V(Y)}}$

$$X \perp Y \implies \rho_{X,Y} = 0 \iff \text{cov}_{X,Y} = 0$$

$$\iff E(XY) = E(X)E(Y)$$

- La **varianza**  $V(X) = E[(X - E(X))^2] = E(X^2) - E(X)^2$ 
  - $V(\lambda X) = \text{cov}(\lambda X, \lambda X) = \lambda^2 V(X)$
  - $V(X + Y) = V(X) + V(Y) + 2\text{cov}_{X,Y}$

## 2.1. Modelos de distribución

- Bernouilli:**  $X \sim \text{Ber}(p)$ :  $x_k \in \{0, 1\}$ ,  $p = P(X = 1)$

$$E(X) = p \quad V(X) = p(1 - p) \quad P(X = k) = p^k (1 - p)^{1-k}$$

- Binomial:**  $X \sim \text{Binom}(n, p)$ . Repetir una Bernouilli  $n$  veces y contar los aciertos.  $x_i = k = 0, \dots, n$

$$P(X = k) = \binom{n}{k} p^k (1 - p)^{n-k}$$

$$E(X) = np \quad V(X) = np(1 - p)$$

- Poisson:**  $X \sim \text{Poisson}(\lambda = np)$ . Binomial con  $n \rightarrow \infty$ ,  $p \rightarrow 0$ .

$$E(X) = \lambda \quad V(X) = \lambda \quad P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$$

- Exponencial:**  $X \sim \text{Exp}(\lambda)$ ,  $E(X) = \frac{1}{\lambda}$ ,  $V(X) = \frac{1}{\lambda^2}$

$$f_X(x) = \lambda e^{-\lambda x} \mathbf{1}_{[0, \infty)} \quad F_X(x) = 1 - e^{-\lambda x}, \quad x \in [0, \infty]$$

- Gamma:**  $X \sim \text{Gamma}(\lambda, t)$

$$\text{función gamma } \Gamma(s) = (s - 1)\Gamma(s - 1) = \int_0^\infty t^{s-1} e^{-t} dt$$

$$f_X(x) = \frac{1}{\Gamma(t)} \lambda^t x^{t-1} e^{-\lambda x}, \quad E(X^k) = \frac{(t + k - 1)!}{\lambda^k (t - 1)!}$$

$$\text{y es de utilidad } \int_0^\infty x^t e^{-bx} dx = \frac{\Gamma(t + 1)}{b^{t+1}} = \frac{t!}{b^{t+1}}$$

$$\text{Propiedad: } U \sim \text{Gamma}(a, b) \perp V \sim \text{Gamma}(a, c) \implies U + V \sim \text{Gamma}(a, b + c)$$

- Normal:**  $X \sim \mathcal{N}(\mu, \sigma^2)$

$$f_X(x) = \frac{1}{\sigma \sqrt{2\pi}} \exp \left\{ -\frac{(x - \mu)^2}{2\sigma^2} \right\}, \quad E(X) = \mu, \quad V(X) = \sigma^2$$

## 2.2. Cambio de variable

Si  $Y = H(X)$  entonces

$$f_Y(H(x)) |H'(x)| = f_X(x)$$

$$f_Y(y) = f_X(H^{-1}(y)) |(H^{-1})'(y)|$$

Si  $Z = X + Y$  entonces  $f_Z(z) = \int_{-\infty}^\infty f_{(X,Y)}(x, z - x) dx$

### 3. Vectores aleatorios

$X_1, \dots, X_n$  v.v. aa.,  $\mathbb{X} = (X_1, \dots, X_n)^t$  vector aleatorio.

- La **función de densidad conjunta**  $f_{(X_1, \dots, X_n)}(x_1, \dots, x_n)$
- $X_1, \dots, X_n$  **indep.**  $\iff P(X_1 \in A_1, \dots, X_n \in A_n) = P(X_1 \in A_1) \cdots P(X_n \in A_n)$
- El **vector de medias**  $E(\mathbb{X}) = (E(X_1), \dots, E(X_n))$
- La **matriz de covarianzas**

$$V(\mathbb{X}) = \text{cov}(\mathbb{X}) = \begin{pmatrix} V(X_1) & \dots & \text{cov}_{X_1, X_n} \\ \vdots & \ddots & \vdots \\ \text{cov}_{X_n, X_1} & \dots & V(X_n) \end{pmatrix}$$

es simétrica y semidefinida positiva.

#### 3.1. Vectores normales

Sea  $\vec{m} \in \mathbb{R}^n$ ,  $V$  matriz simétrica def. positiva de  $n \times n$ . Entonces  $\mathbb{X} \sim \mathcal{N}(\vec{m}, V) \iff$

$$f_{\mathbb{X}}(\vec{x}) = \frac{1}{(\sqrt{2\pi})^n} \frac{1}{\sqrt{\det V}} \exp \left\{ -\frac{1}{2} (\vec{x} - \vec{m})^t V^{-1} (\vec{x} - \vec{m}) \right\}$$

$V$  es def. pos  $\implies V = UU^t$ ,  $\det U \neq 0$  escribimos

$$\mathbb{X} \sim \mathcal{N}(\vec{m}, V) \iff \mathbb{X} = \vec{m} + U\mathbb{Y} \iff f_{\mathbb{X}}(\vec{x}) = \frac{1}{(\sqrt{2\pi})^n} \frac{1}{|\det U|} \exp \left\{ -\frac{1}{2} \|U^{-1}(\vec{x} - \vec{m})\|^2 \right\}$$

Además se tiene

$$E(\mathbb{X}) = \vec{m} \quad \text{cov}(\mathbb{X}) = V \quad X_j \sim \mathcal{N}(m_j, V_{jj})$$

- Si  $\vec{h} \in \mathbb{R}^n$ ,  $B \in \mathbb{R}^{n \times n}$ ,  $\det B \neq 0$ ,  $\mathbb{X} \sim \mathcal{N}(\vec{m}, V)$  vector normal, entonces  $\mathbb{Z} = h + B\mathbb{X} \sim \mathcal{N}(\vec{h} + B\vec{m}, BV B^t)$ .
- Si  $\vec{a} \in \mathbb{R}^n$  entonces la combinación lineal  $\sum_{i=1}^n a_i X_i = \vec{a}^t \mathbb{X}$  es v.a. y  $\sum_{i=1}^n a_i X_i \sim \mathcal{N}(\vec{a}^t \vec{m}, \vec{a}^t V \vec{a})$ .

#### 3.2. Distr. asociadas a vectores normales

- Chi-cuadrado**  $Z \sim \chi_n^2$  con  $n$  grados de libertad.  $Z = \|\mathbb{X}\|^2 = \sum X_i^2$ ,  $X_i \sim \mathcal{N}(0, 1)$ ,  $X_1, \dots, X_n$  indep.

$$E(Z) = n, \quad E(Z^\alpha) = 2^\alpha \frac{\Gamma(n/2 + \alpha)}{\Gamma(n/2)}, \quad V(Z) = 2n$$

$$X_i^2 \sim \text{Gamma}(1/2, 1/2), \quad Z \sim \text{Gamma}(n/2, 1/2)$$

- F de Fischer**  $Z \sim F_{n,m}$ . Si  $U \sim \chi_n^2$ ,  $V \sim \chi_m^2$ ,  $U \perp V$  entonces

$$Z = \frac{U/n}{V/m} \sim F_{n,m}$$

$$E(Z) = \frac{m}{m-2}, \quad V(Z) = \frac{2m^2(m+n-2)}{n(m-4)(m-2)^2}$$

- t de Student**  $Z \sim t_n$ . Si  $Y \sim \mathcal{N}(0, 1)$ ,  $U_n \sim \chi_n^2$ ,  $Y \perp U_n$  entonces

$$Z = \frac{Y}{\sqrt{\frac{U_n}{n}}} \sim t_n, \quad E(Z) = 0, \quad V(Z) = \frac{n}{n-2}$$

### 4. Modelo de muestreo aleatorio

Sea  $X$  una v.a. Sean  $X_1, \dots, X_n$  clones independientes de  $X$  que cumplen

- $X_i, X_j$  independientes si  $i \neq j$
- $X_i \stackrel{d}{=} X$ ,  $i = 1, \dots, n$  ( $X_i$  es igual a  $X$  en distribución)

Un **estadístico**  $T$  es una función  $T = H(X_1, \dots, X_n)$  donde  $H: \mathbb{R}^n \rightarrow \mathbb{R}$ .

#### 4.1. Estadístico media muestral

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i \quad E(\bar{X}) = E(X) \quad V(\bar{X}) = \frac{1}{n} V(X)$$

- Si  $X \sim \text{Ber}(p)$  entonces  $n\bar{X} \sim \text{Binom}(n, p)$ :

$$P(\bar{X} = \frac{k}{n}) = \binom{n}{k} p^k (1-p)^{n-k}$$

- Si  $X \sim \text{Poisson}(\lambda)$

$$P(\bar{X} = \frac{k}{n}) = P(n\bar{X} = k) = e^{-n\lambda} \frac{(n\lambda)^k}{k!}$$

- Si  $X \sim \mathcal{N}(\mu, \sigma^2)$  entonces  $\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n})$

$$P(|\bar{X} - E(X)| \geq \lambda) \leq \frac{V(X)}{n\lambda^2}$$

$$\sqrt{n}(\bar{X} - E(X)) \stackrel{d}{\underset{n \rightarrow \infty}{\rightarrow}} \mathcal{N}(0, V(X))$$

#### 4.2. Estadístico cuasivarianza muestral

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

$$n(n-1)S^2 = (n-1) \sum_{i=1}^n X_i^2 - \sum_{i \neq j} x_i x_j$$

$$E(S^2) = V(X)$$

$$V(S^2) = \frac{1}{n} E((X - E(X))^4) - \frac{n-3}{n(n-1)} V(X)^2$$

Estimación de la dispersión de  $S^2$  con Chebyshev:

$$P(|S^2 - V(X)| > \lambda) \leq \frac{V(S^2)}{\lambda^2} \leq \frac{1}{\lambda n^2} E((X - E(X))^4)$$

**Teorema** (de Fischer-Cochran). Si  $\mathbf{X} \sim \mathcal{N}(\mu, \sigma^2)$ ,  $X_1, \dots, X_n$  clones independientes de  $X$  y  $\bar{X}$ ,  $S^2$  son los estadísticos habituales entonces  $\bar{X}$  y  $S^2$  son independientes. Además

$$\bar{X} \sim \mathcal{N}(\mu, \frac{\sigma^2}{n}), \quad \frac{(n-1)}{\sigma^2} S^2 \sim \chi_{n-1}^2, \quad \frac{\bar{X}}{S/\sqrt{n}} \sim t_{n-1}$$

#### 4.3. Estadísticos máximo y mínimo

$$Mn = \max\{X_1, \dots, X_n\} \quad m_n = \min\{X_1, \dots, X_n\}$$

$$F_{M_n}(t) = F_X(t)^n \quad F_{m_n}(t) = 1 - (1 - F_X(t))^n$$

Para  $\alpha$  **mínimo esencial** ( $F_X(x < \alpha) = 0 \wedge F_X(x > \alpha) > 0$ ) y  $\beta$  **máximo esencial** ( $F_X(x < \beta) < 1 \wedge F_X(x > \beta) = 1$ ) de  $X$

$$P(\alpha \leq x \leq r) \xrightarrow{n \rightarrow \infty} 1, \quad r > \alpha$$

$$P(r \leq x \leq \beta) \xrightarrow{n \rightarrow \infty} 1, \quad r < \beta$$

## 5. Estimación de parámetros

Sea  $X$  v.a. con modelo descrito por  $f(x; \theta)$ ,  $\theta \in \Theta$ . Sea  $x_1, \dots, x_n$  una muestra de  $X$ .

- El **soporte**  $\text{sop}_\theta = \{x \in \mathbb{R} \mid f(x; \theta) > 0\}$
- Un **estimador** es un estadístico  $T = h(x_1, \dots, x_n) = \hat{\theta}$  que dados unos datos da una estimación de  $\theta$  llamada  $\hat{\theta}$   
Sean  $T, T'$  estimadores
  - El **sesgo** de  $T$  es  $\text{sesgo}_\theta(T) = E_\theta(T) - \theta$
  - $T$  es **insesgado**  $\iff E_\theta(T) = \theta, \forall \theta \in \Theta$
  - El **error cuadrático medio** de  $T$  es  $\text{ECM}_\theta(T) = E_\theta((T - \theta)^2)$ 
    - $T$  insesgado  $\implies \text{ECM}_\theta(T) = V_\theta(T)$
- $T$  es más **eficiente** que  $T'$   $\iff \text{ECM}_\theta(T) < \text{ECM}_\theta(T')$  para todo  $\theta \in \Theta$
- $T$  es eficiente  $\iff V_\theta(T) = \frac{1}{nI_X(\theta)}$

### 5.1. Método de momentos

El **estimador por momentos de orden**  $n \in \mathbb{R}$  se obtiene de

$$\overline{x^n} = E_\theta(X^n)$$

- En las distribuciones simétricas, se utilizan momentos de orden par
- Si el momento no depende de  $\theta$  entonces no hay estimador

### 5.2. Método de máxima verosimilitud

Sea  $X$  una v.a. con distribución  $f(x; \theta)$  discreta y finita,  $x_1, \dots, x_n$  muestra de  $X$

- La **función de verosimilitud**  $\text{VERO} : \Theta \rightarrow \mathbb{R}$

$$\text{VERO}(\theta; x_1, \dots, x_n) = \prod_{j=1}^n f(x_j; \theta)$$

- La **estimación de máxima verosimilitud**, si existe, es  $\hat{\theta}$ , el **máximo global (único)** de  $\text{VERO}$  en  $\Theta$ 
  - Para maximizar, es útil tomar  $\log \text{VERO}$

### 5.3. Límites de calidad de estimadores

**Lema** (de Diotivede). Sea  $X$  una v.a. con soporte que no depende de  $\theta$  y algunas hipótesis técnicas más. Entonces  $E_\theta(Y) = 0, \forall \theta \in \Theta$ .

- La **variable de información**

$$Y = \frac{\partial}{\partial \theta} \log f(x; \theta) = \frac{\partial_\theta f(x; \theta)}{f(x; \theta)}$$

- La **cantidad de información**  $I_X(\theta) = V_\theta(Y)$

**Lema** (Cramér-Rao). Para todo estadístico insesgado  $T$  de una v.a.  $X$  se tiene

$$V_\theta(T) \geq \frac{1}{nI_X(\theta)}$$

## 5.4. Comportamiento asintótico

**Teorema** (Método delta). Sea  $Z_n$  una sucesión de v.a. tales que para ciertos  $\alpha \in \mathbb{R}$  y  $\beta > 0$  se cumple

$$\sqrt{n}(Z_n - \alpha) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, \beta)$$

y sea  $g \in C^2(a, b)$  con  $\alpha \in (a, b)$  y  $g'(\alpha) \neq 0$ . Entonces

$$\sqrt{n}(g(Z_n) - g(\alpha)) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, |g'(\alpha)|^2 \beta)$$

Si  $g'(\alpha) = 0$  entonces aplicamos

$$n(g(Z_n) - g(\alpha)) \xrightarrow[n \rightarrow \infty]{d} \frac{g''(\alpha)}{2} \beta \chi_1^2$$

exigiendo que  $g \in C^3(a, b)$  y que  $g''(\alpha) \neq 0$

## 6. Contraste de hipótesis

Para una distribución que depende de un parámetro  $\theta \in \Theta$  definimos una hipótesis nula  $H_0 \equiv \theta \in \Theta_0$

- p-valor** es el valor tal que para  $\alpha > p$  rechazamos
  - se calcula forzando igualdad con  $\alpha = p$  y despejando el percentil  $c_p$

$$F(c_p) = 1 - p \implies p = 1 - F(c_p)$$

en una normal, la simetría nos da

$$\Phi(z_{p/2}) = 1 - \frac{p}{2} \implies p = 2(1 - \Phi(z_{p/2}))$$

- El **error de tipo 1** se da cuando se rechaza algo bueno
- El **error de tipo 2** se da cuando se acepta algo malo
- La **función de potencia**  $\beta(\theta) = P(\text{rechazar})$
- La **significación** es  $\sup_{\theta \in \Theta_0} \beta(\theta)$

### 6.1. Test por razón de verosimilitudes

Dada  $H_0 \equiv \theta \in \Theta_0$  sobre un parámetro  $\theta \in \Theta \dots$

- El Test se diseña haciendo que  $RV < c$  donde  $c \in (0, 1)$  es el calibre.
- Pasos a seguir

- Construir  $\text{VERO}(\theta; x_1, \dots, x_n)$
- Hallar  $\sup_{\theta \in \Theta} \text{VERO}(\theta)$  y  $\sup_{\theta \in \Theta_0} \text{VERO}(\theta)$
- Construir  $RV = \frac{\sup_{\theta \in \Theta_0} \text{VERO}(\theta)}{\sup_{\theta \in \Theta} \text{VERO}(\theta)}$
- Definimos el test: Rechazo  $H_0 \iff RV < c \in (0, 1)$
- Calculamos  $\beta(\theta) = P_\theta(\text{rechazar}) = P_\theta(RV < c)$
- Obtenemos la significación de  $\sup_{\theta \in \Theta_0} \beta(\theta)$

## 7. Caramelos

$$\text{Chebyshev } P(|X - E(X)| > \varepsilon) \leq \frac{V(X)}{\varepsilon^2}$$

$$\text{TCL } \sqrt{n}(\overline{X}_n - E(X)) \xrightarrow[n \rightarrow \infty]{} \mathcal{N}(0, V(X))$$

### 7.1. para ladrillitos

- Cuasivarianza muestral es  $s^2$ , cuasidesviación típica muestral es  $s$ . Son experimentales si son minúsculas, son de modelo teórico de muestreo aleatorio si son mayúsculas.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$
$$\Phi(z_{\alpha/2}) = 1 - \frac{\alpha}{2}$$

### 7.2. Más distribuciones

- **Uniforme:**  $X \sim \text{Unif}([0, a])$

$$f_X(x) = \frac{1}{a} \mathbf{1}_{[0,a]}, \quad E(X) = \frac{a}{2}, \quad V(X) = \frac{a^2}{12}$$

- **Geométrica:**  $X \sim \text{Geo}(p)$

$$f_X(k) = p(1-p)^{k-1}, \quad E(X) = \frac{1}{p}, \quad V(X) = \frac{1-p}{p^2}$$

### 7.3. Las funciones de distribución

- Para las **distribuciones discretas**

$$F_X(k) = P(X \leq k) = \sum_{j \leq k} f_X(j)$$

- **Poisson**  $X \sim \text{Pois}(\lambda = np)$

$$F_X(k) = e^{-\lambda} \sum_{j=0}^k \frac{\lambda^j}{j!}$$

- **Uniforme**  $X \sim \text{Unif}([0, a]), \ a > 0$

$$F_X(x < 0) = 0 \quad F_X(0 \leq x \leq a) = \frac{x}{a} \quad F_X(a < x) = 1$$

- **Exponencial**  $X \sim \text{Exp}(\lambda), \ \lambda > 0$

$$F_X(x < 0) = 0 \quad F_X(0 \leq x) = 1 - e^{-\lambda x}$$

- **Normal**  $X \sim \mathcal{N}(\mu, \sigma^2)$

$$F_X(x) = F_Y\left(\frac{x - \mu}{\sigma}\right)$$
$$F_Y(y < 0) = 1 - \Phi(-y) \quad F_Y(0 < y) = \Phi(y)$$