

CLASE 2 - MÉTODOS DE ESTIMACIÓN - Matías Carrasco

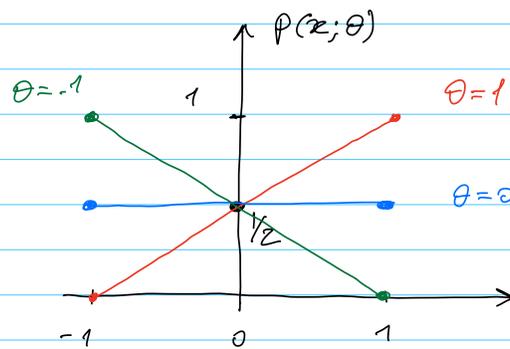
REPASO:

DEF: Un modelo estadístico es un par $(\mathcal{P}, \mathcal{Z})$ con:

- \mathcal{Z} un espacio de observaciones
- \mathcal{P} una familia de distribuciones de probabilidad en \mathcal{Z}

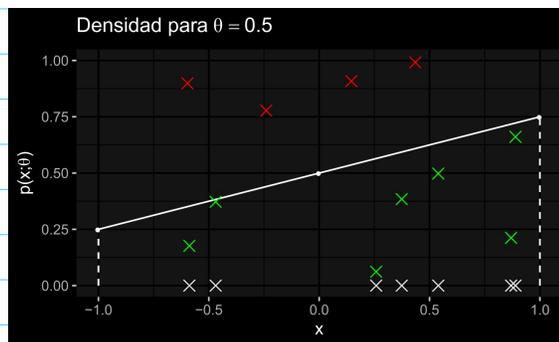
EJEMPLO: Consideremos $\mathcal{Z} = [-1, 1]$ $\mathcal{P} = \{p(x; \theta) : \theta \in [-1, 1]\}$

$$p(x; \theta) = \frac{1}{2} (1 + \theta x) \quad x \in [-1, 1]$$



En la clase pasada obtuvimos una muestra con el método de aceptación-rechazo:

$-0.47, 0.89, 0.26, -0.59, 0.97, 0.54, 0.87$



DEF: Sea $M = (\mathcal{P}, \mathcal{Z})$ modelo estadístico.

Siempre podemos escribir \mathcal{P} de la forma $\mathcal{P} = \{p(z; \theta) : \theta \in \Theta\}$

Es decir, siempre podemos "parametrizar" \mathcal{P} con algún espacio de parámetros Θ .

Dependiendo de la dimensión de Θ decimos que M es paramétrico o no paramétrico:

Si $\dim \Theta$ $\left\{ \begin{array}{l} \text{No depende del tamaño de la muestra} \\ \Rightarrow M \text{ es paramétrico} \\ \text{Depende del tamaño de la muestra} \\ \Rightarrow M \text{ es no paramétrico.} \end{array} \right.$

El ejemplo anterior es un ejemplo de modelo paramétrico.

En breve veremos un ejemplo de modelo no paramétrico.

Un poco de terminología

La muestra anterior consiste de una muestra i.i.d de $p(z; \theta)$

Esto quiere decir que tenemos X_1, X_2, \dots, X_n v.a. independientes todas ellas con distribución $p(z; \theta)$.

La muestra es una realización concreta de las variables:

$$(X_1)_{\text{obs}} = x_1, (X_2)_{\text{obs}} = x_2, \dots, (X_n)_{\text{obs}} = x_n$$

Siempre usamos $\left\{ \begin{array}{l} X_i \text{ para la variable aleatoria} \\ x_i \text{ para la realización } (X_i)_{\text{obs}} \end{array} \right.$

En el ejemplo: $\left\{ \begin{array}{l} x_1 = -0.47, x_2 = 0.89, \dots, x_7 = 0.87 \\ n = 7 \end{array} \right.$

El problema de estimación:

El problema de estimación consiste en **apartir de la muestra observada** $\{(X_i)_{obs}\} = \{x_i\}_{i=1}^n$ **estimar el parámetro** θ .

DEF: Un **estimador** es una **función** que a partir de una muestra genera una estimación para θ :

$$\hat{\theta}_n : \mathbb{R}^n \rightarrow \Theta \\ (x_1 \dots x_n) \mapsto \hat{\theta}_n(x_1 \dots x_n)$$

Podemos considerar $\hat{\theta}_n$ como una **variable aleatoria** si escribimos

$$\hat{\theta}_n(X_1, \dots, X_n)$$

Con este punto de vista un estimador tiene una **distribución asociada**!

Para nosotros será **clave** entender algunas **características** de la distribución de un estimador, por ejemplo su **sesgo** y su **varianza**.

Máxima Verosimilitud

Es el método más popular para **definir estimadores**.

$$M_\theta = (\mathcal{P}, \bar{x}) \text{ modelo estadístico } \mathcal{P} = \{p(x; \theta) : \theta \in \Theta\}$$

Z_1, \dots, Z_n muestra iid de la distribución real θ_0 .

DEF: La **verosimilitud** de un parámetro θ es la **"probabilidad"** de la muestra

$$V_n(\theta) = p(z_1; \theta) p(z_2; \theta) \dots p(z_n; \theta) \\ = \prod_{i=1}^n p(z_i; \theta)$$

El **negativo del logaritmo de la verosimilitud** es:

$$l_n(\theta) = -\ln V_n(\theta) = -\sum_{i=1}^n \ln(p(z_i; \theta))$$

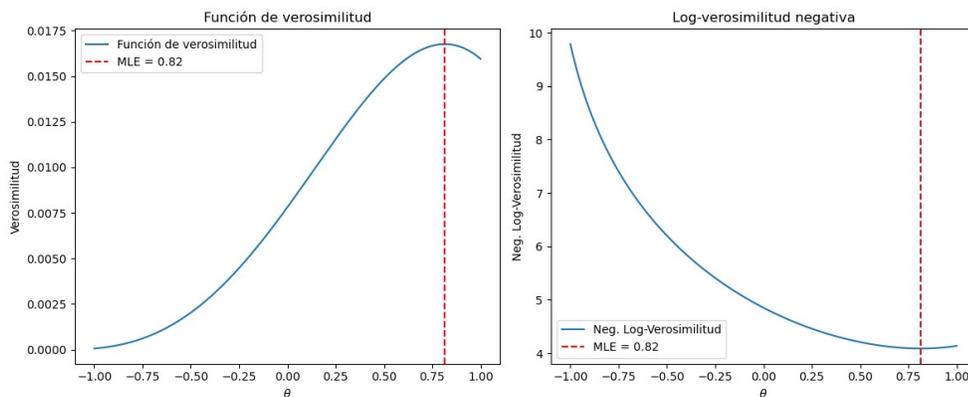
El **estimador de máxima verosimilitud (MLE)** es aquel valor de θ que maximiza $V_n(\theta)$

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} V_n(\theta)$$

Como el logaritmo es **monótono**, esto es equivalente a

$$\hat{\theta}_n = \arg \min_{\theta \in \Theta} l_n(\theta)$$

EJEMPLO: en el ejemplo anterior



$$l(\theta) = -\sum_{i=1}^n \ln\left(\frac{1}{2}(1 + \theta X_i)\right) = n \ln 2 - \sum_{i=1}^n \ln(1 + \theta X_i)$$

Para la muestra dada $\hat{\theta}_{obs} = 0.82$

EJEMPLO: $M = (\mathcal{P}, \mathcal{Z})$ $\mathcal{Z} = \{0, 1\}$
 $\mathcal{P} = \{p(z; \theta) : \theta \in [0, 1]\}$

Bernoulli $p(z; \theta) = \begin{cases} \theta & \text{si } z=1 \\ 1-\theta & \text{si } z=0 \end{cases} = \theta^z (1-\theta)^{1-z}$

Si tenemos una muestra $Z_1 \dots Z_n$:

$$V_n(\theta) = \prod_{i=1}^n p(Z_i; \theta) = \theta^{\sum z_i} (1-\theta)^{n-\sum z_i}$$

$$l_n(\theta) = -(\sum z_i) \ln \theta - (n - \sum z_i) \ln(1-\theta)$$

$$\begin{cases} l'_n(\theta) = -\frac{\sum z_i}{\theta} + \frac{(n - \sum z_i)}{1-\theta} \\ l''_n(\theta) = \frac{\sum z_i}{\theta^2} + \frac{(n - \sum z_i)}{(1-\theta)^2} > 0 \end{cases}$$

$$l'_n(\theta) = 0 \Leftrightarrow \frac{\sum z_i}{\theta} = \frac{(n - \sum z_i)}{1-\theta}$$

$$\Leftrightarrow (1-\theta) \sum z_i = (n - \sum z_i) \theta$$

$$\Leftrightarrow \sum z_i - \theta \sum z_i = (n - \sum z_i) \theta$$

$$\Leftrightarrow \sum z_i = n\theta$$

$$\Leftrightarrow \hat{\theta} = \sum z_i / n = \boxed{\bar{z}_n} \rightarrow \text{promedio}$$

EJEMPLO: $M = (\mathcal{P}, \mathcal{Z})$ $\mathcal{Z} = \mathbb{R}$
 $\mathcal{P} = \{N(x; \mu, \sigma^2) : \theta = (\mu, \sigma^2)\}$

$$N(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x-\mu)^2} \quad \text{Normal}$$

$$V_n(\mu, \sigma^2) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(z_i - \mu)^2}$$

$$\ln(\mu, \sigma^2) = \frac{1}{2} n \ln 2\pi + \frac{1}{2} n \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^n (z_i - \mu)^2$$

$$\begin{cases} \frac{\partial \ln}{\partial \mu} = -\frac{1}{\sigma^2} \sum_{i=1}^n (z_i - \mu) = 0 \\ \frac{\partial \ln}{\partial \sigma^2} = \frac{n}{2\sigma^2} - \frac{1}{2\sigma^4} \sum_{i=1}^n (z_i - \mu)^2 = 0 \end{cases}$$

$$\begin{cases} \hat{\mu} = \bar{z}_n \\ \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (z_i - \bar{z})^2 \end{cases}$$

Sesgo y Varianza de un estimador

Un estimador $\hat{\theta}_n$ depende de la muestra $X_1, \dots, X_n \Rightarrow \hat{\theta}_n$ es una v.a.

$\theta_0 \in \Theta$ es el valor real que queremos estimar

DEF: 1) El error cuadrático medio (MSE) de $\hat{\theta}_n$ es:

$$\text{MSE}(\hat{\theta}_n) = \mathbb{E}[(\hat{\theta}_n - \theta_0)^2]$$

2) El sesgo de $\hat{\theta}_n$ es

$$\text{Sesgo}(\hat{\theta}_n) = \mathbb{E}[\hat{\theta}_n] - \theta_0$$

3) La varianza de $\hat{\theta}_n$ es

$$\text{Var}(\hat{\theta}_n) = \mathbb{E}[(\hat{\theta}_n - \mathbb{E}\hat{\theta}_n)^2]$$

La esperanza es como un promedio en todas las muestras posibles

Por ejemplo:

$$\mathbb{E}[\hat{\theta}_n] = \text{Promedio}[\hat{\theta}_n]_{X_1, \dots, X_n \text{ posibles}}$$

Estas tres cantidades están relacionadas de la siguiente manera:

Proposición:

$$\text{MSE}(\hat{\theta}_n) = \text{Sesgo}(\hat{\theta}_n)^2 + \text{Var}(\hat{\theta}_n)$$

DEMOSTRACIÓN:

$$\text{MSE}(\hat{\theta}_n) \stackrel{\text{DEF}}{=} \mathbb{E}[(\hat{\theta}_n - \theta_0)]^2 \stackrel{\text{DISTRIBUTIVA}}{=} \mathbb{E}[\hat{\theta}_n^2 - 2\hat{\theta}_n\theta_0 + \theta_0^2]$$

$$\stackrel{\text{LINEALIDAD}}{=} \mathbb{E}[\hat{\theta}_n^2] - 2\theta_0 \mathbb{E}[\hat{\theta}_n] + \theta_0^2$$

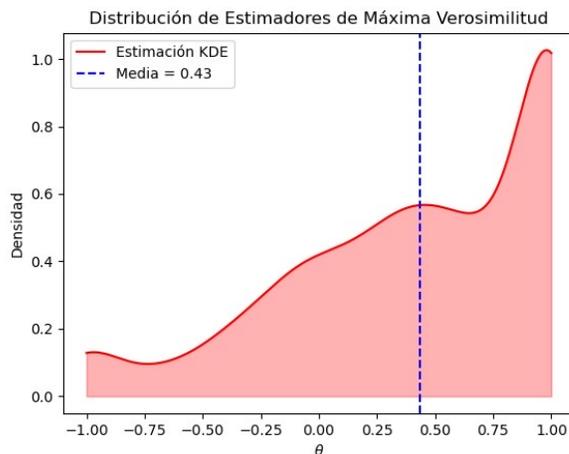
$$= \text{Var}(\hat{\theta}_n) + \mathbb{E}[\hat{\theta}_n]^2 - 2\theta_0 \mathbb{E}[\hat{\theta}_n] + \theta_0^2$$

$$= \text{Var}(\hat{\theta}_n) + \underbrace{[\mathbb{E}[\hat{\theta}_n] - \theta_0]^2}_{\text{Sesgo}(\hat{\theta}_n)} = \text{Var}(\hat{\theta}_n) + \text{Sesgo}(\hat{\theta}_n)^2$$

EJEMPLO: Consideremos nuevamente el ejemplo $\mathcal{X} = [-1, 1]$

$$\mathcal{P} = \{p(x; \theta) : \theta \in [-1, 1]\} \quad p(x; \theta) = \frac{1}{2}(1 + \theta x) \quad -1 \leq x \leq 1$$

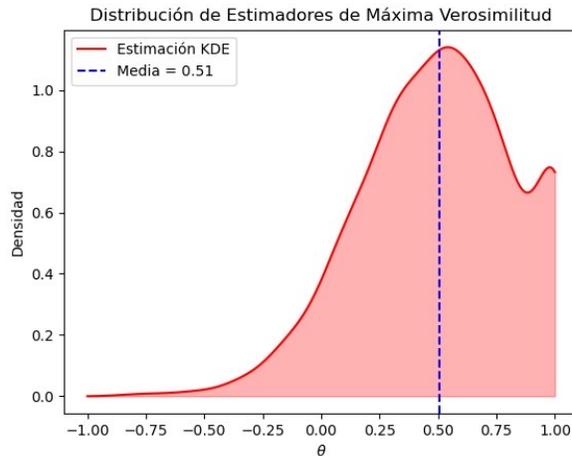
Consideremos ahora muchas ⁽²⁰⁰⁰⁾ muestras de tamaño $n=7$ y grafiquemos la distribución de $\hat{\theta}$ el estimador MLE.



la varianza de $\hat{\theta}_7$
en este caso es

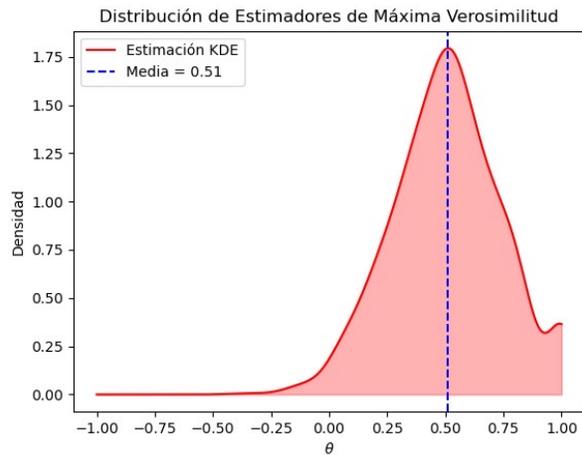
$$\text{Var}(\hat{\theta}_7) = 0.32$$

Observemos ahora el efecto que tiene el tamaño de la muestra en el MLE:



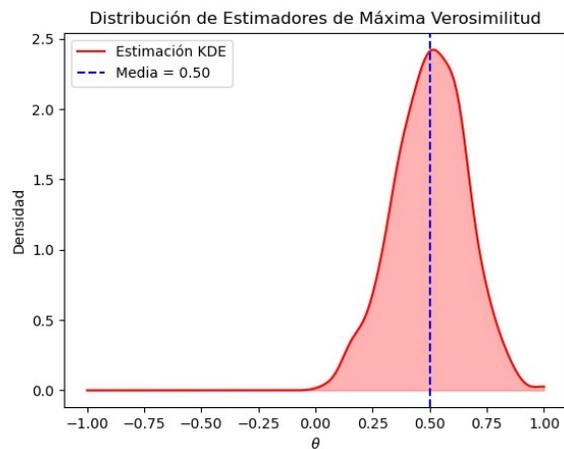
$$n = 25$$

$$\text{Var}(\hat{\theta}_{25}) = 0.10$$



$$n = 50$$

$$\text{Var}(\hat{\theta}_{50}) = 0.06$$



$$n = 100$$

$$\text{Var}(\hat{\theta}_{100}) = 0.02$$

Vemos que al aumentar n :

- 1) Sesgo $(\hat{\theta}_n) \rightarrow 0$ (el MLE es **asintóticamente insesgado**)
- 2) $\text{Var}(\hat{\theta}_n) \rightarrow 0$

1) y 2) juntos implican $\text{MSE}(\hat{\theta}_n) \rightarrow 0$: $\hat{\theta}_n$ es **consistente**

Además vemos que la distribución de $\hat{\theta}_n$ es cada vez más parecida a la distribución normal. El MLE es **asintóticamente normal**.

Cómo se calcula el MLE en la práctica: métodos de optimización

Veremos solo un ejemplo para ilustrar el concepto:

El Algoritmo Golden Search

La proporción áurea:



grande = todo
chico grande

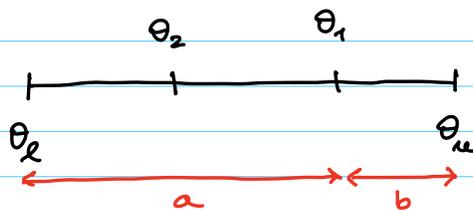
$$\frac{a}{b} = \frac{a+b}{a} = 1 + \left(\frac{a}{b}\right)^{-1}$$

$$\phi = 1 + \frac{1}{\phi}, \quad \phi^2 - \phi - 1 = 0$$

Dado un intervalo $[\theta_L, \theta_U]$

podemos subdividirlo de acuerdo a la proporción áurea desde ambos extremos:

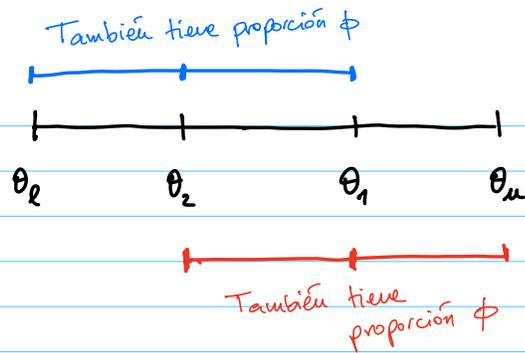
$$\phi = \frac{1 + \sqrt{5}}{2} = 1.618 \dots$$



$$\frac{a}{b} = \frac{a+b}{a}$$

$$\frac{\theta_1 - \theta_L}{\theta_U - \theta_1} = \frac{\theta_U - \theta_L}{\theta_1 - \theta_L} = \phi$$

Hacemos lo mismo con θ_2



El Algoritmo:

1 - Comenzar con dos valores iniciales θ_l y θ_u que encierran al mínimo de la función

2 - Subdividir:

$$\begin{cases} \theta_1 = \theta_l + \frac{\theta_u - \theta_l}{\phi} \\ \theta_2 = \theta_u - \frac{\theta_u - \theta_l}{\phi} \end{cases}$$

3 - Evaluar $l(\theta_1)$ y $l(\theta_2)$

4 - Si: $l(\theta_1) < l(\theta_2)$:

$\hat{\theta} \leftarrow \theta_1$
 $\theta_l \leftarrow \theta_2$
 $\theta_u \leftarrow \theta_u$
 $\theta_2 \leftarrow \theta_1$
 $\theta_1 \leftarrow \theta_l + \frac{\theta_u - \theta_l}{\phi}$

Si: $l(\theta_1) > l(\theta_2)$:

$\hat{\theta} \leftarrow \theta_2$
 $\theta_u \leftarrow \theta_1$
 $\theta_l \leftarrow \theta_l$
 $\theta_1 \leftarrow \theta_2$
 $\theta_2 \leftarrow \theta_u - \frac{\theta_u - \theta_l}{\phi}$

5 - Repetir hasta alcanzar $|\theta_u - \theta_l| < \underbrace{\text{tolerancia}}_{\tau}$

EJERCICIO: Implementarlo en el ejemplo

Un modelo no paramétrico: el histograma

$Z = [0,1]$ y queremos estimar una densidad $p(x)$

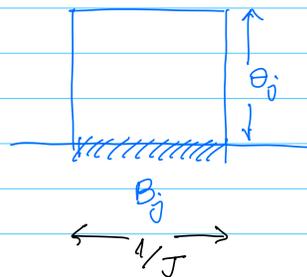
Recordar que para hacer un histograma:

- Dividimos $[0,1]$ en intervalos (bins) $\{B_j\}_{j=1}^J$ de ancho Δ
- Contamos cuantos elementos de la muestra $X_1 \dots X_n$ caen en cada bin B_j .

Supongamos dado el ancho Δ de los bins.

Definimos \mathcal{P} la familia de distribuciones de la siguiente manera:

$$p(x; \theta) = \sum_{j=1}^J \theta_j \mathbb{1}_{\{x \in B_j\}}$$



El parámetro es $\theta = (\theta_1, \dots, \theta_J)$ que tiene dimensión J . ($= 1/\Delta$)

Notar que $\frac{1}{J} \sum_{j=1}^J \theta_j = 1$ pues queremos área = 1.

(Observar también que $1/J = \Delta$)

O sea que $\Theta = \left\{ \theta = (\theta_1, \dots, \theta_J) : \frac{1}{J} \sum_{j=1}^J \theta_j = 1 \right\}$

Se puede ver que el MLE $\hat{\theta}$ es:

$$\hat{\theta}_j = J_x \frac{\#\{i: X_i \in B_j\}}{n}$$

Es decir: el histograma clásico es el MLE.

Para obtener consistencia debemos hacer $\Delta \rightarrow 0$ a medida que $n \rightarrow \infty$, por lo que el modelo es NO PARAMÉTRICO.