

Ciencias de Datos con R: Fundamentos Estadísticos

Ana M. Bianco, Jemina García y Mariela Sued.

Estimación.

Estadística

- Cuándo hacemos estadística no conocemos a F ni el valor del parámetro poblacional de interés.
- Cuándo hacemos estadística queremos hacer una cuenta con la muestra que nos permita estimar el valor θ de interés.

$$\text{Estimador: } \hat{\theta}_n = \hat{\theta}_n(X_1 \dots, X_n)$$

$$\text{Estimación: } \hat{\theta}_{n,\text{obs}} = \hat{\theta}_n(x_1 \dots, x_n),$$

donde x_1, \dots, x_n representan datos. Valores observados.

Muestra - Datos (Observaciones)

- Muestra (aleatoria simple):

X_1, \dots, X_n Variables aleatorias iid.

- Datos - Observaciones - Valores observados x_1, \dots, x_n :
Números.

Datos-Observaciones: son realizaciones de las variables aleatorias

Datos-Observaciones: son los resultados obtenidos al realizar el "experimento"

Estadística

POBLACION $\leftrightarrow F$	MUESTRA X_1, \dots, X_n i.i.d. $X_i \sim F$
Parámetro: Valor asociado de F $\theta = \theta(F)$ θ : valor poblacional	Estimador: estadístico para estimar θ $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ $\hat{\theta}_n$ NUEVA VARIABLE ALEATORIA

Estadística

Parámetro (en palabras)	Parámetro (en matemática)	Estimador Var. Alea	estimación (en R)
tita	θ	$\hat{\theta}_n$	$\hat{\theta}_{n,obs}$
esperanza	$\mu = \mathbb{E}(X)$	$\hat{\mu}_n = \bar{X}_n$	<code>mean(datos)</code>
desvío	$\sigma = \sqrt{V(X)}$	$\sqrt{S_n^2}$	<code>sd(datos)</code>
probabilidad	$p = \mathbb{P}(X \leq 3)$	$\frac{1}{n} \sum_{i=1}^n I_{(X_i \leq 3)}$	<code>mean(datos<=3)</code>
mediana	$F^{-1}(0.5)$	$X_{(n/2)}$	<code>median(datos)</code>
p-quantil	$F^{-1}(p)$	$X_{(pn)}$	<code>quantile(datos, p)</code>

Estimación de $F(t)$ - LA empírica

La LGN demuestra que la frecuencia relativa converge a la probabilidad.

$$(X_i)_{i \geq 1} \text{ i.i.d., } X_i \sim F$$

$$\text{LGN: } \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq 3\}} \xrightarrow{p} \mathbb{E}(I_{\{X_1 \leq 3\}}) = \mathbb{P}(X_1 \leq 3) = F(3).$$

$$\hat{F}_n(3) := \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq 3\}}$$

Estimación de $F(t)$ - LA empírica

La LGN demuestra que la frecuencia relativa converge a la probabilidad.

$$(X_i)_{i \geq 1} \text{ i.i.d., } X_i \sim F$$

$$\text{LGN: } \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq 3\}} \xrightarrow{p} \mathbb{E}(I_{\{X_1 \leq 3\}}) = \mathbb{P}(X_1 \leq 3) = F(3).$$

$$\hat{F}_n(3) := \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq 3\}}$$

$$\hat{F}_n(t) := \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq t\}}$$

$$\text{LGN: } \hat{F}_n(t) = \frac{1}{n} \sum_{i=1}^n I_{\{X_i \leq t\}} \xrightarrow{p} \mathbb{E}(I_{\{X_1 \leq t\}}) = \mathbb{P}(X_1 \leq t) = F(t).$$

Manos a la obra

Mediciones de gas - Equipo 1

Considere $n = 100$ datos obtenidos al utilizar el equipo 1.

1. Realice un histograma.
2. Calcule el promedio de los datos.
3. Calcule el percentil 0.9 de los datos.
4. Estime la probabilidad de que una medición realizada con este equipo diste de 70 en más de 2 unidades.
5. Repita los ítems anteriores utilizando ahora los primeros $n = 5$ y $n = 30$ datos

Sigue amasando, sigue amasando

Duración de lámparas

Considere los datos de la duración de n lámparas (en meses), para $n \in \{5, 30, 100\}$. En cada caso,

1. realice un histograma.
2. calcule el promedio y el percentil 0.9.
3. estime la probabilidad de que la lámpara dure a lo sumo un año (12 meses).

Mediciones de Gas - Con otro lenguaje

Consideremos las mediciones de gas realizadas por el equipo 1. Sea X_i el resultado de la i -ésima medición, para $i = 1, \dots, n$.

Asumiremos que X_1, \dots, X_n son v.a.i.i.d.

1. Indicar cuál cuenta hay que hacer con la muestra (X_1, \dots, X_n) para estimar $\mu = \mathbb{E}(X_1)$. Es decir, proponer un estimador $\hat{\mu}_n$ para μ .
2. Considerar $n = 5$ datos correspondientes al equipo 1 y calcular la estimación de μ correspondiente a estos datos. Repetir considerando $n = 30$ y $n = 100$.
3. Sea $q = F^{-1}(0.9)$ con $X_i \sim F$. Indicar cuál cuenta hay que hacer con la muestra (X_1, \dots, X_n) para estimar q . Es decir, proponer un estimador \hat{q}_n para q .
4. Sea $p = \mathbb{P}(|X - 70| > 2)$ con $X \sim X_1$. Indicar cuál cuenta hay que hacer con la muestra (X_1, \dots, X_n) para estimar p . Es decir, proponer un estimador \hat{p}_n para p .
5. Considerar $n = 5$ datos duraciones de lámparas y calcular la estimación de p correspondiente a estos datos. Repetir considerando $n = 30$ y $n = 100$.

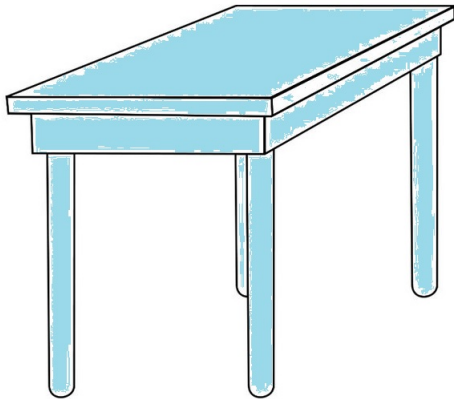
Duración de lámparas - Con otro lenguaje

Consideremos las duraciones de lámparas en meses. Sea X_i la duración de la i -ésima lámpara, para $i = 1, \dots, n$. Asumiremos que X_1, \dots, X_n son v.a.i.i.d.

1. Indicar cuál cuenta hay que hacer con la muestra (X_1, \dots, X_n) para estimar $\mu = \mathbb{E}(X_1)$. Es decir, proponer un estimador $\hat{\mu}_n$ para μ .
2. Considerar $n = 5$ datos de duraciones de lámparas y calcular la estimación de μ correspondiente a estos datos. Repetir considerando $n = 30$ y $n = 100$.
3. Sea $p = \mathbb{P}(X \leq 12)$ con $X \sim X_1$. Indicar cuál cuenta hay que hacer con la muestra (X_1, \dots, X_n) para estimar p . Es decir, proponer un estimador \hat{p}_n para p .
4. Considerar $n = 5$ datos duraciones de lámparas y calcular la estimación de p correspondiente a estos datos. Repetir considerando $n = 30$ y $n = 100$.

Recreo

¿Cuánto mide la mesa?



¿Cuánto mide la mesa?

Estas son las $n = 5$ primeras observaciones:

1.17, 1.36, 0.15, 2.52, 0.21, 1.78, 2.67

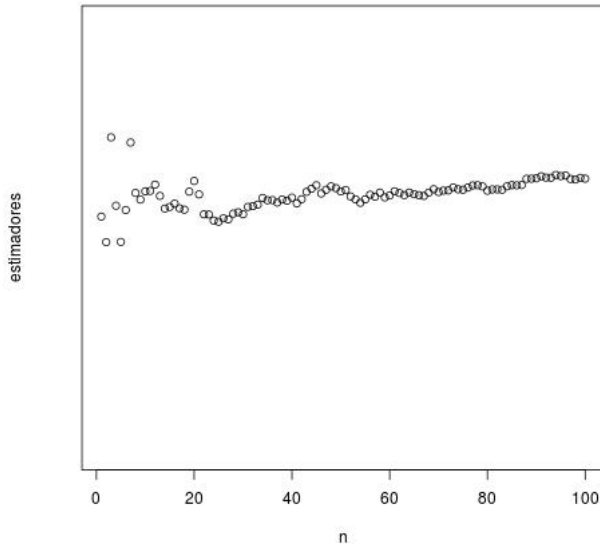
¿Qué hacemos?



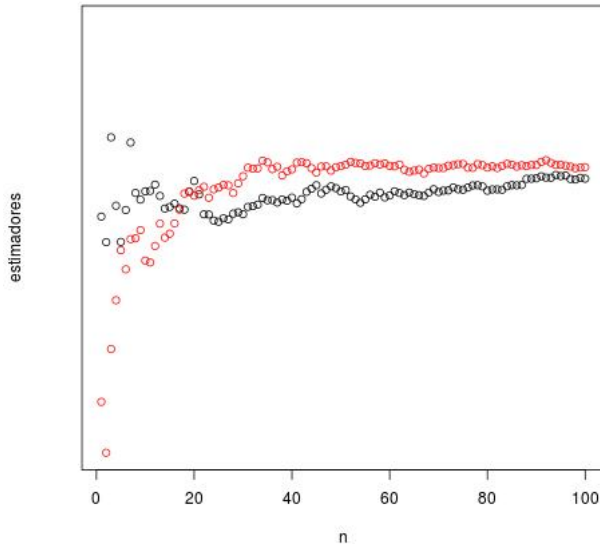
Modelo y Estimadores

Hacemos los puntos 1 y 2 de la Parte 3 de la Guía de actividades

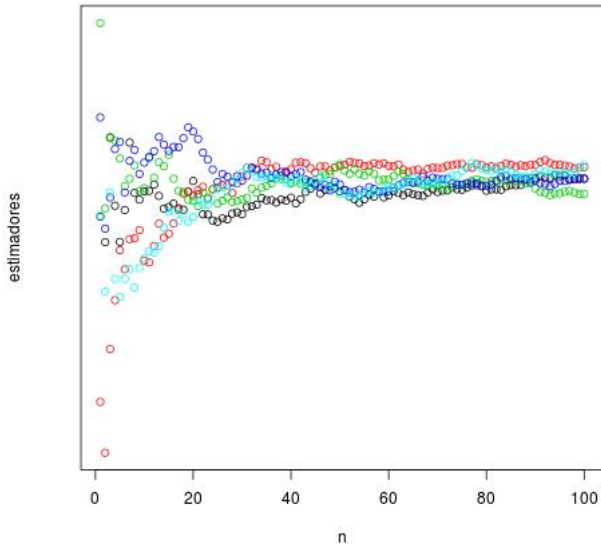
Juan cada vez con más datos. $\hat{\theta}_n = 2\overline{X}_n$



Juan y Andrea, cada vez con más datos. $\hat{\theta}_n = 2\bar{X}_n$



Varios, cada vez con más datos. $\hat{\theta}_n = 2\overline{X}_n$

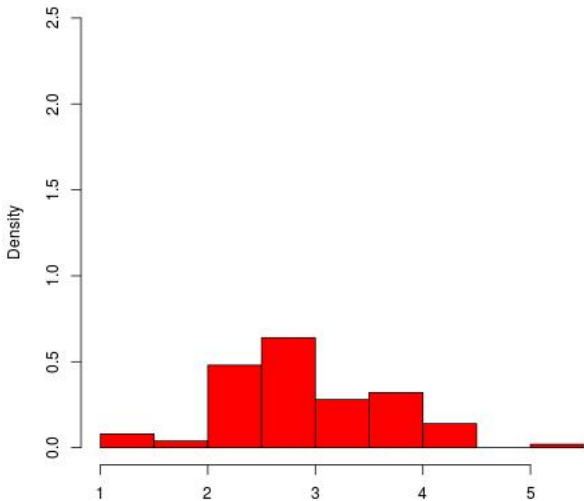


Cada uno con lo suyo. $\hat{\theta}_n = 2\bar{X}_n$

	Nombre	n=5	n=30	n=50
1	Juan	1.08	3.2	2.96
2	Andrea	2.87	2.95	2.88
3	Flor	3.47	3.2	3.18
4	Gonzalo	3.88	3.23	3.18
5	Paula	3.79	2.93	2.81
6	Agustin	3.01	2.9	2.59
7	Julieta	3.55	3.03	3.01
8	Marina	2.09	2.79	3.1
9	Pablo	4.14	3.41	3.01
10	Enrique	2.65	3.29	3.11
.
.
.
.

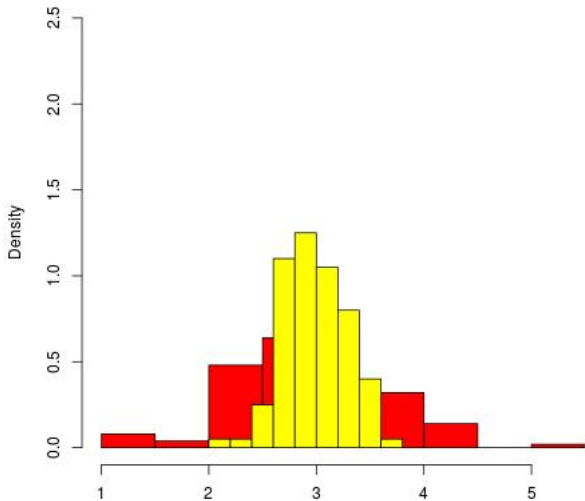
Histogramas de $\hat{\theta}_n = 2\overline{X}_n$

(empirical) Sampling Distribution of $\hat{\theta}_n$



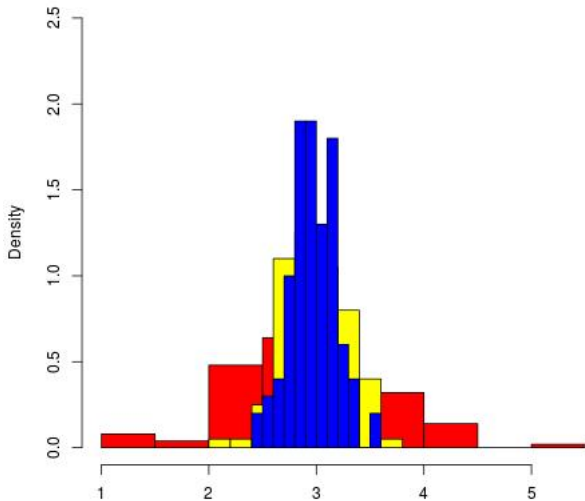
Histogramas de $\hat{\theta}_n = 2\overline{X}_n$

(empirical) Sampling Distribution of $\hat{\theta}_n$



Histogramas de $\hat{\theta}_n = 2\bar{X}_n$

(empirical) Sampling Distribution of $\hat{\theta}_n$



Estimación

Point estimation refers to providing a single “best guess” of some quantity of interest.

All of statistics. Wasserman

Estimación

Point estimation refers to providing a single “best guess” of some quantity of interest.

All of statistics. Wasserman

- translate(some quantity of interest)= Objeto de interés.
- some quantity of interest: largo de la mesa (θ)
- best guess: Estimador: *cuenta hecha con la muestra*
- best guess: Estimador: Función de la muestra

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

Estimación

Point estimation refers to providing a single “best guess” of some quantity of interest.

All of statistics. Wasserman

- translate(some quantity of interest)= Objeto de interés.
- some quantity of interest: largo de la mesa (θ)
- best guess: Estimador: *cuenta hecha con la muestra*
- best guess: Estimador: Función de la muestra

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

- Estimación: Valor del estimador en un conjunto de datos:

$$\hat{\theta}_n(x_1, \dots, x_n)$$

Notemos que el estimador ...

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

- $\hat{\theta}_n$ es una variable aleatoria.
- $\hat{\theta}_n$ tiene distribución (siempre).

Sampling distribution of $\hat{\theta}_n$: $f_{\hat{\theta}_n}$

- $\hat{\theta}_n$ tiene (en general) esperanza: $\mathbb{E}(\hat{\theta}_n)$

Notemos que el estimador ...

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

- $\hat{\theta}_n$ es una variable aleatoria.
- $\hat{\theta}_n$ tiene distribución (siempre).

Sampling distribution of $\hat{\theta}_n$: $f_{\hat{\theta}_n}$

- $\hat{\theta}_n$ tiene (en general) esperanza: $\mathbb{E}(\hat{\theta}_n) = \int u f_{\hat{\theta}_n}(u) du$
- $\hat{\theta}_n$ tiene (en general) varianza: $\mathbb{V}(\hat{\theta}_n)$
- $\hat{\theta}_n$ tiene (en general) desvío estándar.

$$se = se(\hat{\theta}_n) = \sqrt{\mathbb{V}(\hat{\theta}_n)} \quad \text{Standard error of } \hat{\theta}_n.$$

Veamos el shiny

Consistencia

A medida que aumenta el tamaño n de la muestra, el estimador se aproxima al objeto de interés.

$$\hat{\theta}_n \longrightarrow \theta , \quad \text{cuando } n \rightarrow \infty$$

Error cuadrático medio (ECM)

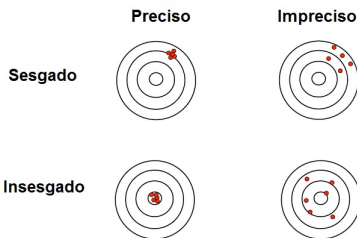
$$\text{ECM} : \mathbb{E} \left\{ (\hat{\theta}_n - \theta)^2 \right\}.$$

Lema: Si el ECM de un estimador converge a cero entonces vale la consistencia:

$$\mathbb{E} \left\{ (\hat{\theta}_n - \theta)^2 \right\} \longrightarrow 0 \quad \text{implica que} \quad \hat{\theta}_n \longrightarrow \theta .$$

Exactitud (In - Sesgado) - Precisión (Varianza)

Sesgo vs. Precisión



Sesgo - Bias

$$\text{Sesgo} : \mathbb{E}(\hat{\theta}_n) - \theta.$$

Insesgado: *El estimador $\hat{\theta}_n$ se dice insesgado si su sesgo vale cero*

$$\text{Insesgado} : \mathbb{E}(\hat{\theta}_n) - \theta = 0$$

En otras palabras, el estimador $\hat{\theta}_n$ se dice insesgado si su esperanza coincide con el valor de interés que queremos estimar:

$$\text{Insesgado} : \mathbb{E}(\hat{\theta}_n) = \theta$$

Propiedades

Lema: El error cuadrático medio de un estimador se descompone de la siguiente manera:

$$\text{ECM}(\hat{\theta}_n) = \mathbb{V}(\hat{\theta}_n) + \left\{ \mathbb{E}(\hat{\theta}_n) - \theta \right\}^2$$

En particular... Si

$$\mathbb{V}(\hat{\theta}_n) \rightarrow 0 \quad \text{y} \quad \mathbb{E}(\hat{\theta}_n) \rightarrow \theta$$

tenemos que ECM converge a cero, y por lo tanto el estimador es consistente:

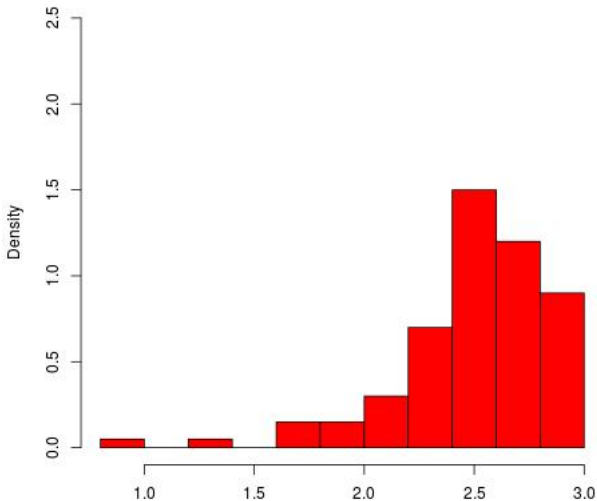
$$\hat{\theta}_n \longrightarrow \theta$$

Miremos todo en el ejemplo: $\hat{\theta}_n = 2\bar{X}_n$

- $(X_i)_{i \geq 1}$ i.i.d., $X_i \sim \mathcal{U}[0, \theta]$
- Objeto de interés: θ
- Estimador: $\hat{\theta}_n = 2\bar{X}_n$
- Distribución de $\hat{\theta}_n$?
- $\mathbb{E}(\hat{\theta}_n) = \theta$: Es insesgado
- $\mathbb{V}(\hat{\theta}_n) = \mathbb{V}(2\bar{X}_n) = 4\mathbb{V}(\bar{X}_n) = 4\frac{\mathbb{V}(X_1)}{n} = 4\frac{\theta^2/12}{n}$
- $\text{ECM}(\hat{\theta}_n) = 0^2 + 4\frac{\theta^2/12}{n}$

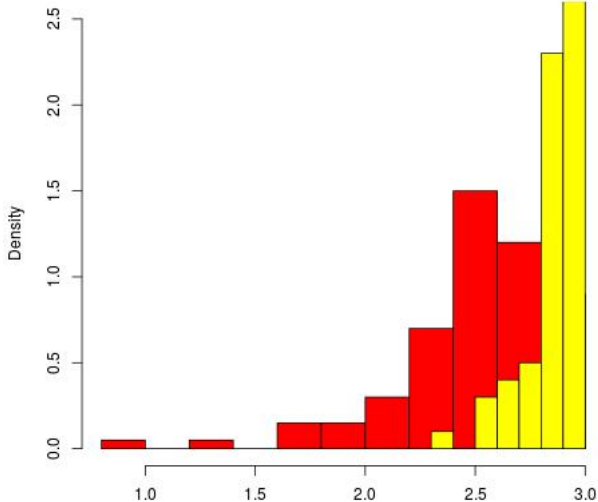
Histogramas de $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$

(empirical) Sampling Distribution of $\tilde{\theta}_n$



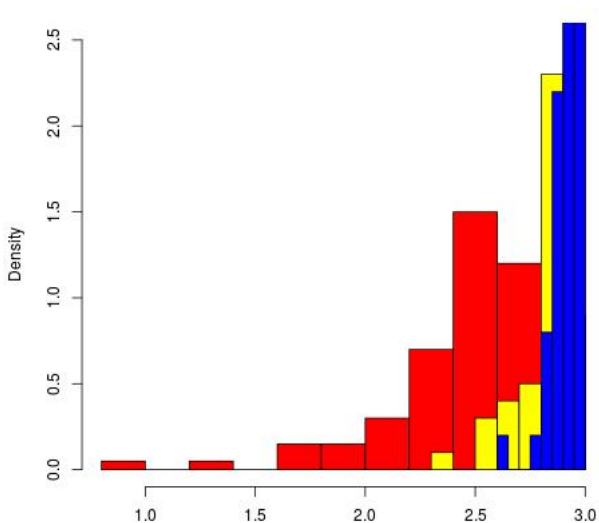
Histogramas de $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$

(empirical) Sampling Distribution of $\tilde{\theta}_n$



Histogramas de $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$

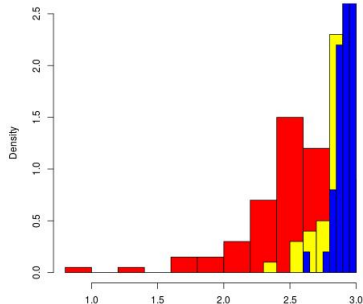
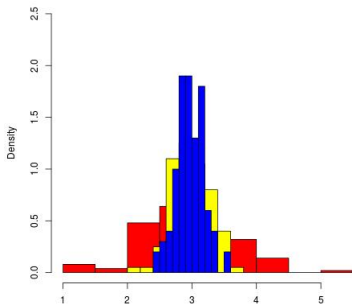
(empirical) Sampling Distribution of $\tilde{\theta}_n$



Miremos todo ahora para $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$

- $(X_i)_{i \geq 1}$ i.i.d., $X_i \sim \mathcal{U}[0, \theta]$
- Objeto de interés: θ
- Estimador: $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$
- Distribución de $\tilde{\theta}_n$?
- $\mathbb{E}(\tilde{\theta}_n)$?
- $\mathbb{V}(\tilde{\theta}_n)$?
- ECM $(\tilde{\theta}_n)$

Histogramas de $\hat{\theta}_n = 2\bar{X}_n$ y de $\tilde{\theta}_n = \max\{X_1, \dots, X_n\}$



Vamos guia.

Estadístico: Cuenta hecha con la muestra

$$h(X_1, \dots, X_n)$$

Estadística

POBLACION $\leftrightarrow F$	MUESTRA X_1, \dots, X_n i.i.d. $X_i \sim F$
Parámetro: Valor asociado de F $\theta = \theta(F)$ θ : valor poblacional	Estimador: estadístico para estimar θ $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$ $\hat{\theta}_n$ NUEVA VARIABLE ALEATORIA

Enfaticemos en que el estimador ... (si, ya la vimos!)

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

- $\hat{\theta}_n$ es una variable aleatoria.
- $\hat{\theta}_n$ tiene distribución (siempre)
- $\hat{\theta}_n$ tiene (en general) esperanza: $\mathbb{E}(\hat{\theta}_n)$
- $\hat{\theta}_n$ tiene (en general) varianza: $\mathbb{V}(\hat{\theta}_n)$

Consistencia

- $(X_i)_{i \geq 1}$ i.i.d., $X_i \sim F$, $F \in \mathcal{F}$
- \mathcal{F} : modelo estadístico.
- $\theta(F)$ objeto de interés definido para cada posible $F \in \mathcal{F}$
- estimador $\hat{\theta}_n = \hat{\theta}_n(X_1, \dots, X_n)$.
- Consistencia:

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

cuando $n \rightarrow \infty$, $X_i \sim F$, cualquiera sea $F \in \mathcal{F}$

A medida que aumenta el tamaño n de la muestra, el estimador se aproxima al objeto de interés.

$$\hat{\theta}_n \longrightarrow \theta, \quad \text{cuando } n \rightarrow \infty$$

Estimación: ejemplo

- X_1, \dots, X_n i.i.d. $X_i \sim X$.
- Parámetro de interés; $\sigma^2 = \mathbb{E}\{(X - \mu)^2\} = \mathbb{V}(X)$
- ¿Estimador?

Estimación: ejemplo

- X_1, \dots, X_n i.i.d. $X_i \sim X$.
- Parámetro de interés; $\sigma^2 = \mathbb{E}\{(X - \mu)^2\} = \mathbb{V}(X)$
- ¿Estimador?

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

Estimación: ejemplo

- X_1, \dots, X_n i.i.d. $X_i \sim X$.
- Parámetro de interés; $\sigma^2 = \mathbb{E}\{(X - \mu)^2\} = \mathbb{V}(X)$
- ¿Estimador?

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

- $\mathbb{E}(\hat{\sigma}_n) = (n - 1)n^{-1}\sigma^2$
- Estimador (insesgado) de la varianza: $S^2 = S_n^2$

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

- $\mathbb{E}[S^2] = \sigma^2$ (insesgado)
- $S^2 \rightarrow \sigma^2$ en probabilidad (consistencia)

Propiedades - si, de nuevo!, pero todas juntas.

- Consistencia

$\hat{\theta}_n(X_1, \dots, X_n) \rightarrow \theta(F)$ en probabilidad, cuando $X_i \sim F$

abreviado: $\hat{\theta} \rightarrow \theta$

- Error cuadrático medio: $\text{ECM} = \mathbb{E}\{(\hat{\theta}_n - \theta)^2\}$
- Lema: Si $\mathbb{E}\{(\hat{\theta}_n - \theta)^2\} \rightarrow 0$, entonces $\hat{\theta}_n \rightarrow \theta$
- Sesgo: $\mathbb{E}(\hat{\theta}_n) - \theta$.
- Estimador insesgado: Sesgo=0: $\mathbb{E}(\hat{\theta}_n) - \theta$
- Lema: $\mathbb{E}\{(\hat{\theta}_n - \theta)^2\} = \mathbb{V}(\hat{\theta}_n) + \left\{\mathbb{E}(\hat{\theta}_n) - \theta\right\}^2$
- Si $\mathbb{V}(\hat{\theta}_n) \rightarrow 0$ y $\mathbb{E}(\hat{\theta}_n) \rightarrow \theta$, entonces

$$\mathbb{E}\{(\hat{\theta}_n - \theta)^2\} \rightarrow 0$$