

# Análisis estadístico de datos simulados

## Estimadores

**Patricia Kisbye**

FaMAF

11 de mayo, 2010

# Análisis estadístico

## Inferencia estadística:

- ▶ Elegir una distribución en base a los datos observados.
- ▶ Estimar los parámetros de la distribución (EMV).
- ▶ Pruebas de bondad de ajuste.

## Estimación de parámetros

- ▶ Varianza del estimador.  $\text{Var}(\hat{\theta})$ .
- ▶ Error cuadrático medio del estimador.  $E[(\hat{\theta} - \theta)^2]$ .
- ▶ Estimadores por intervalo e intervalos de confianza.
- ▶ Pruebas de hipótesis. Nivel de significación  $\alpha$ . Valor  $p$ .

# Media muestral

Dadas  $n$  observaciones:  $X_1, X_2, \dots, X_n$ , con una misma distribución, la media muestral se define por

$$\bar{X}(n) = \frac{X_1 + X_2 + \dots + X_n}{n}.$$

La media muestral se utiliza como un estimador de la media  $\theta$ , es decir, de  $\theta = E[X_i]$ .

Estimador insesgado.

$$E[\bar{X}(n)] = E\left[\sum_{i=1}^n \frac{X_i}{n}\right] = \sum_{i=1}^n \frac{E[X_i]}{n} = \frac{n\theta}{n} = \theta.$$

## Error cuadrático medio

- ▶  $\hat{\theta}$ : estimador del parámetro  $\theta$  de una distribución  $F$
- ▶ Se define el error cuadrático medio (ECM) de  $\hat{\theta}$  con respecto al parámetro  $\theta$  como

$$ECM(\hat{\theta}, \theta) = E[(\hat{\theta} - \theta)^2].$$

$$\begin{aligned} E[(\hat{\theta} - \theta)^2] &= E[(\hat{\theta} - E[\hat{\theta}] + E[\hat{\theta}] - \theta)^2] \\ &= E[(\hat{\theta} - E[\hat{\theta}])^2] + (E[\hat{\theta}] - \theta)^2 \\ &= \text{Var}(\hat{\theta}) + (E(\hat{\theta}) - \theta)^2 \end{aligned}$$

- ▶ El error cuadrático medio de un estimador es igual a su varianza más el sesgo al cuadrado.
- ▶ Si el estimador es insesgado, su ECM es igual a la varianza.

## ECM de la media muestral respecto de la media

Muestra de  $X$ :  $X_1, X_2, \dots, X_n$ ,  $E[X_i] = \theta$

$$\begin{aligned} ECM(\bar{X}(n), \theta) &= E[(\bar{X}(n) - \theta)^2] \\ &= \text{Var}(\bar{X}(n)) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}(X_i) = \frac{\sigma^2}{n} \end{aligned}$$

La media muestral es un buen estimador de  $E[X]$  si  $\sigma/\sqrt{n}$  es pequeño.

- ▶ El ECM depende de la distribución de  $X_i$  y del tamaño de la muestra.
- ▶ Teorema central del límite. Si  $Z \sim N(0, 1)$  y  $n$  es grande:

$$P\left(\frac{|\bar{X}(n) - \theta|}{\sigma/\sqrt{n}} > c\right) \approx P\{|Z| > c\}.$$

# Varianza muestral

El indicador  $\frac{\sigma^2}{n}$  como estimación del error en la media muestral, tiene el inconveniente que  $\sigma$  es en general desconocida.

Para estimar la varianza se utiliza el estimador

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

- ▶ Estimador insesgado de la varianza
- ▶ Fórmula a utilizar:

$$E[S^2(n)] = \text{Var}(X)$$

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

# Varianza muestral

$$E[X_i^2] = \text{Var}(X_i) + (E[X_i])^2 = \sigma^2 + \theta^2.$$

$$E[\bar{X}^2(n)] = \frac{\sigma^2}{n} + \theta^2.$$

$$(n-1)E[S^2(n)] = nE[X_1^2] - nE[\bar{X}^2(n)] = n(\sigma^2 + \theta^2) - n\left(\frac{\sigma^2}{n} + \theta^2\right)$$

$$E[S^2(n)] = \sigma^2$$

Utilizaremos  $S(n) = \sqrt{S^2(n)}$  como estimador de la desviación estándar.

- ▶ Error del estimador  $\bar{X}(n)$ :  $\sigma^2/n$ .
- ▶ Simulación de datos: Si el objetivo es estimar la media, para disminuir el error deben generarse muestras de tamaño  $n$ ,  $n$  grande.

# Media muestral

- ▶ Elegir un valor aceptable  $d$  para la desviación estándar del estimador.
- ▶ Generar  $(n)$  datos hasta que  $\sigma/\sqrt{n} < d$ . ( $S/\sqrt{n} < d$ )
- ▶ Conviene generar al menos 100 datos para:
  - ▶ asegurar normalidad de la distribución de  $\bar{X}(n)$ .
  - ▶ para disminuir la varianza de  $S$ .
- ▶ La estimación de  $\theta$  estará dada por el último valor de  $\bar{X}(n)$ .
- ▶ El algoritmo implica calcular en cada paso  $\bar{X}(n)$  y  $S(n)$ .
- ▶ Es posible calcularlo recursivamente.



# Media muestral

## Cálculo recursivo de $\bar{X}(n)$ y $S^2(n)$

- ▶  $\bar{X}(1) = X_1$ ,
- ▶  $S^2(1) = 0$ .

$$\bar{X}(j+1) = \bar{X}(j) + \frac{X_{j+1} - \bar{X}(j)}{j+1}$$

$$S^2(j+1) = \left(1 - \frac{1}{j}\right) S^2(j) + (j+1)(\bar{X}(j+1) - \bar{X}(j))^2$$

## Estimación de una proporción

El estimador  $\bar{X}(n)$  puede utilizarse también para estimar la proporción de casos en una población.

$$X_i = \begin{cases} 1 & \text{probabilidad } p \\ 0 & \text{probabilidad } 1 - p. \end{cases}$$

- ▶  $\bar{X}(n)$  es un estimador insesgado de  $p$ .
- ▶  $E[(\bar{X}(n) - p)^2] = \text{Var}(\bar{X}(n)) = \frac{p(1 - p)}{n}$
- ▶ En este caso, se estima la varianza del estimador  $\bar{X}(n)$  por:

$$\frac{\bar{X}(n)(1 - \bar{X}(n))}{n}.$$

## Algoritmo: Cálculo de $E[X]$

---

Estimación de la media  $M$  de  $X$  con error  $d$

---

Generar  $X$ ,  $M \leftarrow X$

$$M = \bar{X}(1) = X_1;$$

$S^2 \leftarrow 0$

$$S^2 = S^2(1) = 0;$$

**for**  $1 < j \leq 100$  **do**

    Generar  $X$ ;  $A \leftarrow M$ ;

$M \leftarrow M + (X - M)/j$ ;

$S^2 \leftarrow (1 - 1/(j - 1))S^2 + j(M - A)^2$

**end**

$j \leftarrow 100$ ;

**while**  $\sqrt{S^2/j} > d$  **do**

$j \leftarrow j + 1$ ;

    Generar  $X$ ;

$A \leftarrow M$ ;

$M \leftarrow M + (X - M)/j$ ;

$S^2 \leftarrow (1 - 1/(j - 1))S^2 + j(M - A)^2$

**end**

**return**  $M$

---

# Algoritmo: Cálculo de una probabilidad

---

Estimación de la probabilidad  $p$  de  $X$  con error  $d$

---

Generar  $X$   $X$  es 0 o 1;

$p \leftarrow X$ ;

**for**  $1 < j \leq 100$  **do**

Generar  $X$ ;

$p \leftarrow p + (X - p)/j$

**end**

$j \leftarrow 100$ ;

**while**  $\sqrt{p(1-p)/j} > d$  **do**

$j \leftarrow j + 1$ ;

Generar  $X$ ;

$p \leftarrow p + (X - p)/j$ ;

**end**

**return**  $p$

---

# Estimador por intervalos

Un estimador por intervalo de un parámetro es un intervalo para el que se predice que el parámetro está contenido en él.

La confianza que se da al intervalo es la probabilidad de que el intervalo contenga al parámetro.

## Estimador por intervalo de la media poblacional

- ▶  $\bar{X}(n)$  es un estimador puntual de la media.
- ▶ Si la población es normal con media  $\theta$  y d.s.  $\sigma$ ,

$$\frac{\bar{X}(n) - \theta}{\sigma/\sqrt{n}} \sim Z = N(0, 1)$$

- ▶  $P(Z > z_\alpha) = \alpha$ , para  $0 < \alpha < 1$ .
- ▶ Si el nivel de confianza deseado es  $1 - \alpha$ , utilizamos  $\pm z_{\alpha/2}$ .
- ▶ Ejemplo: nivel de confianza del 95%:  $\alpha = 0.025$ , y  $z_\alpha = 1.96$ .

## Estimador por intervalos

$$P\left(\frac{|\bar{X}(n) - \theta|}{\sigma\sqrt{n}} \leq 1.96\right) = 0.95.$$

$$P\left(\bar{X}(n) - 1.96\frac{\sigma}{\sqrt{n}} \leq \theta \leq \bar{X}(n) + 1.96\frac{\sigma}{\sqrt{n}}\right) = 0.95.$$

- ▶ El intervalo con extremos

$$\bar{X}(n) - 1.96\sigma/\sqrt{n} \quad \text{y} \quad \bar{X}(n) + 1.96\sigma/\sqrt{n}$$

se dice que es un **estimador por intervalo**, con un 95% de confianza para la media  $\theta$ .

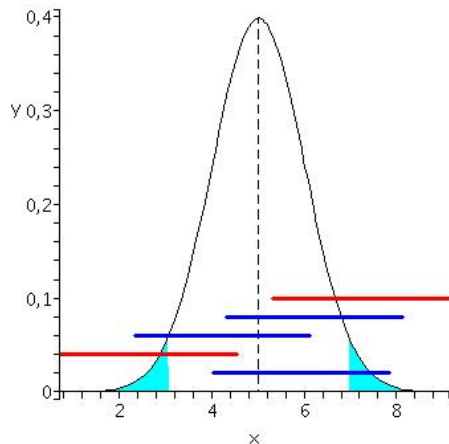
- ▶ Si  $\bar{x}$  es un valor observado de  $\bar{X}(n)$ , el intervalo con extremos

$$\bar{x} - 1.96\sigma/\sqrt{n} \quad \text{y} \quad \bar{x} + 1.96\sigma/\sqrt{n}$$

es el **valor estimado** del estimador por intervalo de  $\theta$ , con un 95% de confianza.

# Estimador por intervalos

Intervalos de confianza  
del 95%



- ▶  $(\bar{X} - \frac{1.96\sigma}{\sqrt{n}}, \bar{X} + \frac{1.96\sigma}{\sqrt{n}})$ .
- ▶  $z_{0.025} = 1.96$ .
- ▶ El 95% de los intervalos cubren la media.

## Estimador por intervalos

- ▶ Si la varianza  $\sigma^2$  es desconocida, utilizamos el estimador  $S^2(n)$ .
- ▶ Para determinar un intervalo de confianza, es necesario conocer la distribución del estadístico:

$$\sqrt{n} \frac{\bar{X}(n) - \theta}{S(n)}$$

### Distribuciones derivadas de la normal

- ▶  $\chi^2$  de Pearson con  $k$  grados de libertad: si  $Z_1, Z_2, \dots, Z_k$  son v.a.  $N(0,1)$ , independientes:

$$\chi_k^2 = Z_1^2 + \dots + Z_k^2$$

- ▶  $T_k$  de Student, con  $k$  grados de libertad: (W. S. Gosset)

$$T_k = \frac{Z}{\sqrt{\frac{\chi_k^2}{k}}}$$



# Intervalos de confianza

- ▶ El estadístico tiene una distribución  $T_{n-1}$ :

$$\sqrt{n} \frac{\bar{X}(n) - \theta}{S(n)} \sim T_{n-1}$$

- ▶ Sea  $t_\alpha$  tal que  $P(|T_{n-1}| > t_\alpha) = 1 - \alpha$ .

$$P\left(\bar{X}(n) - t_{\alpha/2} \frac{S(n)}{\sqrt{n}} \leq \theta \leq \bar{X}(n) + t_{\alpha/2} \frac{S(n)}{\sqrt{n}}\right) = 1 - \alpha.$$

- ▶ Para  $n > 120$ , puede usarse la distribución normal, es decir,  $t_\alpha \approx z_\alpha$ .

# Intervalos de confianza para proporciones

- ▶  $X_1, X_2, \dots, X_n$ : Bernoulli, independientes, con probabilidad  $p$  de éxito.
- ▶ Para  $n$  suficientemente grande tal que  $np$  y  $n(1 - p)$  es mayor que 5,

$$X_1 + \dots + X_n = Bi(n, p) \sim N(np, np(1 - p)).$$

- ▶ Si  $p$  es desconocido, podemos estimar  $p$  con la media muestral:

$$\hat{p} = \bar{X}(n) \quad y \quad \text{Var}(\hat{p}) = \frac{\hat{p}(1 - \hat{p})}{n}.$$

- ▶ Intervalos de confianza del  $100(1 - \alpha)\%$ :

$$\left( \hat{p} - z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}, \hat{p} + z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}} \right)$$

# Longitud del intervalo de confianza

- ▶ Estimación de la media:  $s(n)$ : valor observado de la varianza muestral.

$$2 \frac{z_{\alpha/2} \sigma}{\sqrt{n}} \quad \text{o} \quad 2 \frac{z_{\alpha/2} s(n)}{\sqrt{n}}.$$

- ▶ Estimación de la proporción:

$$2z_{\alpha/2} \sqrt{\frac{\hat{p}(1 - \hat{p})}{n}}$$

- ▶ La longitud del intervalo de confianza al  $100(1 - \alpha)\%$  depende del tamaño de la muestra.