

# Probabilidad y estadística

## Clase 4

# Estadística

# ¿Qué es la estadística?

Ahora lo que tenemos son observaciones (realizaciones) de las variables aleatorias.

El objetivo es poder hacer algún tipo de inferencia a partir de los valores observados.

# ¿Qué es la estadística?

## 6.1 Introduction

Statistical inference, or “learning” as it is called in computer science, is the process of using data to infer the distribution that generated the data. A typical statistical inference question is:

$$\textcircled{H} \ni \emptyset$$

Given a sample  $X_1, \dots, X_n \stackrel{iid}{\sim} F$ , how do we infer  $F$ ?

In some cases, we may want to infer only some feature of  $F$  such as its mean.

# Muestra aleatoria

Supongamos que tenemos un experimento aleatorio relacionado con una v.a.  $X$ . La v.a.  $X$  representa un **observable** del experimento aleatorio.

Los valores de  $X$  son la **población** de estudio, y el objetivo es saber como se comporta esa población.

Una **muestra aleatoria** de tamaño  $n$ , es una sucesión de  $n$  v.a **independientes**  $\underline{X}_n = (X_1, \dots, X_n)$ , tal que  $X_i \stackrel{iid}{\sim} X$

# Estimadores puntuales

# Estimadores puntuales

$$\mathcal{N}(\underbrace{\mu}_0, \underbrace{\sigma^2}_1)$$

$$\hat{\theta}$$

$$\underline{\theta} \in \Theta$$

$$\underline{\theta}: \mathcal{X} \rightarrow \mathcal{R}$$

$$x_i \sim \mathcal{X}$$

**Def:** Dada una muestra aleatoria  $\underline{X}_n$ , un estimador  $\hat{\theta}$  de un parámetro  $\underline{\theta}$  es una función de la muestra aleatoria que provee un valor aproximado del parámetro o característica desconocido.

# Bondades de los estimadores

$$\hat{\theta}_n = \rho(x_n)$$

$$B := E[\hat{\theta} - \theta] = 0 \quad E[\hat{\theta}] - E[\theta] = E[\hat{\theta}] - \theta = 0 \quad \theta \Rightarrow E[\hat{\theta}] = \theta$$

**Def:**  $\hat{\theta}$  es un estimador **insesgado** para  $\theta$  si  $B := \mathbb{E}[\hat{\theta} - \theta] = 0 \quad \forall \theta$ .  
En caso contrario diremos que es **sesgado**. A B se lo conoce como **sesgo**

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

**Def:**  $\hat{\theta}$  es un estimador **asintóticamente insesgado** para  $\theta$  si

$$\lim_{n \rightarrow \infty} \mathbb{E}[\hat{\theta} - \theta] = 0 \quad \forall \theta$$

$$x_n, \quad \hat{\theta}_n \xrightarrow{n \rightarrow \infty} \theta$$



# Bondades de los estimadores

$$\hat{\theta} - \theta$$

$$ECM(\hat{\theta}) =$$

**Def:** El error cuadrático medio (ECM) como  $\mathbb{E}[(\hat{\theta} - \theta)^2]$

**Obs:** El ECM se puede descomponer como:

$$ECM = \underbrace{var(\hat{\theta})}_{\text{varianza}} + \underbrace{B(\hat{\theta})^2}_{\text{sesgo}}$$

donde  $var(\hat{\theta}) = \mathbb{E}[(\hat{\theta} - \mathbb{E}[\hat{\theta}])^2]$  y  $B = \mathbb{E}[\hat{\theta} - \theta]$

**Def:** Un estimador  $\theta^*(\underline{X}_n)$  es **óptimo** (en media cuadrática) si  $ECM(\theta^*) \leq ECM(\hat{\theta})$  para todo  $\hat{\theta}(\underline{X}_n)$

$$\begin{aligned}
 E(M(\hat{\theta})) &= E[(\hat{\theta} - \theta)^2] = E[(\hat{\theta} - \theta - \underbrace{E[\hat{\theta}] + E[\theta]}_{=0})^2] = \\
 &= E\left[\left(\underbrace{(\hat{\theta} - E(\hat{\theta}))}_{\text{var}(\hat{\theta})} + \underbrace{(E(\hat{\theta}) - \theta)}_{E(\hat{\theta}) - E(E(\hat{\theta})) = E(\hat{\theta}) - E(\hat{\theta}) = 0}\right)^2\right] = \\
 &= E[(\hat{\theta} - E(\hat{\theta}))^2] + 2 \underbrace{E[(\hat{\theta} - E(\hat{\theta}))(E(\hat{\theta}) - \theta)]}_{=0} + \\
 &\quad + E[(\underbrace{E(\hat{\theta}) - \theta}_{E[B^2] = B^2})^2] = \text{var}(\hat{\theta}) + B^2(\hat{\theta})
 \end{aligned}$$

# Bondades de los estimadores

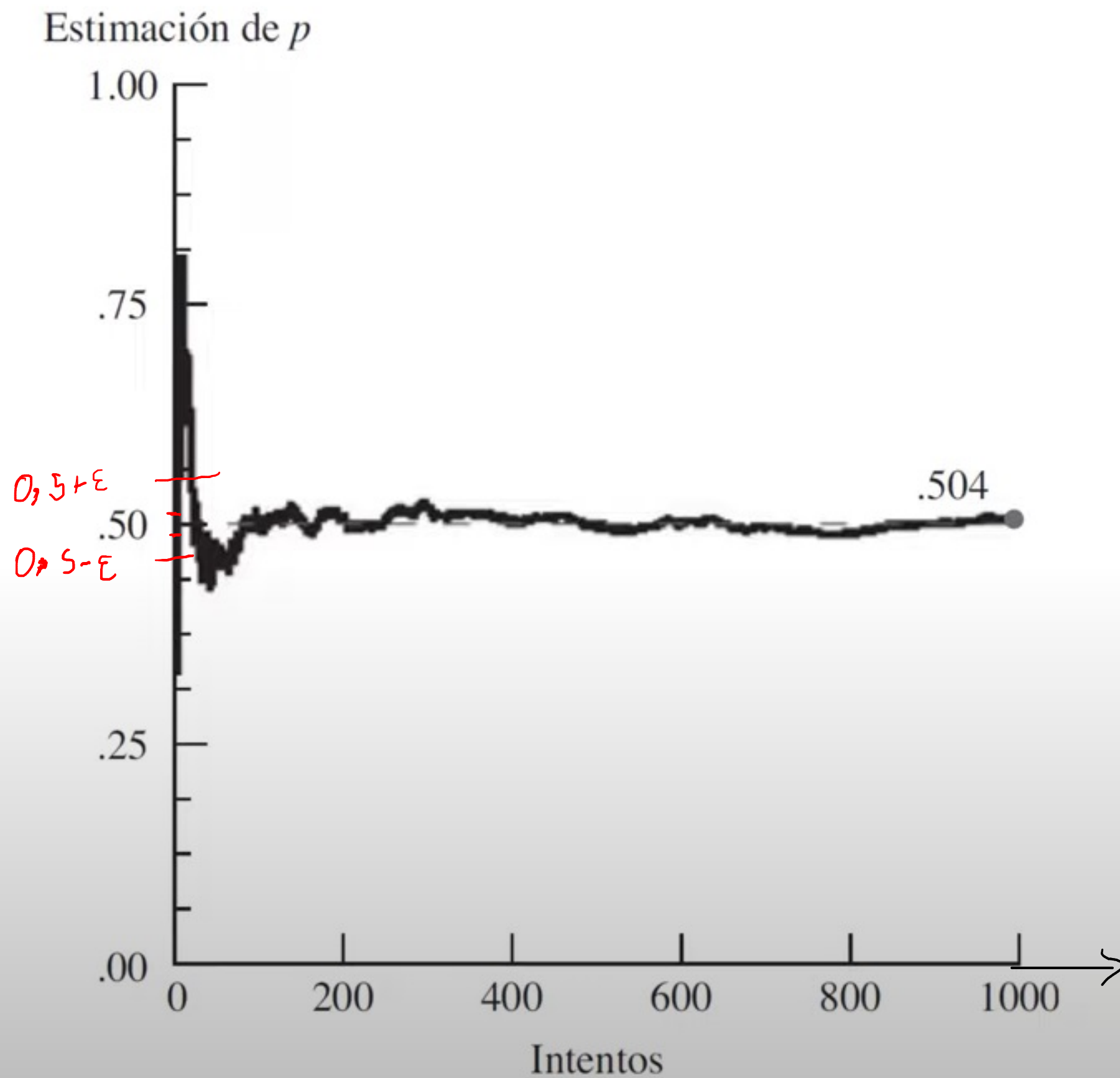
**Def:** Dada una sucesión de estimadores  $\hat{\theta}_n$  de  $\theta$ , diremos que  $T = \hat{\theta}$  es (débilmente) **consistente** si  $\forall \varepsilon > 0, \mathbb{P}(|T - \theta| > \varepsilon) \rightarrow 0$

**Teorema:** Si  $\text{var}(\hat{\theta}) \rightarrow 0$  y  $\mathbb{E}[\hat{\theta}] \rightarrow \theta$ , entonces  $\hat{\theta}$  es consistente.

**Def:** Un estimador es **consistente en media cuadrática** si

$$\lim_{n \rightarrow \infty} ECM(\hat{\theta}_n) = 0, \forall \theta$$





¿Cómo se interpreta este gráfico?

# Ejercicio 1

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$X_i \stackrel{iid}{\sim} X ; E(X_i) = \mu, \forall i$   
 $\mu = 0$   
 $\theta = \theta_1$   
 $X_i \sim X$   
 $X \sim N(\mu, \theta)$

Se desea estimar la media de una variable con distribución  $N(\mu, \theta)$  a partir del promedio de n realizaciones. Analizar las bondades de las que goza dicho estimador.

a)  $E(\bar{X}) \stackrel{\text{linealidad}}{=} \frac{1}{n} \sum_{i=1}^n E(X_i) = \frac{1}{n} \sum_{i=1}^n \mu = \frac{1}{n} (\mu + \mu + \dots + \mu) = \frac{1}{n} n\mu = \mu$

$\theta = 0$

b)  $\text{Var}(\bar{X}) = \text{Var}\left(\frac{1}{n} \sum_{i=1}^n X_i\right) = \frac{1}{n^2} \text{Var}\left(\sum_{i=1}^n X_i\right) \stackrel{\text{Indep.}}{\stackrel{\text{linealidad}}{=}} \frac{1}{n^2} \sum_{i=1}^n \theta = \frac{\theta}{n} \rightarrow 0$

$n \rightarrow \infty$

Por a) y b)  $\bar{X}$  es consistente

# Algunos comportamientos límite

LGN + TCL

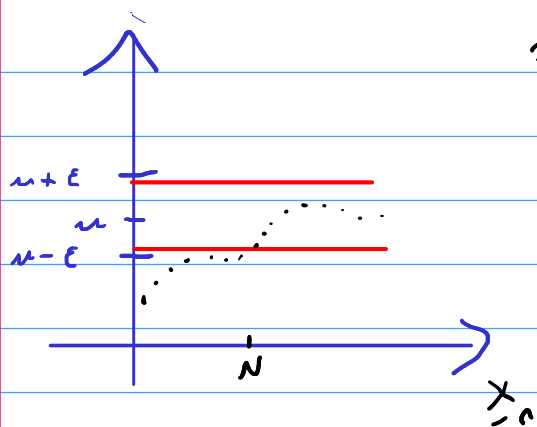
# Ley (débil) de los grandes números

Sea  $\underline{X} = (X_1, \dots, X_n)$  una muestra aleatoria, con  $X_i$  i.i.d. con media  $\mu$  y variancia  $\sigma^2$ . Para  $\epsilon > 0$  se tiene que

$$\mathbb{P}(|\bar{X}_n - \mu| > \epsilon) \xrightarrow{n \rightarrow \infty} 0$$

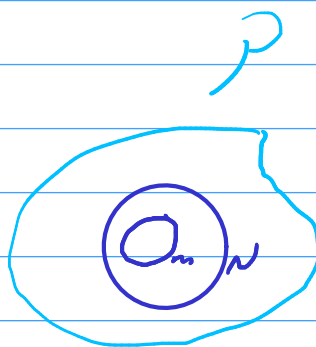
Es decir, a medida que crece la cantidad de muestras, el promedio  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$  tiende a la media real de la distribución.





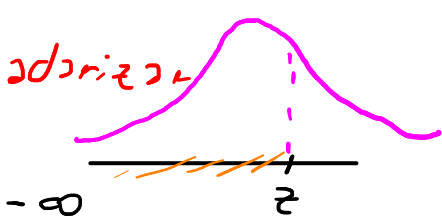
$$\underline{x}_n$$

$$\overline{x} \xrightarrow{n \rightarrow \infty} \mu$$



# Teorema central del límite

Sea  $\underline{X} = (X_1, \dots, X_n)$  una muestra aleatoria, con  $X_i$  i.i.d. con media  $\mu$  y varianza  $\sigma^2 < \infty$ . Sea  $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ , luego

$$Z_n = \frac{\bar{X}_n - \mu}{\sqrt{\mathbb{V}(\bar{X}_n)}} = \frac{\sqrt{n}(\bar{X}_n - \mu)}{\sigma} \rightsquigarrow Z \quad \left. \begin{array}{l} \sim 0 \\ \text{Estad. normal} \end{array} \right\}$$


con  $Z \sim \mathcal{N}(0, 1)$ , o equivalentemente:

$$\lim_{n \rightarrow \infty} \mathbb{P}(Z_n \leq z) = \Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

*f. densidad*

Es decir,  $\bar{X}_n$  tiene una distribución aproximadamente Normal, con media  $\mu$  y variancia  $\sigma^2/n$ . Para  $n$  finito, podemos aproximar

Ejercicio

$$\mathbb{P}\left(\frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n}\sigma} \leq z\right) \approx \Phi(z), \quad \mathbb{P}\left(\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \leq z\right) \approx \Phi(z),$$

# Enfoque frecuentista

# Estadístico suficiente

$$F_{\theta}(x), \theta \in \Theta$$

$$\lambda e^{-\lambda \sum_{i=1}^n x_i} \mathbb{I}\{x_i \geq 0\} \quad \theta = \lambda \quad T = \sum_{i=1}^n x_i \quad \mathbb{I}\{x_i \geq 0\}$$

**Def:** Dada una muestra aleatoria  $\underline{X}_n$ , un **estadístico** es cualquier función  $T_n = T(\underline{X}_n)$

**Def:** Sea una muestra aleatoria  $\underline{X}_n$ , cuya distribución es  $F_{\theta}(x), \theta \in \Theta$ , se dice que  $T = r(\underline{X}_n)$  es un **estadístico suficiente** para  $\theta$  si  $F_{\underline{X}|T=t}(x)$  no depende de  $\theta$ .

**Teorema de factorización:** Diremos que  $T = r(\underline{X}_n)$  es un est. suficiente para  $\theta$  sii existen funciones  $h$  y  $g$  tales que:

$$f_{\theta}(x) = \underbrace{g(r(x), \theta)}_{\text{depende de } \theta} \underbrace{h(x)}_{\text{no depende de } \theta} \quad g(r(x), \theta) = \quad h(x) =$$

# Ejercicio 2

$$X, Y \sim \mathcal{E}(\lambda) ; f_{X,Y}(x,y) =$$

$$X, Y \text{ iid} \quad f_{X,Y} = f_X \cdot f_Y$$

1. Hallar un estadístico suficiente para el parámetro de una v.a. con distribución exponencial
2. Hallar un estadístico suficiente para el parámetro de una v.a. con distribución de Bernoulli
3. Hallar un estadístico suficiente para el parámetro de una v.a. con distribución Uniforme en el intervalo  $(0, \theta)$

$$\begin{aligned} 1) f_{\lambda}(x_1, \dots, x_n) &= \prod_{i=1}^n \lambda e^{-\lambda x_i} \mathbb{I}\{x_i \geq 0\} = \lambda e^{-\lambda x_1} \cdot \lambda e^{-\lambda x_2} \cdot \dots \cdot \lambda e^{-\lambda x_n} \mathbb{I}\{x_i \geq 0, \forall i\} \\ &= \lambda^n e^{-\lambda \sum_{i=1}^n x_i} \prod_{i=1}^n \mathbb{I}\{x_i \geq 0\} \quad (x_i)^n \end{aligned}$$

$$Y \sim \text{Ber}(p) ; \quad p^x (1-p)^{1-x} \mathbb{I}\{x=0, x=1\}$$

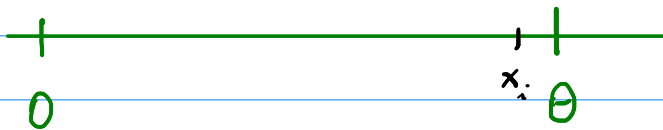
$$2) \quad f_p(\underline{x}) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \mathbb{I}\{x_i \in \{0,1\}\} = \underbrace{p^{\sum_{i=1}^n x_i}}_{\sum_{i=1}^n x_i} (1-p)^{n - \sum_{i=1}^n x_i} \underbrace{\prod_{i=1}^n \mathbb{I}\{x_i \in \{0,1\}\}}_{\prod_{i=1}^n \mathbb{I}\{x_i \in \{0,1\}\}}$$

$$Z \sim \mathcal{U}(0, \theta) , \quad \frac{1}{b-a} \mathbb{I}\{a < x < b\}$$

$$3) \quad f_{\theta}(\underline{x}) = \prod_{i=1}^n \frac{1}{\theta} \cdot \underbrace{\mathbb{I}\{x_i > 0 \wedge x_i < \theta\}}_{\mathbb{I}\{x_i > 0 \wedge x_i < \theta\}} = \frac{1}{\theta^n} \prod_{i=1}^n \mathbb{I}\{x_i < \theta\} \prod_{i=1}^n \mathbb{I}\{x_i > 0\}$$

$$= \frac{1}{\theta^n} \mathbb{I}\{\max_{i=1, \dots, n} x_i < \theta\} \cdot \mathbb{I}\{\min_{i=1, \dots, n} x_i > 0\}$$

$$V = \max(x_i)$$



# Método de máxima verosimilitud

# Método de Máxima Verosimilitud

**Def:** Diremos que  $\hat{\theta}(\underline{X})$  es un Estimador de Máxima Verosimilitud (MLE) si se cumple que:

$$f(\underline{X}, \hat{\theta}) = \max_{\theta} f_{\theta}(\underline{X})$$



La idea es que si observé una determinada muestra, entonces esta debería tener alta probabilidad de ocurrir, por lo tanto busco el  $\theta$  que maximiza esa probabilidad de ocurrencia



# Método de Máxima Verosimilitud

**Def:** Definimos la función de verosimilitud como

$L(\theta) = f(\underline{x}, \theta)$  (vista como función de  $\theta$ ) luego,

$$\hat{\theta} = \arg \max_{\theta \in \Theta} L(\theta).$$

Si el soporte de  $X$  no depende  $\theta$ ,  $\Theta$  es un conjunto abierto y  $f_{\theta}(x)$  es derivable respecto de  $\theta$ , entonces para hallar el EMV puedo hallar  $\theta$  tal que

$$\frac{\partial \ln(L(\theta))}{\partial \theta} = 0$$

# Ejercicio 3

La probabilidad de acertar a un blanco es  $p$ . Se realizan 10 tiros independientes, en los cuales se observaron 4 aciertos. A partir este valor observado estimar el valor de  $p$  por MV.

$$L(p) = f(x, p) = p^4 (1-p)^6$$

$$\ln(L(p)) = 4 \ln p + 6 \ln(1-p) ; \frac{\partial(\ln(L(p)))}{\partial p} = 0 ; \hat{p} = 0,4$$

$$\frac{\partial(\ln(L(p)))}{\partial p} = \frac{4}{p} - \frac{6}{1-p} = 0 \Rightarrow \hat{p} = 0,4$$

# Ejercicio 4

1. Sea  $X \sim U(0, \theta)$ . Hallar el EMV para  $\theta$  basado en una muestra de tamaño  $n$ .
2. Si en una muestra de tamaño 3 se observaron los valores (0.756, 1.1, 3.11), hallar el valor estimado de  $\theta$ .

# Principio de invariancia

Supongamos que ahora queremos estimar por máxima verosimilitud a  $\lambda = q(\theta)$ .

**Teorema:** Si  $\hat{\theta}$  es MLE de  $\theta$ , entonces  $\hat{\lambda} = q(\hat{\theta})$ .

¿Por qué es útil? Por ejemplo, podría querer estimar una probabilidad de la v.a.  $X$ , que en general no puedo porque desconozco el parámetro de la distribución.

# Ejercicio 5

Siguiendo el ejercicio 3, estimar la probabilidad de que se necesiten al menos 2 tiros para observar el primer acierto.

$$Z \sim G(p) \quad , \quad \hat{p} = 0,4$$

$$P(Z=z) = (1-p)^{z-1} p$$

$$= (1 - 0,4) = 0,6$$

# Estimadores de cuadrados mínimos

# Estimador de cuadrados mínimos

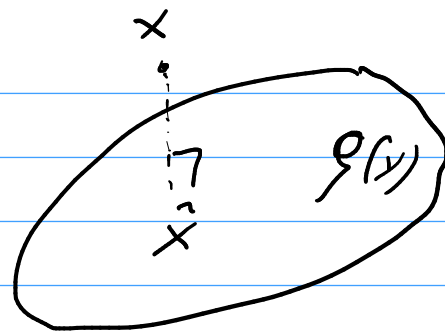
Es común querer estimar el valor de una v.a.  $X$  a partir de una medición  $Y$ . Ejemplo:  $Y$  es una versión ruidosa de  $X$ .

Buscamos un estimador  $\hat{X}(y)$  de  $X$  tal que tenga mínimo error cuadrático medio

$$\text{ECM} = \mathbb{E}[(X - \hat{X})^2 \mid Y]$$

Observar que se corresponde con la distancia asociada al p.i. canónico para v.a.

$$\langle x, x \rangle = \|x - \bar{x}\|^2$$





# Estimador de mínimos cuadrados

En otras palabras, queremos  $\hat{X} = g^*(Y)$  tal que

$$\mathbb{E}[(X - \hat{X})^2 | Y] \leq \mathbb{E}[(X - g(Y))^2 | Y] \quad \forall g(Y) \text{ (medible)}$$

¿Quién era  $\hat{X}$ ?  $\hat{X} = \mathbb{E}[X|Y]$

Idea de demostración: [Ejercicio]

1. Probar que el mejor estimador constante es  $\mathbb{E}[X]$
2. Probar que el mejor estimador condicional es  $\mathbb{E}[X|Y = y]$ .
3. Dejar que  $Y$  tome todos los valores posibles (i.e. reemplazo  $y$  por  $Y$ ), recupero la esperanza condicional.

$$1) E[(x - c)^2 | Y] \Rightarrow \frac{d}{dc} E[(x - c)^2] = 0$$

$$2. E[x - c] = 0$$

$$E[x] = c$$

$$2) E[(x - \hat{x}(y))^2 | Y=y] \Rightarrow E[x | Y=y] = \hat{x}(y)$$

Isol q' ctté

$$3) E[x | Y] = \hat{x}(y)$$

# Mínimos cuadrados: caso lineal

$$X|Y \sim N(E_{X|Y}, \Sigma_{\epsilon})$$

A veces obtener  $\mathbb{E}[X|Y]$  puede ser muy complicado, entonces nos restringimos a los estimadores lineales.

Buscamos  $a, b$  tq  $\mathbb{E}[(X - \underbrace{(aY + b)}_{\hat{x}(Y)})^2 | Y]$  sea mínima.

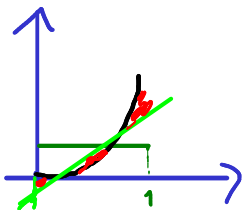
Resulta que  $a = \frac{\text{cov}(X,Y)}{\text{var}(Y)}$  y  $b = \frac{\text{cov}(X,Y)}{\text{var}(Y)} \mathbb{E}[Y] + \mathbb{E}[X]$

**Obs:** Si se asume que  $X$  e  $Y$  son conjuntamente gaussianas, el estimador de mínimos cuadrados (la esperanza condicional) coincide con el estimador de mínimos cuadrados asumiendo un modelo lineal y puede obtenerse a partir del modelo condicional de  $\mathbb{E}[X | Y]$

$$\hat{y} = \frac{\text{cov}(x, y)}{\text{var}(x)} (x - Ex) + Ey$$

# Ejercicio 6

Sea  $X \sim U(0,1)$  e  $Y=X^2$ . Hallar la mejor aproximación lineal de  $Y$  basada en  $X$ . Comparar con la mejor estimación de  $Y$  basada en  $X$ .



$$\hat{Y} = aX - b$$

$$\text{cov}(X, Y) = E[XY] - E[X]E[Y]$$

$$\begin{aligned} E[XY] &= E[E[XY|X]] = E[XE[X^2]] = E[X^3] \\ &= \int_0^1 x^3 dx = \frac{1}{4} \end{aligned}$$

$$EX = \int_0^1 x dx = \frac{1}{2} \quad EY = EX^2 = \int_0^1 x^2 dx = \frac{1}{3}$$

$$\text{Var}[X] = \frac{1}{12} \quad , \quad \text{Cov}(x, y) = \frac{1}{4} - \frac{1}{6} = \frac{1}{12}$$

$$\hat{y} = (x - EX) + EY$$

$$\hat{y} = (x - \frac{1}{2}) + \frac{1}{3} = x - \frac{1}{6}$$

# Regresión lineal (OLS)

$$\langle a, b \rangle = a^T b$$

$$\|a\|^2 = \sum_{i=1}^n a_i^2$$

Contamos con  $n$  observaciones conjuntas  $(y_1, x_1), (y_2, x_2), \dots, (y_n, x_n)$   
el modelo de regresión lineal es

donde  $n \times 1$   $n \times (p+1)$

$$Y = X\beta + \varepsilon \quad \text{con} \quad \mathbb{E}(\varepsilon_i | X_i) = 0 \quad \varepsilon_i \sim \mathcal{N}(0, \sigma^2)$$

$$Y = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}, \quad X = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_n \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

El estimador de cuadrados mínimos de  $\hat{\beta}$  está dado por

$$\hat{\beta} = (X^T X)^{-1} X^T Y$$

**Obs:** Si se asume que  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$  y son i.i.d., este estimador coincide con el estimador de máxima verosimilitud.

# Bibliografía



# Bibliografía

- Notas de estadística. Graciela Boente y Víctor Yohai
- “All of Statistics : A Concise Course in Statistical Inference”, Larry Wasserman, Springer-Verlag New York Inc.
- “Mathematical statistics with applications”, Mendenhall, W., Scheaffer, R. L., & Wackerly, D. D. Boston: Duxbury Press.