

Inferencia Estadística

G1: Distribución en el muestreo y principios de reducción

Gabriel Martos
Matías Pérez

Email: gmartos@utdt.edu
Email: lic.matiasdperez@gmail.com

Listado de ejercicios teórico-prácticos

1. Ejercicio de investigación: Piensa o busca ejemplos (en particular en Economía) de uso de modelos estadísticos paramétricos y no paramétricos. ¿Qué ventajas tienen los modelos no paramétricos por sobre los paramétricos? ¿Se te ocurren ventajas en el sentido contrario?
2. El soporte de $X \sim f(x; \theta)$ es el conjunto definido como:

$$\text{Soporte}(f(x; \theta)) = \{x \in \mathbb{R} \mid f(x; \theta) > 0\}.$$

- (a) ¿El conjunto $\text{Soporte}(f(x; \theta))$ puede depender del parámetro θ si $f(x; \theta)$ pertenece a la familia exponencial? (Formaliza tu respuesta utilizando la definición de familia Exponencial vista en clase).
- (b) Demostrar que la familia exponencial se puede escribir, de manera equivalente a la expresión que dimos en clase, como sigue:

$$f(x; \theta) = \exp \left(w(\theta)t(x) + m(x) + d(\theta) \right).$$

Esta manera de escribir las distribuciones en la familia exponencial resulta práctica cuando se discuten, por ejemplo, los modelos lineales generalizados.

3. Indicar si los siguientes modelos estadísticos pertenecen a la familia exponencial y en caso afirmativo determinar las expresiones analíticas de las funciones h , c , w , t :

- (a) Poisson: $f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$, con $\lambda \in (0, \infty)$ y $x \geq 0$.
- (b) Exponencial: $f(x, \lambda) = \frac{1}{\lambda} e^{-\frac{1}{\lambda} x}$, con $\lambda \in (0, \infty)$ y $x \geq 0$.
- (c) Truncada en θ : $f(x; \theta) = \frac{1}{\theta} e^{1-x/\theta}$ con $0 < \theta < x$.
- (d) Laplace: $f_X(x; \mu, \sigma) = \frac{1}{2\sigma} \exp \left(-\frac{|x-\mu|}{\sigma} \right)$ con $\mu \in \mathbb{R}$, $\sigma > 0$ y $x \in \mathbb{R}$.
- (e) Loc-escala Cauchy: $f(x; \mu, \sigma) = \frac{1}{\pi} \left[\frac{\sigma}{(x-\mu)^2 + \sigma^2} \right]$ con $\mu \in \mathbb{R}$, $\sigma > 0$ y $x \in \mathbb{R}$.
- (f) Gamma: $f(x; \lambda, k) = \lambda e^{-\lambda x} \frac{(\lambda x)^{k-1}}{\Gamma(k)}$, con $\lambda > 0$, $k > 0$ y $x > 0$.
- (g) Beta: $f(x; \beta, \alpha) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha-1} (1-x)^{\beta-1}$, con $\alpha > 0$, $\beta > 0$ y $0 \leq x \leq 1$.

Nota: En (f) y (g) $\Gamma(\cdot)$ denota la función **Gamma**.

$$\underbrace{\left(\frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} \right)}_{c(\theta)} \underbrace{\left[x^{\alpha-1} (1-x)^{\beta-1} \right]}_{k(x)} = c(\theta) k(x) \exp \left\{ \alpha \ln(x) + \beta \ln(1-x) \right\}$$

4. Sea $Z \sim f(z)$; definimos el cuantíl z_p como aquel valor que verifica que:

$$P(Z \leq z_p) = \int_{-\infty}^{z_p} f(z)dz = p.$$

Demostrar que si $X \sim \sigma^{-1}f((x - \mu)/\sigma)$ (es decir, X es una variable aleatoria con una distribución en la familia de localización y escala generada por f) entonces los cuantiles de X y Z están linealmente relacionados: $x_p = \sigma z_p + \mu$ para todo $p \in (0, 1)$.

5. Considere el siguiente *modelo de regresión*:

$$Y = \beta_0 + h(X) + \sigma\varepsilon,$$

donde h es una función conocida y $\varepsilon \sim N(0, 1)$. Identifique el modelo de localización y escala (determine la distribución y los parámetros) que sigue $Y|X$. ¿Cómo se relaciona éste modelo con el modelo lineal habitualmente utilizado en Econometría?

6. Sabiendo que la tasa de desempleo del último trimestre en Argentina fue del $\theta = 7.2\%$ y que pretendes encuestar a 1000 personas de la población económicamente activa del país, se plantean las siguientes cuestiones:

- Llamemos X_i a la v.a. que representa la condición de empleo del encuestado i -ésimo en la muestra; define el modelo, el parámetro y el soporte de la v.a. X_i .
- ¿Qué representa el estadístico $T = \sum_{i=1}^{1000} X_i/1000$?
- ¿Cuántas personas desempleadas esperas encontrar en la muestra aleatoria?
- ¿Cuál es la varianza y el desvío estándar de T ?
- ¿Cuál es la probabilidad de que exactamente 60 personas en la muestra estén desempleados? ¿Cómo cambia esta probabilidad cuando θ se incrementa?
- Utiliza el Teorema del Límite Central (TLC) para aproximar la probabilidad de que por lo menos 40 personas en la muestra aleatoria estén desempleados.

7. Una consultora económica quiere estimar la distribución del ingreso familiar en CABA. Suponga que en la Ciudad de Buenos Aires viven 2 millones de familias.

- Si X es la variable aleatoria ingresos familiares en la población, discuta que modelo estadístico (paramétrico) resulta razonable para X y determine su(s) parámetro(s).
- ¿Cómo estimaría el(los) parámetro(s) que caracterizan la distribución del ingreso? Imagine que usted le debe presentar al directorio de la consultora una justificación de su elección del estadístico a utilizar, que argumentos se le ocurre plantear.
- Suponé ahora que ya realizaste una encuesta a 100 familias elegidos aleatoriamente y que de los datos se observa que el ingreso promedio de las familias es de 35 mil pesos mensuales. Indica (con esta información) tus estimaciones de los cuartiles 1, 2 y 3 de la distribución del ingreso familiar.

8. Sea $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} X$ con $V(X) = \sigma^2 < \infty$, demuestre las siguientes propiedades:

- $E(\bar{X}_n) = \mu$.
- $V(\bar{X}_n) = \sigma^2/n$.
- $E(S_n^2) = \sigma^2$.

9. Sabiendo que $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} N(\mu, \sigma^2)$, definimos las variables aleatorias Y_i como:

$$Y_i = \begin{cases} 1 & \text{si } X_i \geq \mu, \\ 0 & \text{en otro caso.} \end{cases} \quad \text{para } i = 1, \dots, n.$$

Sea $T_n = \sum_{i=1}^n Y_i$, se pide:

- (a) Computa $E(T_n)$ y $V(T_n)$.
 - (b) Determinar la distribución del estadístico T_n (Ayuda: Determine primero como se distribuye Y_1 y tenga en cuenta que $\{Y_1, \dots, Y_n\} \stackrel{iid}{\sim} Y_1$).
10. Sea $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} \text{Unif}(0, 10)$ (i.e. X_1 es uniforme en el intervalo $[0, 10]$):
- (a) Utiliza el TCL para aproximar las siguientes probabilidades: $P(4.5 \leq \bar{X}_n \leq 5.5)$ y $P(|\bar{X}_n - 5| > 1)$ en función del tamaño de la muestra.
 - (b) ¿Qué ocurre con las dos probabilidades anteriores cuando $n \rightarrow \infty$?
11. Sea $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} \text{Unif}(0, 1)$, se pide calcular los siguientes momentos:
- (a) $E(X_{(1)})$ y $V(X_{(1)})$, donde $X_{(1)} = \min\{X_1, \dots, X_n\}$ es el mínimo en la muestra.
 - (b) $E(X_{(n)})$ y $V(X_{(n)})$, donde $X_{(n)} = \max\{X_1, \dots, X_n\}$ es el máximo en la muestra.
 - (c) ¿Cómo cambian las distribuciones del mínimo y el máximo en la muestra aleatoria si se tiene que $X \sim \text{Exp}(\theta)$?
12. Sea $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} X$ con $E(X) = \mu \neq 0$ y $V(X) < \infty$:
- (a) ¿Cuál es la media y varianza (aproximadas) de la variable aleatoria $g(\bar{X}_n) = \bar{X}_n^2$?
 - (b) ¿Cómo se distribuye (aproximadamente) \bar{X}_n^2 ?
 - (c) ¿Cómo se distribuye (aproximadamente) $e^{\bar{X}_n}$?
 - (d) ¿Cómo se distribuye (aproximadamente) \bar{X}_n^2 si $\mu = 0$? (Ayuda: Debes utilizar una aproximación de segundo orden).

13. Sea $\{X_1, \dots, X_n\}$ una muestra aleatoria de una población en donde la variable aleatoria de interés tiene una función de densidad de parámetro $\theta > 0$ dada por

$$f(x; \theta) = \begin{cases} \theta^{-1} e^{-x/\theta} & \text{si } x > 0, \\ 0 & \text{en otro caso.} \end{cases}$$

- (a) Comprobar que el estadístico $T(X_1, \dots, X_n) = \sum_{i=1}^n X_i$ es suficiente para θ determinando de manera explícita las funciones $g(T(\mathbf{x}) = t; \theta)$ y $h(\mathbf{x})$ (Fisher–Neyman).
 - (b) Porque es redundante utilizar el teorema de factorización para probar que T es suficiente para θ ? (Ayuda: ¿A qué familia pertenece $f(x; \theta)$ y quién es T ?).
14. Sea $\{X_1, \dots, X_n\}$ una muestra aleatoria de una población cuya variable aleatoria de interés tiene una densidad como sigue

$$f(x; \theta) = \begin{cases} \theta x^{\theta-1} & \text{si } 0 \leq x \leq 1, \\ 0 & \text{en otro caso,} \end{cases}$$

para $\theta \in (-\infty, 0)$. Comprueba que el estadístico $T(X_1, \dots, X_n) = \prod_{i=1}^n X_i$ es suficiente para θ y determine las expresiones de las funciones $g(T(\mathbf{x}) = t; \theta)$ y $h(\mathbf{x})$.

15. Considere el siguiente modelo estadístico

$$f(x; \theta) = \begin{cases} \left(\frac{\theta}{2}\right)^{|x|} (1 - \theta)^{1-|x|} & \text{si } x \in \{-1, 0, 1\}; \text{ y } 0 < \theta < 1, \\ 0 & \text{en otro caso,} \end{cases}$$

- (a) ¿Pertenece este modelo a la familia exponencial?
- (b) Dada $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} f(x; \theta)$ encuentre un estadístico suficiente para θ .
- (c) ¿Es completo el estadístico que encontraste en el punto anterior?

16. Sea $\{X_1, \dots, X_n\} \stackrel{iid}{\sim} \text{Unif}(0, \theta)$, se pide:

- (a) Verificar que $T(X_1, \dots, X_n) = \max\{X_1, \dots, X_n\}$ es suficiente para θ .
- (b) Demuestre que $E\left(\frac{n+1}{n}T(X_1, \dots, X_n)\right) = \theta$.
- (c) ¿Es $\frac{n+1}{n}T(X_1, \dots, X_n)$ un estadístico suficiente para θ ?

17. Argumentar (utilizando resultados teóricos discutidos en clase) porque cuando se trata de una población normal los estadísticos (\bar{X}_n, S_n^2) son independientes.

18. Consideremos una muestra de tamaño $n = 10$ de una población Poisson (ver 3.a):

- (a) Construye la función de verosimilitud.
- (b) Grafica (en una misma figura) $L(\lambda | \mathbf{x})$ cuando $\sum_{i=1}^{10} x_i \in \{10, 15, 20\}$. ¿Qué cambiaría en estas gráficas, si en vez de graficar $L(\lambda | \mathbf{x})$, graficas $\ell(\lambda | \mathbf{x}) \equiv \log L(\lambda | \mathbf{x})$.
- (c) ¿Si $\sum_{i=1}^{10} x_i = 10$, es más verosímil que la muestra se corresponda con una población Poisson donde $\lambda = 2$ o con una población Poisson donde $\lambda = 1$?
- (d) Repite los puntos (a), (b) y (c), pero asumiendo que la población es exponencial.

Recap:

Ingrediente de la Inferencia Estadística

- m.a: $X = \{X_1, X_2, \dots, X_n\}$
- una muestra: $\{x_1, x_2, \dots, x_n\}$
- un modelo estadístico:

2 supuestos:

* La m.a. es iid de un modelo estadístico con función de distribución (acumulada)

F desconocida ($F = F(x; \theta)$)

* $F \in \mathcal{F} = \{F(x; \theta) \mid \theta \in \Theta\}$

Terminología: . Parámetro $\rightarrow \theta$ es el valor que indexa $\mathcal{F} = \{F(\cdot | \theta) \mid \theta \in \Theta\}$

θ^*

. Parámetro de interés \rightarrow aspecto de la distribución de X s/ el que quiero hacer inferencia. Suele ser una función $g(\theta)$.

. Estimador de $g(\theta)$: es una función de X / $\hat{g} = \delta(X) \rightarrow$ variable observable

. Estimación: estimador evaluado en los datos \rightarrow valor F° jo

. Estadístico: cualquier función de los datos $T(X)$

2. El soporte de $X \sim f(x; \theta)$ es el conjunto definido como:

$$\text{Soporte}(f(x; \theta)) = \{x \in \mathbb{R} \mid f(x; \theta) > 0\}.$$

(a) ¿El conjunto $\text{Soporte}(f(x; \theta))$ puede depender del parámetro θ si $f(x; \theta)$ pertenece a la familia exponencial? (Formaliza tu respuesta utilizando la definición de familia Exponencial vista en clase).

(b) Demostrar que la familia exponencial se puede escribir, de manera equivalente a la expresión que dimos en clase, como sigue:

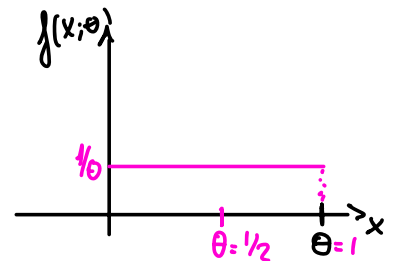
$$f(x; \theta) = \exp(w(\theta)t(x) + m(x) + d(\theta)).$$

Esta manera de escribir las distribuciones en la familia exponencial resulta práctica cuando se discuten, por ejemplo, los modelos lineales generalizados.

(c)

$$f(x; \theta) = c(\theta) k(x) \cdot \exp(w(\theta)^T t(x)) \quad x \in \mathcal{X} \subset [0, \theta]$$

Por definición, el soporte de $f(x; \theta)$ no puede depender de los parámetros $\theta \Rightarrow$ todas las distribuciones en un modelo en particular \mathcal{F} tienen el mismo soporte independiente de cual sea la configuración de los parámetros.



$$\begin{aligned}
 (b) \quad f(x|\theta) &= h(x) c(\theta) \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) \right\}, \quad x \in \mathbb{R}^n \\
 &= \exp[\ln(h(x))] \exp[\ln(c(\theta))] \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) \right\} \\
 &= \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) + \underbrace{\ln[h(x)]}_{m(x)} + \underbrace{\ln[c(\theta)]}_{d(\theta)} \right\} \\
 &= \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) + m(x) + d(\theta) \right\} \\
 &= \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) + \ln[h(x)] - \ln \left[\frac{1}{c(\theta)} \right] \right\} \\
 &= \exp \left\{ \sum_{i=1}^n w_i(\theta) t_i(x) + m(x) - d(\theta) \right\}
 \end{aligned}$$

3. Indicar si los siguientes modelos estadísticos pertenecen a la familia exponencial y en caso afirmativo determinar las expresiones analíticas de las funciones h , c , w , t :

- (a) Poisson: $f(x, \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$, con $\lambda \in (0, \infty)$ y $x \geq 0$.
 (b) Exponencial: $f(x, \lambda) = \frac{1}{\lambda} e^{-\frac{1}{\lambda} x}$, con $\lambda \in (0, \infty)$ y $x \geq 0$.
 (c) Truncada en θ : $f(x; \theta) = \frac{1}{\theta} e^{1-x/\theta}$ con $0 < \theta < x$.
 (d) Laplace: $f_X(x; \mu, \sigma) = \frac{1}{2\sigma} \exp \left(-\frac{|x-\mu|}{\sigma} \right)$ con $\mu \in \mathbb{R}$, $\sigma > 0$ y $x \in \mathbb{R}$.
 (e) Loc-escala Cauchy: $f(x; \mu, \sigma) = \frac{1}{\pi} \left[\frac{\sigma}{(x-\mu)^2 + \sigma^2} \right]$ con $\mu \in \mathbb{R}$, $\sigma > 0$ y $x \in \mathbb{R}$.

Ⓐ Poisson: $f(x; \lambda) = \frac{e^{-\lambda} \lambda^x}{x!}$

• $\text{supp} \{f(x; \lambda)\} = \{x \in \mathbb{Z}_0^+ \mid f(x; \theta) > 0\}$ ✓

• $h(x) = \frac{1}{x!} > 0, \forall x \in \mathbb{Z}_0^+ \rightarrow$ no depende λ

• $c(\theta) = e^{-\lambda} \geq 0, \forall \lambda \in \Theta = \mathbb{R}^+$

• $\lambda^x = \exp \{ \ln[\lambda^x] \} = \exp \{ x \cdot \ln[\lambda] \}$

$f(x; \lambda) = \left(\frac{1}{x!} \right) \cdot e^{-\lambda} \cdot \underbrace{\ln[\lambda]}_{t(x)} \cdot \underbrace{x}_{w(\theta)} \in FE$
 estadístico suficiente

MLE: $\hat{\lambda}$ (tarea)



(b) $f(x; \theta) = \frac{1}{\theta} e^{-\frac{1}{\theta} x}$, $\theta \in (0, +\infty)$ y $x \geq 0$

• $\text{supp} \{f(x; \theta)\} = \mathbb{R}^+$ y no depende de θ

• $f(x; \theta) = \mathbb{1}_{(x \geq 0)} \cdot \frac{1}{\theta} e^{-\frac{1}{\theta} x}$

$h(x) = \mathbb{1}_{(x \geq 0)} \rightarrow$ no depende de θ y $x \geq 0$
 $c(\theta) = \frac{1}{\theta} \rightarrow$ no depende de x y $\theta > 0$
 $w(\theta) = -1/\theta$
 $t(x) = x$

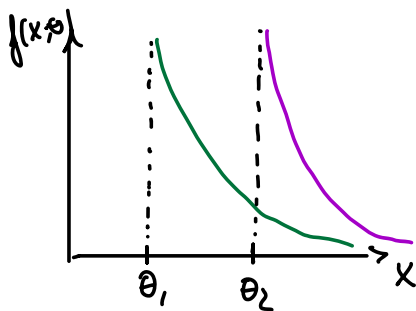
$$\int_0^{+\infty} e^{-\frac{1}{\theta} x} dx = -\theta e^{-\frac{1}{\theta} x} \Big|_0^{+\infty} = 0 - (-\theta) = \theta$$

$$\int_0^{+\infty} \frac{1}{\theta} e^{-\frac{1}{\theta} x} dx = 1$$

$$\int_{-\infty}^{+\infty} \exp \left\{ -\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right\} dx = \sqrt{2\pi\sigma^2}$$

(c) Truncada: $f(x; \theta) = \frac{1}{\theta} \exp \{1 - x/\theta\}$ con $0 < \theta < x$

$\text{supp} \{f(x; \theta)\} = \{x \in \mathbb{R}^+ \mid x > \theta\} \Rightarrow$ el soporte depende de $\theta \Rightarrow \notin FE$



e) Log-scaled Cauchy: $f(x; \mu, \sigma) = \frac{1}{\pi} \left[\frac{\sigma}{(x-\mu)^2 + \sigma^2} \right]$

• $\text{Supp} \{f(x; \mu, \sigma)\} = \mathbb{R} \Rightarrow$ no depende de $\theta = \begin{bmatrix} \mu \\ \sigma \end{bmatrix}$

• $h(x) = \mathbb{1}_{\{x \in \mathbb{R}\}} = 1$

• $c(\theta) = \frac{\sigma}{\pi} \cdot \mu^0$

$\exp \left\{ \ln \left[\frac{1}{(x-\mu)^2 + \sigma^2} \right] \right\} = \exp \left\{ -\ln [(x-\mu)^2 + \sigma^2] \right\}$
 $= \exp \left\{ -\ln [x^2 + \mu^2 - 2x\mu + \sigma^2] \right\}$

$\therefore f(x; \mu, \sigma) \notin FE$

\nrightarrow no puede descomponerse como $w(\theta)t(x)$

v

$$\begin{bmatrix} t_1(X) & t_2(X) \end{bmatrix} \begin{bmatrix} w_1(\theta) \\ w_2(\theta) \end{bmatrix} = t_1(X) w_1(\theta) + t_2(X) w_2(\theta)$$

5. Considere el siguiente *modelo de regresión*:

$$Y = \beta_0 + h(X) + \sigma \varepsilon,$$

donde h es una función conocida y $\varepsilon \sim N(0, 1)$. Identifique el modelo de localización y escala (determine la distribución y los parámetros) que sigue $Y|X$. ¿Cómo se relaciona éste modelo con el modelo lineal habitualmente utilizado en Econometría?

$$\bullet \mathbb{E}_{Y|X}[Y] = \mathbb{E}_{Y|X}[\beta_0 + h(X) + \sigma \varepsilon] = \beta_0 + h(X) + \mathbb{E}_{Y|X}[\sigma \varepsilon] = \beta_0 + h(X)$$

$$\bullet \text{Var}_{Y|X}[Y] = \text{Var}_{Y|X}[\underbrace{\beta_0 + h(X)}_{\text{Fijo}} + \sigma \varepsilon] = \text{Var}_{Y|X}[\sigma \varepsilon] = \sigma^2 \text{Var}_{Y|X}[\varepsilon] = \sigma^2$$

$$\frac{Y - \beta_0 - h(X)}{\sigma} = \varepsilon \Rightarrow \phi(Y) \cup N(0, 1)$$

$$F_\varepsilon(e) = P(\varepsilon \leq e) = P\left(\frac{Y - \beta_0 - h(X)}{\sigma} \leq e\right) = P(Y \leq \beta_0 + h(X) + \sigma e) = F_Y(\beta_0 + h(X) + \sigma e)$$

$$f_\varepsilon(e) = \frac{\partial F_\varepsilon}{\partial e} = \frac{\partial F_Y}{\partial e} = \frac{\partial F_Y}{\partial Y} \cdot \frac{\partial Y}{\partial e} = \dots$$