

EXAMEN LUNES 18 DE DICIEMBRE DE 2017.

Número de Examen	Cédula	Nombre y Apellido

Para uso docente:

Ej. 1	Ej. 2	Ej.3	Ej.4	Ej.5/6	TOTAL

El examen es con material y dura 4 horas. Se deberá entregar la resolución de un máximo de 5 ejercicios.

Ejercicio 1 (15 puntos)

En una ciudad hay 100 taxis, uno azul y 99 verdes. Una noche un taxi atropella a un peatón y se da a la fuga. Un testigo asegura que el taxi era azul, por lo tanto el conductor del taxi azul (que esa noche estaba trabajando) es detenido. Durante el juicio, se contrata a un especialista para evaluar la capacidad del testigo de distinguir si el color era azul o verde en condiciones similares a las de la noche del accidente. Los datos indican que el testigo ve autos azules como azules en el 99% de los casos y que ve autos verdes como azules en el 2% de los casos.

1. Calcular la probabilidad de que el taxi sea efectivamente azul dado que el testigo indica que es azul.
2. Si usted fuera el juez ¿diría que existe una duda razonable sobre la culpabilidad del conductor del taxi azul o no? Justifique su respuesta.

Ejercicio 2 (20 puntos)

Sean X e Y dos variables aleatorias independientes con distribución Bernoulli de parámetro $p = 0.5$. Se definen $S = X + Y$ y $T = X - Y$.

1. Hallar la función de probabilidad puntual conjunta de S y T .
2. ¿Son S y T independientes? Justifique su respuesta.
3. Calcular la probabilidad del suceso $S = T$.

Ejercicio 3 (20 puntos)

Los siguientes datos corresponden al precio en miles de dólares de 10 vehículos usados de la misma marca y año (datos sacados de *Mercado Libre Uruguay*):

i	1	2	3	4	5	6	7	8	9	10
P_i	8.400	10.490	9.200	9.900	8.900	8.500	7.800	9.500	8.990	8.900

(P_i es el precio del auto i).

1. Realizar un histograma de los precios de estos autos. Considere 4 intervalos de igual longitud.
2. Dada una muestra de datos x_1, x_2, \dots, x_n ($x_i \in \mathbb{R}$), defina qué es la *mediana*, *primer cuartil* y *tercer cuartil* de esta muestra.
3. Calcular la media, mediana, primer cuartil y tercer cuartil de P .

- Realizar un diagrama de caja (boxplot) de P . Indicar si existen datos atípicos y estudiar la simetría de la distribución correspondiente a la variable P .

Ejercicio 4 (25 puntos)

Suponga que un investigador toma n medidas independientes de la radiación ambiente en Santiago de Chile (en particular se midieron niveles de $MP_{2,5}$), obteniendo así una muestra $X_1, X_2 \dots X_n$ independiente e idénticamente distribuida. Se asume que estas observaciones siguen una distribución Rayleigh de parámetro $\lambda > 0$, esto es:

$$f(x) = 2\lambda x e^{-\lambda x^2} \quad \forall x \geq 0.$$

- Hallar el estimador por momentos del parámetro λ . Recordar que $\int_0^{\infty} e^{-x^2/2\sigma^2} dx = \sqrt{\frac{\pi}{2}}\sigma$.
- Hallar el estimador de máxima verosimilitud del parámetro λ .
- Se tiene una muestra de 5300 datos (100 muestras de niveles de $MP_{2,5}$ por cada una de las 53 ciudades más pobladas de Chile). Según las mediciones se tiene entonces que:

$$\frac{1}{5300} \sum_{i=1}^{5300} x_i = 20 \quad \sum_{i=1}^{5300} x_i^2 = 2141995,$$

siendo x_i la concentración de $MP_{2,5}$ presente en la muestra i (medida en microgramos por metro cúbico). Calcular los valores del estimador por momentos y del estimador por máxima verosimilitud para esta muestra.

Deberán realizar uno y solo uno de los ejercicios que se enuncian a continuación.

Ejercicio 5 (20 puntos)

Se considera la siguiente muestra de datos aleatorios independientes e idénticamente distribuidos.

35.90	36.55	37.15	37.30	37.35	37.45	37.75	39.90
-------	-------	-------	-------	-------	-------	-------	-------

- Asumiendo que los datos provienen de una variable aleatoria normal con varianza $\sigma^2 = 0.27$:
 - Dé un intervalo de confianza al 95 % para el parámetro μ de la distribución.
 - ¿Es razonable suponer que la muestra proviene de una variable aleatoria con distribución normal con $\mu = 37$ y $\sigma^2 = 0.27$? Decidir a nivel $\alpha = 0.05$.
- Se considera una segunda muestra que se supone aleatoria e independiente de la muestra anterior.

35.95	37.10	37.25	37.35	37.90	40.05	40.10	40.65
-------	-------	-------	-------	-------	-------	-------	-------

Realizar un test de comparación de muestras para decidir si se puede afirmar que estas dos muestras provienen de la misma distribución. Decidir a nivel $\alpha = 0.05$.

Ejercicio 6 (20 puntos)

Se realiza una encuesta para determinar la fracción de la población p que apoyaría un referéndum que requiera que todos los ciudadanos conozcan los principios básicos de la probabilidad y estadística.

- Asumiendo que $p = 0.5$, aproximar a la probabilidad de que en una encuesta a 25 personas al menos 14 apoyen el referéndum.
- Asumiendo ahora que p es desconocido, dar una cota para la mínima cantidad de personas que hay que encuestar para que la fracción de personas encuestadas que apoyan el referéndum difiera del verdadero valor de p en menos de 0.01, con una probabilidad de 0.9. *Sugerencia: puede utilizar que $p(1-p) \leq \frac{1}{4}$ para todo $p \in [0, 1]$.*
- Asumiendo que en una encuesta a 1000 personas, 25 se manifestaron a favor del referéndum, hallar un intervalo de confianza (aproximado) del 95% para el verdadero valor de p .

Solución

- Ej.1** 1. Sean: $AA = \text{“el Auto es Azul”}$, $AV = \text{“el Auto es Verde”}$ y $TA = \text{“el Testigo ve el auto como Azul”}$. Usando Bayes:

$$\mathbb{P}(AA|TA) = \frac{\mathbb{P}(TA|AA) \times \mathbb{P}(AA)}{\mathbb{P}(TA)}$$

Como $\mathbb{P}(TA) = 0.01 \times 0.99 + 0.99 \times 0.02 = 0.0297$, entonces:

$$\mathbb{P}(AA|TA) = \frac{0.99 \times 0.01}{0.0297} = 0.333$$

2. $\mathbb{P}(AV|TA) = 1 - \mathbb{P}(AA|TA) = 1 - 0.333 = 0.667$. Es decir que es más probable que el auto que vio el testigo el día del accidente fuera un auto verde. Es bastante razonable dudar entonces de la culpabilidad del conductor del auto azul.

- Ej.2** 1. Observar que:

X/Y	0	1
0	$S = 0, T = 0$	$S = 1, T = 1$
1	$S = 1, T = -1$	$S = 2, T = 0$

Luego, $Rec(S) = \{0, 1, 2\}$ y $Rec(T) = \{0, 1, -1\}$.

- $\mathbb{P}(S = 0; T = 0) = \mathbb{P}(X = 0; Y = 0) = \frac{1}{4}$
- $\mathbb{P}(S = 1; T = 1) = \mathbb{P}(X = 0; Y = 1) = \frac{1}{4}$
- $\mathbb{P}(S = 1; T = -1) = \mathbb{P}(X = 1; Y = 0) = \frac{1}{4}$
- $\mathbb{P}(S = 2; T = 0) = \mathbb{P}(X = 1; Y = 1) = \frac{1}{4}$

2. Observar que:

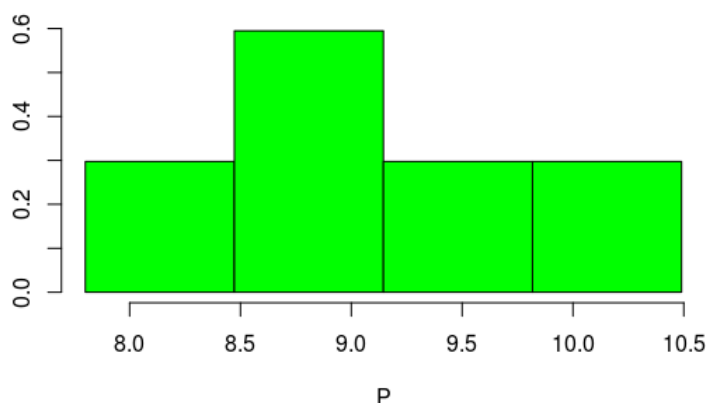
- $\mathbb{P}(S = 0) = \mathbb{P}(S = 0; T = 0) = \frac{1}{4}$
- $\mathbb{P}(T = 0) = \mathbb{P}(S = 0; T = 0) + \mathbb{P}(S = 2; T = 0) = \frac{1}{2}$

Luego, $\mathbb{P}(S = 0)\mathbb{P}(T = 0) = \frac{1}{4} \times \frac{1}{2} \neq \mathbb{P}(S = 0; T = 0)$, entonces S y T no son independientes.

3. $\mathbb{P}(S = T) = \mathbb{P}(S = 0; T = 0) + \mathbb{P}(S = 1; T = 1) = \frac{1}{2}$.

- Ej.3** 1.

Histograma de frecuencias de P



2. Dos definiciones:

- Def.1
- $\hat{m}_X = \inf\{x \in \mathbb{R} : F_n(x) \geq 0.5\} = x_i^*$, siendo x_i^* el primer dato que deja *al menos* el 50% de los datos a su izquierda en la muestra ordenada. (En particular, el ínfimo anterior es un mínimo).
 - $\hat{q}_1 = \inf\{x \in \mathbb{R} : F_n(x) \geq 0.25\} = x_j^*$, siendo x_j^* el primer dato que deja *al menos* el 25% de los datos a su izquierda en la muestra ordenada.

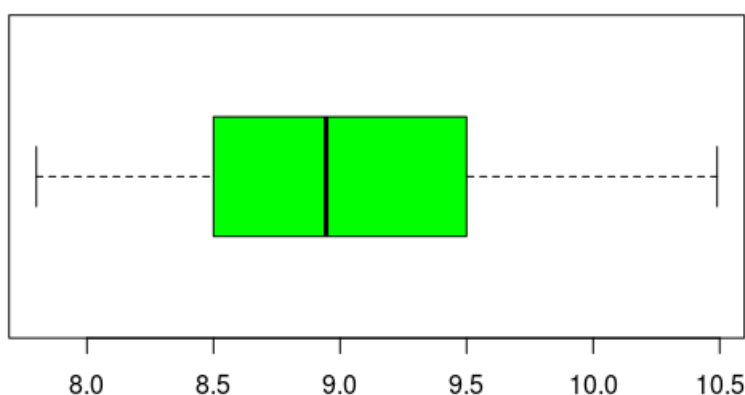
- $\hat{q}_3 = \inf\{x \in \mathbb{R} : F_n(x) \geq 0.75\} = x_j^*$, siendo x_j^* el primer dato que deja *al menos* el 75% de los datos a su izquierda en la muestra ordenada.

Def.2 • Si n es impar, entonces la mediana es el dato del medio: $\hat{m}_X = x_{\frac{n+1}{2}}^*$. Si n es par,

entonces la mediana es el promedio de los datos del medio: $\hat{m}_X = \frac{x_{\frac{n}{2}}^* + x_{\frac{n}{2}+1}^*}{2}$.

- El primer cuartil se define análogamente como la mediana de la primera mitad de los datos ordenados y el tercer cuartil como la mediana de la segunda mitad de los datos ordenados.
3. Usando la primera definición, $q_1 = P_3^* = 8.5$, $m = P_5^* = 8.9$, $q_3 = P_8^* = 9.5$, $\bar{P}_{10} = 10.58$.
Usando la segunda definición, $q_1 = P_3^* = 8.5$, $m = 8.945$, $q_3 = P_8^* = 9.5$, $\bar{P}_{10} = 10.58$.
 4. El límite inferior es: $L_{inf} = \max\{P_1^*, q_1 - 1.5 \times RI\} = 7.8$ y el límite superior es: $L_{sup} = \min\{P_{10}^*, q_3 + 1.5 \times RI\} = 10.49$. Entonces:

Boxplot de P



Ej.4 1. Aplicando integración por partes se tiene que:

$$\mathbb{E}(X) = \int_0^{+\infty} 2\lambda x^2 e^{-\lambda x^2} dx = \int_0^{+\infty} e^{-\lambda x^2} dx = \int_0^{+\infty} e^{-\frac{x^2}{2\sigma^2}} dx,$$

si tomamos $2\sigma^2 = \frac{1}{\lambda}$. Luego, usando la sugerencia:

$$\mathbb{E}(X) = \sqrt{\frac{\pi}{2}} \sigma = \sqrt{\frac{\pi}{2}} \sqrt{\frac{1}{2\lambda}} = \frac{1}{2} \sqrt{\frac{\pi}{\lambda}}.$$

Como $\bar{x}_n \approx \mathbb{E}(X) = \frac{1}{2} \sqrt{\frac{\pi}{\lambda}}$, planteamos:

$$\bar{x}_n = \frac{1}{2} \sqrt{\frac{\pi}{\lambda}} \leftrightarrow \lambda_M = \frac{\pi}{(2\bar{x}_n)^2}.$$

2. La función de verosimilitud es:

$$L(\lambda) = \prod_{i=1}^n f_X(x_i) = (2\lambda)^n (x_1 \dots x_n) e^{-\lambda \sum x_i^2}.$$

Buscamos $\lambda_{MV} = \arg \max_{\lambda > 0} L(\lambda) = \arg \max_{\lambda > 0} \log(L(\lambda))$.

Sea $f(\lambda) = n \log(2\lambda) + \sum \log(x_i) - \lambda \sum x_i^2$, $f'(\lambda) = 0$ si y solo si $\lambda = \lambda_{MV} = \frac{n}{\sum x_i^2}$. En este punto f alcanza el máximo.

3. Según los datos, $n = 5300$, $\bar{x}_n = 20$. Entonces:

- $\lambda_M = \frac{\pi}{(2 \times 20)^2} \approx 0.00196$

- $\lambda_{MV} = \frac{5300}{\sum x_i^2} \approx 0.00247$

Ej.5 1)i. Puesto que los datos provienen de v.a. normales con varianza conocida, un intervalo de confianza a nivel 95% para μ es $I_{0.95}(\mu) = [\bar{x}_n - k, \bar{x}_n + k]$ con $k = \frac{\sigma}{\sqrt{n}} z_{1-\frac{\alpha}{2}}$.
En este caso: $\bar{x}_n = 37.42$, $z_{0.975} = 1.96$ y $k = 0.36$. Entonces:

$$I_{0.95}(\mu) = [37.06; 37.78]$$

1)ii. Realizamos un test de Kolmogorov Smirnov para testear: $\begin{cases} H_0 : X \sim N(\mu = 37, \sigma^2 = 0.27) \\ H_1 : \text{No } H_0 \end{cases}$

El valor más grande de $|\frac{i}{n} - F_0(x_i^*)|$ es 0.239 y el valor más grande de $|\frac{i-1}{n} - F_0(x_i^*)|$ es 0.364, por lo que $d_n = 0.364$. Según la tabla de KS para una muestra:

$$p\text{-valor} = \mathbb{P}(D_n \geq 0.364 | H_0) > 0.2 > \alpha$$

por lo que NO se rechaza H_0 a nivel $\alpha = 0.05$.

2. El estadístico para este test es:

$$d_{n,m} = \max \left\{ \max_{1 \leq i \leq n} |F_n^X(x_i) - F_m^Y(x_i)|; \max_{1 \leq j \leq m} |F_n^X(y_j) - F_m^Y(y_j)| \right\}$$

En este caso: $n = m = 8$ y $d_{n,m} = \max\{0.375; 0.25\} = 0.375$. Según la tabla de KS para dos muestras ($n = m = 8$):

$$p\text{-valor} = \mathbb{P}(D_{n,m} \geq 0,375 | H_0) = \mathbb{P}(nmD_{n,m} \geq nm0,375 | H_0) = \mathbb{P}(nmD_{n,m} \geq 24 | H_0) > \mathbb{P}(nmD_{n,m} \geq 32 | H_0) = 0.283$$

Entonces $p\text{-valor} > 0.283 > \alpha$, por lo tanto NO rechazamos H_0 a nivel $\alpha = 0.05$.

Ej.6 1. Si $X = \#$ de personas que están a favor del referendúm (dentro de las 25 encuestadas), $X \sim Bin(n = 25, p)$. Entonces:

$$\mathbb{P}(X \geq 14) = 1 - \mathbb{P}(X \leq 13) = 1 - \mathbb{P}\left(\bar{X}_{25} \leq \frac{13}{25}\right) \approx 1 - \Phi\left(\frac{\frac{13}{25} - \mu}{\frac{\sigma}{\sqrt{25}}}\right),$$

por el Teorema Central del Límite.

Como $\mu = p = 0.5$ y $\sigma^2 = p(1-p) = 0.25$, entonces:

$$\mathbb{P}(X \geq 14) \approx 1 - \Phi\left(\frac{5(0.52 - 0.5)}{0.5}\right) = 1 - \Phi(0.2) \approx 1 - 0.5793 = 0.4207$$

2. Buscamos n tal que: $\mathbb{P}(|\bar{X}_n - p| < 0.01) = 0.9$. Utilizando nuevamente el TCL:

$$\mathbb{P}(|\bar{X}_n - p| < 0.01) = \mathbb{P}\left(\left|\frac{\sqrt{n}(\bar{X}_n - p)}{\sigma}\right| < \frac{\sqrt{n}0.01}{\sigma}\right) \approx \Phi\left(\frac{\sqrt{n}0.01}{\sigma}\right) - \Phi\left(-\frac{\sqrt{n}0.01}{\sigma}\right) = 2\Phi\left(\frac{\sqrt{n}0.01}{\sigma}\right) - 1$$

Luego,

$$2\Phi\left(\frac{\sqrt{n}0.01}{\sigma}\right) - 1 = 0.9 \Leftrightarrow \frac{\sqrt{n}0.01}{\sigma} = z_{0.95} \Leftrightarrow n = \left(\frac{\sigma}{0.01} z_{0.95}\right)^2 \leq \frac{1}{4} \left(\frac{z_{0.95}}{0.01}\right)^2 = 6806.25$$

Entonces habría que encuestar a lo sumo 6806 personas para saber que la estimación de p difiere del verdadero valor de p en menos de 0.01 con una probabilidad del 90%. Sin importar el tamaño de la población implicada en el referendúm.

3. En este caso $p = \mu = \mathbb{E}(Ber(p))$, entonces un intervalo de confianza aproximado a nivel 95% para p es $I_{0.95}(p) = [\bar{x}_n - k, \bar{x}_n + k]$ con:

$$\bar{x}_n = \frac{25}{1000} = 0.025 \text{ y } k = \frac{\sigma}{\sqrt{n}} z_{0.975}$$

Estimamos σ^2 por $\bar{x}_n(1 - \bar{x}_n) \approx 0.0244$. Luego $k \approx \frac{0.156}{31.6} 1.96 \approx 0.01$ y $I_{0.95}(p) = [0.015, 0.035]$.