

Introducción a la Estadística y Ciencia de Datos

Práctica 1 - Estadística Descriptiva

1. El archivo `Debernardi.csv` contiene los datos referentes a un estudio acerca del cáncer de páncreas (más información en el archivo *Acerca de los datos*, en el Aula Virtual).
 - (a) Construir una tabla con los valores observados para la variable `DIAGNOSIS` y su frecuencia relativa.
 - (b) Realizar un gráfico de barras usando la tabla del ítem anterior.
2. El archivo `datos_titanic.csv` contiene información sobre una muestra seleccionada al azar de las personas, no tripulantes, que viajaban en el barco tristemente célebre *Titanic*, al momento de su hundimiento en el Océano Atlántico (más información en el archivo *Acerca de los datos*, en el Aula Virtual).
 - (a) Estimar la probabilidad de ser mujer sabiendo que sobrevivió y comparar con la estimación de ser mujer a bordo del *Titanic*.
 - (b) Hacer una tabla de contingencia entre las variables categóricas `SURVIVED` y `PCLASS`. A partir de esta tabla estimar la probabilidad de sobrevivir dada la clase para los distintos valores de la variable `PCLASS`.
 - (c) Realizar un gráfico de barras que vincule a las variables categóricas `SURVIVED` y `PCLASS`.
3. En un experimento se midió la temperatura de sublimación del iridio y del rodio. En los archivos `iridio.txt` y `rodio.txt` se encuentran los datos recabados en el experimento.
 - (a) Comparar los dos conjuntos de datos mediante histogramas y boxplots, graficando los boxplots en paralelo.
 - (b) Hallar las medias, las medianas y las medias podadas al 10% y 20% muestrales. Comparar.
 - (c) Hallar los desvíos estándares, las distancias intercuartiles y las MAD muestrales como medidas de dispersión.
 - (d) Hallar los cuantiles 0.90, 0.75, 0.50, 0.25 y 0.10.
4. En un estudio nutricional se consideran las calorías y el contenido de sodio de tres tipos de salchichas y se obtuvieron los datos que se encuentran en los archivos `salchichas_A.txt`, `salchichas_B.txt` y `salchichas_C.txt`.
 - (a) Armar un archivo que se llame `salchichas.txt` que contenga toda la información registrada en los tres archivos mencionados agregando una columna que indique el tipo de salchicha en cada caso.
 - (b) Realizar un histograma para las calorías de cada tipo de salchichas. ¿Observa grupos en algún gráfico? ¿Cuántos grupos observa? ¿Observa algún candidato a dato atípico? ¿Alguno de los histogramas tiene una característica particular?

- (c) Realizar los boxplots paralelos para las calorías. ¿Observa la misma cantidad de grupos que antes? ¿A cuál conclusión llega? De acuerdo con los boxplots graficados, ¿cómo caracterizaría la diferencia entre los tres tipos de salchichas desde el punto de vista de las calorías?
 - (d) Repetir con la cantidad de sodio.
5. El conjunto de datos que figura en el archivo `estudiantes.txt` corresponde a 100 determinaciones repetidas de la concentración de ion nitrato (en $\mu\text{g/l}$), 50 de ellas corresponden a un grupo de estudiantes (Grupo 1) y las restantes 50 a otro grupo (Grupo 2).
- (a) Estudiar si la distribución de los conjuntos de datos para ambos grupos es normal, realizando los correspondientes histogramas y superponiendo la curva normal. Además dibujar los qqplots para cada conjunto de datos superponiendo, en otro color, la recta mediante el comando `qqline`.
 - (b) ¿Le parece a partir de estos datos que ambos grupos están midiendo lo mismo? Responder comparando medidas de centralidad y de dispersión de los datos. Hacer boxplots paralelos.
6. Con la finalidad de incrementar las lluvias en zonas desérticas, se desarrolló un método que consiste en el bombardeo de la nube con átomos. Para evaluar la efectividad del método se realizó el siguiente experimento:
- Para cada nube que se podía bombardear se decidió al azar si se la trataba o no.
 - Las nubes no tratadas fueron denominadas nubes controles.

En el archivo `nubes.txt` se presentan la cantidad de agua caída de 26 nubes tratadas y 26 nubes controles.

- (a) Realizar boxplots paralelos. ¿Le parece que el método produce algún efecto?
 - (b) Analizar la normalidad realizando qqplots e histogramas (de densidad) para ambos conjuntos de datos y superponiendo la curva normal.
 - (c) Realizar la transformación logaritmo natural a los datos (`log` en R) y repetir *b*) para los datos transformados.
 - (d) Realizar boxplots paralelos habiendo transformado las variables con el logaritmo natural. Observar cómo se modificaron los datos atípicos respecto del ítem *a*).
7. El archivo `data_credit_card.csv` tiene información de $n=500$ clientes de un banco, para las siguientes variables: `purchases` es el monto total de compras en el último año, `credit_limit` es el límite de crédito disponible para el cliente, `purchases_freq` es la proporción de semanas del año en las que el cliente realizó compras y `tenure` es la cantidad de meses que restan al cliente para cancelar el crédito. Se pide:
- (a) Para todas las variables, graficar la función de distribución empírica. Discutir sobre el tipo de variable aleatoria que utilizaría para modelar en cada caso.

- (b) Para la variable `credit_limit` hacer un histograma y un gráfico de densidad usando la función `density`, ¿Qué observa? ¿Le parece adecuado realizar estos gráficos para las variables `purchases` y `tenure`?
 - (c) Para la variable `tenure` hacer un *barplot* con las frecuencias relativas de cada valor. ¿Qué observa?
 - (d) Para todas las variables, calcular la media, la mediana y la media α -podada (con $\alpha = 0.1$). Comparar los resultados y justificar. ¿Qué medida de posición del centro de los datos le parece más adecuada en cada caso?
 - (e) Para todas las variables, obtener los cuantiles de nivel 0.25 y 0.75 de los datos. Calcular el rango inter-cuartílico y la MAD muestrales. Graficar *boxplots*. ¿Qué observa?
 - (f) Calcular el desvío estándar, el coeficiente de asimetría y el coeficiente de curtosis muestrales. Interpretar los resultados en relación a las distribuciones vistas.
 - (g) Identificar datos atípicos. ¿Deberían excluirse? ¿Cómo se modifican las medidas obtenidas anteriormente si se los excluye?
8. En el archivo `ciclocombinado.xlsx` hay datos de la potencia entregada por una central térmica de ciclo combinado. Se registraron datos diarios de la potencia máxima entregada (PE, en MW) por la planta funcionando en capacidad máxima. La variable `HighTemp` vale 1 si la temperatura media diaria fue superior a 20°C en el día en el que se tomó el dato y vale 0 en caso contrario.
- (a) Realizar un histograma y un gráfico `density` con los datos de PE, ¿Qué se observa?
 - (b) Clasificar los datos en dos vectores según la variable `HighTemp` y realizar gráficos `density` separados. Visualizar simultáneamente los gráficos en la misma escala. ¿Qué se observa?
 - (c) Estimar $P(\text{PE} < 450 | \text{HighTemp} = 0)$ y $P(\text{PE} < 300 | \text{HighTemp} = 1)$.
 - (d) Estimar $P(\text{PE} < 450)$.
 - (e) Estimar la potencia mínima garantizada con probabilidad 0.9 para un cierto día con `Hightemp = 1`.
 - (f) Estimar la potencia mínima garantizada con probabilidad 0.9 para un cierto día.
9. Considerar nuevamente el conjunto de datos del ejercicio 1.
- (a) Realizar histogramas para la variable `LYVE1` basados en los datos brindados para las observaciones que cumplen `DIAGNOSIS=1`, `DIAGNOSIS=2` y `DIAGNOSIS=3`. Es decir efectuar histogramas según los niveles de la variable factor `DIAGNOSIS`. Indicar las características más sobresalientes de los histogramas y aquellas que los diferencian.
 - (b) Graficar, en distintos colores y superpuestas, las funciones de distribución empíricas de la variable `LYVE1` según los niveles de la variable factor `DIAGNOSIS`. Decidir si la siguiente afirmación es verdadera o falsa y justificar: “los valores de la variable `LYVE1` tienden a ser más altos entre quienes tienen cáncer de páncreas que entre quienes sufren otras enfermedades asociadas al páncreas”.

- (c) Realizar boxplots paralelos para la variable LYVE1 según los niveles de la variable factor DIAGNOSIS, considerando el sexo de los pacientes (variable SEX). Decidir si la siguiente afirmación es verdadera o falsa y justificar: “en términos generales, el sexo del paciente no afecta los niveles de la proteína que se mide en la variable LYVE1”.
- (d) Graficar superpuestas las densidades estimadas, que brinda la función **density**, para la variable LYVE1 según los niveles de la variable factor DIAGNOSIS. Describir las características más sobresalientes de las densidades estimadas y aquellas que las diferencian.
- (e) Repetir a) y d) para el logaritmo de LYVE1.