

Tarea #: 1

Tema: Exploración de datos

Fecha entrega: 11:59 pm Agosto 21 de 2023

Objetivo: Utilizar conceptos estadísticos para entender la relación entre las variables de una base de datos. Adicionalmente, utilizar python como herramienta de exploración de datos y validación de hipótesis.

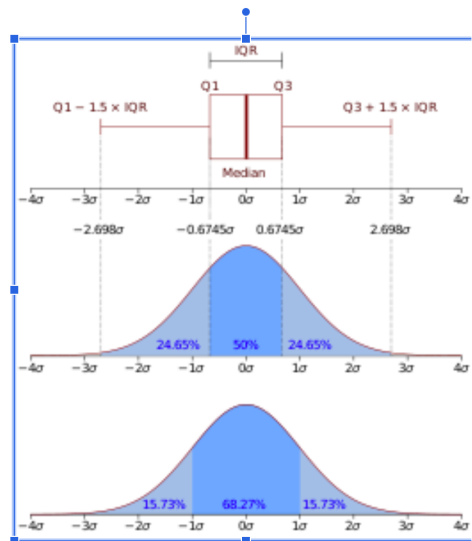
Entrega: Crear un repositorio en su github personal. Dentro del proyecto debe existir una carpeta llamada tarea 1, dentro debe tener una carpeta doc con este documento incluyendo todas las respuestas y los gráficos. Adicionalmente, debe existir una carpeta src con el código del notebook utilizado. Debe adicionar la cuenta jdramirez como colaborador del proyecto y enviar un email antes de q se termine el día indicando el commit desea le sea calificado.

1. Utilizas el siguiente set de datos para calcular paso por paso (mostrar procedimiento y fórmulas):

x1	x2	x3
4	4	28
2	3	24
2	4	30
3	5	32
1	3	18
3	6	41
3	6	44
0	1	5
1	3	18
0	0	1
5	9	62
1	2	17
2	3	24
1	3	19
3	6	42
4	8	56
4	8	56
3	6	44
5	9	64
1	2	17

1	2	17
---	---	----

- 1.1. ¿Cuál es la media, mediana y desviación estándar?, y la moda y los valores repeticiones de la moda para los datos categóricos.
- 1.2. Dibujar un boxplot a mano. Utilizando los datos de la tabla 1 y las siguientes proporciones.



- 1.3. Cual es la covarianza entre las 2 variables X1, X2

$$\text{Cov}(x,y) = \frac{\sum (x_i - \bar{x}) * (y_i - \bar{y})}{N}$$

- 1.4.Cuál es la correlación entre la variable x1 y x2 (Calcularla a mano). Correlación puede ser escrita también como:

$$\text{Cor}(X, Y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}},$$

- 1.5. Explica la relación entre covarianza y correlación.
- 1.6. Calcule el resultado del algoritmo K-means sobre este set de datos a mano como lo hicimos en excel. Vamos a crear 3 grupos, es decir, $k=3$ (clusters).
2. PCA. Utilizar los datos de la tabla 1, para calcular PCA y reducir la dimensionalidad de 2 dimensiones a 1. Para este ejercicio se debe utilizar las variables X_1 , y X_2 y crear un vector con una sola dimensión.
 - 2.1. Cual es la matriz de covarianza
 - 2.2. Cuales son los eigenvalues
 - 2.3. Cuál es la varianza explicada por el eigenvalue.
 - 2.4. Cual es el valor del eigenvector
 - 2.5. Cuál es la matriz proyectada.
 - 2.6. Cual es el error o diferencia entre la matriz proyectada
3. Utilizando el dataset del [proyecto](#) data/CARS.csv crear: **Utilizar la librería de plotly.**
 - 3.1. Distribución de cada variables:
 - 3.1.1. Para las variables categóricas un gráfico de barras. Categoría numero de observaciones.
 - 3.1.2. Para las variables numéricas crear histogramas. Listar los modelos de carros que están más lejos de 5 estándares de desviación, y serían considerados outliers. Hacer test de si es una distribución normal o no.
 - 3.2. Gráfico de la relación de cada variable con respecto a MPG_City:
 - 3.2.1. Variables categóricas debes crear un boxplot. Explique cómo interpreta el gráfico
 - 3.2.2. Variables numéricas vas a crear un scatter plot. Explique cómo interpreta el gráfico
 - 3.3. Matriz de correlación.
 - 3.3.1. Cree la matriz de correlación, cuales son las variables más importantes para explicar la variabilidad de MPG_City. Explique por qué el coeficiente es negativo o positivo.
 - 3.3.2. Cree las dummy variables para todas las variables categóricas y genere la matriz de correlación nuevamente. ¿Cuál es el valor de variable categórica con mayor correlación?

- 3.3.3. Cree la matriz de correlación nuevamente removiendo todas los modelos de carro que fueron catalogados como un outlier. (Puede utilizar `.query('Model in["MDX","TSX 4dr"]')`). Existe alguna variación en la correlación.