

Tarea #: 1

Tema: Exploración de datos

Fecha entrega: 11:59 pm Marzo 06 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.

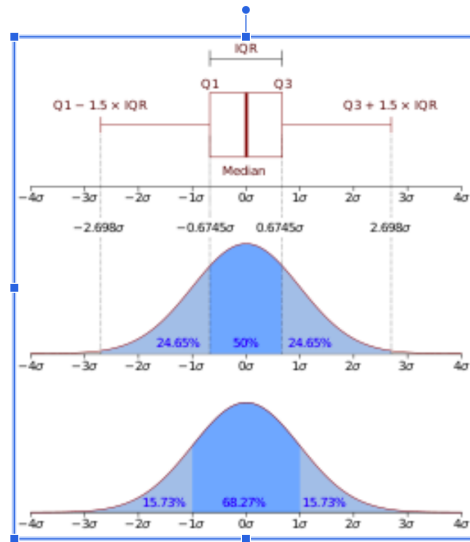
Entrega: Crear un repositorio en su github personal. Dentro del proyecto debe existir una carpeta llamada tarea 1, dentro debe tener una carpeta doc con este documento incluyendo todas las respuestas y los gráficos. Adicionalmente, debe existir una carpeta src con el código del notebook utilizado. Debe 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):

Id	X1	X2
1	1	4
2	1	3
3	0	4
4	5	1
5	6	2
6	4	0

Tabla:1

- 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.
Vamos a crear 2 grupos, es decir, $k=2$ (2 clusters).
2. Utilizando el dataset del [proyecto](#) data/CARS.csv crear:
 - 2.1. Distribución de cada variables:
 - 2.1.1. Para las variables categóricas un gráfico de barras. Categoría numero de observaciones.
 - 2.1.2. Para las variables numéricas crear histogramas. Listar los modelos de carros que están más lejos de 4 estándares de desviación, y serían considerados outliers. Hacer test de si es una distribución normal o no.
 - 2.2. Gráfico de la relación de cada variable con respecto a MPG_City:
 - 2.2.1. Variables categóricas debes crear un boxplot. Explique cómo interpreta el gráfico
 - 2.2.2. Variables numéricas vas a crear un scatter plot. Explique cómo interpreta el gráfico
 - 2.3. Matriz de correlación.
 - 2.3.1. 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.
 - 2.3.2. 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.

7.1

Media:

Formula: $\frac{\sum_{i=1}^n X_i}{n}$

Procedimiento:

$$\bar{X}_1 = \frac{1 + 7 + 0 + 5 + 6 + 4}{6}$$

$$\bar{X}_2 = \frac{4 + 3 + 4 + 7 + 2 + 0}{6}$$

Resultado:

$$X_1 = \frac{17}{6}$$

$$X_2 = \frac{14}{6}$$

Mediana:

Formula:

Para datos ordenados la mediana es el Valor Central

Procedimiento:

$X_1 = [0, 1, 1, 4, 5, 6]$ la mediana es 2.5

$X_2 = [0, 1, 2, 3, 4, 4]$ la mediana es 2.5

(Promedio Valores centrales)

Resultado:

$$X_1 = 4$$

$$X_2 = 2.5$$

Desviación estándar

Formula: $\sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n}}$

Procedimiento:

$$a = \frac{77}{6}$$

$$S(X_1) = \sqrt{\frac{(7-a)^2 + (7-a)^2 + (0-a)^2 + (5-a)^2 + (6-a)^2 + (4-a)^2}{6}}$$

$$b = \frac{74}{6}$$

$$S(X_2) = \sqrt{\frac{(4-b)^2 + (3-b)^2 + (4-b)^2 + (7-b)^2 + (2-b)^2 + (0-b)^2}{6}}$$

Resultado

$$S(X_1) \approx 7.471$$

$$S(X_2) \approx 7.825$$

Moda

Formula:

Valor que aparece con más Frecuencia

Procedimiento y resultado

Moda Para $X_1 = 7$

Moda Para $X_2 = 4$

Boxplots X_1

1. ordenar los datos

$$X_1 = [0, 1, 7, 4, 5, 6]$$

2. (cuartiles)

$$Q_1 = 1 \quad (\text{Valor Posición 2})$$

$$Q_2 = 2.5 \quad (\text{mediana})$$

$$Q_3 = 5$$

$$IQR = Q_3 - Q_1 = 5 - 1 = 4$$

$$\text{Límite inferior} = Q_1 - 1.5 \times IQR$$

$$= 1 - 1.5 \times 4$$

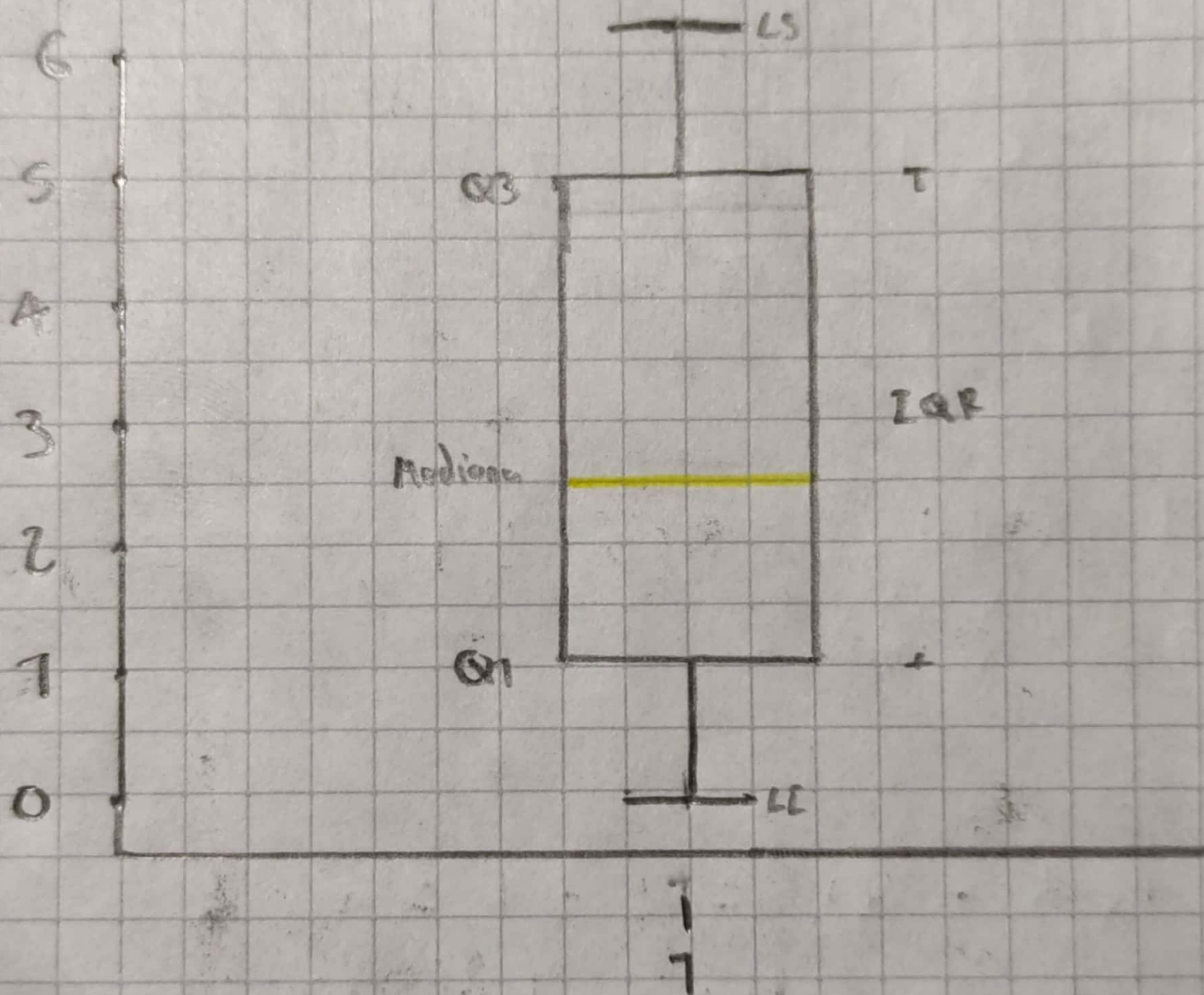
$$= -5.5$$

$$\text{Límite superior} = Q_3 + 1.5 \times IQR$$

$$= 5 + 1.5 \times 4$$

$$= 11.5$$

Box Plot X1



boxplots X_2

1. ordenar datos

$$X_2 = [0, 1, 2, 3, 4, 4]$$

$$Q_1 = X_{2,2} = 1$$

$$Q_2 = 2.5 \quad (\text{mediana})$$

$$Q_3 = X_{2,5} = 4$$

$$IQR = Q_3 - Q_1 = 4 - 1 = 3$$

$$\text{Límite inferior} : Q_1 - 1.5 \times IQR$$

$$= 1 - 1.5 \times 3$$

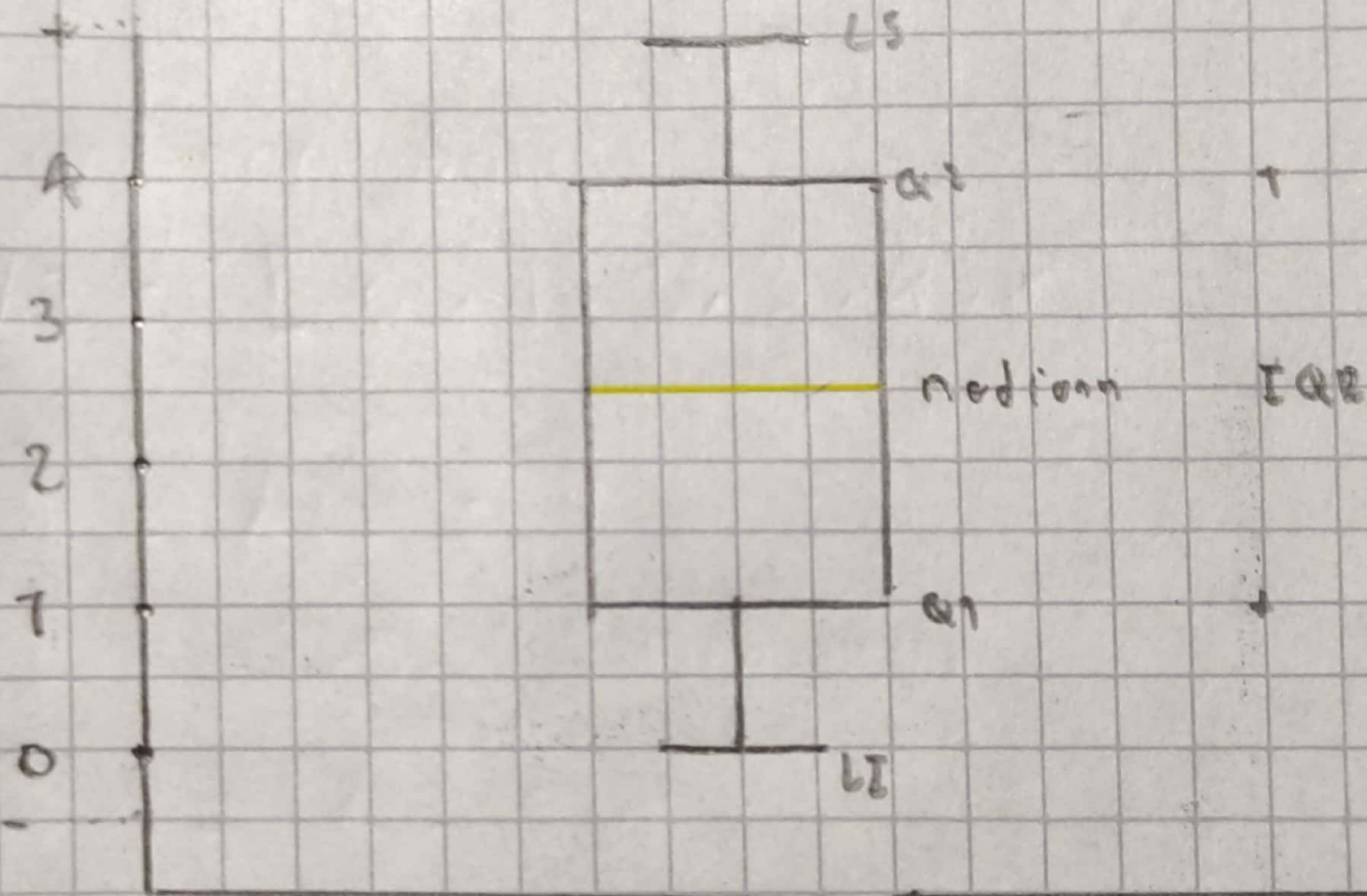
$$= -3.5$$

$$\text{Límite superior} : Q_3 + 1.5 \times IQR$$

$$4 + 1.5 \times 3 = 8.5$$

$$= 8.5$$

Box Plot X1



7-3

Covarianza entre X_1 y X_2

Formula

$$\text{COV}(X_1, X_2) = \frac{\sum_{i=1}^n (X_{1i} - \bar{X}_1)(X_{2i} - \bar{X}_2)}{n}$$

Procedimiento

$$\bar{X}_1 = \frac{12}{6} \text{ (media)} = a$$

$$\bar{X}_2 = \frac{74}{6} \text{ (media)} = b$$

$$\text{COV}(X_1, X_2) =$$

$$= \frac{(0-a)(4-b) + (7-a)(3-b) + (7-a)(4-b) + (4-a)(7-b) + (5-a)(7-b) + (6-a)(0-b)}{6}$$

$$\text{COV}(X_1, X_2) \approx -2.33$$

Correlación X_1 y X_2

$$\text{Corr}(X_1, X_2) = \frac{\text{COV}(X_1, X_2)}{S(X_1) S(X_2)}$$

Procedimiento

$$\text{COV}(X_1, X_2) \approx -2.33 \quad (\text{calcula directamente})$$

$$S(X_1) \approx 1.471 \quad (\text{calculado})$$

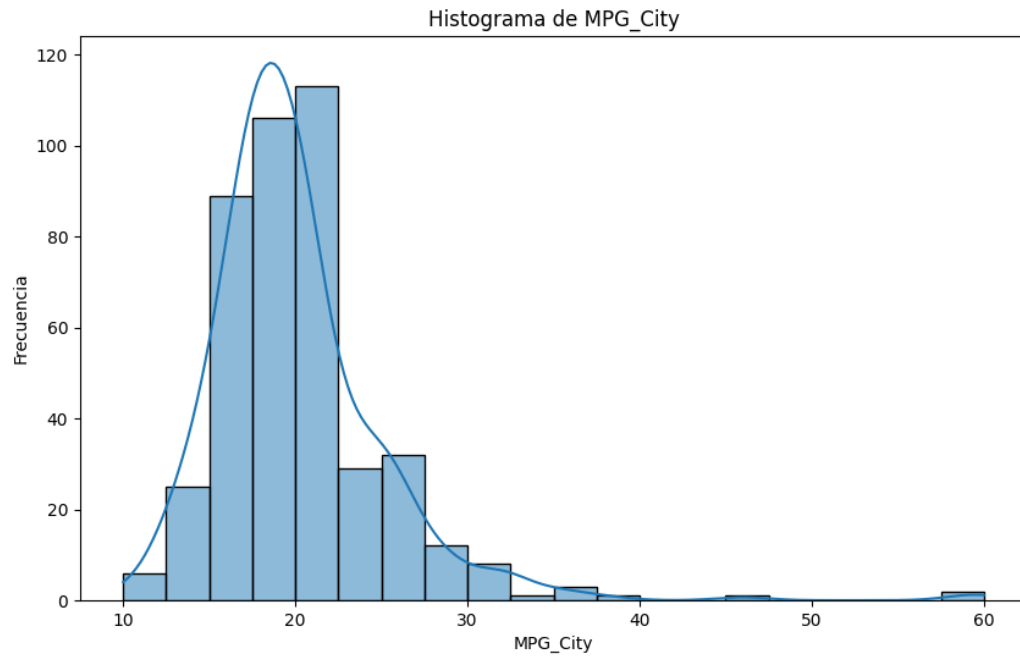
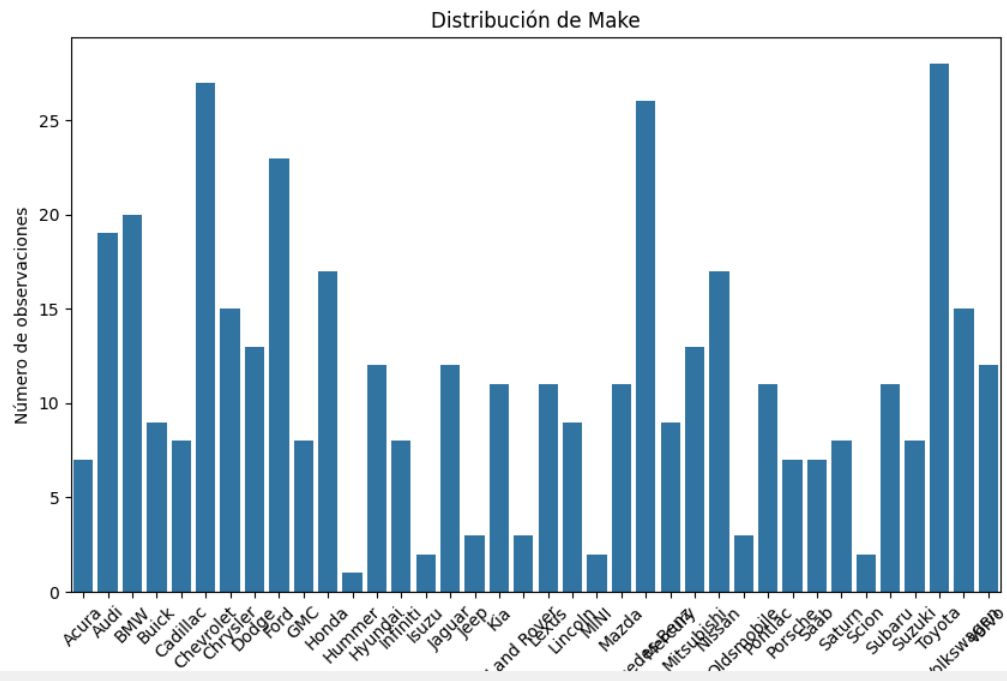
$$S(X_2) \approx 1.825 \quad (\text{calculado})$$

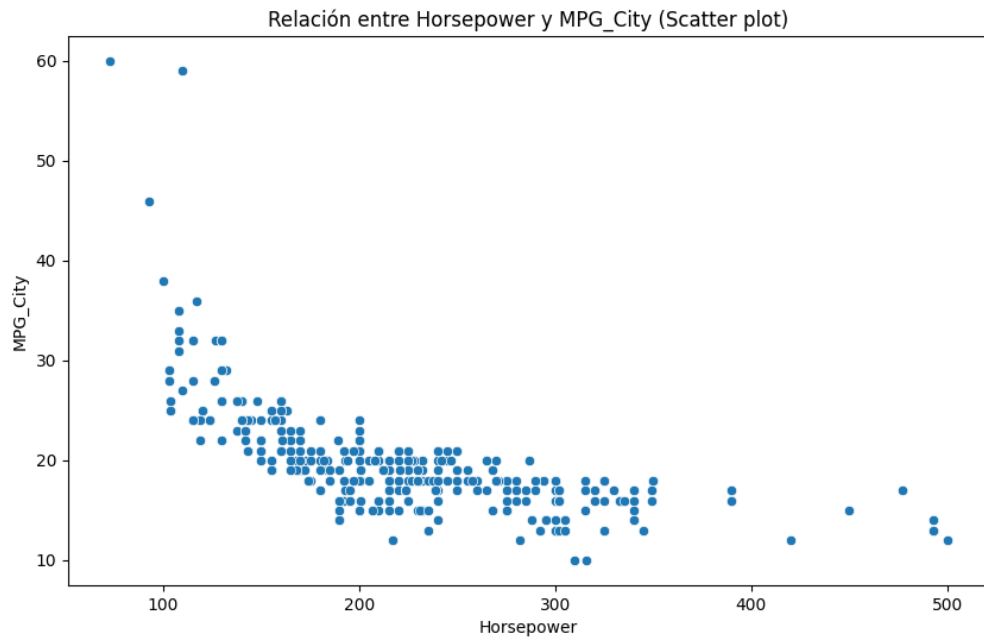
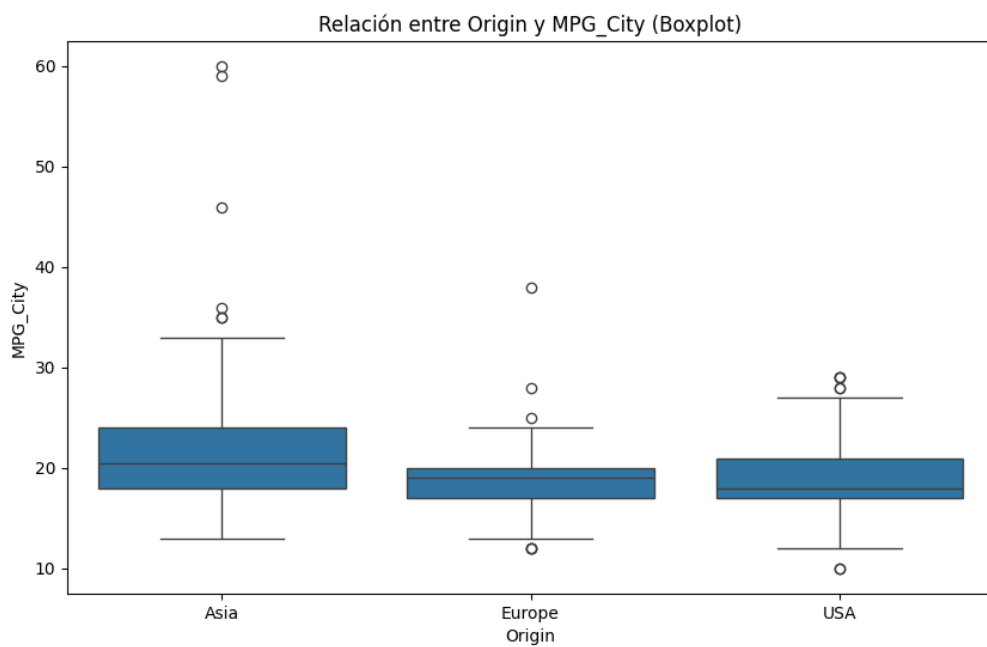
$$\text{Corr}(X_1, X_2) \approx \frac{-2.33}{(1.471)(1.825)}$$

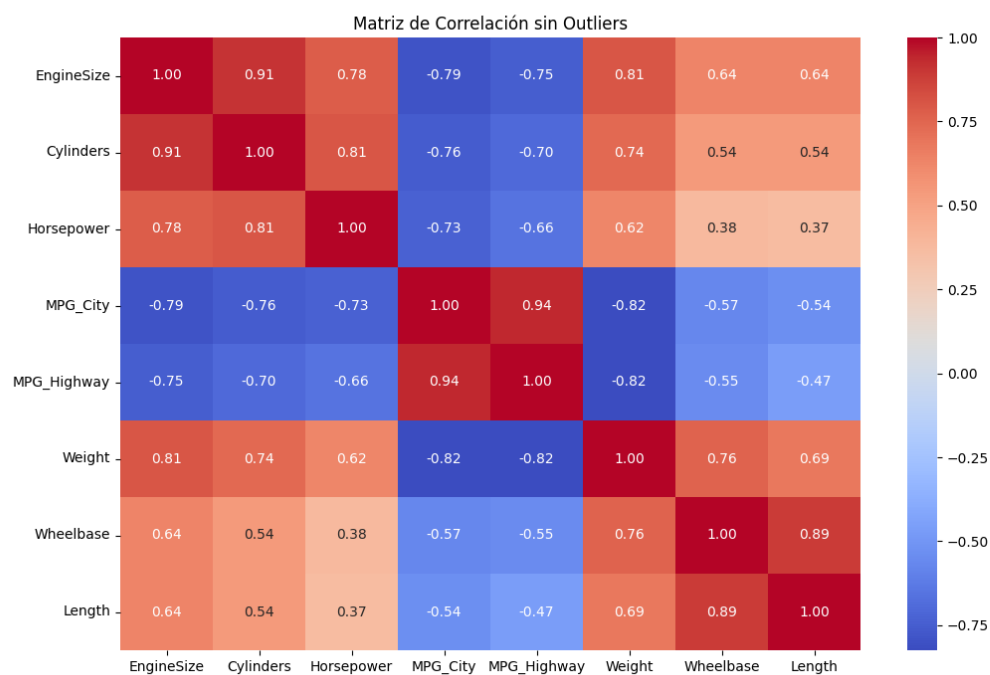
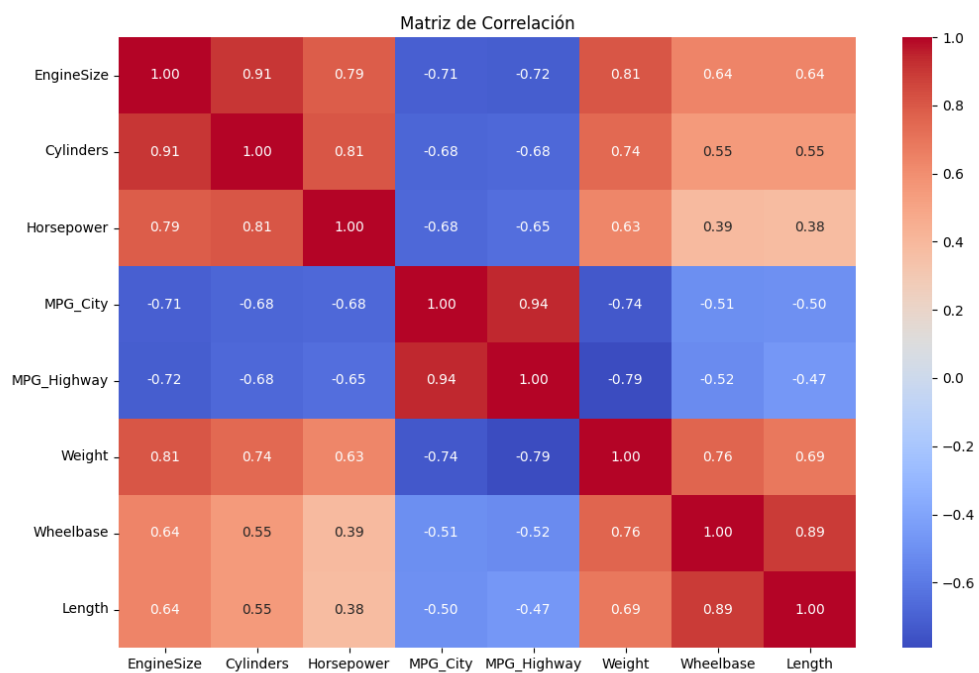
$$\text{Corr}(X_1, X_2) \approx -0.82$$

1.5. Explica la relación entre covarianza y correlación.

La covarianza proporciona información sobre la dirección de la relación, mientras que la correlación además cuantifica la fuerza y normaliza la medida en una escala más fácil de interpretar. Una correlación cercana a 1 o -1 indica una relación más fuerte, mientras que valores cercanos a 0 indican una relación más débil.







```

import pandas as pd
Outliers:
      Model  MPG_City
149 Civic Hybrid 4dr manual (gas/electric)      46
150      Insight 2dr (gas/electric)            60
373      Prius 4dr (gas/electric)              59
Estadístico de prueba: 0.8078395264093217, Valor p: 3.3896695924481945e-22
La distribución no es normal.
Outliers:
      Model  MPG_City
149 Civic Hybrid 4dr manual (gas/electric)      46
150      Insight 2dr (gas/electric)            60
373      Prius 4dr (gas/electric)              59
Estadístico de prueba: 0.8078395264093217, Valor p: 3.3896695924481945e-22
La distribución no es normal.
Variables más importantes:
MPG_Highway    0.941021
Weight         0.737966
EngineSize     0.709471
Cylinders      0.684402
Horsepower     0.676699
Wheelbase      0.507284
Length         0.501526
Name: MPG_City, dtype: float64

```

Por: Juan Pablo Toro Hurtado

CC: 1001477295

Correo: Juanpatoro2011@gmail.com