

Tarea #: 1

Realizado por: Santiago Cárdenas Franco

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 dia 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.

1.1

Media \rightarrow Cantidad de datos = 21 $\Rightarrow \bar{X} = \frac{\sum_{i=1}^n X_i}{N}$

$\bar{X}_1 = \frac{4+2+2+3+1+3+3+0+1+0+5+1+2+1+3+4+4+3+5+1+1}{21}$

$\bar{X}_1 = \frac{49}{21} \approx 2,33$

$\bar{X}_2 = \frac{4+3+4+5+3+6+6+1+3+0+0+2+3+3+6+8+8+6+9+2+2}{21}$

$\bar{X}_2 = \frac{93}{21} \approx 4,428$

$\bar{X}_3 = \frac{28+24+30+32+18+41+44+5+18+1+6+2+17+24+10+42+56+56+44+64+17+19}{21}$

$\bar{X}_3 = \frac{654}{21} \approx 31,143$

Mediana \rightarrow Se ordenan los valores de menor a mayor, Para impares $= Me = \frac{n+1}{2}$

$X_1 = 0, 0, 1, 1, 1, 1, 1, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 4, 5, 5$ $\Rightarrow Me = 11$ (Posición)

$Me_{X_1} = 2$

$X_2 = 0, 1, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 5, 6, 6, 6, 6, 8, 8, 9, 9$

$Me_{X_2} = 4$

$X_3 = 1, 5, 17, 17, 17, 18, 18, 19, 24, 24, 28, 30, 32, 41, 42, 44, 44, 56, 56, 62, 64$

$Me_{X_3} = 28$

Moda \rightarrow La Moda es el dato que más se repite, se usará la organización de los datos anteriores para identificar la moda.

$Mo_{X_1} = 1 \Rightarrow$ Se repite 6 veces

$Mo_{X_2} = 3 \Rightarrow$ Se repite 5 veces

$Mo_{X_3} = 17 \Rightarrow$ Se repite 3 veces

$X_2 = 0, 1, 2, 2, 2, 3, 3, 3, 3, 3, 4, 4, 5, 6, 6, 6, 6, 8, 8, 9, 9$

$Me_{X_2} = 4$

$X_3 = 1, 5, 17, 17, 17, 18, 18, 19, 24, 24, 28, 30, 32, 41, 42, 44, 44, 56, 56, 62, 64$

$Me_{X_3} = 28$

Moda \rightarrow La Moda es el dato que más se repite, se usará la organización de los datos anteriores para identificar la moda.

$Mo_{X_1} = 1 \Rightarrow$ Se repite 6 veces

$Mo_{X_2} = 3 \Rightarrow$ Se repite 5 veces

$Mo_{X_3} = 17 \Rightarrow$ Se repite 3 veces

Para X_1, X_2 y $X_3 \Rightarrow$ Tienen en común el dato 1 una sola vez

Para X_1 y $X_2 \Rightarrow$ Tienen en común los datos 0 (una vez), 1 (una vez), 2 (tres veces), 3 (cuatro veces), 4 (cuatro veces), y el 5 (una vez), la moda entre estos datos es 3 que se repite 5 veces.

Desviación Estándar $\rightarrow S = \sqrt{\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{N}}$

$S_{X_1} = \sqrt{\frac{(1,4)^2 + (-0,33)^2 + (-0,33)^2 + (0,66)^2 + (-1,33)^2 + (0,66)^2 + (0,66)^2 + (-2,33)^2 + (-1,33)^2 + (-2,33)^2 + (2,66)^2 + (-2,33)^2 + (-3,33)^2 + (-3,33)^2}{21}}$

$\dots \sqrt{\frac{(0,66)^2 + (1,66)^2 + (1,66)^2 + (0,66)^2 + (2,66)^2 + (-1,33)^2 + (1,33)^2}{21}} \approx \sqrt{\frac{48,96}{21}} \approx \sqrt{2,331} \approx 1,527$

$$Sx_2 = \sqrt{\frac{(0,428)^2 + (-1,428)^2 + (-0,428)^2 + (0,571)^2 + (-3,428)^2 + (1,571)^2 + (1,571)^2 + (-3,428)^2 + (-1,428)^2}{21}} \dots$$

$$\dots \sqrt{\frac{(0,428)^2 + (1,571)^2 + (-2,428)^2 + (-1,428)^2 + (-1,428)^2 + (1,571)^2 + (3,571)^2 + (3,571)^2 + (1,571)^2 + (0,571)^2}{21}} \dots$$

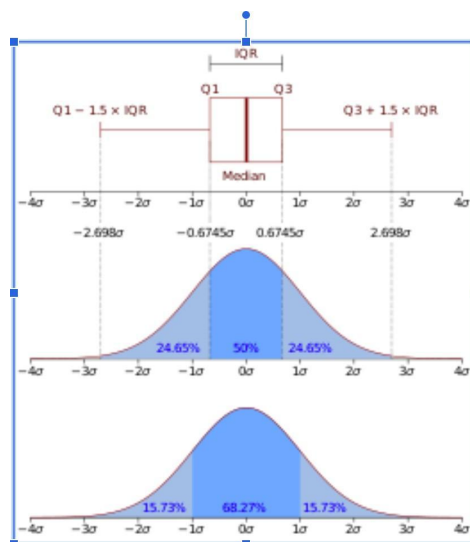
$$\dots \sqrt{\frac{(-2,428)^2 + (-2,428)^2}{21}} \approx \sqrt{\frac{143,93}{21}} \approx \sqrt{6,853} \approx 2,618$$

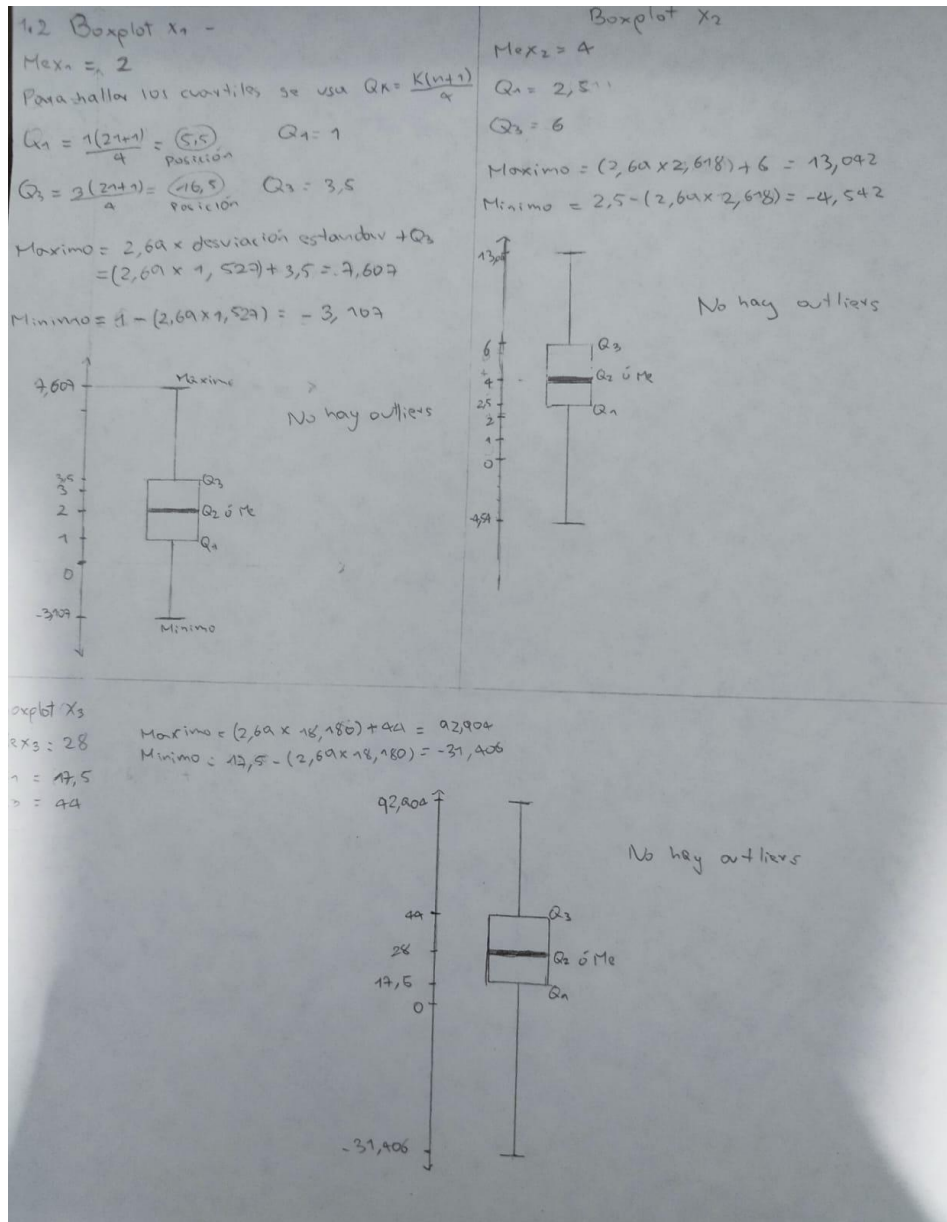
$$Sx_3 = \sqrt{\frac{(-3,380)^2 + (-7,380)^2 + (-1,380)^2 + (0,619)^2 + (-7,380)^2 + (0,619)^2 + (12,619)^2 + (-26,380)^2 + (-13,380)^2 + (-39,380)^2}{21}} \dots$$

$$\dots \sqrt{\frac{(30,619)^2 + (-10,380)^2 + (-7,380)^2 + (-13,380)^2 + (10,619)^2 + (24,619)^2 + (24,619)^2 + (-2,619)^2 + (32,619)^2 + \dots}{21}} \dots$$

$$\dots \sqrt{\frac{(-14,380)^2 + (-14,380)^2}{21}} \approx \sqrt{\frac{6940,26}{21}} \approx \sqrt{330,512} \approx 18,180$$

- 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 X_1, X_2

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

1,3

Covarianza entre $X_1, X_2 \Rightarrow \text{Cov}(X_1, X_2) = \frac{\sum (X_{1i} - \bar{X}_1) * (X_{2i} - \bar{X}_2)}{N}$

$X_1 = 4, 2, 2, 3, 1, 3, 3, 0, 1, 0, 5, 1, 2, 1, 3, 4, 4, 3, 5, 1, 1$
 $X_2 = 4, 3, 4, 5, 3, 6, 6, 1, 3, 0, 9, 2, 3, 3, 6, 8, 8, 6, 9, 2, 2$

$\text{Cov}(X_1, X_2) = \frac{(1,66 * (-0,428)) + (-0,333 * (-1,428)) + (-0,333 * (-0,428)) + (0,667 * 0,571) + (-1,333 * (-1,428)) + \dots}{21}$

$\dots \frac{(0,667 * 1,571) + (0,667 * 1,571) + (-2,333 * (-3,428)) + (-1,333 * (-1,428)) + (-2,333 * (-4,428)) + (2,667 * 4,571) + \dots}{21}$

$\dots \frac{(-1,333 * (-2,428)) + (-0,333 * (-1,428)) + (-1,333 * (-1,428)) + (0,667 * 1,571) + (1,667 * 3,571) + (1,667 * 3,571) + \dots}{21}$

$\dots \frac{(0,667 * 1,571) + (2,667 * 4,571) + (-1,333 * (-2,428)) + (-1,333 * (-2,428))}{21} \approx \frac{75}{21} \approx 3,571$

1.4. Cuál es la correlación entre la variable x_1 y x_2 (Calcularla a mano). Correlación puede ser escrita también como:

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

1,4

Correlación entre X_1 y $X_2 \Rightarrow \text{Cor}(X_1, X_2) = \frac{\sum_{i=1}^n (X_{1i} - \bar{X}_1) * (X_{2i} - \bar{X}_2)}{\sqrt{\sum_{i=1}^n (X_{1i} - \bar{X}_1)^2} * \sqrt{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2}}$

Como la fórmula de covarianza tiene el mismo numerador que correlación se enfocará en el denominador

$\sqrt{\sum_{i=1}^n (X_{1i} - \bar{X}_1)^2} \approx \sqrt{48,96} \rightarrow \text{Es } \sqrt{48,96} \text{ debido al numerador dentro de raíz cuadrada de la desviación estándar}$
 $\rightarrow \approx 6,997$

$\sqrt{\sum_{i=1}^n (X_{2i} - \bar{X}_2)^2} \approx \sqrt{143,93} \approx 11,997$

$\text{Cor}(X_1, X_2) = \frac{75}{6,997 * 11,997} \approx 0,91$

1.5. Explica la relación entre covarianza y correlación.

Respuesta: La covarianza es de 3.571 y la correlación es de 0.9, se puede decir que al ser la correlación tan alta y positiva indican que ambas variables crecen en el mismo sentido, y a la para si una crece la otra también, la covarianza como está de dispersos los datos y con ayuda de la correlación establece que ambas variables tienen una misma dispersión de los datos.

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).

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

Se organiza toda la tabla de acuerdo a la columna Random de manera decreciente

x1	x2	x3	Random
3	6	41	2
0	1	5	2
1	3	18	2
0	0	1	2
2	3	24	2
3	6	42	2
1	2	17	2
1	2	17	2
3	5	32	1
1	3	18	1
3	6	44	1

5	9	62	1
1	2	17	1
4	8	56	1
3	6	44	1
5	9	64	1
4	4	28	0
2	3	24	0
2	4	30	0
1	3	19	0
4	8	56	0

Después de organizar se determina centroides (Para lo centroides se realiza el promedio de acuerdo al grupo) y distancias (Se utiliza la fórmula euclidiana) y se determina el nuevo label de acuerdo a la distancia esta se comparan y la menor determina el grupo

x1	x2	x3	Random	centroide _1	centroide _2	centroide _3	distancia_c1 _2	distancia_c2_ 1	distancia_c3 _0	Nuevo label iteración 1
3	6	41	2	1.375	2.875	20.625	20.67720665	1.131923142	9.74063653	1
0	1	5	2				15.79705273	37.59030793	26.7447191	2
1	3	18	2				2.654595073	24.40350897	13.56760848	2
0	0	1	2				19.88207421	41.67770687	30.82661188	2
2	3	24	2				3.434657916	18.40601125	7.555130707	2
3	6	42	2				21.66326095	0.1767766953	10.72753467	1
1	2	17	2				3.747916088	25.53000685	14.6860478	2
1	2	17	2				3.747916088	25.53000685	14.6860478	2
3	5	32	1	3.125	6	42.125	11.68532734	10.17503071	0.938083152	0
1	3	18	1				2.654595073	24.40350897	13.56760848	2
3	6	44	1				23.63888481	1.879162047	12.70747811	1
5	9	62	1				41.98269733	20.18740325	31.03675241	1
1	2	17	1				3.747916088	25.53000685	14.6860478	2
4	8	56	1				35.84057582	14.04568439	24.90140558	1
3	6	44	1				23.63888481	1.879162047	12.70747811	1
5	9	64	1				43.95505517	22.15922494	33.01030142	1
4	4	28	0	2.6	4.4	31.4	7.90865823	14.29269919	3.698648402	0
2	3	24	0				3.434657916	18.40601125	7.555130707	2
2	4	30	0				9.462921061	12.34022893	1.574801575	0
1	3	19	0				1.67238602	23.41540625	12.58093796	2
4	8	56	0				35.84057582	14.04568439	24.90140558	1

Ahora se cambia la etiqueta Random por el de la iteración 1 y se vuelve a organizar toda la tabla de mayor a menor

x1	x2	x3	Nuevo label iteracción 1
0	1	5	2
1	3	18	2
0	0	1	2

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

Se calcula la iteración 2

x1	x2	x3	Nuevo label iteración 1	centroid e_1	centroid e_2	centroid e_3	distanci a_c1_2	distanci a_c2_1	distanci a_c3_0	Nuevo label iteración 2
0	1	5	2	1	2.2	16	11.11035553	46.6973996	25.3990376	2
1	3	18	2				2.154065923	33.5095602	12.23837317	2
0	0	1	2				15.19341963	150.7852023	29.47503652	2
2	3	24	2				8.10185164	127.51164526	46.227180564	0
1	2	17	2				1.019803903	34.63582863	13.35830994	2
1	2	17	2				1.019803903	34.63582863	13.35830994	2
1	3	18	2				2.154065923	33.5095602	12.23837317	2
1	2	17	2				1.019803903	34.63582863	13.35830994	2
2	3	24	2				8.10185164	127.51164526	46.227180564	0
1	3	19	2				3.104834939	32.52138719	11.25956384	2
3	6	41	1	3.75	7.25	51.125	25.36611914	10.2294003	11.12554618	1
3	6	42	1				26.352229951	29.240704789	12.1151879	1

3	6	44	1				28.327377.27259414.09885	1
							192 104 732	
5	9	62	1				46.6716111.0856032.40027	1
							878 44 435	
4	8	56	1				40.529494.93868626.27630	1
							543 566 957	
3	6	44	1				28.327377.27259414.09885	1
							192 104 732	
5	9	64	1				48.6440113.0533734.37699	1
							299 6 489	
4	8	56	1				40.529494.93868626.27630	1
							543 566 957	
3	5	32	0	3	4.333333	30	16.3658119.271492.108185	0
					333		804 774 107	
4	4	28	0				12.4995923.353592.260776	0
							999 983 661	
2	4	30	0				14.1506121.445061.054092	0
							836 062 553	

Ahora se cambia la etiqueta iteración 1 por la etiqueta iteración 2 y se vuelve a organizar toda la tabla de mayor a menor

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

Se calcula la iteración 3

x1	x2	x3	Nuevo centroide label _1	centroide _2	centroide _3	distancia _c1_2	distancia _c2_1	distancia _c3_0	Nuevo label iteración 3
----	----	----	-----------------------------	-----------------	-----------------	--------------------	--------------------	--------------------	----------------------------

iteración 2										
0	1	5	2	0.75	2	14	9.086390	46.69732	22.92073	2
							923	996	297	
1	3	18	2				4.130677	33.50956	9.765244	2
							91	02	493	
0	0	1	2				13.17431	50.78524	26.99555	2
							213	023	519	
1	2	17	2				3.010398	34.63582	10.87014	2
							645	863	259	
1	2	17	2				3.010398	34.63582	10.87014	2
							645	863	259	
1	3	18	2				4.130677	33.50956	9.765244	2
							91	02	493	
1	2	17	2				3.010398	34.63582	10.87014	2
							645	863	259	
1	3	19	2				5.105144	32.52138	8.784076	2
							464	719	502	
3	6	41	1	3.75	7.25	51.125	27.38726	10.22940	13.58528	1
							894	003	616	
3	6	42	1				28.37362	9.240704	14.57257	1
							331	789	699	
3	6	44	1				30.34901	7.272594	16.55173	1
							152	104	707	
5	9	62	1				48.69355	11.08560	34.87348	1
							707	44	563	
4	8	56	1				42.55070	4.938686	28.74299	1
							505	566	915	
3	6	44	1				30.34901	7.272594	16.55173	1
							152	104	707	
5	9	64	1				50.66618	13.05337	36.84779	1
							695	6	505	
4	8	56	1				42.55070	4.938686	28.74299	1
							505	566	915	
2	3	24	0	2.6	3.8	27.6	10.12731	27.51164	3.736308	0
							455	526	338	
2	3	24	0				10.12731	27.51164	3.736308	0
							455	526	338	
3	5	32	0				18.38647	19.27149	4.578209	0
							601	774	257	
4	4	28	0				14.51077	23.35359	1.469693	0
							186	983	846	
2	4	30	0				16.17289	21.44506	2.481934	0
							399	062	729	

Desde este punto no se realiza más iteraciones debido a que si se ve la iteración 2 ordenada de mayor a menor y la iteración 3, son complementamente iguales, es más si se reemplaza la iteración 2 con la 3 no se tiene que ordenar porque ya están ordenadas, y como se dijo anteriormente la iteración 2 ordenada es igual a la iteración 3 este ciclo se vuelve a repetir infinitamente.

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. PCA entre X_1 y X_2

1) Centrar y estandarizar

$\mu_1 = 2,33$
 $\mu_2 = 4,428$

$\sigma_1 = 1,527$
 $\sigma_2 = 2,618$

La desviación estándar y la media sale del punto 1.1

$X = \begin{bmatrix} X_1 & X_2 \\ 4 & 4 \\ 2 & 3 \\ 2 & 4 \\ 3 & 5 \\ 1 & 3 \\ 3 & 6 \\ 3 & 6 \\ 0 & 1 \\ 1 & 3 \\ 0 & 0 \\ 5 & 9 \\ 2 & 2 \\ 1 & 3 \\ 3 & 6 \\ 4 & 8 \\ 4 & 8 \\ 3 & 6 \\ 5 & 9 \\ 1 & 2 \\ 1 & 2 \end{bmatrix}$

$X_{\text{estandarizada}} = \begin{bmatrix} X_{1n} & X_{2n} \\ 1,091 & -0,163 \\ -0,218 & -0,545 \\ -0,218 & -0,163 \\ 0,436 & 0,218 \\ -0,872 & -0,545 \\ 0,436 & 0,600 \\ 0,436 & 0,600 \\ -1,527 & -1,309 \\ -0,892 & -0,545 \\ -1,527 & -1,691 \\ 1,745 & 1,245 \\ -0,872 & -0,872 \\ -0,218 & -0,545 \\ -0,872 & -0,545 \\ 0,436 & 0,600 \\ 1,091 & 1,363 \\ 1,091 & 1,363 \\ 0,436 & 0,600 \\ 1,745 & 1,745 \\ -0,872 & -0,872 \\ -0,872 & -0,872 \end{bmatrix}$

$\mu_{1n} = 0$
 $\mu_{2n} = 0$

$\sigma_{1n} = 1$
 $\sigma_{2n} = 1$

2.1. Cual es la matriz de covarianza

2) Matriz de covarianza

$\Sigma = \begin{bmatrix} \sigma_{X_1}^2 & \text{Cov}(X_1, X_2) \\ \text{Cov}(X_2, X_1) & \sigma_{X_2}^2 \end{bmatrix} = \begin{bmatrix} 1 & 0,892 \\ 0,892 & 1 \end{bmatrix}$

$\text{Cov}(X_1, X_2) = 0,892$
 $\text{Cov}(X_2, X_1) = 0,892$

2.2. Cuales son los eigenvalues

3) Eigen values $\rightarrow \det(\Sigma - \lambda I) = 0$

$$\det \left(\begin{bmatrix} 1 & 0,892 \\ 0,892 & 1 \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right) = \det \begin{pmatrix} 1-\lambda & -0,892 \\ 0,892 & 1-\lambda \end{pmatrix} \Rightarrow (1-\lambda)^2 - 0,892^2 \Rightarrow \dots$$

$$\dots \Rightarrow 1 - 2\lambda + \lambda^2 - 0,892^2 = 0 \Rightarrow 0,203 = 2\lambda + \lambda^2$$

$$\lambda = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \Rightarrow \lambda = \frac{+2 \pm \sqrt{4 - 4 \times 0,203 \times 1}}{2 \times 1} \Rightarrow \lambda = \frac{2 \pm \sqrt{4 - 0,812}}{2} \Rightarrow \lambda = \frac{2 \pm 1,785}{2}$$

$\lambda_1 = 1,892$
 $\lambda_2 = 0,107$ Eigen values.

2.3. Cuál es la varianza explicada por el eigenvalue.

$\lambda_1 = 1,892$
 $\lambda_2 = 0,107$ Eigen values.

Para calcular las varianzas de los eigen values, se determina la suma total de los eigen values

$$\lambda_T = \lambda_1 + \lambda_2 \Rightarrow \lambda_T = 1,892 + 0,107 \Rightarrow \lambda_T = 1,999$$

Varianza explicada para $\lambda_1 = (1,892/1,999) \times 100\% = 94,65\%$
 $(1,892/2) \times 100\% = 94,64\%$

$\lambda_2 = (0,107/1,999) \times 100\% = 5,35\%$
 $(0,107/2) \times 100\% = 5,36\%$

2.4. Cual es el valor del eigenvector

Eigen vector

Se tomará el eigenvalue λ_1 debido a que es el mayor, tomando un 94,64% de los datos

$$\begin{bmatrix} 1 & 0,892 \\ 0,892 & 1 \end{bmatrix} \cdot \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \lambda_1 \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

① $x_1 + 0,892x_2 = 1,892x_1 \Rightarrow 0,892x_2 = 0,892x_1$
 ② $0,892x_1 + x_2 = 1,892x_2 \Rightarrow 0,892x_1 = 0,892x_2$

\Downarrow
 $x_1 = x_2$
 $x_2 = x_1$

$V_1 = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \Rightarrow$ se verifica $\begin{bmatrix} 1 & 0,892 \\ 0,892 & 1 \end{bmatrix} \begin{bmatrix} 1 \\ 1 \end{bmatrix} = 1,892 \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 1,892 \\ 1,892 \end{bmatrix} = 1,892 \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$\|V_1\| = \sqrt{x_1^2 + x_2^2} = 1 \Rightarrow \sqrt{2x_1^2} \Rightarrow \sqrt{\frac{1}{2}} = x_1 \Rightarrow x_1 = 0,707$

$V_1 = [0,707 \quad 0,707]$

2.5. Cuál es la matriz proyectada.

4) Cálculo de matriz proyectada en una dimensión

$$Z = \begin{bmatrix} 1,091 & -0,163 \\ -0,218 & -0,545 \\ -0,218 & -0,163 \\ 0,436 & 0,218 \\ -0,872 & -0,545 \\ 0,436 & 0,600 \\ 0,436 & 0,600 \\ -1,527 & -1,300 \\ -0,872 & -0,545 \\ -1,527 & -1,691 \\ 1,745 & 1,745 \\ -0,872 & -0,927 \\ -0,218 & -0,545 \\ -0,872 & -0,545 \\ 0,436 & 0,600 \\ 1,091 & 1,363 \\ 1,091 & 1,363 \\ 0,436 & 0,600 \\ 1,745 & 1,745 \\ -0,872 & -0,927 \\ -0,872 & -0,927 \end{bmatrix} \cdot \begin{bmatrix} 0,707 \\ 0,707 \end{bmatrix} = \begin{bmatrix} 0,671 \\ -0,553 \\ -0,27 \\ 0,47 \\ -1,02 \\ 0,75 \\ 0,75 \\ -2,05 \\ -1,02 \\ -2,33 \\ 2,52 \\ -1,30 \\ -0,553 \\ -1,02 \\ 0,75 \\ 1,77 \\ 1,77 \\ 0,75 \\ 2,52 \\ -1,30 \\ -1,30 \end{bmatrix}$$

2.6. Cual es el error o diferencia entre la matriz proyectada

5) Error de la matriz proyectada

Para el error se resta la matriz estandarizada con la conseguida

$$\bar{X} = \begin{bmatrix} 0,671 \\ -0,553 \\ -0,27 \\ 0,47 \\ -1,02 \\ 0,75 \\ 0,75 \\ -2,05 \\ -1,02 \\ -2,33 \\ 2,52 \\ -1,30 \\ -0,553 \\ -1,02 \\ 0,75 \\ 1,77 \\ 1,77 \\ 0,75 \\ 2,52 \\ -1,30 \\ -1,30 \end{bmatrix} \cdot \begin{bmatrix} 0,707 & 0,707 \end{bmatrix} = \begin{bmatrix} 0,671 & 0,671 \\ -0,391 & -0,391 \\ -0,190 & -0,190 \\ 0,332 & 0,332 \\ -0,721 & -0,721 \\ 0,530 & 0,530 \\ 0,530 & 0,530 \\ -1,449 & -1,449 \\ -0,721 & -0,721 \\ -1,647 & -1,647 \\ 1,781 & 1,781 \\ -0,919 & -0,919 \\ -0,388 & -0,388 \\ -0,721 & -0,721 \\ 0,530 & 0,530 \\ 1,251 & 1,251 \\ 1,251 & 1,251 \\ 0,530 & 0,530 \\ 1,781 & 1,781 \\ -0,919 & -0,919 \\ -0,919 & -0,919 \end{bmatrix}$$

$$Error = \begin{bmatrix} 0,616 & -0,638 \\ 0,192 & -0,154 \\ -0,079 & 0,027 \\ 0,104 & -0,114 \\ -0,151 & 0,175 \\ -0,093 & 0,069 \\ -0,093 & 0,069 \\ -0,079 & 0,120 \\ -0,151 & -0,036 \\ -0,151 & 0,175 \\ 0,120 & -0,093 \\ -0,036 & -0,008 \\ 0,006 & -0,156 \\ 0,170 & 0,175 \\ -0,151 & 0,069 \\ -0,093 & 0,112 \\ -0,160 & 0,112 \\ -0,093 & 0,069 \\ -0,036 & -0,036 \\ 0,006 & -0,08 \\ 0,006 & -0,08 \end{bmatrix}$$

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.

Se realizan las gráficas en el notebook en la celda 9.

- 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.

Se realizan las gráficas en el notebook en la celda 10.

Se realiza el test de distribución normal para cada columna en la celda 11 del notebook.

Se listan los modelos de carros que están más lejos de 5 desviaciones estándar en la celda 12 del notebook.

- 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

Se realizan los gráficos (boxplot) con respecto MPG_City y sus respectivos análisis en la celda 14 del notebook.

- 3.2.2. Variables numéricas vas a crear un scatter plot. Explique cómo interpreta el gráfico

Se realizan los gráficos (scatter) con respecto MPG_City y sus respectivos análisis en la celda 15 del notebook.

- 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.

Se realiza la matriz de correlación y se analizan la variables más importantes para MPG_City en la celda 16 del notebook

Como se ve en la matriz de correlación se identifica que la variable con mayor relación es MPG_Highway y tiene sentido debido a que la gráfica los datos crecen de manera similar, esto también significa que los vehículos que tienen un buen rendimiento en consumo de combustible en la ciudad tienden a tener también un buen rendimiento en consumo de combustible en autopistas o en carreteras.

Vemos también que MPG_Highway tiene una correlación fuertes con EngineSize(-0.717302), Horsepower(-0.647195), Cylinders(-0.676100) y Weight(-0.790989) al ser MPG_Highway muy similar en comportamiento con MPG_City ambos tienen correlaciones similares obviamente no tienen la misma proporción, pero se pueden sacar las mismas conclusiones para ambas.

Otra correlación bastante fuerte es MPG_City con EngineSize, es una correlación negativa que indica que entre menor sea el tamaño del motor, más combustible ahorra.

Otra correlación bastante fuerte es MPG_City con Cylinders, es una correlación negativa que indica que entre menor sea la cantidad de cilindros, menos combustible gasta.

Otra correlación bastante fuerte es MPG_City con Horsepower, es una correlación negativa que indica que si la potencia del motor aumenta, el consumo de combustible en la ciudad tiende a disminuir, es bastante lógico, ya que entre más potente es el motor generalmente consumen más combustible.

Otra correlación bastante fuerte es MPG_City con Weight, es una correlación negativa que indica que si un auto pesa mucho, su consumo de combustible se ve bastante afectado de manera negativa, es decir consume mucho más debido a su peso.

En si la correlación con mayor valor es MPG_Highway debido a que son muy similares, pero también se tiene EngineSize, Cylinders, Horsepower y Weight.

- 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?

Se generan las dummy variables en la celda 17 del notebook.

Se realiza la matriz de correlación y se analiza la variable categórica más importante para MPG_City en la celda 18 del notebook.

Se observa que la variable categórica con mayor relación a MPG_City es Type_Hybrid con una correlación del 0.561053, es decir que los autos híbridos tienden a consumir menos gasolina, también se puede observar que esta afirmación es correcta con respecto a la gráfica realizada anteriormente.

- 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.

Se excluyen los outliers de la data de carros y se realiza la matriz de correlación en la celda 19 del notebook.

Se logra observar que los datos se alteraron, si volvemos a comparar MPG_Highway que tiene una correlación de 0.944083 con MPG_City, antes de eliminar los outliers se observa 0.941021, no es una alteración tan grande, pero para otros datos si puede llegar a modificarlos bastante como a Weight que pasó de -0.737966 a -0.802215, es decir que los outliers estaban actuando como un disminuidor de fuerzas.